**now**

the essence of knowledge

# A Primer on Reproducing Kernel Hilbert Spaces

Jonathan H. Manton
The University of Melbourne
Victoria 3010 Australia
j.manton@ieee.org

Pierre-Olivier Amblard
CNRS
Grenoble 38402 France
pierre-olivier.amblard@gipsa-lab.grenoble-inp.fr

# Contents

## Abstract

Reproducing kernel Hilbert spaces are elucidated without assuming prior familiarity with Hilbert spaces. Compared with extant pedagogic material, greater care is placed on motivating the definition of reproducing kernel Hilbert spaces and explaining when and why these spaces are efficacious. The novel viewpoint is that reproducing kernel Hilbert space theory studies extrinsic geometry, associating with each geometric configuration a canonical overdetermined coordinate system. This coordinate system varies continuously with changing geometric configurations, making it well-suited for studying problems whose solutions also vary continuously with changing geometry. This primer can also serve as an introduction to infinite-dimensional linear algebra because reproducing kernel Hilbert spaces have more properties in common with Euclidean spaces than do more general Hilbert spaces.

# 1

## Introduction

Hilbert space theory is a prime example in mathematics of a beautiful synergy between symbolic manipulation and visual reasoning. Two-dimensional and three-dimensional pictures can be used to reason about infinite-dimensional Hilbert spaces, with symbolic manipulations subsequently verifying the soundness of this reasoning, or suggesting modifications and refinements. Visualising a problem is especially beneficial because over half the human brain is involved to some extent with visual processing. Hilbert space theory is an invaluable tool in numerous signal processing and systems theory applications [62, 12, 10].

Hilbert spaces satisfying certain additional properties are known as Reproducing Kernel Hilbert Spaces (RKHSs), and RKHS theory is normally described as a transform theory between Reproducing Kernel Hilbert Spaces and positive semi-definite functions, called kernels: every RKHS has a unique kernel, and certain problems posed in RKHSs are more easily solved by involving the kernel. However, this description hides the crucial aspect that the kernel captures not just intrinsic properties of the Hilbert space but also how the Hilbert space is embedded in a larger function space, which is referred to here as its extrinsic geometry. A novel feature of this primer is drawing attention to this

extrinsic geometry, and using it to explain why certain problems can be solved more efficiently in terms of the kernel than the space itself.

Another novel feature of this primer is that it motivates and develops RKHS theory in finite dimensions before considering infinite dimensions. RKHS theory is ingenious; the underlying definitions are simple but powerful and broadly applicable. These aspects are best brought out in the finite-dimensional case, free from the distraction of infinite-dimensional technicalities. Essentially all of the finite-dimensional results carry over to the infinite-dimensional setting.

This primer ultimately aims to empower readers to recognise when and how RKHS theory can profit them in their own work. The following are three of the known uses of RKHS theory.

1. If a problem involves a subspace of a function space, and if the subspace (or its completion) is a RKHS, then the additional properties enjoyed by RKHSs may help solve the problem. (Explicitly computing limits of sequences in Hilbert spaces can be difficult, but in a RKHS the limit can be found pointwise.)

2. Certain classes of problems involving positive semi-definite functions can be solved by introducing an associated RKHS whose kernel is precisely the positive semi-definite function of interest. A classic example, due to Parzen, is associating a RKHS with a stochastic process, where the kernel of the RKHS is the covariance function of the stochastic process (see §7.2).

3. Given a set of points and a function specifying the desired distances between points, the points can be embedded in a RKHS with the distances between points being precisely as prescribed; see §5. (Support vector machines use this to convert certain nonlinear problems into linear problems.)

In several contexts, RKHS methods have been described as providing a unified framework [77, 32, 46, 59]; although a subclass of problems was solved earlier by other techniques, a RKHS approach was found to be more elegant, have broader applicability, or offer new insight for obtaining actual solutions, either in closed form or numerically. Parzen

describes RKHS theory as facilitating a coordinate-free approach [46]. While the underlying Hilbert space certainly allows for coordinate-free expressions, the power of a RKHS beyond that of a Hilbert space is the presence of two coordinate systems: the pointwise coordinate system coming from the RKHS being a function space, and a canonical (but overdetermined) coordinate system coming from the kernel. The pointwise coordinate system facilitates taking limits while a number of geometric problems have solutions conveniently expressed in terms of what we define to be the canonical coordinate system. (Geometers may wish to think of a RKHS as a subspace $V \subset \mathbb{R}^X$ with pointwise coordinates being the extrinsic coordinates coming from $\mathbb{R}^X$ while the canonical coordinates are intrinsic coordinates on $V$ relating directly to the inner product structure on $V$.)

The body of the primer elaborates on all of the points mentioned above and provides simple but illuminating examples to ruminate on. Parenthetical remarks are used to provide greater technical detail that some readers may welcome. They may be ignored without compromising the cohesion of the primer. Proofs are there for those wishing to gain experience at working with RKHSs; simple proofs are preferred to short, clever, but otherwise uninformative proofs. Italicised comments appearing in proofs provide intuition or orientation or both.

This primer is neither a review nor a historical survey, and as such, many classic works have not been discussed, including those by leading pioneers such as Wahba [73, 72].

**Contributions**   This primer is effectively in two parts. The first part (§1–§7), written by the first author, gives a gentle and novel introduction to RKHS theory. It also presents several classical applications. The second part (§8–§9), with §8 written jointly and §9 written by the second author, focuses on recent developments in the machine learning literature concerning embeddings of random variables.

## 1.1 Assumed Knowledge

Basic familiarity with concepts from finite-dimensional linear algebra is assumed: vector space, norm, inner product, linear independence, basis, orthonormal basis, matrix manipulations and so forth.

Given an inner product $\langle \cdot, \cdot \rangle$, the induced norm is $\|x\| = \sqrt{\langle x, x \rangle}$. Not every norm comes from an inner product, meaning some norms cannot be written in this form. If a norm does come from an inner product, the inner product can be uniquely determined from the norm by the polarisation identity $4\langle x, y \rangle = \|x + y\|^2 - \|x - y\|^2$. (A corresponding formula exists for complex-valued vector spaces.)

A metric $d(\cdot, \cdot)$ is a "distance function" describing the distance between two points in a metric space. To be a valid metric, it must satisfy several axioms, including the triangle inequality. A normed space is automatically a metric space by the correspondence $d(x, y) = \|x - y\|$.

## 1.2 Extrinsic Geometry and a Motivating Example

Differential geometry groups geometric properties into two kinds: intrinsic and extrinsic. Intrinsic properties depend only on the space itself, while extrinsic properties depend on precisely how the space is embedded in a larger space. A simple example in linear algebra is that the orientation of a straight line passing through the origin in $\mathbb{R}^2$ describes the *extrinsic* geometry of the line.

The following observation helps motivate the development of finite-dimensional RKHS theory in §2. Let

$$L(\theta) = \{(t \cos \theta, t \sin \theta) \mid t \in \mathbb{R}\} \subset \mathbb{R}^2 \tag{1.1}$$

denote a straight line in $\mathbb{R}^2$ passing through the origin and intersecting the horizontal axis at an angle of $\theta$ radians; it is a one-dimensional subspace of $\mathbb{R}^2$. Fix an arbitrary point $p = (p_1, p_2) \in \mathbb{R}^2$ and define $f(\theta)$ to be the point on $L(\theta)$ closest to $p$ with respect to the Euclidean metric. It can be shown that

$$f(\theta) = (r(\theta) \cos \theta, r(\theta) \sin \theta), \qquad r(\theta) = p_1 \cos \theta + p_2 \sin \theta. \tag{1.2}$$

Visualising $f(\theta)$ as the projection of $p$ onto $L(\theta)$ shows that $f(\theta)$ depends *continuously* on the orientation of the line. While (1.2) veri-

fies this continuous dependence, it resorted to introducing an *ad hoc* parametrisation $\theta$, and different values of $\theta$ (e.g., $\theta$, $\pi + \theta$ and $2\pi + \theta$) can describe the same line.

> Is there a more natural way of representing $L(\theta)$ and $f(\theta)$, using linear algebra?

A first attempt might involve using an orthonormal basis vector to represent $L(\theta)$. However, there is no *continuous* map from the line $L(\theta)$ to an orthonormal basis vector $v(\theta) \in L(\theta)$. (This should be self-evident with some thought, and follows rigorously from the Borsuk-Ulam theorem.) Note that $\theta \mapsto (\cos \theta, \sin \theta)$ is not a well-defined map from $L(\theta)$ to $\mathbb{R}^2$ because $L(0)$ and $L(\pi)$ represent the same line yet $(\cos 0, \sin 0) \neq (\cos \pi, \sin \pi)$.

RKHS theory uses not one but two vectors to represent $L(\theta)$. Specifically, it turns out that the kernel of $L(\theta)$, in matrix form, is

$$K(\theta) = \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix} \begin{bmatrix} \cos \theta & \sin \theta \end{bmatrix} = \begin{bmatrix} \cos^2 \theta & \sin \theta \cos \theta \\ \sin \theta \cos \theta & \sin^2 \theta \end{bmatrix}. \qquad (1.3)$$

The columns of $K(\theta)$ are

$$k_1(\theta) = \cos \theta \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix}, \qquad k_2(\theta) = \sin \theta \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix}. \qquad (1.4)$$

Note that $L(\theta)$ is spanned by $k_1(\theta)$ and $k_2(\theta)$, and moreover, both $k_1$ and $k_2$ are well-defined (and continuous) functions of $L$; if $L(\theta) = L(\phi)$ then $k_1(\theta) = k_1(\phi)$ and $k_2(\theta) = k_2(\phi)$. To emphasise, although $\theta$ is used here for convenience to describe the construction, RKHS theory defines a map from $L$ to $k_1$ and $k_2$ that does not depend on any *ad hoc* choice of parametrisation. It is valid to write $k_1(L)$ and $k_2(L)$ to show they are functions of $L$ alone.

Interestingly, $f$ has a simple representation in terms of the kernel:

$$f(L) = p_1 \, k_1(L) + p_2 \, k_2(L). \qquad (1.5)$$

Compared with (1.2), this is both simple and natural, and does not depend on any *ad hoc* parametrisation $\theta$ of the line $L$. In summary,

- the kernel represents a vector subspace by a possibly overdetermined (i.e., linearly dependent) ordered set of vectors, and the correspondence from a subspace to this ordered set is continuous;

- this *continuous* correspondence cannot be achieved with an ordered set of basis (i.e., linearly independent) vectors;

- certain problems have solutions that depend continuously on the subspace and can be written elegantly in terms of the kernel:

$$\text{subspace} \rightarrow \text{kernel} \rightarrow \text{solution.} \tag{1.6}$$

The above will be described in greater detail in §2.

**Remark**   The above example was chosen for its simplicity. Ironically, the general problem of projecting a point onto a subspace is not well-suited to the RKHS framework for several reasons, including that RKHS theory assumes there is a norm only on the subspace; if there is a norm on the larger space in which the subspace sits then it is ignored. A more typical optimisation problem benefitting from RKHS theory is finding the minimum-norm function passing through a finite number of given points; minimising the norm acts to regularise this interpolation problem; see §6.1.

## 1.3   Pointwise Coordinates and Canonical Coordinates

Aimed at readers already familiar with Hilbert space theory, this section motivates and defines two coordinate systems on a RKHS.

A separable Hilbert space $\mathcal{H}$ possesses an orthonormal basis $e_1, e_2, \cdots \in \mathcal{H}$. An arbitrary element $v \in \mathcal{H}$ can be expressed as an infinite series $v = \sum_{i=0}^{\infty} \alpha_i e_i$ where the "coordinates" $\alpha_i$ are given by $\alpha_i = \langle v, e_i \rangle$. A classic example is using a Fourier series to represent a periodic function. The utility of such a construction is that an arbitrary element of $\mathcal{H}$ can be written as the limit of a linear combination of a manageable set of fixed vectors.

RKHS theory generalises this ability of writing an arbitrary element of a Hilbert space as the limit of a linear combination of a manageable

set of fixed vectors. If $\mathcal{H} \subset \mathbb{R}^T$ is a (not necessarily separable) RKHS then an arbitrary element $v \in \mathcal{H}$ can be expressed as the limit of a sequence $v_1, v_2, \cdots \in \mathcal{H}$ of vectors, each of which is a finite linear combination of the vectors $\{K(\cdot, t) \mid t \in T\}$, where $K \colon T \times T \to \mathbb{R}$ is the kernel of $\mathcal{H}$. It is this ability to represent an arbitrary element of a RKHS $\mathcal{H}$ as the limit of a linear combination of the $K(\cdot, t)$ that, for brevity, we refer to as the presence of a canonical coordinate system. The utility of this canonical coordinate system was hinted at in §1.2.

There is another natural coordinate system: since an element $v$ of a RKHS $\mathcal{H} \subset \mathbb{R}^T$ is a function from $T$ to $\mathbb{R}$, its $t$th coordinate can be thought of as $v(t)$. The relationship between this pointwise coordinate system and the aforementioned canonical coordinates is that $v(t) = \langle v, K(\cdot, t) \rangle$. Note though that whereas an arbitrary linear combination of the $K(\cdot, t)$ is guaranteed to be an element of $\mathcal{H}$, assigning values arbitrarily to the $v(t)$, i.e., writing down an arbitrary function $v$, may not yield an element of $\mathcal{H}$; canonical coordinates are intrinsic whereas pointwise coordinates are extrinsic. The utility of the pointwise coordinate system is that limits in a RKHS can be determined pointwise: if $v_k$ is a Cauchy sequence, implying there exists a $v$ satisfying $\|v_k - v\| \to 0$, then $v$ is fully determined by $v(t) = \lim_k v_k(t)$.

# 2

---

## Finite-dimensional RKHSs

---

Pedagogic material on RKHSs generally considers infinite-dimensional spaces from the outset[1]. Starting with finite-dimensional spaces though is advantageous because the remarkable aspects of RKHS theory are already evident in finite dimensions, where they are clearer to see and easier to study. The infinite-dimensional theory is a conceptually straightforward extension of the finite-dimensional theory.

Although elements of a RKHS must be functions, there is a canonical correspondence between $\mathbb{R}^n$ and real-valued functions on $\{1, 2, \cdots, n\}$ given by the rule that $(x_1, \cdots, x_n) \in \mathbb{R}^n$ is equivalent to the function $f \colon \{1, 2, \cdots, n\} \to \mathbb{R}$ satisfying $f(i) = x_i$ for $i = 1, \cdots, n$. To exemplify, the point $(5, 8, 4) \in \mathbb{R}^3$ is canonically equivalent to the function $f \colon \{1, 2, 3\} \to \mathbb{R}$ given by $f(1) = 5$, $f(2) = 8$ and $f(3) = 4$. For simplicity and clarity, this primer initially works with $\mathbb{R}^n$.

**Remark** The term finite-dimensional RKHS is potentially ambiguous; is merely the RKHS itself finite-dimensional, or must the embedding

---

[1]The tutorial [66] considers finite-dimensional spaces but differs markedly from the presentation here. An aim of [66] is explaining the types of regularisation achievable by minimising a RKHS norm.

space also be finite dimensional? For convenience, we adopt the latter interpretation. While this interpretation is not standard, it is consistent with our emphasis on the embedding space playing a significant role in RKHS theory. Precisely, we say a RKHS is finite dimensional if the set $X$ in §4 has finite cardinality.

## 2.1 The Kernel of an Inner Product Subspace

Central to RKHS theory is the following question. Let $V \subset \mathbb{R}^n$ be endowed with an inner product. How can this configuration be described efficiently? The configuration involves three aspects: the vector space $V$, the orientation of $V$ in $\mathbb{R}^n$, and the inner product on $V$. Importantly, the inner product is not defined on the whole of $\mathbb{R}^n$, unless $V = \mathbb{R}^n$. (An alternative viewpoint that will emerge later is that RKHS theory studies possibly degenerate inner products on $\mathbb{R}^n$, where $V$ represents the largest subspace on which the inner product is not degenerate.)

One way of describing the configuration is by writing down a basis $\{v_1, \cdots, v_r\} \subset V \subset \mathbb{R}^n$ for $V$ and the corresponding Gram matrix $G$ whose $ij$th element is $G_{ij} = \langle v_i, v_j \rangle$. Alternatively, the configuration is completely described by giving an orthonormal basis $\{u_1, \cdots, u_r\} \subset V \subset \mathbb{R}^n$; the Gram matrix associated with an orthonormal basis is the identity matrix and does not need to be stated explicitly. However, these representations do not satisfy the following requirements because there is no unique choice for a basis, even an orthonormal basis.

**One-to-one** The relationship between a subspace $V$ of $\mathbb{R}^n$ and its representation should be one-to-one.

**Respect Topology** If $(V_1, \langle \cdot, \cdot \rangle_1), (V_2, \langle \cdot, \cdot \rangle_2), \cdots$ is a sequence of inner product spaces "converging" to $(V, \langle \cdot, \cdot \rangle)$ then the representations of $(V_1, \langle \cdot, \cdot \rangle_1), (V_2, \langle \cdot, \cdot \rangle_2), \cdots$ should "converge" to the representation of $(V, \langle \cdot, \cdot \rangle)$, and *vice versa*.

**Straightforward Inverse** It should be easy to deduce $V$ and its inner product from its representation.

The lack of a canonical basis for $V$ can be overcome by considering spanning sets instead. Whereas a basis induces a coordinate system, a

spanning set that is not a basis induces an overdetermined coordinate system due to linear dependence.

RKHS theory produces a unique ordered spanning set $k_1, \cdots, k_n$ for $V \subset \mathbb{R}^n$ by the rule that $k_i$ is the unique vector in $V$ satisfying $\langle v, k_i \rangle = e_i^\top v$ for all $v \in V$. Here and throughout, $e_i$ denotes the vector whose elements are zero except for the $i$th which is unity; its dimension should always be clear from the context. In words, taking the inner product with $k_i$ extracts the $i$th element of a vector. This representation looks *ad hoc* yet has proven to be remarkably useful, in part because it unites three different perspectives given by the three distinct but equivalent definitions below.

**Definition 2.1.** Let $V \subset \mathbb{R}^n$ be an inner product space. The *kernel* of $V$ is the unique matrix $K = [k_1 \ k_2 \ \cdots \ k_n] \in \mathbb{R}^{n \times n}$ determined by any of the following three equivalent definitions.

1. $K$ is such that each $k_i$ is in $V$ and $\langle v, k_i \rangle = e_i^\top v$ for all $v \in V$.

2. $K = u_1 u_1^\top + \cdots + u_r u_r^\top$ where $u_1, \cdots, u_r$ is an orthonormal basis for $V$.

3. $K$ is such that the $k_i$ span $V$ and $\langle k_j, k_i \rangle = K_{ij}$.

The third definition is remarkable despite following easily from the first definition because it shows that $K \in \mathbb{R}^{n \times n}$ is the Gram matrix corresponding to the vectors $k_1, \cdots, k_n$. The kernel simultaneously encodes a set of vectors that span $V$ and the corresponding Gram matrix for those vectors!

The name *reproducing kernel* can be attributed either to the kernel $K$ reproducing itself in that its $ij$th element $K_{ij}$ is the inner product of its $j$th and $i$th columns $\langle k_j, k_i \rangle$, or to the vectors $k_i$ reproducing an arbitrary vector $v$ by virtue of $v = (\langle v, k_1 \rangle, \cdots, \langle v, k_n \rangle)$.

**Example 2.1.** Equip $V = \mathbb{R}^n$ with an inner product $\langle u, v \rangle = v^\top Q u$ where $Q$ is symmetric (and positive definite). The equation $\langle v, k_i \rangle = e_i^\top v$ implies $k_i = Q^{-1} e_i$, that is, $K = Q^{-1}$. Alternatively, an eigendecomposition $Q = X D X^\top$ yields an orthonormal basis for $V$ given by the scaled columns of $X$, namely, $\{ X D^{-\frac{1}{2}} e_1, \cdots, X D^{-\frac{1}{2}} e_n \}$. Therefore

$K = \sum_{i=1}^{n}(XD^{-\frac{1}{2}}e_i)(XD^{-\frac{1}{2}}e_i)^\top = XD^{-1}X^\top = Q^{-1}$. Observe that $\langle k_j, k_i \rangle = \langle Q^{-1}e_j, Q^{-1}e_i \rangle = e_i^\top Q^{-1}e_j = e_i^\top K e_j = K_{ij}$.

**Example 2.2.** Let $V \subset \mathbb{R}^2$ be the subspace spanned by the vector $(1,1)$ and endowed with the inner product giving the vector $(1,1)$ unit norm. Then $\{(1,1)\}$ is an orthonormal basis for $V$ and therefore $K = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$.

**Example 2.3.** The configuration having $K = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$ as its kernel can be found as follows. Since $K$ is two-by-two, $V$ must sit inside $\mathbb{R}^2$. Since $V$ is the span of $k_1 = [1\ 0]^\top$ and $k_2 = [0\ 0]^\top$, $V = \mathbb{R} \times \{0\}$. The vector $k_1$ has unit norm in $V$ because $\langle k_1, k_1 \rangle = K_{11} = 1$.

Before proving the three definitions of $K$ are equivalent, several remarks are made. Existence and uniqueness of $K$ follows most easily from definition one because the defining equations are linear. From definition one alone it would seem remarkable that $K$ is always positive semi-definite, denoted $K \geq 0$, yet this fact follows immediately from definition two. Recall that a positive semi-definite matrix is a symmetric matrix whose eigenvalues are greater than or equal to zero. That $K$ is unique is not clear from definition two alone. Definition three gives perhaps the most important characterisation of $K$, yet from definition three alone it is not at all obvious that such a $K$ exists. Definition three also implies that any element of $V$ can be written as a linear combination $K\alpha$ of the columns of $K$, and $\langle K\alpha, K\beta \rangle = \beta^\top K\alpha$. The reader is invited to verify this by writing $K\beta$ as $\beta_1 k_1 + \cdots + \beta_n k_n$.

**Lemma 2.1.** Given an inner product space $V \subset \mathbb{R}^n$, there is precisely one $K = [k_1, \cdots, k_n] \in \mathbb{R}^{n \times n}$ for which each $k_i$ is in $V$ and satisfies $\langle v, k_i \rangle = e_i^\top v$ for all $v \in V$.

*Proof. This is an exercise in abstract linear algebra that is the generalisation of the linear equation $Ax = b$ having a unique solution if $A$ is square and non-singular. Let $v_1, \cdots, v_r$ be a basis for $V$. Linearity implies the constraints on $k_i$ are equivalent to requiring $\langle v_j, k_i \rangle = e_i^\top v_j$ for $j = 1, \cdots, r$. Let $L$ denote the linear operator taking $k \in V$ to*

$(\langle v_1, k \rangle, \langle v_2, k \rangle, \cdots, \langle v_r, k \rangle)$. Both the domain $V$ and the range $\mathbb{R}^r$ of $L$ have dimension $r$. *The linear operator $L$ is equivalent to a square matrix.* Now, $L$ is injective because if $L(k) = L(\tilde{k})$ then $L(k - \tilde{k}) = 0$, that is, $\langle v, k - \tilde{k} \rangle = 0$ for all $v \in V$ (because $v_1, \cdots, v_r$ is a basis), and in particular, $\langle k - \tilde{k}, k - \tilde{k} \rangle = 0$, implying $k = \tilde{k}$. *The linear operator $L$ is non-singular.* Therefore, $L$ is also surjective, and in particular, $L(k_i) = (e_i^\top v_1, \cdots, e_i^\top v_r)$ has one and only one solution. □

Definition three gives a non-linear characterisation of $K$ requiring $k_i$ to extract the $i$th element of the vectors $k_j$ whereas definition one requires $k_i$ to extract the $i$th element of any vector $v$. However, definition three also requires the $k_j$ to span $V$, therefore, by writing an arbitrary vector $v$ as a linear combination of the $v_j$, it becomes evident that definition three satisfies definition one.

**Lemma 2.2.** If $K$ is such that the $k_i$ span $V$ and $\langle k_j, k_i \rangle = K_{ij}$ then $\langle v, k_i \rangle = e_i^\top v$ for all $v \in V$.

*Proof.* Fix a $v \in V$. Since $v$ is in $V$ and the $k_i$ span $V$, there is a vector $\alpha$ such that $v = K\alpha$. Also, $k_i = Ke_i$. Therefore, $\langle v, k_i \rangle = \langle K\alpha, Ke_i \rangle = e_i^\top K\alpha = e_i^\top v$, as required. (Recall from earlier that $\langle k_j, k_i \rangle = K_{ij}$ implies $\langle K\alpha, K\beta \rangle = \beta^\top K\alpha$.) □

The existence of a $K$ satisfying definition three is implied by the existence of a $K$ satisfying definition one.

**Lemma 2.3.** If $K$ is such that each $k_i$ is in $V$ and $\langle v, k_i \rangle = e_i^\top v$ for all $v \in V$ then the $k_i$ span $V$ and $\langle k_j, k_i \rangle = K_{ij}$.

*Proof.* That $\langle k_j, k_i \rangle = K_{ij}$ follows immediately from $K_{ij}$ being the $i$th element of $k_j$, namely, $e_i^\top k_j$. If the $k_i$ do not span $V$ then there is a non-zero $k \in V$ which is orthogonal to each and every $k_i$, that is, $\langle k, k_i \rangle = 0$. Yet this implies $e_i^\top k = 0$ for all $i$, that is, $k = 0$, a contradiction. □

If the columns of $U \in \mathbb{R}^{n \times r}$ form an orthonormal basis for $V \subset \mathbb{R}^n$ then an arbitrary vector in $V$ can be written as $U\alpha$, and moreover, $\langle U\alpha, U\beta \rangle = \beta^\top \alpha$. Referring to definition two, if the $u_i$ are the columns of $U$, then $u_1 u_1^\top + \cdots + u_r u_r^\top = UU^\top$. Armed with these facts, showing definition two satisfies definition one becomes straightforward.

**Lemma 2.4.** If $u_1, \cdots, u_r$ is an orthonormal basis for $V$ then $K = u_1 u_1^\top + \cdots + u_r u_r^\top$ satisfies $\langle v, k_i \rangle = e_i^\top v$ for all $v \in V$.

*Proof.* Let $U = [u_1, \cdots, u_r]$, so that $K = UU^\top$ and $k_i = Ke_i = UU^\top e_i$. An arbitrary $v \in V$ can be represented as $v = U\alpha$. Therefore $\langle v, k_i \rangle = \langle U\alpha, UU^\top e_i \rangle = e_i^\top U\alpha = e_i^\top v$, as required.    $\square$

Since definition one produces a unique $K$, and definitions two and three both produce a $K$ satisfying definition one, the equivalence of the three definitions has been proven.

The next pleasant surprise is that for every positive semi-definite matrix $K \in \mathbb{R}^{n \times n}$ there is an inner product space $V \subset \mathbb{R}^n$ whose kernel is $K$. The proof relies on the following.

**Lemma 2.5.** If $K \in \mathbb{R}^{n \times n}$ is positive semi-definite then $\alpha^\top K \alpha = 0$ implies $K\alpha = 0$, where $\alpha$ is a vector in $\mathbb{R}^n$.

*Proof.* Eigendecompose $K$ as $K = UDU^\top$. Let $\beta = U^\top \alpha$. Then $\alpha^\top K \alpha = 0$ implies $\beta^\top D \beta = 0$. Since $D \geq 0$, $D\beta$ must be zero. Therefore $K\alpha = UDU^\top \alpha = UD\beta = 0$.    $\square$

**Lemma 2.6.** Let $V = \mathrm{span}\{k_1, \cdots, k_n\}$ be the space spanned by the columns $k_1, \cdots, k_n$ of a positive semi-definite matrix $K \in \mathbb{R}^{n \times n}$. There exists an inner product on $V$ satisfying $\langle k_j, k_i \rangle = K_{ij}$.

*Proof.* Let $u, v \in V$. Then there exist vectors $\alpha, \beta$ such that $u = K\alpha$ and $v = K\beta$. Define $\langle u, v \rangle$ to be $\beta^\top K \alpha$. To show this is well-defined, let $u = K\tilde{\alpha}$ and $v = K\tilde{\beta}$ be possibly different representations. Then $\beta^\top K \alpha - \tilde{\beta}^\top K \tilde{\alpha} = (\beta - \tilde{\beta})^\top K \alpha + \tilde{\beta}^\top K (\alpha - \tilde{\alpha}) = 0$ because $K\alpha = K\tilde{\alpha}$, $K\beta = K\tilde{\beta}$ and $K = K^\top$. Clearly $\langle \cdot, \cdot \rangle$ so defined is bilinear. To prove $\langle \cdot, \cdot \rangle$ is positive definite, assume $\langle K\alpha, K\alpha \rangle = \alpha^\top K \alpha = 0$. By Lemma 2.5 this implies $K\alpha = 0$, as required.    $\square$

Given $K$, there is a *unique* inner product space whose kernel is $K$.

**Lemma 2.7.** Let $V_1 \subset \mathbb{R}^n$ and $V_2 \subset \mathbb{R}^n$ be two inner product spaces having the same kernel $K$. Then $V_1$ and $V_2$ are identical spaces: $V_1 = V_2$ and their inner products are the same.

*Proof.* The columns of $K$ span both $V_1$ and $V_2$, hence $V_1 = V_2$. For the same reason, the inner products on $V_1$ and $V_2$ are uniquely determined from the Gram matrix $K$ corresponding to $k_1, \cdots, k_n$. Since $V_1$ and $V_2$ have the same Gram matrix, their inner products are identical too. (Indeed, the inner product must be given by $\langle K\alpha, K\beta \rangle = \beta^\top K\alpha$.) $\square$

To summarise, for a fixed $n$, there is a bijective correspondence between inner product spaces $V \subset \mathbb{R}^n$ and positive semi-definite matrices $K \in \mathbb{R}^{n \times n}$. Given $V \subset \mathbb{R}^n$ there is precisely one kernel $K \geq 0$ (Definition 2.1 and Lemma 2.1). Given a $K \geq 0$, there is precisely one inner product space $V \subset \mathbb{R}^n$ for which $K$ is its kernel (Lemmata 2.6 and 2.7). Recalling the criteria listed earlier, this correspondence is One-to-one and has a Straightforward Inverse. That it Respects Topology is discussed next.

## 2.2 Sequences of Inner Product Spaces

A sequence of one-dimensional vector spaces in $\mathbb{R}^2$ can be visualised as a collection of lines passing through the origin and numbered $1, 2, \cdots$. It can be recognised visually when such a sequence converges. (This corresponds to the topology of the Grassmann manifold.) The kernel representation of an inner product space induces a concept of convergence for the broader situation of a sequence of *inner product* subspaces, possibly of *differing dimensions*, by declaring that the limit $V_\infty \subset \mathbb{R}^n$ of a sequence of inner product spaces $V_1, V_2, \cdots \subset \mathbb{R}^n$ is the space whose kernel is $K_\infty = \lim_{n \to \infty} K_n$, if the limit exists. This makes the kernel representation of a subspace Respect Topology. The pertinent question is whether the induced topology is sensible in practice.

**Example 2.4.** Define on $\mathbb{R}^2$ a sequence of inner products $\langle u, v \rangle_n = v^\top Q_n u$ where $Q_n = \text{diag}\{1, n^2\}$ is a diagonal matrix with entries 1 and $n^2$. An orthonormal basis for the $n$th space is $\{(1, 0), (0, n^{-1})\}$. It seems acceptable, at least in certain situations, to agree that the limit of this sequence of two-dimensional spaces is the one-dimensional subspace of $\mathbb{R}^2$ spanned by $(1, 0)$. This accords with defining convergence via kernels. Let $K_n$ be the kernel of the $n$th space: $K_n = \text{diag}\{1, n^{-2}\}$.

Then $K_\infty = \lim_{n\to\infty} K_n = \operatorname{diag}\{1,0\}$. The space $V_\infty$ having kernel $K_\infty$ is the subspace of $\mathbb{R}^2$ spanned by $(1,0)$, where the inner product is such that $(1,0)$ has unit norm.

The example above illustrates a general principle: since $K$ can be defined in terms of an orthonormal basis as $K = u_1 u_1^\top + \cdots + u_r u_r^\top$, a sequence of higher-dimensional inner product spaces can converge to a lower-dimensional inner product space if one or more of the orthonormal basis vectors approaches the zero vector. This need not be the only sensible topology, but it is a topology that finds practical use, as later verified by examples.

**Remark**   Recall from Example 2.1 that if $V = \mathbb{R}^n$ then $K = Q^{-1}$. If $Q_n$ is a sequence of positive definite matrices becoming degenerate in the limit, meaning one or more eigenvalues goes to infinity, then $Q_n$ does not have a limit. However, the corresponding eigenvalues of $K_n = Q_n^{-1}$ go to zero, hence the kernel $K_n$ may have a well-defined limit. In this sense, RKHS theory encompasses degenerate inner products. A way of visualising this is presented in §2.4.

## 2.3   Extrinsic Geometry and Interpolation

Finite-dimensional RKHS theory studies subspaces $V$ of $\mathbb{R}^n$ rather than abstract vector spaces. If it is only known that $Z$ is an abstract vector space then there is no way of knowing what the elements of $Z$ look like. The best that can be done is assert the existence of a basis $\{b_1, \cdots, b_r\}$. By comparison, knowing $V \subset \mathbb{R}^n$ allows working with $V$ using *extrinsic coordinates* by writing an element of $V$ as a vector in $\mathbb{R}^n$.

Writing $V \subset \mathbb{R}^n$ may also signify the importance of the orientation of $V$ inside $\mathbb{R}^n$. If $V$ and $W$ are $r$-dimensional linear subspaces of $\mathbb{R}^n$ then their *intrinsic geometry* is the same but their *extrinsic geometry* may differ; $V$ and $W$ are equivalent (precisely, isomorphic) as vector spaces but they lie inside $\mathbb{R}^n$ differently unless $V = W$.

The usefulness of extrinsic geometry is exemplified by considering the interpolation problem of finding a vector $x \in V \subset \mathbb{R}^n$ of smallest norm and some of whose coordinates $e_i^\top x$ are specified. This problem is

posed using extrinsic geometry. It can be rewritten intrinsically, without reference to $\mathbb{R}^n$, by using the linear operator $L_i \colon V \to \mathbb{R}$ defined to be the restriction of $x \mapsto e_i^\top x$. However, the extrinsic formulation is the more amenable to a unified approach.

**Example 2.5.** Endow $V \subset \mathbb{R}^n$ with an inner product. Fix an $i$ and consider how to find $x \in V$ satisfying $e_i^\top x = 1$ and having the smallest norm. Geometrically, such an $x$ must be orthogonal to any vector $v \in V$ satisfying $e_i^\top v = 0$, for otherwise its norm could be decreased. This approach allows $x$ to be found by solving a system of linear equations, however, going further seems to require choosing a basis for $V$.

RKHS theory replaces an *ad hoc* choice of basis for $V \subset \mathbb{R}^n$ by the particular choice $k_1, \cdots, k_n$ of spanning vectors for $V$. Because the kernel representation Respects Topology, the vectors $k_1, \cdots, k_n$ vary continuously with changing $V$. The solution to the interpolation problem should also vary continuously as $V$ changes. There is thus a chance that the solution can be written elegantly in terms of $k_1, \cdots, k_n$.

**Example 2.6.** Continuing the example above, let $K = [k_1, \cdots, k_n]$ be the kernel of $V$. Let $x = K\alpha$ and $v = K\beta$. The constraints $e_i^\top x = 1$ and $\langle x, v \rangle = 0$ whenever $e_i^\top v = 0$ become $k_i^\top \alpha = 1$ and $\beta^\top K \alpha = 0$ whenever $\beta^\top K e_i = 0$. (Recall that $i$ is fixed.) The latter requirement is satisfied by $\alpha = c e_i$ where $c \in \mathbb{R}$ is a scalar. Solving $k_i^\top \alpha = 1$ implies $c = K_{ii}^{-1}$. Therefore, $x = K_{ii}^{-1} k_i$. Note $\|x\|^2 = \langle x, x \rangle = K_{ii}^{-1}$.

In the above example, as $V$ changes, both the kernel $K$ and the solution $x$ change. Yet the relationship between $x$ and $K$ remains constant: $x = K_{ii}^{-1} k_i$. This is further evidence that the representation $K$ of the extrinsic geometry is a useful representation.

There is a geometric explanation for the columns of $K$ solving the single-point interpolation problem. Let $L_i \colon V \to \mathbb{R}$ denote the $i$th coordinate function: $L_i(v) = e_i^\top v$. That $\langle v, k_i \rangle = L_i(v)$ means $k_i$ is the gradient of $L_i$. In particular, the line determined by $k_i$ meets the level set $\{v \mid L_i(v) = 1\}$ at right angles, showing that $k_i$ meets the orthogonality conditions for optimality. Turning this around leads to yet another definition of $K$.

**Lemma 2.8.** The following is a geometric definition of the $k_i$ that is equivalent to Definition 2.1. Let $H_i = \{z \in \mathbb{R}^n \mid e_i^\top z = 1\}$ be the hyperplane consisting of all vectors whose $i$th coordinate is unity. If $V \cap H_i$ is empty then define $k_i = 0$. Otherwise, let $\tilde{k}_i$ be the point in the intersection $V \cap H_i$ that is closest to the origin. Define $k_i$ to be $k_i = c\tilde{k}_i$ where $c = \langle \tilde{k}_i, \tilde{k}_i \rangle^{-1}$.

*Proof.* Since $K$ is unique (Lemma 2.1), it suffices to prove the $k_i$ defined here satisfy definition one of Definition 2.1. If $V \cap H_i$ is empty then $e_i^\top v = 0$ for all $v \in V$ and $k_i = 0$ satisfies definition one. Assume then that $V \cap H_i$ is non-empty. It has a closest point, $\tilde{k}_i$, to the origin. As $e_i^\top \tilde{k}_i = 1$, $\tilde{k}_i$ is non-zero and $c$ is well-defined. Also, $\langle w, \tilde{k}_i \rangle = 0$ whenever $w \in V$ satisfies $e_i^\top w = 0$, for otherwise $\tilde{k}_i$ would not have the smallest norm. Let $v \in V$ be arbitrary and define $w = v - a\tilde{k}_i$ where $a = e_i^\top v$. Then $e_i^\top w = e_i^\top v - a = 0$. Thus $\langle v, k_i \rangle = \langle w + a\tilde{k}_i, c\tilde{k}_i \rangle = ac\langle \tilde{k}_i, \tilde{k}_i \rangle = a = e_i^\top v$, as required.                                     $\square$

The $c$ in the lemma scales $\tilde{k}_i$ so that $\langle k_i, k_i \rangle = e_i^\top k_i$. It is also clear that if $k_i \neq 0$ then $K_{ii} \neq 0$, or in other words, if $K_{ii} = 0$ then the $i$th coordinate of every vector $v$ in $V$ is zero and the interpolation problem has no solution.

## 2.4   Visualising RKHSs

It is more expedient to understand RKHS theory by focusing on the columns $k_i$ of the kernel $K$ rather than on the matrix $K$ itself; $K$ being a positive semi-definite matrix corresponding to the Gram matrix of the $k_i$ is a wonderful bonus. The indexed set of vectors $k_1, \cdots, k_n$ changes in a desirable way as $V \subset \mathbb{R}^n$ changes. This can be visualised explicitly for low-dimensional examples.

The inner product on a two-dimensional vector space $V$ can be depicted by drawing the ellipse $\{v \in V \mid \langle v, v \rangle = 1\}$ on $V$ because the inner product is uniquely determined from the norm, and the ellipse uniquely determines the norm. Based on Lemma 2.8, the columns $k_1$ and $k_2$ of the kernel $K$ corresponding to $V \subset \mathbb{R}^2$ can be found geometrically, as explained in the caption of Figure 2.1.

Figures 2.2 and 2.3 reveal how $k_1$ and $k_2$ vary as the inner product on $V = \mathbb{R}^2$ changes. Figure 2.3 shows that $k_1$ and $k_2$ vary smoothly as the inner product changes. By comparison, an orthonormal basis formed from the principal and minor axes of the ellipse must have a discontinuity somewhere, because after a 180 degree rotation, the ellipse returns to its original shape yet, for example, $\{e_1, \frac{1}{2}e_2\}$ rotated 180 degrees is $\{-e_1, -\frac{1}{2}e_2\}$, which differs from $\{e_1, \frac{1}{2}e_2\}$.

Figure 2.4 illustrates the kernel of various one-dimensional subspaces in $\mathbb{R}^2$. As $V$ is one-dimensional, $k_1$ and $k_2$ must be linearly dependent since they span $V$. Together, $k_1$ and $k_2$ describe not only $V$ and its inner product, but also the orientation of $V$ in $\mathbb{R}^2$.

Figure 2.5 portrays how a sequence of two-dimensional subspaces can converge to a one-dimensional subspace. From the perspective of the inner product, the convergence is visualised by an ellipse degenerating to a line segment. From the perspective of the kernel $K = [k_1 \; k_2]$, the convergence is visualised by $k_1$ and $k_2$ individually converging to vectors lying in the same subspace.

The figures convey the message that the $k_i$ are a spanning set for $V$ that vary continuously with changing $V$. Precisely how they change is of more algebraic than geometric importance; the property $\langle k_j, k_i \rangle = K_{ij}$ is very convenient to work with algebraically.

## 2.5 RKHSs over the Complex Field

RKHS theory extends naturally to complex-valued vector spaces. In finite dimensions, this means considering subspaces $V \subset \mathbb{C}^n$ where $V$ is equipped with a complex-valued inner product $\langle \cdot, \cdot \rangle \colon V \times V \to \mathbb{C}$. The only change to Definition 2.1 is replacing $K = u_1 u_1^\top + \cdots + u_r u_r^\top$ by $K = u_1 u_1^H + \cdots + u_r u_r^H$ where $H$ denotes Hermitian transpose.

The kernel $K$ will always be Hermitian ($K = K^H$) and positive semi-definite. Swapping the order of an inner product introduces a conjugation: $\langle u, v \rangle = \overline{\langle v, u \rangle}$. Beyond such minor details, the real-valued and complex-valued theories are essentially the same.

**Figure 2.1:** The ellipse comprises all points one unit from the origin. It determines the chosen inner product on $\mathbb{R}^2$. The vector $\tilde{k}_1$ is the closest point to the origin on the vertical line $x = 1$. It can be found by enlarging the ellipse until it first touches the line $x = 1$, or equivalently, as illustrated, it can be found by shifting the line $x = 1$ horizontally until it meets the ellipse tangentially, represented by the dashed vertical line, then travelling radially outwards from the point of intersection until reaching the line $x = 1$. The vector $k_1$ is a scaled version of $\tilde{k}_1$. If the dashed vertical line intersects the $x$-axis at $c$ then $k_1 = c^2 \tilde{k}_1$. Equivalently, $k_1$ is such that its tip intersects the line $x = c^2$. The determination of $k_2$ is analogous but with respect to the horizontal line $y = 1$.

**Figure 2.2:** Shown are the vectors $k_1$ (red) and $k_2$ (blue) corresponding to rotated versions of the inner product $\langle u, v \rangle = v^\top Q u$ where $Q = \text{diag}\{1, 4\}$. The magnitude and angle of $k_1$ and $k_2$ are plotted in Figure 2.3.

**Figure 2.3:** Plotted are the magnitude and angle of $k_1$ (red) and $k_2$ (blue) corresponding to rotated versions of the inner product $\langle u, v \rangle = v^\top Q u$ where $Q = \text{diag}\{1, 4\}$, as in Figure 2.2.

**Figure 2.4:** Shown are $k_1$ (red) and $k_2$ (blue) for various one-dimensional subspaces (black) of $\mathbb{R}^2$. In all cases, the inner product is the standard Euclidean inner product. Although not shown, $k_2$ is zero when $V$ is horizontal, and $k_1$ is zero when $V$ is vertical. Therefore, the magnitude of the red vectors increases from zero to a maximum then decreases back to zero. The same occurs for the blue vectors.

**Figure 2.5:** Illustration of how, as the ellipse gets narrower, the two-dimensional inner product space $V = \mathbb{R}^2$ converges to a one-dimensional inner product space. The kernels of the subspaces are represented by red ($k_1$) and blue ($k_2$) vectors.

# 3

---

## Function Spaces

---

Typifying infinite-dimensional vector spaces are function spaces. Given an arbitrary set $X$, let $\mathbb{R}^X = \{f\colon X \to \mathbb{R}\}$ denote the set of all functions from $X$ to $\mathbb{R}$. It is usually given a vector space structure whose simplicity belies its usefulness: vector-space operations are defined *pointwise*, meaning scalar multiplication $\alpha \cdot f$ sends the function $f$ to the function $g$ given by $g(x) = \alpha\,f(x)$, while vector addition sends $f + g$ to the function $h$ defined by $h(x) = f(x) + g(x)$.

The space $\mathbb{R}^X$ is often too large to be useful on its own, but it contains useful subspaces, such as the spaces of all continuous, smooth or analytic functions. (This necessitates $X$ having an appropriate structure; for example, it does not make sense for a function $f\colon X \to \mathbb{R}$ to be continuous unless $X$ is a topological space.) Note that to be a subspace, as opposed to merely a subset, the restrictions placed on the functions must be closed with respect to vector-space operations: if $f$ and $g$ are in the subset and $\alpha$ is an arbitrary scalar then $f + g$ and $\alpha\,f$ must also be in the subset. While the set of all continuous functions $f\colon [0,1] \to \mathbb{R}$ on the interval $[0,1] \subset \mathbb{R}$ is a vector subspace of $\mathbb{R}^{[0,1]}$, the set of all functions $f\colon [0,1] \to \mathbb{R}$ satisfying $f(1) \leq 1$ is not.

## 3.1   Function Approximation

Function approximation is an insightful example of how function spaces facilitate geometric reasoning for solving algebraic problems [41].

**Example 3.1.** Denote by $C_2[0,1]$ the space of continuous functions $f\colon [0,1] \to \mathbb{R}$ equipped with the inner product $\langle f, g \rangle = \int_0^1 f(x)g(x)\,dx$. The square of the induced norm is $\|f\|^2 = \langle f, f \rangle = \int_0^1 [f(x)]^2\,dx$. Let $V$ denote the subspace of $C_2[0,1]$ spanned by the functions $g_0(x) = 1$, $g_1(x) = x$ and $g_2(x) = x^2$. In words, $V$ is the space of all polynomials of degree at most two. Fix an $f$, say, $f(x) = e^x - 1$. Consider finding $g \in V$ that minimises $\|f - g\|^2$. Figure 3.1 visualises this in two ways. Graphing $f$ does not help in finding $g$, whereas treating $f$ and $g$ as elements of a function space suggests $g$ can be found by solving the linear equations $\langle f - g, v \rangle = 0$ for all $v \in V$.

   An orthogonal basis for $V$ is $p_0(x) = 1$, $p_1(x) = 2x - 1$ and $p_2(x) = 6x^2 - 6x + 1$. (These are the shifted Legendre polynomials.) Note $\|p_i\|^2 = (2i + 1)^{-1}$. Write $g$ as $g = \alpha_0 p_0 + \alpha_1 p_1 + \alpha_2 p_2$. Attempting to solve $\langle f - g, v \rangle = 0$ for $v \in \{p_0, p_1, p_2\}$ leads to the equations $\langle \alpha_i p_i, p_i \rangle = \langle f, p_i \rangle$. Since $\langle f, p_0 \rangle = e - 2$, $\langle f, p_1 \rangle = 3 - e$ and $\langle f, p_2 \rangle = 7e - 19$, the coefficients $\alpha_i$ are found to be $\alpha_0 = e - 2$, $\alpha_1 = 9 - 3e$ and $\alpha_2 = 35e - 95$.

   The simple optimisation problem in Example 3.1 could have been solved using calculus by setting to zero the derivatives of $\|f - \sum_i \alpha_i p_i\|^2$ with respect to the $\alpha_i$. Importantly though, setting the derivatives to zero in general does not guarantee finding the optimal solution because the optimal solution need not exist[1]. When applicable, vector space methods go beyond calculus in that they can guarantee the optimal solution has been found. (Calculus examines local properties whereas being a global optimum is a global property.)

   Infinite-dimensional inner product spaces can exhibit behaviour not found in finite-dimensional spaces. The projection of $f$ onto $V$ shown in Figure 3.1 suggesting every minimum-norm problem has a solution is not necessarily true in infinite dimensions, as Example 3.2 will exemplify. Conditions must be imposed before the geometric picture in

---

[1]See [70] for a historical account relating to existence of optimal solutions.
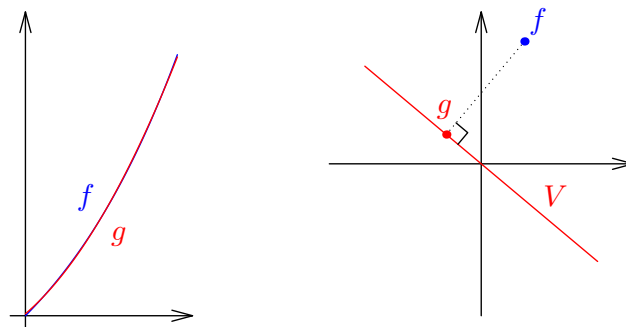
**Figure 3.1:** On the left is the direct way of visualising the function approximation problem: given a function $f$ (blue), find a function $g$ (red) that best approximates $f$, where $g$ is constrained to a subclass of functions. Shown here is the second-order polynomial approximation of $f(x) = e^x - 1$ found in Example 3.1. It approximates $f$ well, in that the graphs of $f$ and $g$ overlap. On the right is a geometric visualisation of the function approximation problem, made possible by representing functions not by their graphs (left) but as elements of a vector space (right). The subspace $V$ is the subclass of functions in which $g$ must lie. Minimising the square of the norm $\|f - g\|$ means finding the point on $V$ that is closest to $f$. As the norm comes from an inner product, the point $g$ must be such that $f - g$ is perpendicular to every function in $V$.

Figure 3.1 accurately describes the infinite-dimensional function approximation problem.

When coming to terms with infinite-dimensional spaces, one of the simplest (Hilbert) spaces to contemplate is $l^2$, the space of square-summable sequences. Elements of $l^2$ are those sequences of real numbers for which the sum of the squares of each term of the sequence is finite. In symbols, $x = (x_1, x_2, \cdots)$ is in $l^2$ if and only if $\sum_{i=1}^{\infty} x_i^2 < \infty$. For example, $(1, \frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \cdots)$ is in $l^2$. The inner product on $l^2$ is implicitly taken to be $\langle x, y \rangle = \sum_{i=1}^{\infty} x_i y_i$. Note that the square of the induced norm is $\|x\|^2 = \sum_{i=1}^{\infty} x_i^2$ and hence the condition for $x$ to be in $l^2$ is precisely the condition that it have finite $l^2$-norm.

Approximating an element of $l^2$ by an element of a subspace $V \subset l^2$ can be considered a function approximation problem because there is a canonical correspondence between $l^2$ and a subclass of functions from $\{1, 2, \cdots\}$ to $\mathbb{R}$, just as there is a canonical correspondence between $R^n$ and functions from $\{1, \cdots, n\}$ to $\mathbb{R}$. (The subclass comprises those

functions $f$ for which $\sum_{i=1}^{\infty} f(i)^2 < \infty$.)

**Example 3.2.** Let $V \subset l^2$ be the vector space of sequences with only a finite number of non-zero terms, thought of as a subspace of $l^2$. There is no element $g$ of $V$ that minimises $\|f - g\|^2$ when $f = (1, \frac{1}{2}, \frac{1}{4}, \cdots) \in l^2$. Given any approximation $g$ of $f$, it can always be improved by changing a zero term in $g$ to make it match the corresponding term in $f$. A rigorous proof goes as follows. Assume to the contrary there exists a $g \in V$ such that $\|f - \tilde{g}\| \geq \|f - g\|$ for all $\tilde{g} \in V$. Let $i$ be the smallest integer such that $g_i = 0$. It exists because only a finite number of terms of $g$ are non-zero. Let $\tilde{g}$ be a replica of $g$ except for setting its $i$th term to $\tilde{g}_i = f_i$. Let $e$ be the zero sequence except for the $i$th term which is unity. Then $f - g = (f - \tilde{g}) + f_i e$. Since $\langle f - \tilde{g}, e \rangle = 0$, $\|f - g\|^2 = \|f - \tilde{g}\|^2 + f_i^2 \|e\|^2$, implying $\|f - \tilde{g}\| < \|f - g\|$ and contradicting the existence of an optimal element $g$.

Reconciling the *algebra* in Example 3.2 proving the non-existence of an optimal approximation, with the *geometry* depicted in Figure 3.1 suggesting otherwise, requires understanding the interplay between vector space structures and topological structures induced by the norm.

Before delving into that topic, it is penetrating to ascertain whether RKHS theory may be beneficially applied to the function approximation problem. A novel feature of this primer is considering RKHS theory from the "dynamic" perspective of how $K$ changes in response to changing $V$. The solution of the interpolation problem in Examples 2.5 and 2.6 is continuous in $K$, and in particular, a solution can be found even for a lower-dimensional $V \subset \mathbb{R}^n$ by working in $\mathbb{R}^n$ and taking limits as the inner product becomes degenerate, much like in Figure 2.5. The function approximation problem lacks this property because enlarging $V$ to include $f$ will mean the optimal solution is $g = f$ no matter how close the inner product is to being degenerate. Since the solution of the function approximation problem does not Respect Topology it is unlikely RKHS theory will be of assistance. (It is irrelevant that the space used in Example 3.1 is not a RKHS; the function approximation problem could have been posed equally well in a genuine RKHS.)

No theory is a panacea for all problems. RKHS theory has broad but not universal applicability.

## 3.2   Topological Aspects of Inner Product Spaces

Reproducing kernel Hilbert spaces have additional properties that more general Hilbert spaces and other inner product spaces do not necessarily enjoy. These properties were chosen to make RKHSs behave more like finite-dimensional spaces. Learning what can "go wrong" in infinite dimensions helps put RKHS theory in perspective.

Some subspaces should not be drawn as sharp subspaces, as in Figure 3.1, but as blurred subspaces denoting the existence of points not lying in the subspace yet no more than zero distance away from it. The blur represents an infinitesimally small region as measured by the norm, but need not be small in terms of the vector space structure. Every point in $l^2$ is contained in the blur around the $V$ in Example 3.2 despite $V$ being a very small subspace of $l^2$ in terms of cardinality of basis vectors. (A Hamel basis for $l^2$ is uncountable while a Hamel basis for $V$ is countable. Note that in an infinite-dimensional Hilbert space, every Hamel basis is uncountable.)

Since the blur is there to represent the geometry coming from the norm, the blur should be drawn as an infinitesimally small region around $V$. An accurate way of depicting $V \subset l^2$ in Example 3.2 is by using a horizontal plane for $V$ and adding a blur of infinitesimal height to capture the rest of $l^2$. This visually implies correctly that there is no vector in $l^2$ that is orthogonal to every vector in $V$. Any $f \in l^2$ that is not in $V$ makes an infinitesimally small angle with $V$, for otherwise the vectors $f, 2f, 3f, \cdots$ would move further and further away from $V$.

As these facts are likely to generate more questions than answers, it is prudent to return to the beginning and work slowly towards deriving these facts from first principles.

An inner product induces a norm, and a norm induces a topology. It suffices here to understand a topology as a rule for determining which sequences converge to which points. A norm coming from an inner product determines the inner product uniquely. It is pragmatic to think of the norm as the more dominant structure. The presence of an inner product means the square of the norm is "quadratic", a very convenient property.

Inner product spaces therefore have two structures: a vector space structure and a norm whose square is "quadratic". The axioms of a normed space ensure these two structures are compatible. For example, the triangle inequality $\|x + y\| \leq \|x\| + \|y\|$ involves both the norm and vector addition. Nevertheless, the two structures capture different aspects that are conflated in finite dimensions.

A norm induces a topology, and in particular, a norm determines what sequences have what limit points. Topology is a weaker concept than that of norm because different norms may give rise to the same topology. In fact, in *finite* dimensions, every norm gives rise to the same topology as every other norm: if $x_k \to x$ with respect to one norm $\|\cdot\|$, meaning $\|x_k - x\| \to 0$, then $x_k \to x$ with respect to any other norm.

The first difference then in infinite dimensions is that different norms can induce different topologies.

**Example 3.3.** Let $\|\cdot\|_2$ denote the standard norm on $l^2$ and let $\|x\|_\infty = \sup_i |x_i|$ be an alternative norm. Let $x^{(i)}$ denote the sequence $(\frac{1}{i}, \cdots, \frac{1}{i}, 0, 0, \cdots)$ whose first $i^2$ terms are $\frac{1}{i}$ and the rest zero. Then $x^{(1)}, x^{(2)}, \cdots$ is a sequence of elements of $l^2$. Since $\|x^{(i)}\|_\infty = \frac{1}{i} \to 0$, the sequence converges to the zero vector with respect to the norm $\|\cdot\|_\infty$. However, $\|x^{(i)}\|_2 = 1$ does not converge to zero and hence $x^{(i)}$ does not converge to the zero vector with respect to the norm $\|\cdot\|_2$. The two norms induce different topologies.

A second difference is that an infinite dimensional subspace need not be topologically *closed*. A subspace $V$ of a normed space $W$ is closed in $W$ if, for any convergent sequence $v_1, v_2, \cdots \to w$, where the $v_i$ are in $V$ but $w$ need only be in $W$, it is nevertheless the case that $w$ is in $V$. (Simply put, $V$ is closed if every limit point of every sequence in $V$ is also in $V$.) The subspace $V$ in Example 3.2 is not closed because the sequence $(1, 0, \cdots), (1, \frac{1}{2}, 0, \cdots), (1, \frac{1}{2}, \frac{1}{4}, 0, \cdots)$ is in $V$ but its limit $(1, \frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \cdots) \in l^2$ is not in $V$.

A closed subspace can be drawn as a sharp subspace but any other subspace should have an infinitesimal blur around it denoting the existence of elements not in the subspace but which are infinitesimally close to it, a consequence of $v_1, v_2, \cdots \to w$ meaning $\|v_i - w\| \to 0$.

In finite dimensions, if $V \subset \mathbb{R}^n$ and $v_1, v_2, \cdots$ is a sequence of

elements of $V$ converging to a point $v \in \mathbb{R}^n$, it is visually clear from a diagram that the limit point $v$ must lie in $V$.

**Lemma 3.1.** Every finite-dimensional subspace $V \subset W$ of an inner product space $W$ is closed.

*Proof.* Let $v_1, v_2, \cdots \to w$ be a convergent sequence, with $v_i$ in $V$ and $w$ in $W$. Let $u_1, \cdots, u_r$ be an orthonormal basis for $V$. Let $\tilde{w} = w - \langle w, u_1 \rangle u_1 - \cdots - \langle w, u_r \rangle u_r$. *Here, $w - \tilde{w}$ is the projection of $w$ onto $V$. In finite dimensions, the projection is always well-defined.* Then $\tilde{w}$ is orthogonal to $u_1, \cdots, u_r$. By Pythagoras' theorem, $\|\tilde{w}\| \leq \|w - v_i\|$. By definition, $v_i \to w$ implies $\|w - v_i\| \to 0$, therefore, $\|\tilde{w}\|$ must equal zero, implying $w$ lies in $V$. □

A key ingredient in the above proof is that a vector $w \notin V$ can be projected onto $V$ and hence is a non-zero distance away from $V$. In Example 3.2 though, there is no orthogonal projection of $f$ onto $V$.

**Lemma 3.2.** Define $f$ and $V$ as in Example 3.2. There is no vector $\tilde{f} \in V$ such that $\langle f - \tilde{f}, v \rangle = 0$ for all $v \in V$.

*Proof.* Let $\tilde{f} \in V$ be arbitrary. Let $i$ be the smallest integer such that the $i$th term of $\tilde{f}$ is zero. Then the $i$th term of $f - \tilde{f}$ is non-zero and the inner product of $f - \tilde{f}$ with the element $(0, \cdots, 0, 1, 0, \cdots, 0) \in V$ having 1 for its $i$th term, is non-zero. □

The space $l^2$, as a vector space, is larger than it looks. The sequences $x^{(i)} = (0, \cdots, 0, 1, 0, \cdots, 0)$, with the unit appearing in the $i$th term, do not form a (Hamel) basis for $l^2$. The $x^{(i)}$ span $V$ in Example 3.2 but do not come close to spanning $l^2$. The $x^{(i)}$ are countable in number whereas a basis for $l^2$ must necessarily be uncountable. This stems from the span of a set of vectors being the vectors representable by a *finite* linear combination of vectors in the set. Elements in $l^2$ having infinitely many non-zero terms cannot be written as finite linear combinations of the $x^{(i)}$. (Since infinite summations involve taking limits, only finite linear combinations are available in spaces with only a vector space structure. A prime mathematical reason for introducing a topology, or *a fortiori* a norm, is to permit the use of infinite summations.)

From the norm's perspective though, $l^2$ is not large: any point in $l^2$ is reachable as a limit of a sequence of points in the span of the $x^{(i)}$. Precisely, $(f_1, f_2, \cdots) \in l^2$ is the limit of $(f_1, 0, \cdots), (f_1, f_2, 0, \cdots), \cdots$. This shows *the norm is coarser than the vector space structure.* The vector space structure in Example 3.2 places $f$ away from $V$ yet the norm places $f$ no more than zero distance from $V$.

To demonstrate the phenomenon in Example 3.2 is not esoteric, the following is a more familiar manifestation of it. Once again, the chosen subspace $V$ is not closed.

**Example 3.4.** Define $C_2[0, 1]$ as in Example 3.1, but this time take $V$ to be the subspace of all polynomials on the unit interval $[0, 1]$. Then there is no polynomial $g \in V$ that is closest to $f(x) = \sin(x)$ because $f$ can be approximated arbitrarily accurately on $[0, 1]$ by including more terms of the Taylor series approximation $\sin(x) = 1 - \frac{x^3}{3!} + \frac{x^5}{5!} - \cdots$.

To summarise, it is sometimes possible to extend an inner product on an infinite dimensional space $V$ to an inner product on $V \oplus \mathrm{span}\{f\}$ so that there is a sequence in $V$ converging to $f$. Here, $f$ is a new vector being added to form a larger space. The affine spaces $cf + V$ for $c \in \mathbb{R}$ would then be infinitesimally close to each other, crammed one on top of the other. Similarly, $l^2$ is obtained from the $V$ in Example 3.2 by adding basis vectors that are all crammed flat against $V$, hence the earlier suggestion of visualising $l^2$ as an infinitesimal blur around $V$.

## 3.3 Evaluation Functionals

Evaluation functionals play a crucial role in RKHS theory. In finite dimensions, they went unnoticed as the linear functionals $v \mapsto e_i^\top v$.

Let $V$ be a subspace of $\mathbb{R}^X$ where $X$ is an arbitrary set. An element $f$ of $V$ is therefore a function $f \colon X \to \mathbb{R}$ and can be evaluated pointwise. This means mathematically that for any $x \in X$, there is a function $l_x \colon V \to \mathbb{R}$ given by $l_x(f) = f(x)$. To become familiar with this definition, the reader is invited to verify $l_x$ is linear.

If there is a norm on $V$ then it can be asked whether the norm "controls" the pointwise values of elements of $V$. The following examples are intended to convey the concept prior to a rigorous definition.

**Example 3.5.** Define $C_2[0,1]$ as in Example 3.1. Let $f(x)$ be the tent function given by $f(x) = 8x + 4$ for $x \in [-\frac{1}{2}, 0]$, $f(x) = 4 - 8x$ for $x \in [0, \frac{1}{2}]$ and $f(x) = 0$ otherwise. Let $f_i(x) = i\,f\left(i^2(x - \frac{1}{2})\right)$ for $x \in [0,1]$. Then $\|f_i\| = \frac{1}{i}$ decreases to zero yet $f_i(\frac{1}{2})$ diverges to infinity. No matter how small the norm is, there are elements of $C_2[0,1]$ taking on arbitrarily large values pointwise.

The next example uses a version of the Cauchy-Schwarz inequality:

$$\left(\int_a^b g(t)\,dt\right)^2 \le (b - a) \int_a^b g(t)^2\,dt. \qquad (3.1)$$

It is valid for any real-valued $a$ and $b$, not just when $a < b$.

**Example 3.6.** Let $C^1[0,1]$ be the space of continuously differentiable functions. (If $f \in C^1[0,1]$ then $f'(x)$ exists and is continuous.) Define a norm by $\|f\|^2 = \int_0^1 f(x)^2\,dx + \int_0^1 f'(x)^2\,dx$. It is shown below that $\max_x |f(x)| \le 2\|f\|$. The norm controls the values of $f(x)$ for all $x$.

Let $c = \max_x |f(x)|$. By considering $-f$ instead of $f$ if necessary, it may be assumed that $c = \max_x f(x)$. If $f(x) \ge \frac{c}{2}$ for all $x$ then $\|f\| \ge \frac{c}{2}$ and thus $\max_x |f(x)| \le 2\|f\|$. Alternatively, there must exist a $y \in [0,1]$ such that $f(y) = \frac{c}{2}$. Let $x \in [0,1]$ be such that $f(x) = c$. Then $2\int_y^x f'(t)\,dt = c$ because $f(x) = f(y) + \int_y^x f'(t)\,dt$. Applying (3.1) with $g(t) = f'(t)$ yields $c^2 \le 4(x - y) \int_y^x f'(t)^2\,dt$, from which it follows that $c^2 \le 4 \int_0^1 f'(t)^2\,dt$, that is, $\|f\| \ge \frac{c}{2}$. Thus, $\max_x |f(x)| \le 2\|f\|$.

Linear functionals, such as $l_x$, fall into two categories, bounded and unbounded, based on whether $c_x = \sup_{f \in V} |l_x(f)| \cdot \|f\|^{-1}$ is finite or infinite. If $l_x$ is bounded then $|l_x(f)| \le c_x \|f\|$ for all $f$, demonstrating that the norm being small implies $f(x) = l_x(f)$ is small in magnitude.

A key requirement for $V \subset \mathbb{R}^X$ to be a RKHS is for the $l_x$ to be bounded for all $x$. This implies that if $f_n \in V$ is a sequence converging in norm, meaning $\|f_n - f\| \to 0$, then the sequence also converges pointwise, meaning $f_n(x) \to f(x)$. This requirement was not mentioned earlier because $l_x$ is automatically bounded when $V$ is finite dimensional.

# 4

---

## Infinite-dimensional RKHSs

---

The finite-dimensional RKHS theory introduced earlier suggests the general theory should take the following form. Fix an arbitrary set $X$. Endow a subspace $V$ of the function space $\mathbb{R}^X$ with an inner product. RKHS theory should associate with the tuple $(V, \langle \cdot, \cdot \rangle, X)$ a set of vectors in $\mathbb{R}^X$ that span $V$, that vary continuously as $(V, \langle \cdot, \cdot \rangle)$ changes, and that encode the inner product on $V$. Furthermore, by analogy with Definition 2.1, an appropriate set should be the vectors $k_x$ satisfying $\langle f, k_x \rangle = f(x)$ for all $f \in V$ and $x \in X$.

The above can be achieved with only several minor adjustments. Requiring the span of the $k_x$ to be $V$ is too demanding. For example, a basis for $l^2$ would necessarily be uncountable and cannot even be constructed explicitly, whereas the countable set of vectors $(1, 0, \cdots)$, $(0, 1, 0, \cdots)$, $\cdots$ suffices for working with $l^2$ if taking limits is acceptable. Although seemingly wanting to reconstruct $V$ by taking the topological closure of the span of the $k_x$, this would necessitate endowing $\mathbb{R}^X$ with a topology, and there might not be a topology on $\mathbb{R}^X$ compatible with all possible choices of $(V, \langle \cdot, \cdot \rangle)$. Instead, the mathematical process of completion (§4.1) can be used to achieve the same aim of adding to the span of the $k_x$ any additional vectors that "should be"

infinitesimally close to the $k_x$. A Hilbert space is an inner product space that is complete, hence the above can be summarised by saying RKHS theory requires $(V, \langle \cdot, \cdot \rangle)$ to be a Hilbert space.

There may not exist a $k_x$ such that $\langle f, k_x \rangle = f(x)$ for all $f \in V$ because the linear functional $f \mapsto \langle f, k_x \rangle$ is always bounded whereas Example 3.5 showed $f \mapsto f(x) = l_x(f)$ need not be bounded. Another advantage of requiring $V$ to be a Hilbert space is that any bounded linear functional can be written as an inner product. (This is the Riesz representation theorem.) Therefore, a necessary and sufficient condition for $\langle f, k_x \rangle = f(x)$ to have a solution $k_x$ is for $l_x$ to be bounded. Requiring $l_x$ to be bounded for all $x \in X$ is the second requirement RKHS theory places on $(V, \langle \cdot, \cdot \rangle)$.

RKHSs are similar to Euclidean spaces. A feature of Euclidean space is the presence of coordinate functions $\pi_i \colon \mathbb{R}^n \to \mathbb{R}$ sending $(x_1, \cdots, x_n)$ to $x_i$. These coordinate functions are continuous. A RKHS replicates this: if $V \subset \mathbb{R}^X$ is a RKHS, the coordinate functions $\pi_x \colon V \to \mathbb{R}$ sending $f$ to $f(x)$ are continuous by definition. (Recall that a linear functional is continuous if and only if it is bounded.)

## 4.1 Completions and Hilbert Spaces

If $V \subset W$ is not closed (§3.2) then there are points in $W$ not in $V$ but infinitesimally close to $V$. The closure $\bar{V}$ of $V$ in $W$ is the union of these infinitesimally close points and $V$. Perhaps a new $W$ can be found though for which $\bar{V} \subset W$ is not closed?

Given a normed space $V$, there is a unique (up to an isometric isomorphism) normed space $\hat{V} \supset V$, called the *completion* of $V$, such that every point in $\hat{V}$ is infinitesimally close to $V$ and there is no normed space $W$ for which $\hat{V} \subset W$ is not closed. (Here, the norm on $W$ must agree with the norm on $\hat{V}$ which must agree with the norm on $V$.) This means there is one and only one way to enlarge $V$ maximally by adding infinitesimally close points. Once this is done, no more points can be added infinitesimally close to the enlarged space $\hat{V}$.

This fact permits a refinement of the previous visualisation of a non-closed subspace as being blurred. If $V \subset W$ is not closed then there are

dents on the surface of $V$. Pedantically, these dents are straight line scratches since $V$ is a vector space. Filling in these scratches closes $V$ in $W$. Choosing a different $W$ cannot re-open these scratches but may reveal scratches elsewhere. The completion $\hat{V}$ results from filling in all possible scratches.

Scratches can be found without having to construct a $W$. A sequence $v_n$ in a normed space $V$ is *Cauchy* if, for all $\epsilon > 0$, there exists a positive integer $N$ such that $\|v_n - v_m\| < \epsilon$ whenever $n, m \geq N$. Every convergent sequence is a Cauchy sequence, and a sequence that is Cauchy "should" converge; if a Cauchy sequence in $V$ does not converge then it "points" to a scratch.

**Remark**   The distance between successive terms of a convergent sequence must go to zero: $1, \frac{1}{2}, \frac{1}{3}, \cdots \to 0$ implies $|\frac{1}{n} - \frac{1}{n+1}| \to 0$. The converse is false though: successive terms of $s_n = \sum_{i=1}^{n} i^{-1}$ become infinitesimally close but too slowly to prevent $s_n$ from diverging to infinity. Hence Cauchy imposed the stronger requirement of non-successive terms becoming infinitesimally close.

To demonstrate that Cauchy sequences point to scratches, let $v_n \in V$ be a non-convergent Cauchy sequence. Enlarge $V$ to $V \oplus \mathrm{span}\{f\}$ where $f$ is a new vector. The aim is to construct a norm on $V \oplus \mathrm{span}\{f\}$ agreeing with the norm on $V$ and placing $f$ precisely at the limit of the Cauchy sequence, that is, $v_n \to f$. Extending the norm requires defining $\|cf + g\|$ for all $c \in \mathbb{R}$ and $g \in V$. Since the aim is for $v_n \to f$, the obvious choice to try is $\|cf + g\| = \lim_{n \to \infty} \|cv_n + g\|$. It can be shown the proposed norm really is a norm; it satisfies the requisite axioms. That it extends the original norm is clear; when $c = 0$, the new and old norms agree. Finally, $v_n \to f$ because $\lim_n \|f - v_n\| = \lim_n \lim_m \|v_m - v_n\| = 0$, the last equality due to $v_n$ being Cauchy.

Once a scratch is filled in, it cannot be re-filled due to uniqueness of limits in normed spaces. If an attempt was made to place $g \in V \oplus \mathrm{span}\{f, g\}$ in the location pointed to by $v_n$, then $v_n \to g$, yet from above, $v_n \to f$, hence $f$ and $g$ must be the same point.

Textbooks explain how $V$ is completed by creating a new space comprising all Cauchy sequences in $V$ then quotienting by declaring two

Cauchy sequences as equivalent if they converge to the same limit. This is not explained here since the description above conveys adequately the underlying principles.

If the norm on $V$ comes from an inner product then the norm on the completion $\hat{V}$ also comes from an inner product, that is, the completion of an inner product space is itself an inner product space.

An inner product space that is complete (with respect to the norm induced by the inner product) is called a *Hilbert space.*

A completion is needed for reconstructing $V \subset \mathbb{R}^X$ from its kernel $K$ because in general the kernel can describe only a dense subspace of the original space. Any space containing this dense subspace and contained in its completion would have the same kernel. A unique correspondence is achieved by insisting $V$ is complete; given $K$, the space $V$ is the completion of the subspace described by $K$.

It would be remiss not to mention another reason Cauchy sequences and completions are important. Completeness ensures existence of solutions to certain classes of problems by preventing the solution from having been accidentally or deliberately removed from the space.

**Example 4.1.** Let $x^{(i)} = (0, \cdots, 0, 1, 0, \cdots)$ be the element of $l^2$ having unity as its $i$th term. Let $f = (1, \frac{1}{2}, \frac{1}{3}, \cdots) \in l^2$. Then $f$ is the unique solution in $l^2$ to the system of equations $\langle x^{(i)}, f \rangle = i^{-1}$ for $i = 1, 2, \cdots$. Let $V$ be a subspace of $l^2$ omitting $f$ but containing all the $x^{(i)}$, such as the $V$ in Example 3.2. Treating $V$ as a vector space, there is no solution $\tilde{f} \in V$ to $\langle x^{(i)}, \tilde{f} \rangle = i^{-1}$ for $i = 1, 2, \cdots$.

If a space is complete, an existence proof of a solution typically unfolds as follows. Assume the problem is to prove the existence of a solution $f$ to a differential equation. Construct a sequence of approximate solutions $f_n$ by mimicking how differential equations are solved numerically, with decreasing step size. The hope is that $f_n \to f$, but since $f$ cannot be exhibited explicity, it is not possible to prove $f_n \to f$ directly. Instead, $f_n$ is shown to be Cauchy and therefore, by completeness, has a limit $\tilde{f}$. Lastly, $\tilde{f}$ is verified by some limit argument to be the solution, thus proving the existence of $f = \tilde{f}$. For details, see [11].

Two examples of determining whether a space is complete are now given. The proof that $l^2$ is complete is quite standard and proceeds

along the following lines. If $f_1, f_2, \cdots \in l^2$ is Cauchy then the $i$th term of each $f_j$ form a Cauchy sequence of real numbers. Since the real numbers are complete — they are the completion of the rational numbers — the $f_j$ converge pointwise to a limit $f$. The norm of $f$ can be shown to be finite and hence $f$ is in $l^2$. Finally, it can be shown that $\|f_i - f\| \to 0$, proving the $f_i$ have a limit in $l^2$.

The space $C_2[0, 1]$ in Example 3.1 is not complete. This can be seen by considering a sequence of continuous functions that better and better approximate a square wave. The limit "should" exist but a square wave is not a continuous function and hence does not lie in $C_2[0, 1]$. The completion of $C_2[0, 1]$ is the space known as $L_2[0, 1]$. Although $L_2$ spaces are often treated as if they were function spaces, technically they are not. A meaning cannot be ascribed to the pointwise evaluation of a function in $L_2$ because an element of $L_2$ is actually an equivalence class of functions, any two members of which may differ pointwise on a set of measure zero. Conditions for when the completion of a function space remains a function space are given in §4.5.

## 4.2 Definition of a RKHS

The following definition of a RKHS differs in two inconsequential ways from other definitions in the literature. Normally, RKHSs over the complex field are studied because every real-valued RKHS extends canonically to a complex-valued RKHS. Often, the evaluation functionals are required to be continuous whereas here they are required to be bounded. A linear operator is bounded if and only if it is continuous.

**Definition 4.1.** Let $X$ be an arbitrary set and denote by $\mathbb{R}^X$ the vector space of all functions $f \colon X \to \mathbb{R}$ equipped with pointwise operations. A subspace $V \subset \mathbb{R}^X$ endowed with an inner product is a *reproducing kernel Hilbert space (RKHS)* if $V$ is complete and, for every $x \in X$, the evaluation functional $f \mapsto f(x) = l_x(f)$ on $V$ is bounded.

The kernel of a RKHS exists and is unique.

**Definition 4.2.** If $V \subset \mathbb{R}^X$ is a RKHS then its *kernel* is the function $K \colon X \times X \to \mathbb{R}$ satisfying $\langle f, K(\cdot, y) \rangle = f(y)$ for all $f \in V$ and $y \in X$. Here, $K(\cdot, y)$ denotes the function $x \mapsto K(x, y)$ and is an element of $V$.

If $V$ is a subspace of $\mathbb{R}^X$ where $X = \{1, \cdots, n\}$ then finite dimensionality implies it is complete and its evaluation functionals bounded. Thus, $V$ is a RKHS by Definition 4.1. The kernel in Definitions 4.2 and 2.1 are equivalent: $K(i, j) = K_{ij}$.

**Remark** The set $X$ in Definition 4.1 is arbitrary. It need not have a topology and, in particular, there are no continuity requirements on the elements of $V$ or the kernel $K$.

## 4.3 Examples of RKHSs

Examples of RKHSs can be found in [32, 10, 3, 30], among other places.

**Example 4.2** (Paley-Wiener Space). Let $V$ consist of bandlimited functions $f \colon \mathbb{R} \to \mathbb{R}$ expressible as $f(t) = \frac{1}{2\pi} \int_{-a}^{a} F(\omega) e^{j\omega t} \, d\omega$ where $F(\omega)$ is square-integrable. Endow $V$ with the inner product $\langle f, g \rangle = \int_{-\infty}^{\infty} f(t) g(t) \, dt$. Then $V$ is a RKHS with kernel [77, 38]

$$K(t, \tau) = \frac{\sin(a(t - \tau))}{\pi(t - \tau)}. \tag{4.1}$$

Using the Fourier transform and its inverse, the kernel (4.1) can be derived as follows from the requirement $\langle f, K(\cdot, \tau) \rangle = f(\tau)$.

$$
\begin{aligned}
f(\tau) &= \frac{1}{2\pi} \int_{-a}^{a} \left( \int_{-\infty}^{\infty} f(t) e^{-j\omega t} \, dt \right) e^{j\omega\tau} \, d\omega \\
&= \int_{-\infty}^{\infty} f(t) \left( \frac{1}{2\pi} \int_{-a}^{a} e^{j\omega(\tau - t)} \, d\omega \right) dt \\
&= \int_{-\infty}^{\infty} f(t) \left( \frac{1}{2\pi} \frac{2}{\tau - t} \sin(a(\tau - t)) \right) dt \\
&= \int_{-\infty}^{\infty} f(t) K(t, \tau) \, dt.
\end{aligned}
$$

Completeness of $V$ in Example 4.2 can be proven essentially as follows. If $f_n \in V$ is a Cauchy sequence then its Fourier transform $F_n$ is Cauchy too. Therefore, $F_n$ converges to a square-integrable function $F$ and $f(t) = \frac{1}{2\pi} \int_{-a}^{a} F(\omega) e^{j\omega t} \, d\omega$ is in $V$. It can be verified $f_n \to f$.

**Example 4.3.** Let $V$ consist of functions of the form $f\colon [0,1] \to \mathbb{R}$ where $f(t)$ is absolutely continuous, its derivative $f'(t)$ (which exists almost everywhere) is square-integrable, and $f(0) = 0$. Then $V$ is a RKHS when endowed with the inner product $\langle f, g \rangle = \int_0^1 f'(t)g'(t)\,dt$. The associated kernel is $K(t,s) = \min\{t,s\}$.

Readers may recognise $K(t,s) = \min\{t,s\}$ as the covariance function of a Wiener process. If $g(t) = K(t,s)$ then $g'(t) = 1$ when $t < s$ and $g'(t) = 0$ when $t > s$. Therefore,

$$\langle f, K(\cdot, s) \rangle = \int_0^s f'(t)\,dt = f(s).$$

Despite this primer focusing on real-valued functions, the following classic example of a RKHS comprises complex-valued functions.

**Example 4.4** (Bergman Space)**.** Let $S = \{z \in \mathbb{C} \mid |z| < 1\}$ denote the unit disk in the complex plane. Let $V$ be the vector space of all analytic and square-integrable functions $f\colon S \to \mathbb{C}$. Equip $V$ with the inner product $\langle f, g \rangle = \int_S f(z)\overline{g(z)}\,dz$. Then $V$ is a RKHS with kernel

$$K(z,w) = \frac{1}{\pi}\frac{1}{(1 - z\bar{w})^2}. \tag{4.2}$$

The above is an example of a Bergman space and its associated kernel. Choosing domains other than the unit disk produce other Bergman spaces and associated Bergman kernels [8]. (Generally the Bergman kernel cannot be determined explicitly though.)

## 4.4   Basic Properties

The fundamental property of RKHSs is the bijective correspondence between RKHSs and positive semi-definite functions: the kernel of a RKHS is positive semi-definite and every positive semi-definite function is the kernel of a unique RKHS [3].

A symmetric function $K\colon X \times X \to \mathbb{R}$ is positive semi-definite (equivalently, of positive type [45]) if, $\forall r \geq 1$, $\forall c_1, \cdots, c_r \in \mathbb{R}$, $\forall x_1, \cdots, x_r \in X$,

$$\sum_{i=1}^r \sum_{j=1}^r c_i c_j K(x_i, x_j) \geq 0. \tag{4.3}$$

Unlike in the complex-valued case when the complex version of (4.3) forces $K$ to be symmetric [10], in the real-valued case, it is necessary to insist explicitly that $K$ be symmetric: $K(x, y) = K(y, x)$.

The kernel of a RKHS is automatically symmetric; choosing $f(y) = K(y, x)$ in Definition 4.2 shows

$$\langle K(\cdot, x), K(\cdot, y) \rangle = K(y, x), \qquad x, y \in X. \tag{4.4}$$

The symmetry of inner products implies $K(x, y) = K(y, x)$, as claimed.

Confusingly, positive semi-definite functions are often referred to in the literature as positive definite functions. The terminology adopted here agrees with the finite-dimensional case when $K$ is a matrix.

A number of different topologies are commonly placed on function spaces. The strong topology comes from the norm: the sequence $f_n$ converges strongly to $f$ if $\|f_n - f\| \to 0$. It converges weakly if $\langle f_n, g \rangle \to \langle f, g \rangle$ for all $g \in V$. (This is weaker because the "rate of convergence" can be different for different $g$.) Pointwise convergence is when $f_n(x)$ converges to $f(x)$ for all $x \in X$. In general, the pointwise topology is unrelated to the strong and weak topologies. In a RKHS however, strong convergence implies pointwise convergence, as does weak convergence. Writing $f_n \to f$ refers to strong convergence.

If $K$ is the kernel of $V \subset \mathbb{R}^X$, let $V_0 = \text{span}\{x \mapsto K(x, y) \mid y \in X\}$. In words, $V_0$ is the space spanned by the functions $K(\cdot, y)$ as $y$ ranges over $X$. Clearly $V_0 \subset V$. In the finite-dimensional case, $V_0$ would equal $V$. In general, $V_0$ is only dense in $V$. Any $f \in V$ can be approximated arbitrarily accurately by finite linear combinations of the $K(\cdot, y)$. Moreover, $V$ is the completion of $V_0$ and hence can be recovered uniquely from $V_0$, where the inner product on $V_0$ is uniquely determined by (4.4). Any Cauchy sequence in $V_0$ converges *pointwise* to an element of the RKHS $V$. (It converges strongly by definition of completion, and strong convergence in a RKHS implies pointwise convergence.)

**Remark**   Proofs have been omitted partly because they can be found in more traditional introductory material, including [3, 30, 59], and partly because the above results should not be surprising once finite-dimensional RKHSs are understood.

## 4.5   Completing a Function Space

If $V$ is incomplete but meets the other conditions in Definition 4.1, a naive hope is for the completion of $V$ to be a RKHS. This hope is close to the mark: the completion of $V$ might not be a function space on $X$ but can be made into a function space on a set larger than $X$.

First a subtlety about completions. Although it is common to speak of *the* completion $\hat{V}$ of a normed space $V$, and to think of $V$ as a subset of $\hat{V}$, it is more convenient to refer to any normed space $W$ as a version of the completion of $V$, or simply, the completion of $V$, provided three conditions are met: $V$ can be identified with a subspace $V' \subset W$ (that is, $V$ and $V'$ are isometrically isomorphic); every Cauchy sequence in $V'$ converges to an element of $W$; every element of $W$ can be obtained in this way. The usual construction of $\hat{V}$ produces a quotient space of Cauchy sequences, where $V$ is identified with the set $V' \subset \hat{V}$ consisting of (equivalence classes of) Cauchy sequences converging to an element of $V$. The consequence for RKHS theory is that even though $V \subset \mathbb{R}^X$ is a function space, the usual construction produces a completion $\hat{V}$ whose elements are not functions.

Whether there exists a version of the completion that is a subset of $\mathbb{R}^X$ will be addressed in two parts. Assuming a RKHS completion exists, a means for constructing it will be derived, then necessary and sufficient conditions will be found for the RKHS completion to exist.

Assume $V_0 \subset \mathbb{R}^X$ has a completion $V \subset \mathbb{R}^X$ that is a RKHS containing $V_0$. Let $f_n$ be a Cauchy sequence in $V_0$. Its limit $f$ can be found pointwise: $f(x) = \langle f, K(\cdot, x) \rangle = \lim_n \langle f_n, K(\cdot, x) \rangle = \lim_n f_n(x)$. Hence $V$ can be reconstructed as the set of all functions on $X$ that are pointwise limits of Cauchy sequences in $V_0$.

This construction can fail in two ways to produce a completion of an arbitrary $V_0 \subset \mathbb{R}^X$. The pointwise limit might not exist, or the pointwise limit might be the same for Cauchy sequences with distinct limits. Being in a vector space, the latter is equivalent to the existence of a Cauchy sequence $f_n$ not converging to zero, meaning $\lim_n \|f_n\| \neq 0$, but whose pointwise limit is zero: $f_n(x) \to 0$ for $x \in X$.

**Proposition 4.1.** An arbitrary inner product space $V_0 \subset \mathbb{R}^X$ has a

RKHS completion $V$, where $V_0 \subset V \subset \mathbb{R}^X$, if and only if

1. the evaluation functionals on $V_0$ are bounded;

2. if $f_n \in V_0$ is a Cauchy sequence converging pointwise to zero then $\|f_n\| \to 0$.

*Proof.* Assume $V$ exists. Condition (1) holds because the evaluation functionals on $V$ are bounded. If $f_n \in V_0$ is Cauchy then $f_n \to f$ for some $f$ in $V$. Since strong convergence implies pointwise convergence, $f(x) = \lim_n f_n(x)$ for all $x \in X$. Thus, if $f_n$ converges pointwise to zero then its strong limit $f$ must be zero, in which case $\|f_n\| = \|f - f_n\| \to 0$, showing condition (2) holds.

Conversely, assume (1) and (2) hold. Let $f_n \in V_0$ be Cauchy. The evaluation functionals being bounded implies $f_n(x)$ is Cauchy for each $x \in X$, hence $f(x) = \lim_n f_n(x)$ exists. Let $V$ be the set of all such pointwise limits $f$. Consider endowing $V$ with a norm satisfying $\|f\| = \lim_n \|f_n\|$ where $f_n$ is a Cauchy sequence converging pointwise to $f$. Condition (2) ensures different Cauchy sequences converging pointwise to the same limit yield the same norm: if $f_n(x) \to f(x)$ and $g_n(x) \to f(x)$ then $f_n(x) - g_n(x) \to 0$, hence $|\|f_n\| - \|g_n\|| \leq \|f_n - g_n\| \to 0$. Omitted are the routine proofs showing $\|\cdot\|$ satisfies the axioms of a norm and the parallelogram law, hence giving $V$ an inner product agreeing with the original inner product on $V_0 \subset V$.

To verify $V$ is the completion of $V_0$, let $f_n \in V_0$ be Cauchy, with pointwise limit $f \in V$. The sequence $f_m - f_n$ in $m$ is a Cauchy sequence in $V_0$ converging pointwise to $f - f_n$. Therefore, $\lim_n \|f - f_n\| = \lim_n \lim_m \|f_m - f_n\| = 0$, proving the pointwise limit of a Cauchy sequence in $V_0$ is also the strong limit. Therefore, any $f \in V$ can be written as a pointwise limit and hence as a strong limit of a Cauchy sequence in $V_0$, and any Cauchy sequence in $V_0$ converges pointwise and hence in the strong limit to an element of $V$. $\square$

The following example is taken from [3, pp. 349–350] and shows condition (2) in the above proposition is not automatically true. Complex-valued functions are used because analytic functions are better behaved than real-analytic functions.

Define $S$ and $\langle \cdot, \cdot \rangle$ as in Example 4.4 and let $\|\cdot\|$ denote the norm induced from the inner product. The sequence $a_n = 1 - (\frac{1}{2})^n$ belongs to $S$. The function (known as a Blaschke product) $f(z) = \prod_{n=1}^{\infty} \frac{a_n - z}{1 - a_n z}$ is analytic and bounded on $S$. It satisfies $f(a_n) = 0$ for all $n$. (An analytic function vanishing on a convergent set of points need itself only be zero if the limit point is within the domain of definition. Here, the domain is $S$ and $a_n \to 1 \notin S$.) There exists a sequence of polynomials $f_n(z)$ for which $\lim_{n \to \infty} \|f_n - f\| = 0$.

These ingredients are assembled into an example by defining $X = \{a_1, a_2, \cdots\} \subset S$ and taking $V_0$ to be the restriction to $X$ of all polynomials $p \colon S \to \mathbb{C}$, keeping the above norm. The $f_n$ above are a Cauchy sequence of polynomials on $S$ and *a fortiori* on $X$. Pointwise they converge to zero: $f_n(x) \to f(x) = 0$ for $x \in S$. However, $\|f_n\| \to \|f\| \neq 0$. Condition (2) in Proposition 4.1 is not met.

Put simply, $X$ did not contain enough points to distinguish the strong limit of the Cauchy sequence $f_n$ from the zero function. In this particular example, enlarging $X$ to contain a sequence and its limit point would ensure condition (2) is met.

This principle holds generally. Assume $V_0 \subset \mathbb{R}^X$ fails to meet condition (2). Augment $X$ to become $X' = X \sqcup \hat{V}$, the disjoint union of $X$ and the completion $\hat{V}$ of $V_0$. Extend each $f \in V_0$ to a function $\tilde{f} \in \mathbb{R}^{X'}$ by declaring $\tilde{f}(x) = f(x)$ if $x \in X$ and $\tilde{f}(x) = \langle f, x \rangle$ if $x \in \hat{V}$. (The inner product is the one on $\hat{V}$, and $f \in V_0$ is identified with its corresponding element in $\hat{V}$.) Let $V_0'$ be the image of this linear embedding $f \mapsto \tilde{f}$ of $V_0$ into $\mathbb{R}^{X'}$. The inner product on $V_0$ gets pushed forwards to an inner product on $V_0'$, that is, $\langle \tilde{f}, \tilde{g} \rangle$ is defined as $\langle f, g \rangle$. To show condition (2) of Proposition 4.1 is met, let $\tilde{f}_n \in V_0'$ be a Cauchy sequence converging pointwise to zero. Then $\tilde{f}_n(x) \to 0$ for all $x \in X'$, and in particular, $\langle f_n, g \rangle \to 0$ for all $g \in \hat{V}$. This implies $f = 0$ where $f \in \hat{V}$ is the limit of $f_n$. Therefore, $\|f_n\| \to \|f\| = 0$, as required.

## 4.6  Joint Properties of a RKHS and its Kernel

The correspondence between a RKHS $V \subset \mathbb{R}^X$ and its kernel $K$ induces a correspondence between certain properties of $V$ and of $K$. Modifying

or combining RKHSs to form new ones can translate into familiar operations on the corresponding kernels. A small selection of examples is presented below.

### 4.6.1 Continuity

If $X$ is a metric space then it can be asked whether all the elements of a RKHS $V \subset \mathbb{R}^X$ are continuous. Since $K(\cdot, y)$ is an element of $V$, a necessary condition is for the functions $x \mapsto K(x, y)$ to be continuous for all $y \in X$. An arbitrary element of $V$ though is a limit of a Cauchy sequence of finite linear combinations of the $K(\cdot, y)$. An extra condition is required to ensure this limit is continuous.

**Proposition 4.2.** The elements of a RKHS $V \subset \mathbb{R}^X$ are continuous, where $X$ is a metric space, if and only if

1. $x \mapsto K(x, y)$ is continuous for all $y \in X$; and

2. for every $x \in X$ there is a radius $r > 0$ such that $y \mapsto K(y, y)$ is bounded on the open ball $B(x, r)$.

*Proof.* See [10, Theorem 17]. □

The necessity of the second condition is proved by assuming there is an $x \in X$ and a sequence $x_n \in B(x, \frac{1}{n})$ satisfying $K(x_n, x_n) \geq n$. The functions $K(\cdot, x_n)$ grow without limit: $\|K(\cdot, x_n)\| = K(x_n, x_n) \to \infty$. By a corollary of the Banach-Steinhaus theorem, a weakly convergent sequence is bounded in norm, therefore, there exists a $g \in V$ such that $\langle g, K(\cdot, x_n) \rangle \nrightarrow \langle g, K(\cdot, x) \rangle$. This $g$ is not continuous: $g(x_n) \nrightarrow g(x)$.

### 4.6.2 Invertibility of Matrices

Equation (4.3) asserts the matrix $A_{ij} = K(x_i, x_j)$ is positive semi-definite. If this matrix $A$ is singular then there exist constants $c_i$, not

all zero, such that (4.3) is zero. Since

$$\sum_{i=1}^{r}\sum_{j=1}^{r} c_i c_j K(x_i, x_j) = \sum_{i=1}^{r}\sum_{j=1}^{r} c_i c_j \langle K(\cdot, x_j), K(\cdot, x_i) \rangle$$

$$= \left\langle \sum_{j=1}^{r} c_j K(\cdot, x_j), \sum_{i=1}^{r} c_i K(\cdot, x_i) \right\rangle$$

$$= \left\| \sum_{i=1}^{r} c_i K(\cdot, x_i) \right\|^2,$$

the matrix $A$ is non-singular if and only if $\sum_{i=1}^{r} c_i K(\cdot, x_i) = 0$ implies $c_1 = \cdots = c_r = 0$. Now, $\sum_{i=1}^{r} c_i K(\cdot, x_i) = 0$ if and only if $\sum_{i=1}^{r} c_i f(x_i) = 0$ for all $f \in V$, where $V$ is the RKHS whose kernel is $K$. Whenever $V$ is sufficiently rich in functions, there will be no non-trivial solution to $\sum_{i=1}^{r} c_i f(x_i) = 0$ and every matrix of the form $A_{ij} = K(x_i, x_j)$ is guaranteed to be non-singular. For example, since the corresponding RKHS contains all polynomials, the kernel (4.2) is strictly positive, and for any collection of points $x_1, \cdots, x_r$ in the unit disk, the matrix $A_{ij} = (1 - x_i \bar{x}_j)^{-2}$ is non-singular.

### 4.6.3   Restriction of the Index Set

If $V$ is a RKHS of functions on the interval $[0, 2]$ then a new space $V'$ of functions on $[0, 1]$ is obtained by restricting attention to the values of the functions on $[0, 1]$. Since two or more functions in $V$ might collapse to the same function in $V'$, the norm on $V$ does not immediately induce a norm on $V'$. Since a kernel $K$ on $X$ restricts to a kernel $K'$ on $X'$, it is natural to ask if there is a norm on $V'$ such that $K'$ is the kernel of $V'$. (The norm would have to come from an inner product.)

When interpolation problems are studied later, it will be seen that $K(\cdot, x)$ solves the one-point interpolation problem; it is the function of smallest norm satisfying $f(x) = K(x, x)$. This suggests the norm of $f$ should be the smallest norm of all possible functions in $V$ collapsing to $f$. This is precisely the case and is now formalised.

If $X' \subset X$ then any $f \colon X \to \mathbb{R}$ restricts to a function $f|_{X'} \colon X' \to \mathbb{R}$ given by $f|_{X'}(x) = f(x)$ for $x \in X'$. If $V \subset \mathbb{R}^X$ is a RKHS then a new space $V' \subset \mathbb{R}^{X'}$ results from restricting each element of $V$ to $X'$.

Precisely, $V' = \{f|_{X'} \mid f \in V\}$. Define on $V'$ the norm

$$\|f\| = \inf_{\substack{g \in V \\ g|_{X'}=f}} \|g\|. \tag{4.5}$$

**Remark** In (4.5), inf can be replaced by min because $V$ is complete.

A kernel $K \colon X \times X \to \mathbb{R}$ restricts to a kernel $K' \colon X' \times X' \to \mathbb{R}$ where $K'(x,y) = K(x,y)$ for $x, y \in X'$. As proved in [3, p. 351] and [10, Theorem 6], the kernel of $V'$ is $K'$.

### 4.6.4 Sums of Kernels

If $K_1$ and $K_2$ are kernels on the same set $X$ then $K = K_1 + K_2$ is positive semi-definite and hence a kernel of a RKHS $V \subset \mathbb{R}^X$. There is a relatively straightforward expression for $V$ in terms of the RKHSs $V_1$ and $V_2$ whose kernels are $K_1$ and $K_2$ respectively.

The space itself is $V = V_1 \oplus V_2$, that is,

$$V = \{f_1 + f_2 \mid f_1 \in V_1, \ f_2 \in V_2\}. \tag{4.6}$$

The norm on $V$ is given by a minimisation:

$$\|f\|^2 = \inf_{\substack{f_1 \in V_1 \\ f_2 \in V_2 \\ f_1 + f_2 = f}} \|f_1\|^2 + \|f_2\|^2. \tag{4.7}$$

The three norms in (4.7) are defined on $V$, $V_1$ and $V_2$, in that order. Proofs of the above can be found in [3, p. 353] and [10, Theorem 5].

If there is no non-zero function belonging to both $V_1$ and $V_2$ then $f \in V$ uniquely decomposes as $f = f_1 + f_2$. To see this, assume $f_1 + f_2 = f_1' + f_2'$ with $f_1, f_1' \in V_1$ and $f_2, f_2' \in V_2$. Then $f_1 - f_1' = f_2 - f_2'$. As the left-hand side belongs to $V_1$ and the right-hand side to $V_2$, the assertion follows. In this case, (4.7) becomes $\|f\|^2 = \|f_1\|^2 + \|f_2\|^2$, implying from Pythagoras' theorem that $V_1$ and $V_2$ have been placed at right-angles to each other in $V$.

# 5

## Geometry by Design

The delightful book [44] exemplifies the advantages of giving a geometry to an otherwise arbitrary collection of points $\{p_t \mid t \in T\}$. Although this can be accomplished by a recipe for evaluating the distance $d(p_t, p_\tau)$ between pairs of points, a metric space lacks the richer structure of Euclidean space; linear operations are not defined, and there is no inner product or coordinate system.

### 5.1 Embedding Points into a RKHS

Assume a recipe has been given for evaluating $\langle p_t, p_\tau \rangle$ that satisfies the axioms of an inner product. If the $p_t$ are random variables then the recipe might be $\langle p_t, p_\tau \rangle = E[p_t p_\tau]$. It is not entirely trivial to construct an inner product space and arrange correctly the $p_t$ in it because the recipe for $\langle p_t, p_\tau \rangle$ may imply linear dependencies among the $p_t$. It is not possible to take by default $\{p_t \mid t \in T\}$ as a basis for the space.

RKHSs have the advantage of not requiring the kernel to be non-singular. The set $\{p_t \mid t \in T\}$ can serve directly as an "overdetermined basis", as now demonstrated. Note $T$ need not be a finite set.

Define $K(\tau, t) = \langle p_t, p_\tau \rangle$. If the inner product satisfies the axioms

required of an inner product then $K \colon T \times T \to \mathbb{R}$ is positive semi-definite and therefore is the kernel of a RKHS $V \subset \mathbb{R}^T$. Just like an element of $\mathbb{R}^2$ can be drawn as a point and given a label, interpret the element $K(\cdot, t)$ of $V$ as a point in $\mathbb{R}^T$ labelled $p_t$. In this way, the $p_t$ have been arranged in $V$. They are arranged correctly because

$$\langle K(\cdot, t), K(\cdot, \tau) \rangle = K(\tau, t) = \langle p_t, p_\tau \rangle. \tag{5.1}$$

Not only have the points $p_t$ been arranged correctly in the RKHS, they have been given coordinates: since $p_\tau$ is represented by $K(\cdot, \tau)$, its $t$th coordinate is $\langle K(\cdot, \tau), K(\cdot, t) \rangle = K(t, \tau)$. This would not have occurred had the points been placed in an arbitrary Hilbert space, and even if the Hilbert space carried a coordinate system, the assignment of coordinates to points $p_t$ would have been *ad hoc* and lacking the reproducing property.

Even if the $p_t$ already belong to a Hilbert space, embedding them into a RKHS will give them a pointwise coordinate system that is consistent with taking limits (i.e., strong convergence implies pointwise convergence). See §6.2 and §7 for examples of why this is beneficial.

Let $\mathcal{H}$ be a Hilbert space, $\{p_t \mid t \in T\} \subset \mathcal{H}$ a collection of vectors and $V \subset \mathbb{R}^T$ the RKHS whose kernel is $K(\tau, t) = \langle p_t, p_\tau \rangle$. Then $V$ is a replica of the closure $U$ of the subspace in $\mathcal{H}$ spanned by the $p_t$; there is a unique isometric isomorphism $\phi \colon U \to V$ satisfying $\phi(p_t) = K(\cdot, t)$. An isometric isomorphism preserves the linear structure and the inner product. In particular, $\langle \phi(p), \phi(q) \rangle = \langle p, q \rangle$. The $t$th coordinate of $u \in U$ is defined to be the $t$th coordinate of $\phi(u)$, namely, $\langle \phi(u), K(\cdot, t) \rangle$.

## 5.2   The Gaussian Kernel

Points in $\mathbb{R}^n$ can be mapped to an infinite-dimensional RKHS by declaring the inner product between $x, y \in \mathbb{R}^n$ to be

$$\langle x, y \rangle = \exp \left\{ -\frac{1}{2} \|x - y\|^2 \right\}. \tag{5.2}$$

As in §5.1, the RKHS is defined via its kernel $K(x, y) = \langle x, y \rangle$, called a Gaussian kernel. The point $x \in \mathbb{R}^n$ is represented by the element $K(\cdot, x)$ in the RKHS, a Gaussian function centred at $x$.

The geometry described by an inner product is often easier to understand once the distance between points has been deduced. Since $\langle K(\cdot, x), K(\cdot, x) \rangle = K(x, x) = 1$, every point $x \in \mathbb{R}^n$ is mapped to a unit length vector in the RKHS. The squared distance between $x, y \in \mathbb{R}^n$ is

$$\langle K(\cdot, x) - K(\cdot, y), K(\cdot, x) - K(\cdot, y) \rangle$$
$$= \langle K(\cdot, x), K(\cdot, x) \rangle - 2\langle K(\cdot, x), K(\cdot, y) \rangle + \langle K(\cdot, y), K(\cdot, y) \rangle$$
$$= K(x, x) - 2K(x, y) + K(y, y) \tag{5.3}$$
$$= 2\left( 1 - \exp\left\{ -\frac{1}{2} \|x - y\|^2 \right\} \right). \tag{5.4}$$

As $y$ moves further away from $x$ in $\mathbb{R}^n$, their representations in the infinite-dimensional RKHS also move apart, but once $\|x - y\| > 3$, their representations are close to being as far apart as possible, namely, the representations are separated by a distance close to 2.

Although the maximum straight-line distance between two representations is 2, a straight line in $\mathbb{R}^n$ does not get mapped to a straight line in the RKHS. The original geometry of $\mathbb{R}^n$, including its linear structure, is completely replaced by (5.2). The ray $t \mapsto tx$, $t \geq 0$, is mapped to the curve $t \mapsto K(\cdot, tx)$. The curve length between $K(\cdot, 0)$ and $K(\cdot, tx)$ is

$$\lim_{N \to \infty} \sum_{i=0}^{N-1} \left\| K\left(\cdot, \frac{i+1}{N} tx\right) - K\left(\cdot, \frac{i}{N} tx\right) \right\|$$
$$= \lim_{N \to \infty} \sum_{i=0}^{N-1} \sqrt{2\left(1 - \exp\left\{-\frac{1}{2N^2} \|tx\|^2\right\}\right)}$$
$$= \lim_{N \to \infty} \sqrt{2N^2 \left(1 - \left[1 - \frac{1}{2N^2} \|tx\|^2 + \cdots\right]\right)}$$
$$= t\|x\|.$$

Calculations similar to the above can be carried out for any given kernel $K$, at least in theory.

# 6

## Applications to Linear Equations and Optimisation

Linearly constrained norm-minimisation problems in Hilbert spaces benefit from the norm coming from an inner product. The inner product serves as the derivative of the cost function, allowing the minimum-norm solution to be expressed as the solution to a linear equation. This equation, known as the normal equation, can be written down directly as an orthogonality constraint, as shown on the right in Figure 3.1. Whenever the constraint set forms a *closed* subspace, a minimum-norm solution is guaranteed to exist [41].

Reproducing kernel Hilbert spaces have additional benefits when the constraints are of the form $f(x) = c$ because the intersection of a finite or even infinite number of such pointwise constraints is a *closed* affine space. The minimum-norm solution is guaranteed to exist and is expressible using the kernel of the RKHS.

Example 2.6 showed mathematically that the kernel features in minimum-norm problems. A simple geometric explanation comes from the defining equation $\langle f, K(\cdot, x) \rangle = f(x)$ implying $K(\cdot, x)$ is orthogonal to every function $f$ having $f(x) = 0$. This is precisely the geometric constraint for $K(\cdot, x)$ to be a function of smallest norm satisfying $f(x) = c$. The actual value of $c$ is determined by substituting $f = K(\cdot, x)$ into

the defining equation, yielding the self-referential $f(x) = K(x, x)$.

## 6.1  Interpolation Problems

The interpolation problem asks for the function $f \in V \subset \mathbb{R}^X$ of minimum norm satisfying the finite number of constraints $f(x_i) = c_i$ for $i = 1, \cdots, r$. Assuming for the moment there is at least one function in $V$ satisfying the constraints, if $V$ is a RKHS with kernel $K$ then a minimum-norm solution exists.

Let $U = \{f \in V \mid f(x_i) = c_i, \ i = 1, \cdots, r\}$ and $W = \{g \in V \mid g(x_i) = 0, \ i = 1, \cdots, r\}$. The subspace $W$ is closed because $g(x) = 0$ is equivalent to $\langle g, K(\cdot, x) \rangle = 0$ and the collection of vectors orthogonal to a particular vector forms a closed subspace. (Equivalently, the boundedness of evaluation functionals implies they are continuous and the inverse image of the closed set $\{0\}$ under a continuous function is closed.) Being the intersection of a collection of closed sets, $W$ is closed. Thus $U$ is closed since it is a translation of $W$, and a minimum-norm solution $f \in U$ exists because $U$ is closed (§4.1).

A minimum-norm $f \in U$ must be orthogonal to all $g \in W$, for else $f + \epsilon g \in U$ would have a smaller norm for some $\epsilon \neq 0$. Since $W$ is the orthogonal complement of $O = \text{span}\{K(\cdot, x_1), \cdots, K(\cdot, x_r)\}$, and $O$ is closed, $O$ is the orthogonal complement of $W$. The minimum-norm solution must lie in $O$.

The minimum-norm $f$ can be found by writing $f = \sum_{j=1}^r \alpha_j K(\cdot, x_j)$ and solving for the scalars $\alpha_j$ making $f(x_i) = c_i$. If no such $\alpha_j$ exist then $U$ must be empty. Since $f(x_i) = \sum_{j=1}^r \alpha_j K(x_i, x_j)$, the resulting linear equations can be written in matrix form as $A\alpha = c$ where $A_{ij} = K(x_i, x_j)$ and the $\alpha$ and $c$ are vectors containing the $\alpha_j$ and $c_i$ in order. A sufficient condition for $A$ to be non-singular is for the kernel to be strictly positive definite (§4.6.2). See also [38, 77].

It is remarked that the Representer Theorem in statistical learning theory is a generalisation of this basic result.

Example 4.2 demonstrates that certain spaces of bandlimited signals are reproducing kernel Hilbert spaces. Reconstructing a bandlimited signal from discrete-time samples can be posed as an interpolation

problem [77].

## 6.2 Solving Linear Equations

If the linear equation $Ax = b$ has more than one solution, the Moore-Penrose pseudo-inverse $A^+$ of $A$ finds the solution $x = A^+b$ of minimum norm. The involvement of norm minimisation suggests RKHS theory might have a role in solving infinite-dimensional linear equations.

The minimum-norm solution is that which is orthogonal to the null-space of $A$. The $r$ rows of $A$ are orthogonal to the null-space, therefore, the required solution is that which is expressible as a linear combination of the rows of $A$. This motivates writing $Ax = b$ as $\langle x, a_i \rangle = b_i$, $i = 1, \cdots, r$, where $a_i$ is the $i$th row of $A$, and looking for a solution $x$ lying in $U = \text{span}\{a_1, \cdots, a_r\}$.

If $U$ was a RKHS and the $a_i$ of the form $K(\cdot, i)$ then $\langle x, a_i \rangle = b_i$ would be an interpolation problem (§6.1). This form can always be achieved by embedding $U$ into a RKHS (§5.1) because the embedding $\phi$ sends $a_i$ to $K(\cdot, i)$ and preserves the inner product: $\langle x, a_i \rangle = \langle \phi(x), \phi(a_i) \rangle = \langle \phi(x), K(\cdot, i) \rangle$. The solution to $\langle \phi(x), K(\cdot, i) \rangle = b_i$ is $\phi(x) = b$ and hence the solution to $\langle x, a_i \rangle = b_i$ is $x = \phi^{-1}(b)$.

**Theorem 6.1.** Let $\mathcal{H}$ be a Hilbert space and $\{a_t \mid t \in T\} \subset \mathcal{H}$ a collection of vectors in $\mathcal{H}$. Denote the closure of the span of the $a_t$ by $U$. The equations $\langle x, a_t \rangle = b_t$ for $t \in T$ have a solution $x \in U$ if and only if the function $t \mapsto b_t$ is an element of the RKHS $V \subset \mathbb{R}^T$ whose kernel $K \colon T \times T \to \mathbb{R}$ is $K(s, t) = \langle a_t, a_s \rangle$. If a solution $x$ in $U$ exists then it is unique and has the smallest norm out of all solutions $x$ in $\mathcal{H}$. Moreover, $x = \phi^{-1}(b)$ where $b(t) = b_t$ and $\phi \colon U \to V$ is the unique isometric isomorphism sending $a_t$ to $K(\cdot, t)$ for $t \in T$.

*Proof.* Follows from the discussion above; see also [10, Theorem 42]. $\square$

The $\mathcal{H}$-norm of the solution $x$ equals the $V$-norm of $b$ because $\phi$ in Theorem 6.1 preserves the norm. Sometimes this knowledge suffices. Otherwise, in principle, $x$ can be found by expressing $b$ as a linear combination of the $K(\cdot, t)$. If $b = \sum_k \alpha_k K(\cdot, t_k)$ then $x = \phi^{-1}(b) =$

$\sum_k \alpha_k a_{t_k}$ by linearity of $\phi$. More generally, limits of linear combinations can be used, including integrals and derivatives of $K(\cdot, t)$; see [32, Section III-B] or [10, Theorem 43]. For an alternative, see Theorem 6.2.

**Example 6.1.** The equation $Ax = b$ can be written as $\langle x, a_i \rangle = b_i$ with respect to the Euclidean inner product. Then $K(i, j) = \langle a_j, a_i \rangle$. In matrix form, $K = AA^\top$. The map $\phi$ sending $a_i = A^\top e_i$ to $K(\cdot, i) = AA^\top e_i$ is $\phi(x) = Ax$. It is more useful though to think of $\phi$ as sending $A^\top v$ to $AA^\top v$ for arbitrary $v$ because then $\phi^{-1}$ is seen to send $AA^\top v$ to $A^\top v$. Expressing $b$ as a linear combination of the $K(\cdot, i)$ means finding the vector $\alpha$ such that $b = AA^\top \alpha$. Assuming $AA^\top$ is invertible, $\alpha = (AA^\top)^{-1}b$. As $b = (AA^\top)(AA^\top)^{-1}\alpha$ and $\phi^{-1}$ changes the leading $AA^\top$ to $A^\top$, $x = \phi^{-1}(b) = A^\top(AA^\top)^{-1}b$.

The kernel $K = AA^\top$ in the above example corresponds to the inner product $\langle b, c \rangle = c^\top(AA^\top)^{-1}b$ by the observation in Example 2.1. The solution $x = A^\top(AA^\top)^{-1}b$ also involves $(AA^\top)^{-1}$, suggesting $x$ can be written using the inner product of the RKHS. This is pursued in §6.2.1.

**Example 6.2.** Consider the integral transform

$$f(t) = \int_Z F(z)h(z, t)\, dz, \qquad t \in T, \tag{6.1}$$

taking a function $F \in \mathbb{R}^Z$ and returning a function $f \in \mathbb{R}^\top$. Regularity conditions are necessary for the integral to exist: assume $F(\cdot)$ and $h(\cdot, t)$ for all $t \in T$ are square-integrable, as in [59, Chapter 6]. Let $\mathcal{H}$ be the Hilbert space $L_2(Z)$ of square-integrable functions on $Z$. Then (6.1) becomes $\langle F, h(\cdot, t) \rangle = f(t)$ for $t \in T$. Define $K(s, t) = \langle h(\cdot, t), h(\cdot, s) \rangle = \int_Z h(z, s)h(z, t)\, dz$ and let $V \subset \mathbb{R}^\top$ be the associated RKHS. Then (6.1) has a solution $F$ if and only if $f \in V$.

### 6.2.1  A Closed-form Solution

The solution $x = A^\top(AA^\top)^{-1}b$ to $Ax = b$ in Example 6.1 can be written as $x_j = \langle b, c_j \rangle$ where $c_j$ is the $j$th column of $A$ and the inner product comes from the kernel $K = AA^\top$, that is, $\langle b, c \rangle = c^\top(AA^\top)^{-1}b$. Indeed, $\langle b, c_j \rangle = \langle b, Ae_j \rangle = e_j^\top A^\top(AA^\top)^{-1}b = e_j^\top x = x_j$, as claimed.

Similarly, under suitable regularity conditions, a closed-form solution to (6.1) can be found as follows. Let $V \subset \mathbb{R}^X$ be the RKHS with $K(s,t) = \int h(z,t)h(z,s)\,dz$ as its kernel. Any $f \in V$ therefore satisfies

$$
\begin{aligned}
f(t) &= \langle f, K(\cdot, t) \rangle \\
&= \left\langle f, \int h(z,t)h(z,\cdot)\,dz \right\rangle \\
&= \int h(z,t) \langle f, h(z,\cdot) \rangle \, dz \\
&= \int h(z,t)F(z)\,dz,
\end{aligned}
\tag{6.2}
$$

showing the integral equation $f(t) = \int F(z)h(z,t)\,dz$ has $F(z) = \langle f, h(z,\cdot) \rangle$ as a solution. This derivation does not address whether the solution has minimum norm; sometimes it does [59, Chapter 6].

The following novel derivation builds on the idea in (6.2) and culminates in Theorem 6.2. It relies on the concept of a Parseval frame [28, 52]. A collection of vectors $c_s \in V$ is a Parseval frame for $V \subset \mathbb{R}^T$ if and only if $b = \sum_s \langle b, c_s \rangle c_s$ for all $b \in V$. By [52, Theorem 3.12], $c_s$ is a Parseval frame for $V$ if and only if $K = \sum_s c_s c_s^\top$. This condition extends the second definition in Definition 2.1 because an orthonormal basis is *a fortiori* a Parseval frame. Note that a solution to $\langle x, a_t \rangle = b_t$ can be found analogous to (6.2) by substituting $K = \sum_s c_s c_s^\top$ into $b_t = \langle b, Ke_t \rangle$ in the particular case when $T = \{1, 2, \cdots, n\}$. A more general strategy is sought though.

The difficulty of solving $\phi(x) = b$ in Theorem 6.1 using the kernel directly is having to express $b$ as a linear combination of the $K(\cdot, t)$. Expressing $b$ as a linear combination of a Parseval frame $c_s$ though is trivial: $b = \sum_s \langle b, c_s \rangle c_s$. For this to be useful, it must be possible to evaluate $\phi^{-1}(c_s)$. It turns out that the choice $c_s = \psi(v_s)$ achieves both objectives, where $v_s$ is an orthonormal basis for $\mathcal{H}$ and $\psi \colon \mathcal{H} \to \mathbb{R}^T$ is the map sending $x$ to $b$ where $b_t = \langle x, a_t \rangle$.

The map $\phi \colon U \to V$ in Theorem 6.1 is the restriction of $\psi$ to $U$. Denote by $P$ the orthogonal projection onto $U$. If $b = \psi(x)$ for some $x \in \mathcal{H}$ then the minimum-norm solution of $\psi(\hat{x}) = b$ is $\hat{x} = P(x)$. Equivalently, $\phi^{-1}(\psi(x)) = P(x)$. Moreover, $\phi(P(x)) = \psi(x)$.

The inverse image $\phi^{-1}(c_s)$ of $c_s = \psi(v_s)$ can be deduced from the

following calculation, valid for any $x \in \mathcal{H}$.

$$
\begin{aligned}
\phi^{-1}(\psi(x)) &= \sum_s \langle P(x), v_s \rangle v_s \\
&= \sum_s \langle P(x), P(v_s) \rangle v_s \\
&= \sum_s \langle \phi(P(x)), \phi(P(v_s)) \rangle v_s \\
&= \sum_s \langle \psi(x), \psi(v_s) \rangle v_s.
\end{aligned}
\tag{6.3}
$$

In fact, this shows not only how $\phi^{-1}(c_s)$ can be found, but that the minimum-norm solution $\hat{x}$ to $\psi(x) = b$ is given by

$$
\hat{x} = \sum_s \langle b, c_s \rangle v_s, \quad c_s = \psi(v_s).
\tag{6.4}
$$

Equivalently, $x$ is found "elementwise" by $\langle x, v_s \rangle = \langle b, c_s \rangle$. Note that $b = \sum_s \langle b, c_s \rangle c_s$, therefore, $x$ is obtained by using the same linear combination but with $c_s$ replaced by its preimage $v_s$. Even though the $v_s$ need not lie in $U$, the linear combination of them does.

That the $c_s$ form a Parseval frame only enters indirectly to guide the derivation (6.3). It helps explain why it works. For completeness, the following verifies the $c_s$ form a Parseval frame.

$$
\begin{aligned}
e_t^\top \left( \sum_s c_s c_s^\top \right) e_\tau &= \sum_s (e_t^\top c_s)(e_\tau^\top c_s) \\
&= \sum_s \langle v_s, a_t \rangle \langle v_s, a_\tau \rangle \\
&= \left\langle \sum_s \langle a_t, v_s \rangle v_s, a_\tau \right\rangle \\
&= \langle a_t, a_\tau \rangle = \langle a_\tau, a_t \rangle = K(t, \tau).
\end{aligned}
$$

The above leads to the following elementary but new result.

**Theorem 6.2.** With notation as in Theorem 6.1, let $\{v_s \mid s \in S\}$ be an orthonormal basis for $\mathcal{H}$. Define $c_s(t) = \langle v_s, a_t \rangle$. The minimum-norm solution $x$ in Theorem 6.1 is the unique solution to $\langle x, v_s \rangle = \langle b, c_s \rangle$.

*Proof.* As $v_s$ is an orthonormal basis, the solution $x$ to $\langle x, v_s \rangle = \langle b, c_s \rangle$ is unique. It therefore suffices to show $\langle x, v_s \rangle = \langle \phi(x), c_s \rangle$ for $x \in U$. This

follows from $\langle x, v_s \rangle = \langle P(x), v_s \rangle = \langle x, P(v_s) \rangle = \langle \phi(x), \phi(P(v_s)) \rangle = \langle \phi(x), c_s \rangle$ where $P$ is projection onto $U$. That $c_s = \phi(P(v_s))$ can be derived from $c_s(t) = \langle v_s, P(a_t) \rangle = \langle P(v_s), a_t \rangle$. $\qquad\square$

Theorem 6.2 can be used to solve numerically integral equations via series expansions (Example 6.2). If $x = \sum_s \alpha_s v_s$ then each $\alpha_s$ is found by evaluating numerically the RKHS inner product $\langle b, c_s \rangle$. Here, $c_s$ can be found by evaluating an integral numerically.

# 7

---

# Applications to Stochastic Processes

---

Parzen is credited with introducing RKHS theory into statistics [48, 49, 50, 51]. He learnt about RKHSs by necessity, as a graduate student, back when the theory was in its infancy, because it was his prescribed minor thesis topic [46]. A decade later, Kailath demonstrated to the signal processing community the relevance of RKHS theory in a series of lucid papers concerning detection, filtering and parameter estimation [32, 33, 16, 17, 35]. Kailath and Parzen were colleagues at Stanford [46].

Despite Hilbert spaces being linear spaces, RKHS theory is not confined to studying *linear* estimation theory; a clever use of characteristic functions leads to a RKHS theory for optimal nonlinear estimators [29, 36, 16, 68]. Robust estimation theory has also been considered [7].

A RKHS approach usually entails finding an expression for the kernel. If this cannot be done, Kailath observes it may suffice to use a numerical scheme for computing the norms of functions on a RKHS [34].

The ideas underlying the application of RKHS theory to stochastic processes are presented without regard to technical considerations (such as measurability).

## 7.1 Detection in Finite Dimensions

Let $m \in \mathbb{R}^n$ be a known vector called the message. Let $w \in \mathbb{R}^n$ be a realisation of Gaussian white noise, that is, $w \sim N(0, I)$. The detection problem is deciding whether the received signal $y \in \mathbb{R}^n$ represents only noise, $y = w$, or the message plus noise, $y = m + w$.

The standard mechanism for making such a decision is the likelihood ratio test. The likelihood of $y$ being just noise is proportional to $\exp\{-\frac{1}{2}\|y\|^2\}$ while the likelihood of $y$ being the message plus noise is proportional to $\exp\{-\frac{1}{2}\|y - m\|^2\}$, where the constant of proportionality is the same in both cases. The message is therefore deemed to be present if and only if $\|y - m\| < \|y\|$.

The test $\|y - m\| \lessgtr \|y\|$ has a simple geometric explanation. Interpret the case $y = w$ as sending the zero message: $y = 0 + w$. In $\mathbb{R}^n$, label the point $m$ as the message, label the origin as the zero message and label the point $y$ as the received signal. Since the distribution of $w$ is radially symmetric, the likelihood of $y = s + w$ is proportional to $\|y - s\|$ and decays monotonically. Deciding between $s = 0$ and $s = m$ comes down to deciding whether $y$ is closer to $0$ or $m$.

The geometry suggests $\langle y, m \rangle - \frac{1}{2}\langle m, m \rangle \lessgtr 0$ as an equivalent formulation of the likelihood ratio test, corresponding to projecting the signal $y$ onto the "message space" $m$.

Generalising to coloured noise $N(0, \Sigma)$ is achieved by replacing $w$ by $\Sigma^{\frac{1}{2}}w$. Testing between $y = \Sigma^{\frac{1}{2}}w$ and $y = m + \Sigma^{\frac{1}{2}}w$ is achieved by whitening the received signal, thus deciding between $\Sigma^{-\frac{1}{2}}y = w$ and $\Sigma^{-\frac{1}{2}}y = \Sigma^{-\frac{1}{2}}m + w$. Specifically, the test is

$$\langle \Sigma^{-\frac{1}{2}}y, \Sigma^{-\frac{1}{2}}m \rangle - \frac{1}{2}\langle \Sigma^{-\frac{1}{2}}m, \Sigma^{-\frac{1}{2}}m \rangle \lessgtr 0, \tag{7.1}$$

assuming of course $\Sigma$ is non-singular.

Intriguingly, (7.1) can be written as

$$\langle y, m \rangle_K - \frac{1}{2}\langle m, m \rangle_K \lessgtr 0, \tag{7.2}$$

where $\langle \cdot, \cdot \rangle_K$ is the inner product on the RKHS whose kernel is $K = \Sigma$. Indeed, from Example 2.1, the RKHS whose kernel is $K$ is $\mathbb{R}^n$ equipped with the inner product $\langle x, y \rangle_K = y^T \Sigma^{-1} x$.

The advantage of (7.2) is it can be shown to remain valid when $\Sigma$ is singular. Importantly, note that if $m$ does not lie in range($\Sigma^{\frac{1}{2}}$) = range($\Sigma$) then the detection problem can be solved deterministically. Otherwise, if $m$ is in range($\Sigma$) then $y$ will also lie in range($\Sigma$), hence (7.2) is meaningful because both $m$ and $y$ will be elements of the RKHS.

The RKHS structure gives a useful link between the geometry $\langle \cdot, \cdot \rangle_K$ and the statistics of the noise $\Sigma^{\frac{1}{2}}w$. This manifests itself in several ways, one of which is that the distribution of $\langle \Sigma^{\frac{1}{2}}w, z \rangle_K$ is zero-mean Gaussian with variance $\langle z, z \rangle_K$, and moreover, the correlation between any two such Gaussian random variables is given by

$$E[\,\langle \Sigma^{\frac{1}{2}}w, x \rangle_K \, \langle \Sigma^{\frac{1}{2}}w, z \rangle_K\,] = \langle x, z \rangle_K. \qquad (7.3)$$

In particular, with respect to $\langle \cdot, \cdot \rangle_K$, the distribution of $\Sigma^{\frac{1}{2}}w$ is spherically symmetric, hence the ratio test reduces to determining whether $y$ is closer to 0 or $m$ with respect to the RKHS norm $\|\cdot\|_K$.

## 7.2   The RKHS Associated with a Stochastic Process

A stochastic process $X$ is a collection of random variables $X(t)$ indexed by a parameter $t \in T$ often taken to be time. An archetypal process is the Wiener process, described from first principles in [43], among other introductory material. A sample path is a realisation $t \mapsto X(t)$ of $X$, also known as a randomly chosen signal or waveform.

Despite RKHSs being function spaces and sample paths being functions, a RKHS structure is not given to the space of all possible sample paths. Since linear filters form linear combinations of the $X(t)$, the Hilbert space approach [62] to filtering endows the space spanned by the individual $X(t)$ with the inner product $\langle X(s), X(t) \rangle = E[X(s)X(t)]$. This inner product equates the statistical concept of conditional expectation with the simpler geometric concept of projection. For example, uncorrelated random variables are orthogonal to each other.

The RKHS approach goes one step further. The space spanned by the $X(t)$ has a useful geometry but no convenient coordinate system. A point in that space is nothing more than a point representing a random variable, whereas a point in $\mathbb{R}^n$ comes with $n$ coordinates describing the point's precise location. In hindsight, the RKHS approach can be

understood as giving a convenient coordinate system to the space of random variables. (While other authors have considered a RKHS approach to be coordinate free, we argue $\langle f, K(\cdot, x) \rangle$ is precisely the $x$th coordinate of $f$. RKHS theory is a blend of coordinate-free geometric reasoning with algebraic manipulations in pointwise coordinates.)

The coordinate system introduced by the RKHS approach assigns to each random variable $U$ the function $t \mapsto E[UX(t)]$. The $t$th coordinate is $E[UX(t)]$. This provides explicitly a wealth of information about $U$ in a form compatible with the geometry of the space.

This coordinate system comes from taking the covariance function $R(t, s) = E[X(t)X(s)]$, which is always positive semi-definite, to be the kernel of a RKHS $V \subset \mathbb{R}^T$. Recall from §5 that this construction arranges the individual points $X(t)$ in $V$ according to the geometry $\langle X(s), X(t) \rangle = E[X(s)X(t)]$. Precisely, the element $R(\cdot, s) \in V$ represents the random variable $X(s)$. The space $V$ is the completion of the space spanned by the $R(\cdot, s)$ for $s \in T$. Any (mean-square) limit $U$ of finite linear combinations of the $X(t)$ can therefore be represented in $V$ by a limit $E[UX(t)]$ of corresponding finite linear combinations of the $R(\cdot, s)$. This map $U \mapsto E[UX(t)]$ corresponds to the isometric isometry $\phi$ in §5.1. In other words, associated with the random variable $U$ are the coordinates $t \mapsto E[UX(t)]$ describing the location of the embedding of $U$ in the RKHS $V$.

While the elements of $V$ and sample paths are both functions on $T$, the former represents the coordinates $t \mapsto E[UX(t)]$ of a random variable $U$ while the latter represents a realisation of $t \mapsto X(t)$. Especially since the former is deterministic and the latter stochastic, there is not necessarily any relationship between the two. Surprisingly then, several relationships have been found for Gaussian processes. In typical cases, the sample paths are elements of the RKHS with probability zero or one, depending on the process itself [15, 42]. A detection problem is non-singular if and only if the signal belongs to $V$ [32]. (This is a manifestation of the Cameron-Martin Theorem, with the RKHS $V$ being known as the Cameron-Martin space [13].)

## 7.3   Signal Detection

Salient features of the RKHS approach to signal detection are summarised below, the purpose being to add substance to the higher level description above. Greater detail can be found in [32].

The signal detection problem is a hypothesis test for deciding between the null hypothesis that the signal $X(t)$ is composed of just noise — $X(t) = N(t)$ — versus the hypothesis that a known signal $m(t)$ is present — $X(t) = m(t) + N(t)$. For concreteness, $t$ is assumed to belong to the interval $T = [0, 1]$. The noise process $N(t)$ is assumed to be a zero-mean second-order Gaussian process. Its covariance function $R(t, s) = E[N(s)N(t)]$ is assumed to be continuous (and hence bounded) on $T \times T$. (This is equivalent to $N(t)$ being mean-square continuous.) Example calculations will take $N(t)$ to be the standard Wiener process [43] with $N(0) = 0$ and covariance function $R(t, s) = \min\{s, t\}$.

Hypothesis testing generally reduces to comparing the likelihood ratio against a threshold [56]. Two methods for determining the likelihood ratio are discussed below.

### 7.3.1   The Karhunen-Loéve Approach

The continuous-time detection problem has a different flavour from the conceptually simpler discrete-time detection problem when $T$ is a finite set. If $m(t)$ varies faster than a typical sample path of $N(t)$, and especially if $m(t)$ contains a jump, correct detection is possible with probability one [56]. Otherwise, if exact detection is not possible, the continuous-time detection problem can be solved as a limit of discrete-time detection problems; being continuous, a sample path of $N(t)$ is fully defined once its values at rational times $t$ are known.

A more practical representation of $N(t)$ using a countable number of random variables is the Karhunen-Loéve expansion. Provided $R(t, s)$ is continuous on the unit square $T \times T$, it has an eigendecomposition

$$R(t, s) = \sum_{k=1}^{\infty} \lambda_k \psi_k(t) \psi_k(s) \tag{7.4}$$

where the $\lambda_k \geq 0$ are the eigenvalues and the $\psi_k$ are the orthonormal

eigenfunctions defined by

$$\int_0^1 R(t,s)\psi(s)\,ds = \lambda\psi(t). \tag{7.5}$$

Orthonormality imposes the extra constraint

$$\int_0^1 \psi_n(t)\psi_m(t)\,dt = \begin{cases} 1 & m=n; \\ 0 & \text{otherwise.} \end{cases}$$

This is a generalisation of the eigendecomposition of a positive semi-definite matrix and is known as Mercer's theorem [45].

If $N(t)$ is the standard Wiener process then $R(t,s) = \min\{s,t\}$ and

$$\psi_k(t) = \sqrt{2}\sin\left(\left(k-\tfrac{1}{2}\right)\pi t\right),$$

$$\lambda_k = \frac{1}{\left(k-\tfrac{1}{2}\right)^2 \pi^2},$$

as verified in part by

$$\int_0^1 \min\{s,t\}\sqrt{2}\sin\left(\left(k-\tfrac{1}{2}\right)\pi s\right)\,ds$$

$$= \sqrt{2}\left[\int_0^t s\sin\left(\left(k-\tfrac{1}{2}\right)\pi s\right)\,ds + t\int_t^1 \sin\left(\left(k-\tfrac{1}{2}\right)\pi s\right)\,ds\right]$$

$$= \lambda_k\sqrt{2}\sin\left(\left(k-\tfrac{1}{2}\right)\pi t\right).$$

Note the $\lambda_k$ decay to zero. Define

$$N_k = \int_0^1 N(t)\psi_k(t)\,dt.$$

The orthonormality of the $\psi_k$ and $N(t)$ having zero mean imply $E[N_k] = 0$, $E[N_k^2] = \lambda_k$ and $E[N_n N_m] = 0$ for $n \neq m$. Indeed,

$$E[N_n N_m] = \int_0^1 \int_0^1 R(t,s)\psi_n(s)\psi_m(t)\,ds\,dt$$

$$= \int_0^1 \lambda_n\psi_n(t)\psi_m(t)\,dt$$

$$= \begin{cases} \lambda_n & m=n; \\ 0 & \text{otherwise.} \end{cases}$$

Since $N(t)$ is a Gaussian process, the $N_k$ are actually independent Gaussian random variables with zero mean and variance $\lambda_k$.

The Karhunen-Loéve expansion of $N(t)$ is

$$N(t) = \sum_{k=1}^{\infty} N_k \psi_k(t) \qquad (7.6)$$

where it is simplest to interpret the right-hand side as a way of generating a Wiener process: generate $N_k \sim N(0, \lambda_k)$ at random then form the summation. See [76] for a proof.

Assume the signal $m(t)$ looks sufficiently like the noise that it too can be represented as $m(t) = \sum_{k=1}^{\infty} m_k \psi_k(t)$. Substituting this expansion of $m(t)$ into $\int_0^1 m(t) \psi_n(t)\, dt$, and recalling the orthonormality of the $\psi_k$, shows $m_k = \int_0^1 m(t) \psi_k(t)\, dt$.

Transforming the received signal $X(t)$ in the same way leads to a corresponding hypothesis test for distinguishing between $X_k = N_k$ and $X_k = m_k + N_k$, where $X_k = \int_0^1 X(t) \psi_k(t)\, dt$. Since the $\lambda_k$ decrease to zero, the likelihood ratio can be approximated by using only a finite number of the $X_k$. For details, see [56].

The use of (Lebesgue measure) *ds* in (7.5) is *ad hoc* [32]. Changing it will result in a different discrete-time approximation of the detection problem. An equivalent hypothesis test will still result, nevertheless, avoiding an *ad hoc* choice is preferable. Moreover, computing an eigendecomposition (7.5) of $R(t, s)$ can be difficult [32]. The RKHS approach avoids the need for computing an eigendecomposition of $R(t, s)$.

### 7.3.2   The RKHS Approach

The logarithm of the likelihood ratio of $X(t) = m(t) + N(t)$ to $X(t) = N(t)$ when $N(t)$ is a Gaussian process takes the form

$$\int_0^1 X(t)\, dH(t) - \frac{1}{2} \int_0^1 m(t)\, dH(t) \qquad (7.7)$$

if there exists an $H(t)$ (of bounded variation) satisfying

$$\int_0^1 R(t, s) dH(s) = m(t). \qquad (7.8)$$

This is a consequence of a martingale representation theorem; see [55, Theorem 2.3], [56, Proposition VI.C.2] and [76, Section 7.2] for details. In other words, the linear operator $X \mapsto \int_0^1 X(t)\,dH(t)$ reduces the signal detection problem to the one-dimensional problem of testing whether a realisation of $\int_0^1 X(t)\,dH(t)$ came from $\int_0^1 N(t)\,dH(t)$ or $\int_0^1 m(t) + N(t)\,dH(t)$.

Directly obtaining a closed-form solution to (7.8) can be difficult. The RKHS approach replaces this problem by the sometimes simpler problem of computing the inner product of a particular RKHS; as discussed later, the RKHS approach leads to the mnemonic

$$\langle X, m \rangle - \frac{1}{2}\langle m, m \rangle \tag{7.9}$$

for the likelihood ratio, where the inner product is that of the RKHS whose kernel is the covariance $R(t, s)$. (See Example 7.2.)

Let $U = \int_0^1 N(t)\,dH(t)$. Being the limit of certain finite linear combinations of the $N(t)$, it belongs to the Hilbert space (of square-integrable random variables) generated by the $N(t)$. Since each $N(t)$ is zero-mean Gaussian, so is $U$. Its distribution is thus fully determined by its "coordinates" $t \mapsto E[UN(t)]$. Since

$$E[UN(t)] = \int_0^1 E[N(s)N(t)]\,dH(s)$$
$$= \int_0^1 R(t, s)\,dH(s),$$

if the requirement (7.8) holds then $U = \int_0^1 N(t)\,dH(t)$ satisfies

$$E[UN(t)] = m(t). \tag{7.10}$$

The existence of an $H(t)$ satisfying (7.8) is only a sufficient condition for the signal detection problem to be non-singular. The existence of a random variable $U$, belonging to the Hilbert space generated by the $N(t)$ and satisfying (7.10), is both necessary and sufficient [32]. The discrepancy arises because a $U$ satisfying (7.10) need not be expressible in the form $U = \int_0^1 N(t)\,dH(t)$; see [55, p. 44] for an example.

Let $V \subset \mathbb{R}^T$ be the RKHS whose kernel is the noise covariance function $R(t, s) = E[N(s)N(t)]$. Recall that the Hilbert space generated by the $N(t)$ is isometrically isomorphic to $V$, with $U \mapsto E[UN(\cdot)]$

an isometric isomorphism. Therefore, (7.10) has a solution if and only if $m$ is an element of the RKHS $V$.

Condition (7.8) requires expressing $m(t)$ as a linear combination of the $R(\cdot, s)$, but with the caveat that the linear combination take the form of a particular integral. The requirement that $m$ be an element of $V$ is that $m$ is the limit of finite linear combinations of the $R(\cdot, s)$ where the limit is with respect to the norm on $V$, defined implicitly by $\langle R(\cdot, s), R(\cdot, t) \rangle = R(t, s)$.

If $U$ satisfies (7.10) then the likelihood ratio can be found as follows. Write $U$ as $U(N)$ to emphasise that $U$ can be considered to be an almost surely linear function of the $N(t)$. The likelihood ratio is $U(X) - \frac{1}{2}U(m)$. This agrees with (7.7) if $U(N) = \int_0^1 N(t)\, dH(t)$.

**Example 7.1.** With notation as above, assume $m(t) = \sum_i \alpha_i R(\cdot, s_i)$. A solution to (7.8) is obtained by choosing $H(s)$ to have increments at the $s_i$ of height $\alpha_i$. This leads to $\int_0^1 X(t)\, dH(t) = \sum_i \alpha_i X(s_i)$ and $\int_0^1 m(t)\, dH(t) = \sum_i \alpha_i m(s_i)$. Alternatively, let $\phi$ be the linear isomorphism $U \mapsto E[UN(\cdot)]$, so that (7.10) becomes $\phi(U) = m$. Then $U(N) = \phi^{-1}(m) = \sum_i \alpha_i \phi^{-1}(R(\cdot, s_i)) = \sum_i \alpha_i N(s_i)$. Thus $U(X) = \sum_i \alpha_i X(s_i)$ and $U(m) = \sum_i \alpha_i m(s_i)$, resulting in the same likelihood ratio.

A more explicit expression can be given for the second term of the likelihood ratio (7.7) by using the inner product on the RKHS $V$. Let $U(f) = \int_0^1 f(t)\, dH(t)$ where $H(t)$ satisfies (7.8), that is, $U(R(t, \cdot)) = m(t)$. Assume $U(f)$ for $f \in V$ can be expressed as $\langle f, g \rangle$ where $g \in V$ is to be determined. The requirement (7.8) becomes $\langle R(t, \cdot), g \rangle = m(t)$, whose unique solution is $g = m$. In particular, the second term of the likelihood ratio is proportional to

$$\int_0^1 m(t)\, dH(t) = \langle m, m \rangle = \|m\|^2. \qquad (7.11)$$

This suggests, and it turns out to be true, that $U(m) = \langle m, m \rangle$ in general, justifying the second part of (7.9). (Even if $U$ cannot be written in the form $\int \cdot\, dH$, it can be approximated arbitrarily accurately by such an expression [55].)

Equation (7.11) can be verified to hold in Example 7.1. The square of the RKHS norm of $m(t)$ is

$$\|m\|^2 = \left\langle \sum_i \alpha_i R(\cdot, s_i), \sum_j \alpha_j R(\cdot, s_j) \right\rangle$$
$$= \sum_{i,j} \alpha_i \alpha_j R(s_j, s_i).$$

The integral can be written as

$$\int_0^1 m(t)\, dH(t) = \sum_j \alpha_j m(s_j)$$
$$= \sum_j \alpha_j \sum_i \alpha_i R(s_j, s_i).$$

These two expressions are the same, as claimed in (7.11).

The first term of (7.7) is more subtle. It naively equals $\langle X, m \rangle$, leading to (7.9). The difficulty is that a realisation of $X(t)$ is almost surely not an element of the RKHS $V$, making $\langle X, m \rangle$ technically meaningless. As advocated in [32] though, the mnemonic $\langle X, m \rangle$ is useful for guessing the solution then verifying its correctness by showing it satisfies (7.10). Moreover, it is intuitive to interpret the detection problem as $X \mapsto \langle X, m \rangle$ since this projects the received signal $X$ onto the signal $m$ to be detected. The RKHS inner product gives the optimal projection.

**Example 7.2.** Assume $N(t)$ is the standard Wiener process with $N(0) = 0$ and $R(t, s) = \min\{s, t\}$. The associated RKHS $V$ is given in Example 4.3. The detection problem is non-singular if and only if $m \in V$. Assuming $m \in V$, the second term of the likelihood ratio (7.7) is $\frac{1}{2}\|m\|^2 = \frac{1}{2}\int_0^1 (m'(s))^2\, ds$. The mnemonic $\langle X, m \rangle$ suggests the first term of (7.7) is $\int_0^1 X'(s)m'(s)\, ds$. This is technically incorrect because a Wiener process is nowhere differentiable, that is, $X'(t)$ is not defined. However, assuming $m$ is sufficiently smooth and applying integration by parts leads to the alternative expression $m'(1)X(1) - \int_0^1 X(s)m''(s)\, ds$. If $X$ lies in $V$ then $\int_0^1 X'(s)m'(s)\, ds$ and $m'(1)X(1) - \int_0^1 X(s)m''(s)\, ds$ are equal. The advantage of the latter is it makes sense when $X$ is a realisation of a Wiener process. To verify this is correct, let

$U = m'(1)N(1) - \int_0^1 N(s)m''(s)\,ds$. Then

$$\begin{aligned}
E[UN(t)] &= m'(1)R(t,1) - \int_0^1 R(t,s)m''(s)\,ds \\
&= m'(1)t - \int_0^t s\,m''(s)\,ds - t\int_t^1 m''(s)\,ds \\
&= m'(1)t - (tm'(t) - m(t)) - t(m'(1) - m'(t)) \\
&= m(t),
\end{aligned}$$

verifying the correct linear operator has been found. Here, integration by parts has been used, as has the fact that $m(0) = 0$ because $m$ is assumed to be in $V$. (If $m$ is not sufficiently smooth to justify integration by parts, it is possible to write $U$ in terms of a stochastic integral.)

It remains to demystify how the technically nonsensical $\langle X, m \rangle$ can lead to the correct linear operator being found. At the heart of the matter is the concept of an abstract Wiener space. For the Wiener process in Example 7.2 the setting is the following. Let $\mathcal{B}$ be the Banach space of continuous functions $f : [0,1] \to \mathbb{R}$ satisfying $f(0) = 0$, equipped with the sup norm $\|f\| = \sup_{t \in [0,1]} |f(t)|$. Let $\mathcal{H}$ be the (reproducing kernel) Hilbert space $V$ from Example 7.2. Treated as sets, observe that $\mathcal{H} \subset \mathcal{B}$. The inclusion map $\iota : \mathcal{H} \to \mathcal{B}$ is continuous because the norm on $\mathcal{B}$ is weaker than the norm on $\mathcal{H}$. Furthermore, $\iota(\mathcal{H})$ is dense in $\mathcal{B}$. The dual map $\iota^* : \mathcal{B}^* \to \mathcal{H}^*$ is injective and dense; here, $\mathcal{B}^*$ is the space of all bounded linear functionals on $\mathcal{B}$, and $\iota^*$ takes a functional $\mathcal{B} \mapsto \mathbb{R}$ and returns a functional $\mathcal{H} \mapsto \mathbb{R}$ by composing the original functional with $\iota$. Put simply, a linear functional on $\mathcal{B}$ is *a fortiori* a linear functional on $\mathcal{H}$ because $\mathcal{H}$ is a subset of $\mathcal{B}$. (The functions $\iota$ and $\iota^*$ are used to account for the different topologies on the two spaces.)

Every realisation of the noise process $N$ lies in $\mathcal{B}$. If $g \in \mathcal{B}^*$ is a linear functional then, as in the finite dimensional case (§7.1), $g(N)$ is a zero-mean Gaussian random variable. For technical reasons, the geometry of $\mathcal{B}$ cannot correspond to the statistics of $N$. Instead, it is the geometry of $\mathcal{H}$ that corresponds to the statistics of $N$. The variance of $g(N)$ is $\|\iota^*(g)\|_{\mathcal{H}}$. In particular, the Gaussian distribution of $N$ in $\mathcal{B}$ is radially symmetric with respect to the norm $g \mapsto \|\iota^*(g)\|_{\mathcal{H}}$. Believably then, the detection problem can be solved by projecting $X$ onto the

one-dimensional space spanned by the message $m$, as in §7.1.

If there exists a $g \in \mathcal{B}^*$ such that $\iota^*(g)$ is the linear operator $\cdot \mapsto \langle \cdot, m \rangle$ then $g(X)$ is the precise definition of the mnemonic $\langle X, m \rangle$ in (7.9). This corresponds to the case when $m$ was assumed to be sufficiently smooth in Example 7.2. Otherwise, the extension of $\cdot \mapsto \langle \cdot, m \rangle$ from $\mathcal{H}$ to $\mathcal{B}$ is not a *bounded* linear functional. It turns out though that there exists a subspace $\mathcal{B}'$ of $\mathcal{B}$ such that $\langle \cdot, m \rangle$ extends to $\mathcal{B}'$ and the probability of $N$ lying in $\mathcal{B}'$ is unity. In particular, $\langle X, m \rangle$ can be defined to be the extension of $\langle \cdot, m \rangle$ by continuity: if $X_k \to X$ in $\mathcal{B}$, where the $X_k$ belong to $\mathcal{H}$, then $\langle X, m \rangle$ is defined to be $\lim_k \langle X_k, m \rangle$. This limit exists with probability one. (Mathematically, $\mathcal{B}^*$ can be treated as a subset of $L^2(\mathcal{B}; \mathbb{R})$, where the $L^2$-norm of $g \colon \mathcal{B} \to \mathbb{R}$ is the square-root of $E[g(N)^2]$. Then $(\iota^*)^{-1} \colon \iota^*(\mathcal{B}^*) \subset \mathcal{H}^* \to \mathcal{B}^* \subset L^2(\mathcal{B}; \mathbb{R})$ extends by continuity to a linear isometry from $\mathcal{H}^*$ to $L^2(\mathcal{B}; \mathbb{R})$. The image of $\langle \cdot, m \rangle$ under this extension is the true definition of $X \mapsto \langle X, m \rangle$.)

**Remark**   Here is another perspective. By [55, Lemma 1.2], the process $N(t)$ lies in $\mathcal{H} = L^2([0,1])$ with probability one: $\int_0^1 N^2(t) \, dt < \infty$ almost surely. Furthermore, the same lemma proves that if the detection problem is non-singular then $m(t)$ lies in $\mathcal{H}$. Therefore, without loss of generality, the sample space can be taken to be $\mathcal{H}$. As in [13], a Gaussian measure $N_{\mu, Q}$ can be placed on the separable Hilbert space $\mathcal{H}$, where $\mu \in \mathcal{H}$ and $Q \colon \mathcal{H} \to \mathcal{H}$ represent the mean and covariance of the Gaussian measure. The signal detection problem is deciding between the probability measures $N_{0,Q}$ and $N_{m,Q}$. The Cameron-Martin formula [13, Theorem 2.8] yields the likelihood ratio. Furthermore, the image of $\mathcal{H}$ under the square-root of $Q$, denoted $Q^{\frac{1}{2}}(\mathcal{H})$, is precisely the RKHS $V$ discussed earlier, and goes by the name Cameron-Martin space.

# 8

## Embeddings of Random Realisations

Random variables are normally taken to be real-valued [75]. Random vectors and random processes are collections of random variables. Random variables can be geometrically arranged in a RKHS according to pairwise correlations, as described in §7.2. A random variable is mapped to a single point in the RKHS.

A different concept is embedding the *realisations* of a random variable into a RKHS. Let $X$ be an $\mathbb{X}$-valued random variable encoding the original random variables of interest. For example, the random variables $Y$ and $Z$ are encoded by $X = (Y, Z)$ and $\mathbb{X} = \mathbb{R}^2$. Let $V \subset \mathbb{R}^{\mathbb{X}}$ be a RKHS whose elements are real-valued functions on $\mathbb{X}$. The realisation $x \in \mathbb{X}$ of $X$ is mapped to the element $K(\cdot, x)$ of $V$, where $K$ is the kernel of $V$.

The kernel does not encode statistical information as in §7.2. Its purpose is "pulling apart" the realisations of $X$. A classic example is embedding the realisations of $X = (Y, Z) \in \mathbb{R}^2$ into an infinite-dimensional RKHS by using the Gaussian kernel (§5.2). *Linear* techniques for classification and regression applied to the RKHS correspond to *nonlinear* techniques in the original space $\mathbb{R}^2$.

If the realisations are pulled sufficiently far apart then a probability

distribution on $X$ can be represented uniquely by a point in the RKHS. The kernel is then said to be characteristic. Thinking of distributions as points in space can be beneficial for filtering, hypothesis testing and density estimation.

**Remark:** Information geometry also treats distributions as points in space. There though, the differential geometry of the space encodes intrinsic statistical information. The relative locations of distributions in the RKHS depend on the chosen kernel.

RKHS was first presented in the finite dimensional setting, and this choice is done here as well, for presenting the interest of embedding random variables in a RKHS. In section 8.2, the theory of random variables with values in Hilbert space is outlined as a preparation. Next, embeddings of random variables in RKHS are presented and studied from the probabilistic as well as the empirical points of view. Some generalizations and considerations needed in chapter 9 close the chapter.

## 8.1 Finite Embeddings

Definitions and concepts concerning the embedding of random realisations in a RKHS are more easily grasped when the number of possible outcomes is finite.

Let $X$ be an $\mathbb{X}$-valued random variable where $\mathbb{X} = \{1, 2, \cdots, n\}$. Its distribution is described by $p_1, \cdots, p_n$ where $p_x$ is the probability of $X$ taking the value $x$. Choose the RKHS to be Euclidean space $\mathbb{R}^n$. Its kernel is the identity matrix. The realisation $x \in \mathbb{X}$ is therefore embedded as the unit vector $e_x$. Effectively, the elements $e_1, \cdots, e_n$ of the RKHS are chosen at random with respective probabilities $p_1, p_2, \cdots, p_n$. This choice of kernel leads to the embedded points being uniformly spaced apart from each other: $\|e_i - e_j\| = \sqrt{2}$ for $i \neq j$.

**Remark:** This embedding is used implicitly in [18] to great advantage. The elements of $\mathbb{X} = \{1, \cdots, n\}$ are the possible states of a Markov chain. Representing the $i$th state by the vector $e_i$ converts nonlinear functions on $\mathbb{X}$ into linear functions on $\mathbb{R}^n$, as will be seen presently.

When $n > 1$, the distribution of $X$ cannot be recovered from its expectation $E[X] = p_1 + 2p_2 + \cdots + np_n$. However, if $\check{X}$ is the embedding

of $X$ then $E[\check{X}] = p_1 e_1 + p_2 e_2 + \cdots + p_n e_n$. The distribution of $X$ is fully determined by the expectation of its embedding! *Points in the RKHS can represent distributions.*

This has interesting consequences, the first being that the distribution of $X$ can be estimated from observations $\check{x}_1, \cdots, \check{x}_M \in \mathbb{R}^n$ by approximating $E[\check{X}]$ by the "sample mean" $\frac{1}{M} \sum_{m=1}^{M} \check{x}_m$. Hypothesis testing (§7.3) reduces to testing which side of a hyperplane the sample mean lies: the two distributions under test are thought of as two points in the RKHS. An observation, or the average of multiple observations, is a third point corresponding to the estimated distribution. Intuitively, the test is based on determining which of the two distributions is closer to the observation. The following lemma validates this. The probabilities $p_i$ are temporarily written as $p(i)$.

**Lemma 8.1.** Let $p$ and $q$ be two probability distributions on the $\mathbb{X}$-valued random variable $X$, where $\mathbb{X} = \{1, \cdots, n\}$. Given $M$ observations $x_1, \cdots, x_M$ of $X$, the log likelihood ratio $\ln \left\{ \prod_{m=1}^{M} \frac{p(x_m)}{q(x_m)} \right\}$ equals $\langle v, \sum_{m=1}^{M} \check{x}_m \rangle$ where $v_i = \ln \frac{p(i)}{q(i)}$.

*Proof.* If $x_m = i$ then $\check{x}_m = e_i$ and $\langle v, \check{x}_m \rangle = v_i = \ln \frac{p(i)}{q(i)} = \ln \frac{p(x_m)}{q(x_m)}$. Therefore $\ln \left\{ \prod \frac{p(x_m)}{q(x_m)} \right\} = \sum \ln \frac{p(x_m)}{q(x_m)} = \sum \langle v, \check{x}_m \rangle$. □

More advanced probability texts tend to favour working with expectations (and conditional expectations) over probabilities (and conditional probabilities). This generally involves expressions of the form $E[f(X)]$ where $f \colon \mathbb{X} \to \mathbb{R}$ is a "test" function; if $E[f(X)]$ is known for sufficiently many functions $f$ then the distribution of $X$ can be inferred.

An arbitrary function $f \colon \mathbb{X} \to \mathbb{R}$ is fully determined by its values $f(1), \cdots, f(n)$. Let $\check{f} \colon \mathbb{R}^n \to \mathbb{R}$ be the unique *linear* function on $\mathbb{R}^n$ satisfying $\check{f}(e_x) = f(x)$ for $x \in \mathbb{X}$. Then

$$E[\check{f}(\check{X})] = \sum_x p_x \check{f}(e_x) = \sum_x p_x f(x) = E[f(X)]. \tag{8.1}$$

The advantage of $E[\check{f}(\check{X})]$ over $E[f(X)]$ is that $E[\check{f}(\check{X})] = \check{f}(E[\check{X}])$ because $\check{f}$ is linear. Once $E[\check{X}]$ is known, $E[f(X)] = \check{f}(E[\check{X}])$ is effectively known for any $f$. By comparison, $E[f(X)]$ generally cannot be deduced from $E[X]$.

Correlations take the form $E[f(X)\,g(X)]$. As above, this is equivalent to $E[\check{f}(\check{X})\,\check{g}(\check{X})]$. Being linear, $\tilde{f}(\check{X})$ and $\tilde{g}(\check{X})$ can be rewritten using the Euclidean inner product: there exist $v_f, v_g \in \mathbb{R}^n$ such that $\check{f}(\check{X}) = \langle \check{X}, v_f \rangle$ and $\check{g}(\check{X}) = \langle \check{X}, v_g \rangle$. Moreover, $v_f$ and $v_g$ can be expressed as linear combinations $\sum_i \alpha_i e_i$ and $\sum_j \beta_j e_j$ of the unit basis vectors $e_1, \cdots, e_n$. In particular,

$$E[f(X)\,g(X)] = E[\langle \check{X}, \sum_i \alpha_i e_i \rangle \langle \check{X}, \sum_j \beta_j e_j \rangle] \qquad (8.2)$$

$$= \sum_{i,j} \alpha_i \beta_j E[\langle \check{X}, e_i \rangle \langle \check{X}, e_j \rangle]. \qquad (8.3)$$

Once the "correlation" matrix $C$ given by $C_{ij} = E[\langle \check{X}, e_i \rangle \langle \check{X}, e_j \rangle]$ is known, $E[f(X)\,g(X)]$ can be determined from (8.3).

## 8.2 Random elements in Hilbert spaces

Prior studying the case of the embedding in a RKHS, let us consider the case of random variables which take values in a (separable) Hilbert space $V$ [9, 21]. A lemma states that a function from a probability space $(\Omega, \mathcal{F}, P)$ to $V$ is a random variable with values in $V$ if and only if $x^*(X)$ is a real random variable for any linear form $x^* \in V^*$, the dual of $V$. Since the dual of a Hilbert space can be identified to itself, the linear form simply writes $x^*(X) = \langle x, X \rangle$ where $x \in V$.

The linear form on $V$ defined by $\ell_X(x) = E[\langle x, X \rangle]$ is bounded whenever $E[\|X\|] < +\infty$. Indeed, $|\ell_X(x)| \le \|x\|E\|X\|$ thanks to the Cauchy-Schwartz inequality. Thus, Riesz representation theorem [14] shows the existence of a unique element $m_X$ of $V$ such that $E[\langle x, X \rangle] = \langle x, m_X \rangle$. $m_X$ is the mean element and is denoted as $E[X]$.

Denote the space of square integrable random elements of $V$ as $L_V^2(P)$ (a short notation for $L_V^2(\Omega, \mathcal{F}, P)$.) It is the space of $V$ valued random variables on $(\Omega, \mathcal{F}, P)$ such that $E\|X\|^2 < +\infty$. When equipped with $\langle X, Y \rangle_{L^2} := E[\langle X, Y \rangle_V]$, $L_V^2(P)$ is itself a Hilbert space.

The covariance operator is a linear operator from $V$ to $V$ defined by $\Sigma_X : x \longmapsto \Sigma_X(x) := E[\langle x, X - m_X \rangle (X - m_X)]$. It is bounded whenever $X \in L_V^2(P)$. To see this, suppose for the sake of simplicity that $EX = 0$. Recall that that the operator norm is defined as $\|\Sigma_X\| =$

$\sup_{\|x\| \le 1} \|\Sigma_X(x)\|$ and that $\|\Sigma_X\| = \sup_{\|x\| \le 1, \|y\| \le 1} |\langle y, \Sigma_X(x) \rangle|$. But applying Cauchy-Schwartz inequality leads to

$$
\begin{aligned}
|\langle y, \Sigma_X(x) \rangle| &= |E[\langle x, X \rangle \langle y | X \rangle]| \\
&\le E[\|x\| \|y\| \|X\|^2] = \|x\| \|y\| E[\|X\|^2]
\end{aligned}
$$

which shows that $\|\Sigma_X\| < +\infty$ whenever $X \in L_V^2(P)$ .

Likewise, we can define a cross-covariance operator between two elements $X, Y$ of $L_V^2(P)$ by the bounded linear operator from $V$ to itself defined by $\Sigma_{YX}(x) := E[\langle x, (X - m_X) \rangle (Y - m_Y)]$. The adjoint operator defined by $\langle \Sigma_{YX}^*(y), x \rangle = \langle y, \Sigma_{YX}(x) \rangle$ is then $\Sigma_{XY}$ since by definition $\Sigma_{XY}(y) = E[\langle y, (Y - m_Y) \rangle (X - m_X)]$. The two operators are completely defined by $\langle y, \Sigma_{XY}(x) \rangle = E[\langle x, X \rangle \langle y, Y \rangle]$ (if the mean elements are assumed to be equal to zero.) The cross-covariance can even be generalized to the case of two different Hilbert spaces. Consider two random variables $X$ and $Y$ defined on a common probability space $(\Omega, \mathcal{F}, P)$, but taking values in two different Hilbert spaces $V_x$ and $V_y$. The cross-covariance operator has the same definition, but $\Sigma_{YX}$ has a domain of definition included in $V_x$ and a range included in $V_y$.

Covariance and cross-covariance operators have furthermore the properties to be nuclear as well as Hilbert-Schmidt operators. A Hilbert-Schmidt operator $\Sigma$ from a Hilbert space $V_1$ to another $V_2$ is such that $\sum_i \|\Sigma e_i\|^2 = \sum_{ij} \langle f_j, \Sigma e_i \rangle^2 < +\infty$ where $\{e_i\}, \{f_j\}$ are orthornomal bases of respectively $V_1$ and $V_2$. $\|\Sigma\|_{HS} = \sum_i \|\Sigma e_i\|^2 = \sum_{i,j} \langle f_j, \Sigma e_i \rangle^2$ is the Hilbert-Schmidt norm (it can be shown independent of the choice of the bases.)

$\Sigma$ is nuclear if there exist orthornomal bases $\{e_i\}, \{f_j\}$ and a sequence $\{\lambda_i\}$ verifying $\sum_i |\lambda_i| < +\infty$ such that $\Sigma = \sum_i \lambda_i e_i \otimes f_i$, where the tensorial product is defined as $(e_i \otimes f_i)(x) = \langle x, e_i \rangle f_i$. Then $\|\Sigma\|_N = \sum_i |\lambda_i|$ is the nuclear norm. It also holds $\|\Sigma\|_{HS} = \sum_j |\lambda_i|^2$. Furthermore, the three norms of operators so far introduced satisfy the inequalities $\|\Sigma\| \le \|\Sigma\|_{HS} \le \|\Sigma\|_N$.

## 8.3 Random elements in reproducing kernel Hilbert spaces

The theory of random elements in Hilbert space is specialized in the sequel to the case where $V$ is a reproducing kernel Hilbert space. Specif-

ically, the embedding of an $\mathbb{X}$ valued random variable is developed.

Consider the embedding of an $\mathbb{X}$ valued random variable $X$ into a RKHS $V \subset \mathbb{R}^{\mathbb{X}}$ defined by its kernel $K \colon \mathbb{X} \times \mathbb{X} \to \mathbb{R}$. This mean that each realisation $x \in X$ is mapped to the element $K(\cdot, x) \in V$. This is equivalent to defining a random variable $\check{X} = K(\cdot, X)$ with values in the RKHS $V$.

The expected value of $\check{X} = K(\cdot, X)$ was defined coordinate-wise in §8.1. This generalises immediately to defining $E[\check{X}]$ by

$$\langle E[\check{X}], K(\cdot, z) \rangle = E[\langle \check{X}, K(\cdot, z) \rangle], \qquad z \in \mathbb{X}, \tag{8.4}$$

where the left-hand side is the $z$th coordinate of $E[\check{X}]$ and the right-hand side is the expected value of the $z$th coordinate of $\check{X}$. Evaluating $E[\check{X}]$ is straightforward in principle: if $E[X] = \int x \, d\mu(x)$ then $E[\langle \check{X}, K(\cdot, z) \rangle] = \int \langle K(\cdot, x), K(\cdot, z) \rangle \, d\mu(x) = \int K(z, x) \, d\mu(x)$.

**Remark:** For $E[\check{X}]$ to be well-defined, $x \mapsto \langle K(\cdot, x), K(\cdot, z) \rangle = K(z, x)$ must be measurable for every $z \in \mathbb{X}$. By symmetry of the kernel, this is the same as $K(\cdot, x)$ being measurable for every $x \in \mathbb{X}$. By [10, Theorem 90], this is equivalent to every element of the RKHS being measurable.

While (8.4) defines $E[\check{X}]$ coordinate-wise, a subtlety is whether $E[\check{X}] \in \mathbb{R}^{\mathbb{X}}$ belongs to the RKHS $V$. It is usually desirable to require $E[\|\check{X}\|]$ to be finite, and this suffices for $E[\check{X}]$ to be an element of the RKHS. Rather than prove this directly, an equivalent definition of $E[\check{X}]$ is given, using the theory outlined in the preceding section.

For $v \in V$, $|\langle \check{X}, v \rangle| \leq \|\check{X}\| \|v\|$ by the Cauchy-Schwartz inequality. Therefore, $E[|\langle \check{X}, v \rangle|] \leq E[\|\check{X}\|] \|v\|$. In particular, if $E[\|\check{X}\|] < \infty$ then $E[|\langle \check{X}, v \rangle|] < \infty$ and $E[\langle \check{X}, v \rangle]$ is well-defined. Moreover, $v \mapsto E[\langle \check{X}, v \rangle]$ is a bounded linear function. The Riesz representation theorem implies the existence of an element $m_X \in V$ such that this linear map is given by $v \mapsto \langle m_X, v \rangle$. Then $E[\check{X}]$ is defined to be $m_X$.

Henceforth, this definition of $E[\check{X}]$ will be adopted, so that $E[\check{X}]$ is an element of the RKHS by definition. The notations $m_X$ and $E[\check{X}]$ will be used interchangeably and will be called the *mean element* of $X$.

**Remark:** Condition (8.4) only required $\langle E[\check{X}], v \rangle = E[\langle \check{X}, v \rangle]$ to hold for $v = K(\cdot, z)$, $z \in \mathbb{X}$. However, the $K(\cdot, z)$ are dense in $V$, hence the two definitions agree provided $E[\|\check{X}\|] < \infty$.

Finally, let $f \in V$. The mean element completely determine $E[f(X)]$. This is easy to show because $E[f(X)] = E[\langle K(., X), f \rangle] = E[\langle \check{X}, f \rangle] = \langle m_X, f \rangle$. Therefore, knowing $m_x$ allows to evaluate the mean of any transformation of $X$, provided the transformation belongs to $V$.

It is convenient to complement $m_X$ with a covariance operator $\Sigma_X$ capturing $E[(f(X) - E[f(X)])(g(X) - E[g(X)])]$ for $f, g \in V$. Observe

$$E[(f(X) - E[f(X)])(g(X) - E[g(X)])]$$
$$= E[\langle f, \check{X} - m_X \rangle \langle g, \check{X} - m_X \rangle]$$
$$= E[\langle f, \langle g, \check{X} - m_X \rangle (\check{X} - m_X) \rangle]$$
$$= \langle f, E[\langle g, \check{X} - m_X \rangle (\check{X} - m_X)] \rangle$$
$$= \langle f, \Sigma_X(g) \rangle$$

where the linear operator $\Sigma_X \colon V \to V$ is given by

$$\Sigma_X(g) = E[\langle g, \check{X} - m_X \rangle (\check{X} - m_X)]. \tag{8.5}$$

Provided $E[\|\check{X}\|^2] < \infty$, $\Sigma_X$ is well-defined and its operator norm is finite: $\|\Sigma_X\| = \sup_{\|g\|=1} \|\Sigma_X(g)\| = \sup_{\|f\|=\|g\|=1} \langle f, \Sigma_X(g) \rangle$.

**Remark:** If $E[\|\check{X}\|^2] = E[K(X, X)] < \infty$ then $E[\|\check{X}\|] < \infty$, a consequence of Jensen's inequality [75, Section 6.7]. In particular, $E[\|\check{X}\|^2] < \infty$ ensures both $m_X$ and $\Sigma_X$ exist. See also [10, §4.5].

## 8.4   Universal and characteristic kernels

An aim of embedding points of a space into a RKHS is to reveal features that are difficult to study or see in the initial space. Therefore, the RKHS has to be sufficiently rich to be useful in a given application.

For example, the fact that the mean element reveals $E[f(X)]$ simply as $\langle m_x, f \rangle$ is an interesting property essentially if this can be applied for a wide variety of functions $f$. It is thus desirable that $V \subset \mathbb{R}^{\mathbb{X}}$ is a sufficiently rich function space. The richness of the RKHS is obviously provided by some properties of the kernel, since a kernel gives rise to a unique RKHS.

Two important notions occur: the notion of universality, defined initially by Steinwart [67] and the notion of characteristic kernel.

Universality is linked to the denseness of the RKHS into a target space of functions and is therefore linked to the ability of the functions of the RKHS to approximate functions in the target space. Since it depends on the kernel, on the initial space and on the target space, there exist different definitions of universality. We will need in the sequel kernel universal in the sense that their reproducing kernel Hilbert space is dense in the space of continuous functions. In some applications, universality refers to denseness in $L^p$ spaces. This is particularly important when dealing with embedding of square integrable random variables. Thus the following definition is considered.

**Definition 8.1.** Let $\mathbb{X}$ be a locally compact Hausdorff space (such as Euclidean space). Denote by $C_0(\mathbb{X})$ the class of real-valued continuous functions on $\mathbb{X}$ vanishing at infinity. Let $K \colon \mathbb{X} \times \mathbb{X} \to \mathbb{R}$ be a bounded kernel for which $K(\cdot, x) \in C_0(\mathbb{X})$ for all $x \in \mathbb{X}$. Then $K$, or its corresponding RKHS $V$, is *universal* if $V$ is dense in $C_0(\mathbb{X})$ with respect to the uniform norm. (Actually, universality has several different definitions depending on the class of functions of most interest [65].)

The notion of characteristic kernel is useful when embedding probability measures and is linked to the ability to discriminate two original measures in the RKHS. Let $K : \mathbb{X} \times \mathbb{X} \longrightarrow \mathbb{R}$ a kernel on $\mathbb{X}$. Let $\mathcal{P}$ the set of probability measures on $\mathbb{X}$ equipped with its Borel $\sigma$ algebra. Then $K$ (and thus $V$) is characteristic if and only if the map from $\mathcal{P}$ to $V$ defined by $P \mapsto E_P[K(., X)]$ is injective, where $X$ is any random variable distributed under $P$.

There exist links between the two notions which are studied in [65]. For this paper and applications in signal processing, it suffices to know that some well-known and often used kernels are universal and charac-

teristic. Importantly, it is shown in [65], Proposition 5, that for radial kernels on $\mathbb{R}^d$, kernels are universal (in many diffferent sense) if and only if they are characteristic. Furthermore the proposition gives other necessary and sufficient condition for universality. For example, strict positive-definiteness of the radial kernel insures it is characteristic.

Thus, important kernels such as the Gaussian kernel, $K(x,y) = \exp(-\sigma\|x-y\|^2)$, the inverse multiquadrics $K(x,y) = (c+\|x-y\|^2)^{-\beta}, \beta > d/2)$ are characteristic and universal.

## 8.5 Empirical estimates of the mean elements and covariance operators

The mean element $m_X$ can be estimated by the sample mean

$$\hat{m}_X^N = \frac{1}{N}\sum_{i=1}^{N} K(\cdot, x_i) \tag{8.6}$$

where $x_1, \cdots, x_N$ are independent realisations of $X$. The law of large numbers implies $\hat{m}_X^N$ converges to $m_X$ almost surely [31]. It also converges in quadratic mean:

$$E[\|\hat{m}_X^N - m_X\|^2] = \frac{1}{N^2} E\left[\left\langle \sum_{i=1}^{N} K(\cdot, X_i) - m_X, \sum_{j=1}^{N} K(\cdot, X_j) - m_X \right\rangle\right]$$

$$= \frac{1}{N^2}\sum_{i,j=1}^{N} E\left[K(X_j, X_i) - m_X(\check{X}_j) - m_X(\check{X}_i) + \|m_X\|^2\right]$$

$$= \frac{1}{N^2}\sum_{i,j=1}^{N} E\left[K(X_j, X_i) - \|m_X\|^2\right].$$

If $i \neq j$ then $E[K(X_j, X_i)] = \|m_X\|^2$ by independence:

$$E[K(X_j, X_i)] = E[\langle K(\cdot, X_i), K(\cdot, X_j)\rangle]$$
$$= \langle E[K(\cdot, X_k)], E[K(\cdot, X_j)]\rangle$$
$$= \langle m_X, m_X\rangle.$$

Therefore,

$$E[\|\hat{m}_X^N - m_X\|^2] = \frac{1}{N}\left(E[K(X,X)] - \|m_X\|^2\right) \to 0$$

as $N \to \infty$, proving convergence in quadratic mean.

A central limit theorem also exists. For arbitrary $f \in V$, $\langle f, \hat{m}_X^N \rangle = \frac{1}{N} \sum_{i=1}^N f(x_i)$ is the sample mean of real-valued independent random variables $f(x_i)$ and therefore $\sqrt{N} \langle f, \hat{m}_X^N - m_X \rangle$ converges to a Gaussian random variable having zero mean and the same variance as $f(X)$. This implies $\sqrt{N}(\hat{m}_X^N - m_X)$ converges weakly on $V$ to a Gaussian distribution [10, Theorem 108, full proof and statement].

Estimating the covariance of $\check{X}$ in $V$ can be reduced to estimating the mean of $(\check{X} - m_X) \otimes (\check{X} - m_X)$ in $V \otimes V$, as now explained.

Given $f, g \in V$, define $f \otimes g$ to be the bilinear continuous functional $V \times V \to \mathbb{R}$ sending $(u, v)$ to $\langle f, u \rangle \langle g, v \rangle$. That is,

$$(f \otimes g)(u, v) = \langle f, u \rangle \langle g, v \rangle. \tag{8.7}$$

The space generated by linear combinations of such functionals is the ordinary tensor product of $V$ and $V$. Completing this space results in the Hilbert space tensor product $V \otimes V$. The completion is with respect to the inner product defined by

$$\langle f_1 \otimes g_1, f_2 \otimes g_2 \rangle = \langle f_1, f_2 \rangle \langle g_1, g_2 \rangle. \tag{8.8}$$

If $V$ is a RKHS then $V \otimes V$ is also a RKHS.

**Theorem 8.2.** Let $V_1$ and $V_2$ be two RKHSs with kernels $K_1$ and $K_2$. Then the Hilbert space tensor product $V = V_1 \otimes V_2$ is a RKHS with kernel $K((x_1, x_2), (y_1, y_2)) = K_1(x_1, y_1) K_2(x_2, y_2)$.

*Proof.* This is Theorem 13 of [10]. $\qquad\square$

Define $\Sigma_{XX} = E[(\check{X} - m_X) \otimes (\check{X} - m_X)]$. It encodes the same information as $\Sigma_X$ because $\Sigma_{XX}(f, g) = E[\langle f, \check{X} - m_X \rangle \langle g, \check{X} - m_X \rangle]$. Being a mean, it can be estimated by

$$\hat{m}_{XX}^N = \frac{1}{N} \sum_{i=1}^N (K(\cdot, x_i) - \hat{m}_X^N) \otimes (K(\cdot, x_i) - \hat{m}_X^N). \tag{8.9}$$

## 8.6 Generalisations and further considerations

A covariance operator between two random variables $X$ and $Y$ defined on a common probability space can be defined. The variables are

embedded in RKHSs $V_x$ and $V_y$ using kernels $K_x$ and $K_y$, and the co-variance operator $\Sigma_{YX}$ is the linear operator from $V_x$ to $V_y$ defined as $\Sigma_{YX} f = E[\langle f, \check{X} - m_X \rangle (\check{Y} - m_Y)]$. It also has the tensor representation $\Sigma_{YX} = E[(\check{X} - m_X) \otimes (\check{Y} - m_Y)]$, is a mean in the tensor product space and can thus be estimated by

$$\hat{m}_{YX}^N = \frac{1}{N} \sum_{i=1}^{N} (K_x(\cdot, x_i) - \hat{m}_X^N) \otimes (K_y(\cdot, x_i) - \hat{m}_Y^N). \qquad (8.10)$$

The covariance operators are nuclear operators and Hilbert-Schmidt operators (see section 8.2.) This fact allows to construct measures of independence and of conditional independence in an elegant and efficient way, as developed in §9.3. As nuclear operators, they admit a nuclear decomposition as

$$\Sigma = \sum_i \lambda_i e_i \otimes f_i$$

where $\{e_i\}, \{f_j\}$ are orthonormal bases of the spaces of interest.

If the operator is the covariance operator of a random variable embedded in the space, it is a positive operator. The $\lambda_i$ are the eigenvalues and are positive or zero, and the $\{e_i\}$ are the eigenfunctions of the operator. In this case, the Hilbert-Schmidt theorem [14] states that any element of the space admits the decomposition

$$x = \sum_i \lambda_i \langle x, e_i \rangle e_i + x'$$

where $x' \in \mathcal{N}(\Sigma)$ is in the null space of the covariance operator. The range of the operator $\mathcal{R}(\Sigma)$ is spanned by the eigen vectors $\{e_i\}$. Thus the range of $\Sigma$ is the subset of $V$ of those functions $f \in V$ not in the null space of $\Sigma$ verifying $\sum_i \lambda_i^2 \langle f, e_i \rangle^2 < +\infty$. Restricting the domain of $\Sigma$ to the space of function of $V$ that can be written $\sum_i \langle f, e_i \rangle / \lambda_i$ for some $f$ in the range of $\Sigma$, we define a bijective restriction, and the inverse is unambiguously defined as $\Sigma^{-1} f = \sum_i \langle f, e_i \rangle / \lambda_i$ for $f \in \mathcal{R}(\Sigma)$. The operator is then said invertible on its range. Note that since $\Sigma$ is positive the $\lambda_i$ considered are strictly positive, $\lambda_i = 0$ characterizing members of the null space.

Hilbert-Schmidt operators are compact. If a compact operator is invertible, then necessarily its inverse is unbounded, otherwise $AA^{-1} = I$

would be compact. But $I$ is not compact in infinite dimension. A compact operator transforms a bounded set into a precompact set (every sequence contains a convergent subsequence). But the unit ball is not precompact in infinite dimension. One can construct a sequence in it from which no subsequence is convergent [14].

The problem of unboundedness of the inverse of Hilbert-Schmidt operators is at the root of ill-posed inverse problems, since unboundedness implies non continuity. Thus, two arbitrary close measurements in the output space may be created by two inputs separated by a large distance in the input space. This is at the heart of regularization theory. In §9, inverting covariance operators is needed in almost all developments in estimation or detection. In theoretical development, we will assume the covariance operators are invertible, at least on their range (see above). However though, the inverse may be unbounded, and regularizing the inversion is needed. Practically, Tikhonov approach will systematically be called for regularization [71]. Typically, the inverse of $\Sigma$ will be replaced by $(\Sigma + \lambda I)^{-1}$, where $I$ is the identity operator, or even by $(\Sigma + \lambda I)^{-2}\Sigma$. Parameter $\lambda$ is the regularization parameter, which can be chosen in some problems by either inforcing solutions (which depend on the inverse) to satisfy some constraints, or by studying convergence of solutions as the number of data grows to infinity.

# 9

## Applications of Embeddings

In this chapter, the use of embeddings in RKHS for signal processing applications is illustrated. Many applications exist in signal processing and machine learning. We made here arbitrary choices that cover topics in signal processing mainly. The aim is not only to cover some topics, but also to provide some practical developments which will allow the reader to implement some algorithms.

Since practical developments deal with finite amount of data, only finite dimensional subspaces of possibly infinite dimensional RKHS are used empirically. This is explained in a first paragraph and used to develop matrix representations for the mean elements and the operators presented in the previous chapter (such as covariance operators.) The following paragraphs illustrate applications of embeddings. We first discuss embeddings for application of statistical testing in signal processing: Comments on the use of the so-called deflection criterion for signal detection are made; The design of independence and conditional independence measures are then presented. Filtering is discussed next. A first approach elaborate on the embedding of Bayes rule into a RKHS, while a second approach directly deals with the embeddings of the realisations of random signals, and how they are used for optimal

filtering.

## 9.1 Matrix representations of mean elements and covariance operators

Practically, mean elements and covariance operators are used by applying them to functions in the RKHS. Furthermore, the operators estimated are usually of finite rank. In inference problems, and when dealing with a finite number of observations $\{x_i\}_{i=1,\ldots,N}$ the representer theorem [37, 60] states that the optimizers of some empirical risk function are to be searched for in the subspace $W_x$ of $V_x$ generated by $\{K_x(.,x_i)\}_{i=1,\ldots,N}$ (from now on, a kernel is indexed to stress different variables and different kernels and spaces jointly.) This is a direct consequence of the reproducing property. Any empirical risk to be minimized with respect to functions $f \in V_x$ is evaluated at the points $x_i$. Let $f(.) = f_{\|}(.) + f_{\perp}(.)$, where $f_{\|}(.) \in W_x$ and $f_{\perp}(.) \in W_x^{\perp}$. Then $f(x_i) = \langle f, K_x(.,x_i) \rangle = f_{\|}(x_i)$ and $f_{\perp}(x_i) = 0$. The useful subspace of $V_x$ is $W_X$, and practically, the empirical mean element or operators are applied to functions in the form

$$f(.) = \sum_{i=1}^{N} \alpha_i K_x(.,x_i) \tag{9.1}$$

Consider first the action of the mean element $m = \langle f, \widehat{m}_X^N \rangle$. The reproducing property leads to

$$
\begin{aligned}
m &= \langle \sum_{i=1}^{N} \alpha_i K_x(.,x_i), \frac{1}{N} \sum_{i=1}^{N} K_x(.,x_i) \rangle \\
&= \frac{1}{N} \sum_{i,j=1}^{N} \alpha_i K_x(x_i,x_j) \quad = \quad \mathbf{1}_N^{\top} \boldsymbol{K}_x \boldsymbol{\alpha} \tag{9.2}
\end{aligned}
$$

where the Gram matrix $\boldsymbol{K}_x$ has entries $\boldsymbol{K}_{x,ij} = K_x(x_i,x_j)$, and where we introduced the vectors $\boldsymbol{\alpha} = (\alpha_1,\ldots,\alpha_N)^{\top}$ and $\mathbf{1}_N = \mathbf{1}/N = (1/N,\ldots,1/N)^{\top}$. This is precisely the inner product of two elements in the finite dimensional RKHS with kernel $\boldsymbol{K}_x$ and internal representations $\mathbf{1}_N$ and $\boldsymbol{\alpha}$. Furthermore, this formula leads to $\langle K(.,x_i), \widehat{m}_X^N \rangle = \delta_i^{\top} \boldsymbol{K}_x \mathbf{1}_N$ where $\delta_i$ is a vector of zeros except a 1 at the $i$th position.

Hence, the mean can be calculated without embedding explicitly the data into the RKHS, but just by using the kernel evaluated at the data points.

To get the analytic form for the application of the covariance operator to functions of $V_X$, consider first the interpretation of the covariance as the mean element in the tensor product $V_x \otimes V_y$, and evaluate

$$
\begin{aligned}
\langle f(.), \widehat{m}_{XY}^N(., v) \rangle_{V_x} \\
= & \ \langle \sum_{i=1}^N \alpha_i K_x(., x_i), \frac{1}{N} \sum_{i=1}^N (K_x(., x_i) - \widehat{m}_X)(K_y(v, y_i) - \widehat{m}_Y(v)) \rangle_{V_x} \\
= & \ \frac{1}{N} \sum_{j=1}^N ((\boldsymbol{\alpha}^\top \boldsymbol{K}_x)_j - \boldsymbol{\alpha}^\top \boldsymbol{K}_x \boldsymbol{1}_N)(K_y(v, y_j) - \widehat{m}_Y(v)) \qquad (9.3)
\end{aligned}
$$

Applying his result to a function $g \in V_y$ allows to obtain

$$
\begin{aligned}
\langle g, \widehat{\Sigma}_{YX} f \rangle = & \ \mathrm{Cov}\,[g(Y), f(X)] \\
= & \ \frac{1}{N} \sum_{j=1}^N ((\boldsymbol{\alpha}^\top \boldsymbol{K}_x)_j - \boldsymbol{\alpha}^\top \boldsymbol{K}_x \boldsymbol{1}_N)\langle \sum_{i=1}^N \beta_i K_y(v, y_i), K_y(v, y_j) - \widehat{m}_Y(v) \rangle \\
= & \ \frac{1}{N} \boldsymbol{\alpha}^\top \boldsymbol{K}_x (\boldsymbol{I} - \frac{1}{N} \boldsymbol{1}\boldsymbol{1}^\top) \boldsymbol{K}_y \boldsymbol{\beta} \qquad (9.4)
\end{aligned}
$$

The matrix $\boldsymbol{C}_N = \boldsymbol{I} - \frac{1}{N} \boldsymbol{1}\boldsymbol{1}^\top$ is the so-called centering matrix. If the mean elements $\widehat{m}_X$ and $\widehat{m}_Y$ are known to be identically equal to zero, then the centering matrix does not appear in the calculations. Note that $\boldsymbol{C}_N$ is idempotent, $\boldsymbol{C}_N^2 = \boldsymbol{C}_N$.

To directly study how the empirical covariance operator acts, the alternative definition of the covariance operator as a linear operator is called for. Let $f(.) = \sum_{i=1}^N \alpha_i K_x(., x_i)$ in $V_x$ and $g(.) = \sum_{i=1}^N \beta_i K_y(., y_i)$ in $V_y$ such that $g = \widehat{\Sigma}_{YX} f$.

Assuming $\widehat{m}_X = 0$ and $\widehat{m}_Y = 0$ for the sake of simplicity, the empirical covariance operator can be written as

$$
\widehat{\Sigma}_{YX} f = \frac{1}{N} \sum_i \langle f(.), K_x(., x_i) \rangle K_y(v, y_i) \qquad (9.5)
$$

Then evaluating $g = \widehat{\Sigma}_{YX} f$ at the data points gives

$$
\begin{aligned}
g(y_k) &= \sum_j \beta_j K_y(y_k, y_j) &=& \quad (\boldsymbol{K}_y \boldsymbol{\beta})_k \\
&= \frac{1}{N} \sum_i \langle f(.), K_x(., x_i) \rangle K_y(y_k, y_i) \\
&= \frac{1}{N} \sum_{i,j} \alpha_j K_x(x_i, x_j) K_y(y_k, y_i) &=& \quad \frac{1}{N} (\boldsymbol{K}_y \boldsymbol{K}_x \boldsymbol{\alpha})_k (9.6)
\end{aligned}
$$

Thus $\boldsymbol{\beta} = N^{-1} \boldsymbol{K}_x \boldsymbol{\alpha}$ and $\langle g, \widehat{\Sigma}_{YX} f \rangle = N^{-1} \boldsymbol{\beta}^\top \boldsymbol{K}_y \boldsymbol{K}_x \boldsymbol{\alpha}$ is recovered.

Finally, the application of the inverse of a covariance operator is important to study. As discussed earlier, the effect of the regularised version of the operator $\widehat{\Sigma}_{r,XX} = \widehat{\Sigma}_{XX} + \lambda \boldsymbol{I}$ is studied, $\boldsymbol{I}$ being the identity (it also denotes the identity matrix in finite dimension.) Let $f = \widehat{\Sigma}_{r,XX}^{-1} g$, or $g = \widehat{\Sigma}_{r,XX} f$ where $f$ and $g$ are in the RKHS. Using the decomposition $f = f_\parallel + f_\perp$ recalled earlier, $f_\parallel(.) = \sum_i \alpha_i K_x(., x_i)$, $f(x_k) = f_\parallel(x_k)$ and $f_\perp(x_k) = 0$. Thus

$$
\begin{aligned}
g(.) &= \frac{1}{N} \sum_i f(x_i) K_x(., x_i) + \lambda f(.) \\
&= \frac{1}{N} \sum_{i,j} (\boldsymbol{K}\boldsymbol{\alpha})_i K_x(., x_i) + \lambda \sum_i \alpha_i K_x(., x_i) + \lambda f_\perp(.) \quad (9.7)
\end{aligned}
$$

Then, evaluating $g(.)$ at all $x_i$,

$$
\boldsymbol{K}_x \boldsymbol{\beta} = \frac{1}{N} (\boldsymbol{K}_x + N\lambda I) \boldsymbol{K}_x \boldsymbol{\alpha} \tag{9.8}
$$

or solving, the action of the regularized inverse is obtained as

$$
\boldsymbol{\alpha} = N(\boldsymbol{K}_x + N\lambda I)^{-1} \boldsymbol{\beta} \tag{9.9}
$$

## 9.2 Signal detection and the deflection criterion

Signal detection is usually modeled as a binary testing problem: Based on the observation of a signal, a detector has to decide which hypothesis among $H_0$ or $H_1$ is true. The detector is in general a functional of the observation, denoted here as a filter. In [53, 54], Picinbono&Duvaut developed the theory of Volterra filters for signal detection. Recall that

Volterra filters are polynomial filters. If $x(n)$ is an input signal, the output of a $M$th order Volterra filter reads

$$y(n) = h_0 + \sum_{i=1}^{M} \sum_{j_1,\ldots,j_i} h_i(j_1,\ldots,j_i) x(n-j_1) \times \ldots \times x(n-j_i) \quad (9.10)$$

where the functions $h_i$ satisfy some hypothesis ensuring the existence of the output $y(n)$. The first term $i=1$ is nothing but a linear filter; the term $i=2$ is called a quadratic filter, and so on. If the range of summation of the $j_i$'s is finite for all $i$, the filter is said to be of finite memory.

In signal detection theory, if the detection problem is set up as a binary hypothesis testing problem, different approaches exist to design an efficient detector. For example, in the Neyman-Pearson approach, the detector that maximizes the probability of detection subject to a maximal given probability of false alarm is sought for. Recall for instance that the optimal detector in this approach of a known signal in Gaussian noise is the matched filter, which is a linear filter, and thus a first order Volterra filter.

**Deflection.** A simpler approach relies on the so-called deflection criterion. This criterion does not require the full modeling of the probability laws under each hypothesis, as is the case for Neyman-Pearson approach. The deflection is a measure that quantifies a contrast (or distance) between the two hypotheses for a particular detector. The greater the deflection the easier the detection, because the greater the separation between the two hypotheses. Let an observation $x$ be either distributed according to $P_0$ under hypothesis $H_0$ or according to $P_1$ under hypothesis $H_1$. Let $y(x)$ be a test designed to decide whether $H_0$ or $H_1$ is true. The deflection is defined as

$$d(y) = \frac{(E_1[y] - E_0[y])^2}{\text{Var}_0[y]} \quad (9.11)$$

where the subscript 0 or 1 corresponds to the distribution under which averages are evaluated. As mentionned above, the deflection quantifies the ability of the test to separate the two hypotheses.

In general the structure of the detector is imposed, and the best constrained structure which maximises the deflection is sought for. If

a linear structure is chosen, such as $y(x) = h^\top x$ in finite dimension, the matched filter is recovered. Picinbono&Duvaut studied the optimal detector according to the deflection when it is constrained to be a Volterra filter of the observation. They particularly develop the geometry of the filter, recoursing to the Hilbert space underlying Volterra filters of finite orders. An interesting fact outlooked in [53, 54] is that this Hilbert space is the finite dimensional reproducing kernel Hilbert space generated by the kernel $K(x, y) = (1 + x^\top y)^M$. This is easily seen using a very simple example for $M = 2$. Consider the map defined by

$$
\begin{aligned}
\Phi \quad &: \quad \mathbb{R}^2 \longrightarrow V \\
&\phantom{:\quad} \boldsymbol{x} \longmapsto \left(1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, \sqrt{2}x_1 x_2, x_2^2\right)^\top
\end{aligned}
\tag{9.12}
$$

which embeds $\mathbb{R}^2$ into a subspace $V$ of $\mathbb{R}^6$. $V$ is a reproducing kernel Hilbert space whose kernel is $K(x, y) = (1 + \boldsymbol{x}^\top \boldsymbol{y})^2$. Indeed, a direct calculation shows that $\Phi(\boldsymbol{x})^\top \Phi(\boldsymbol{y}) = (1 + \boldsymbol{x}^\top \boldsymbol{y})^2$. Now, consider $y(n) = \boldsymbol{H}^\top \Phi\big((x(n), x(n-1))\big)$ where $\boldsymbol{H} \in \mathbb{R}^6$. $y(n)$ is then the output of a linear-quadratic Volterra filter, with $h_0 = H_1, h_1(1) = H_2/\sqrt{2}, h_1(2) = H_3/\sqrt{2}, h_2(1, 1) = H_4, h_2(1, 2) = H_5/\sqrt{2}, h_2(2, 2) = H_6$, and all other parameters set to zeros. This simple example generalizes to any $M$ and any finite memory.

The aim here is to analyse the deflection for filters living in arbitrary RKHS.

**Detecting in a RKHS.** Consider $V$ to be the RKHS associated with kernel $K$. We assume the kernel $K$ to be characteristic (*i.e.* the mapping sending a probability to the mean element in $V$ is injective.) Let $\mu_i$ be the embedding of $P_i$, or $\mu_i(.) = E_i[K(., x)]$. Let $\Sigma_0 : V \to V$ be the covariance operator of a random variable distributed under $P_0$. The detector is sought for as a function $f$ in the RKHS that maximizes the deflection. Thus $y = f(x) = \langle f(.), K(., x) \rangle$. Then $E_i[y] = E_i[f(x)] = \langle \mu_i, f \rangle$ (definition of the mean element.) Furthermore, the definition of the covariance operator gives $\text{Var}_0[y] = \langle f, \Sigma_0 f \rangle$. The deflection for $y$ thus reads

$$
d(y) = \frac{\langle \mu_1 - \mu_0, f \rangle^2}{\langle f, \Sigma_0 f \rangle}
\tag{9.13}
$$

For ease of discussion, the covariance operator is assumed invertible.

Since $\Sigma_0$ is positive definite, so is its inverse, and a new inner product in $V$ can be defined by $\langle f, g \rangle_0 = \langle f, \Sigma_0^{-1} g \rangle$. The deflection of $y$ then writes

$$d(y) = \frac{\langle \mu_1 - \mu_0, \Sigma_0 f \rangle_0^2}{\langle f, \Sigma_0 f \rangle} \tag{9.14}$$

and Schwartz inequality offers an easy means to maximize the deflection. Indeed, the following holds

$$
\begin{aligned}
d(y) &\leq \frac{\langle \mu_1 - \mu_0, \mu_1 - \mu_0 \rangle_0 \langle \Sigma_0 f, \Sigma_0 f \rangle_0}{\langle f, \Sigma_0 f \rangle} \\
&= \langle \mu_1 - \mu_0, \Sigma_0^{-1}(\mu_1 - \mu_0) \rangle
\end{aligned} \tag{9.15}
$$

and the maximum is attained when vectors $\mu_1 - \mu_0$ and $\Sigma_0 f$ are proportional. The constant of proportionality is not important, since this amounts to scale function $f$ and does not change the optimal deflection. The constant of proportionality is thus chosen to be 1. The optimal function hence satisfies $\mu_1 - \mu_0 = \Sigma_0 f$, or $f = \Sigma_0^{-1}(\mu_1 - \mu_0)$.

The assumption of the invertibility of the covariance operator is not fundamental in the derivation. If it is not invertible, the derivative of the deflection (in some functional sense such as the Gateaux derivative) is used to show that the maximum is obtained when $\mu_1 - \mu_0$ and $\Sigma_0 f$ are proportional.

If no structure is imposed to the filter, and if a Neyman-Pearson approach is taken to solve the detection problem, the best strategy is to compare the likelihood ratio to a threshold, chosen in order to satisfy a constraint on the error of the first kind. A link between the optimal detector in the deflection sense and the likelihood ratio can be established. Precisely let $r(x) = l(x) - 1 = p_1(x)/p_0(x) - 1$. In the following, $r(x)$ is assumed square integrable under hypothesis $H_0$.

First note the following. Any $f$ for which $f(x) \in L^2$ satisfies $E_0[f(x)r(x)] = E_1[f(x)] - E_0[f(x)]$. Thus any $f \in V$ such that $f(x) \in L^2$ satisfies $E_0[f(x)r(x)] = \langle f, \mu_1 - \mu_0 \rangle$. In particular, $E_0[K(u,x)r(x)] = \langle K(u,.), \mu_1 - \mu_0 \rangle = (\mu_1 - \mu_0)(u)$. If we seek for the best estimator of the likelihood ratio in $V$ minimizing the mean square error under $H_0$, we must center it and then estimate $r(x)$, since

$E_0[p_1(x)/p_0(x)] = 1$. Then the optimal estimator will lead to an error orthogonal (in the $L^2$ sense) to $V$, and thus the following equation has to be solved

$$E_0[g(x)(r(x) - \widehat{r}(x))] = 0, \quad \forall g \in V \qquad (9.16)$$

But using the reproducing property this last expectation can be written as

$$
\begin{aligned}
E_0[g(x)(r(x) - \widehat{r}(x))] &= E_0\Big[\langle g, K(.,x)\rangle (r(x) - \widehat{r}(x))\Big] \qquad (9.17)\\
&= \Big\langle g \Big| E_0\big[K(.,x)r(x)\big] - E_0\big[K(.,x)\langle \widehat{r}, K(.,x)\rangle\big]\Big\rangle \qquad (9.18)\\
&= \langle g, \mu_1 - \mu_0 - \Sigma_0\widehat{r}\rangle \qquad (9.19)
\end{aligned}
$$

Since the last result is equal to 0 for any $g \in V$, the solution necessarily satisfies the equation $\mu_1 - \mu_0 - \Sigma_0\widehat{r} = 0$, hence proving that $\widehat{r}$ is the detector that maximizes the deflection.

The optimal element in the RKHS which maximises the deflection is therefore also the closest in the mean square sense to the likelihood ratio.

## 9.3 Testing for independence

Recent works in machine learning and/or signal processing use RKHS mainly for the possibility offered to unfold complicated data in a larger space. In classification for example, data nonlinearly separable in the physical space may become linearly separable in a RKHS. Testing independence by embedding data into a RKHS relies in some way on the same idea.

The argument is to use covariance (in the RKHS) to assess independence. It is well known that no correlation between two variables does not imply independence between these variables. However, an intuitive idea is that no correlation between any nonlinear transformation of two variables may reveal independence. This simple idea was at the root of the celebrated Jutten-Herault algorithm, the first device to perform blind source separation. In fact it was studied as early as 1959 by Rényi [57]. He showed $X$ and $Y$, two variables defined on some

common probability space, are independent if and only if the so-called maximal correlation $\sup_{f,g} \text{Cov} \left[ f(X), g(Y) \right]$ is equal to zero, where $f, g$ are continuous bounded functions.

Recently, this result was revisited by Bach, Gretton and co-workers through the lense of RKHS [5, 27]. The maximal correlation as used by Rényi is too complicated to be practically evaluated, because the space over which the supremum has to be calculated is far too big. The idea is then to look for the supremum in a space in which, firstly the maximum correlation can be more easily calculated, secondly, in which Rényi's result remains valid.

Gretton showed in [27] that $X$ and $Y$ are independent if and only if $\sup_{f,g} \text{Cov} \left[ f(X), g(Y) \right]$ is equal to zero, where $f, g$ are living in a RKHS (precisely its unit ball) generated by a universal kernel. Recall that universality is understood here as the denseness of the RKHS into the space of bounded continuous functions. The link between Gretton's result and Rényi's is then intuitive, since under universality, any continuous bounded function may be approximated as closely as needed by a function in the RKHS. The second requirement above is thus satisfied. The first one is also verified, and this is the magical part of Gretton's approach. The maximal correlation evaluated in the unit ball of the universal RKHS

$$\sup_{f \in \mathcal{U}_x, g \in \mathcal{U}_y} \text{Cov} \left[ f(X), g(Y) \right] \tag{9.20}$$

where $\mathcal{U} = \{ f \in V \backslash \| f \| = 1 \}$, is nothing but the operator norm of the covariance operator $\Sigma_{XY}$. Indeed, by definition,

$$
\begin{aligned}
\| \Sigma_{XY} \| &= \sup_{g \in \mathcal{U}_y} \| \Sigma_{XY} g \|_{V_x} &= \sup_{f \in \mathcal{U}_x, g \in \mathcal{U}_y} \left| \langle f, \Sigma_{XY} g \rangle \right| \\
&= \sup_{f \in \mathcal{U}_x, g \in \mathcal{U}_y} \left| \text{Cov} \left[ f(X), g(Y) \right] \right| &\tag{9.21}
\end{aligned}
$$

As shown by Bach [5] and Gretton [27] this quantity can be efficiently evaluated from a finite amount of data. Let $(x_i, y_i)_{i=1,\dots,N}$ be $N$ independent and identically distributed copies of $X$ and $Y$, then $\langle f \,|\, \Sigma_{XY} g \rangle$ approximated by $\left\langle f \,|\, \widehat{\Sigma}_{XY}^N g \right\rangle$ is given by

$$\left\langle f \,|\, \widehat{\Sigma}_{XY}^N g \right\rangle = \frac{1}{N} \boldsymbol{\alpha}^\top \boldsymbol{K}_x \boldsymbol{C}_N \boldsymbol{K}_y \boldsymbol{\beta} \tag{9.22}$$

where the $\alpha$s (respectively $\beta$s) are the coefficients of the expansion of $f$ (respectively $g$) in the subspace of $V_x$ (respectively $V_y$) spanned by $K_x(.,x_i), i = 1, \ldots, N$ (respectively $K_y(.,y_i), i = 1, \ldots, N$). The norm of $f$ is given by $\boldsymbol{\alpha}^\top \boldsymbol{K}_x \boldsymbol{\alpha}$ and the norm of $g$ by $\boldsymbol{\beta}^\top \boldsymbol{K}_y \boldsymbol{\beta}$. Therefore the maximal correlation is approximated by

$$\sup_{\boldsymbol{\alpha}^\top \boldsymbol{K}_x \boldsymbol{\alpha}=1, \boldsymbol{\beta}^\top \boldsymbol{K}_x \boldsymbol{\beta}=1} \frac{1}{N} \boldsymbol{\alpha}^\top \boldsymbol{K}_x \boldsymbol{C}_N \boldsymbol{K}_y \boldsymbol{\beta} = \frac{1}{N} \|\boldsymbol{K}_x^{1/2} \boldsymbol{C}_N \boldsymbol{K}_y^{1/2}\|_2 \quad (9.23)$$

where $\|A\|_2 = \sqrt{\lambda_M(A^\top A)}$ is the usual spectral norm of matrix $A$, that is the square root of its maximal singular value.

The Gram matrices appearing in the estimation of the maximal correlation are of size $N \times N$. Therefore, the calculation of the maximal correlation can be very time consuming for large data sets. A simple idea allows to end up with a measure easier to evaluate. It relies on Hilbert-Schmidt norms. It is known that for any Hilbert-Schmidt operator, the operator norm is lower or equal than the Hilbert-Schmidt norm. Thus Gretton's result remains true if the operator norm is replaced with he Hilbert-Schmidt norm: $X$ and $Y$ are independent if and only if $\|\Sigma_{XY}\|_{HS} = 0$.

Furthermore, when dealing with the data set $(x_i, y_i), i = 1, \ldots, N$, the Hilbert-Schmidt norm of the covariance operator may be approximated by the the Hilbert-Schmidt norm of the empirical estimates of the operator. It is then not difficult to show that

$$\|\widehat{\Sigma}_{XY}^N\|_{HS}^2 = \frac{1}{N^2} \mathrm{Tr}(\boldsymbol{C}_N \boldsymbol{K}_x \boldsymbol{C}_N \boldsymbol{K}_y) \quad (9.24)$$

For independent and identically distributed data, this estimator satisfies a central limit theorem. It is asymptotically unbiased, and its variance can be explicitly written. An unbiased version also exists which in the spirit of $k$-statistics eliminates the $1/N$ bias [64]. Both versions satisfy some concentration inequalities and a central limit theorem (at the usual rate) [27, 64].

This measure has been called HSIC for Hilbert-Schmidt Independence Criterion. It can obviously be used for testing independence between two samples, but has also been used for feature selection [64]. Its evaluation requires $O(N^2)$ operations. This complexity can be lowered by using approximation to Gram matrices such as the incomplete

Cholesky factorization [20] or other types of approximation [60]. As an alternative, we have designed a recursive implementation of HSIC which is exact, of course still requiring $O(N^2)$ operations, but which only manipulates vectors. This algorithm can be used on very long data sets in reasonable time. However, for huge data sets it is still impractical. But its recursive structure allows an efficient approximation leading to the possibility of calculating HSIC on-line or whatever the size of the data set. We just give the forms of the second algorithm, the derivation being developed in [1, 2]. In this procedure, a dictionary is built recursively, the elements of which are used to evaluate HSIC on-line. New data are included in the dictionary if they are sufficiently incoherent with the members of the dictionary. This procedure was proposed in a filtering context in [58] as a simplification of the approximate linear dependence (ALD) criterion proposed in [19]. Coherence is measured in the tensor product $V_x \otimes V_y$. Thanks to the reproducing property, the coherence between two data $(x_\alpha, y_\alpha)$ and $(x_n, y_n)$ is evaluated as $|K_x(x_n, x_\alpha) K_y(y_n, y_\alpha)|$ (assuming the kernels are normalized). The dictionary $\mathcal{D}_n^\mu$ contains the index of the data retained up to time $n$, initializing it with $\mathcal{D}_1^\mu = \{1\}$. It is updated according to

$$\mathcal{D}_n^\mu = \begin{cases} \mathcal{D}_{n-1}^\mu \cup \{n\} & \text{if } \sup_{\alpha \in \mathcal{D}_{n-1}^\mu} |K_x(x_n, x_\alpha) K_y(y_n, y_\alpha)| \le \mu \\ \mathcal{D}_{n-1}^\mu & \text{otherwise} \end{cases} \tag{9.25}$$

Parameter $\mu$ is in $(0, 1]$. If $\mu = 1$ all the new data are aggregated to the dictionary and the following algorithm exactly delivers HSIC. If $\mu < 1$, the new data are added to the dictionary if it is sufficiently incoherent with all the members of the dictionary.

To describe the algorithm, some more notations are needed. Let $\kappa_x^n$ be the norm $K_x(x_n, x_n)$. Let $\boldsymbol{\pi}_n$ be a $|\mathcal{D}_n|$ dimensional vector whose entries $\pi_n(\alpha)$ count the number of times the element $\alpha$ of the dictionary has been chosen by the rule (9.25). $\boldsymbol{\pi}_n$ is initialised by $\pi_1(1) = 1$. Let $\boldsymbol{k}_x^n$ contain the $K_x(x_n, x_\alpha), \forall \alpha \in \mathcal{D}_{n-1}^\mu$. Vectors $\boldsymbol{v}^n$ appearing in the algorithm below are initialised as $\boldsymbol{v}^0 = 1$. Finally, $\circ$ denotes the Hadamard product, *i.e.* the entrywise product for vectors or matrices. Equipped with all this, the algorithm is the following:

Sparse HSIC :

$$\widehat{H}_n^\mu = \|M_{yx}^n\|^2 + \|m_x^n\|^2 \|m_y^n\|^2 - 2c_{yx}^n \tag{9.26}$$

$$\|M_{yx}^n\|^2 = \frac{(n-1)^2}{n^2}\|M_{yx}^{n-1}\|^2 + \frac{2}{n^2}\boldsymbol{\pi}_{n-1}^\top \boldsymbol{k}_x^n \circ \boldsymbol{k}_y^n + \frac{\kappa_x^n \kappa_y^n}{n^2} \tag{9.27}$$

$$\|m_x^n\|^2 = \frac{(n-1)^2}{n^2}\|m_x^{n-1}\|^2 + \frac{2}{n^2}\boldsymbol{\pi}_{n-1}^\top \boldsymbol{k}_x^n + \frac{\kappa_x^n}{n^2} \tag{9.28}$$

$$c_{yx}^n = \frac{1}{n^3}\boldsymbol{\pi}_n^\top \boldsymbol{v}_x^n \circ \boldsymbol{v}_y^n \text{ where:} \tag{9.29}$$

**1.** Si $\mathcal{D}_n = \mathcal{D}_{n-1} \cup \{n\}$

$$\boldsymbol{v}_x^n = \begin{pmatrix} \boldsymbol{v}_x^{n-1} + \boldsymbol{k}_x^n \\ \boldsymbol{\pi}_{n-1}^\top \boldsymbol{k}_x^n + \kappa_x^n \end{pmatrix} \quad \text{and} \quad \boldsymbol{\pi}_n = \begin{pmatrix} \boldsymbol{\pi}_{n-1} \\ 1 \end{pmatrix} \tag{9.30}$$

**2.** Si $\mathcal{D}_n = \mathcal{D}_{n-1} \quad : \quad a = \arg\max_{\alpha \in \mathcal{D}_{n-1}} |K_x(x_n, x_\alpha)K_y(y_n, y_\alpha)|,$

$$\boldsymbol{v}_x^n = \boldsymbol{v}_x^{n-1} + \boldsymbol{k}_x^n \quad \text{and} \quad \boldsymbol{\pi}_n = \boldsymbol{\pi}_{n-1} + \delta_{a\alpha} \tag{9.31}$$

Sparse HSIC is illustrated on the following example, taken from [24]. Consider the couple of independent variables $(X, Y)$ where $X$ is uniformly distributed on $[-a, a]$ and $Y$ is uniformly distibuted on $[-c, -b] \cup [b, c]$, $a, b, c$ being real positive constants. We choose $a, b, c$ to ensure that $X$ and $Y$ have the same variance. Let $Z_\theta$ the random vector obtained by rotating vector $(X, Y)$ by an angle $\theta$. The components of $Z_\theta$ are uncorrelated whatever $\theta$ but are independent if and only if $\theta = \pi/2 \times \mathbb{Z}$. $N$ i.i.d. samples of $Z_\theta$ are simulated. For $\theta$ varying in $[0, \pi/2]$ the final value of sparse HSIC $\widehat{H}_N^\mu$ is plotted in figure (9.1). The right plot displays $\widehat{H}_N^\mu$ for the five values $\mu = 0.8, 0.85, 0.9, 0.95, 1$. The kernel used is the Gaussian kernel $\exp(-\|\boldsymbol{x} - \boldsymbol{y}\|^2)$. As can be seen, the distorsion due to sparsity is very small. Interestingly, the size of the dictionary (left plot) is very small as soon as $\mu < 0.95$ in this example, thus implying a dramatic decrease in computation time and memory load. Moreover, the form of the algorithm can be simply turned into an adaptive estimator by changing decreasing step-sizes into constant step sizes. This allows to track changes into the dependence structure of two variables [1]. The size of the dictionary is probably linked to the speed of decrease of the spectrum of the Gram matrix. In general,
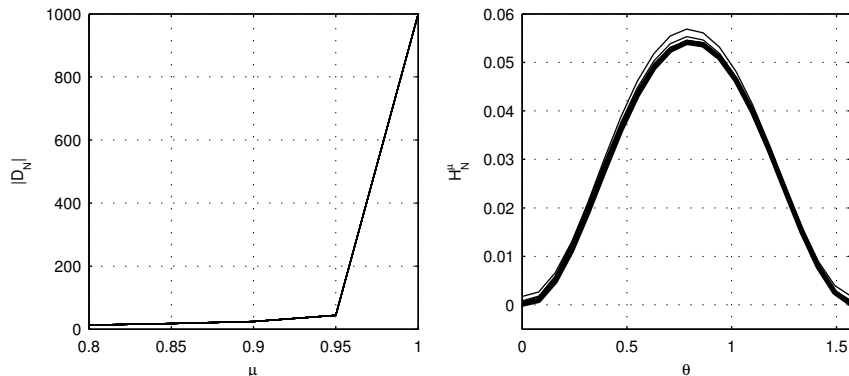
**Figure 9.1:** Right plot: HSIC as a function of $\theta$ in the example in the text. The thick line corresponds to $\mu = 1$ or equivalently to a non sparse evaluation of HSIC. The other lines correspond to $\mu = 0.8, 0.85, 0.9, 0.95$. Left plot: Size of the dictionary achieved after $N$ iterations as a function of $mu$.

when the Gram matrix can be nicely approximated by a low rank matrix, we observed that the size of the dictionary obtained is small. A general study of the approximation of the Gram matrix by the Gram matrix issued from the coherence dictionary remains to be done. Note however that some information are given for the related ALD criterion in [69] (see also [4] for related materials on low rank approximations in a regression context). We will meet again the coherence-based sparsification procedure in the application of RKHS to on-line nonlinear filtering.

**Maximum Mean Discrepancy (MMD).** Another approach to testing independence using embeddings into RKHS relies on the maximum mean discrepancy measures or MMD [26]. MMD measures a disparity between two probability measures $P$ and $Q$. Let $X$ and $Y$ be two random variables taking values in a space $\mathbb{X}$ and respectively distributed according to $P$ and $Q$. Let $\mathcal{F}$ be a function space on $\mathcal{X}$. MMD is defined as

$$MMD(P, Q; \mathcal{F}) = \sup_{f \in \mathcal{F}} \Big( E_P[f(X)] - E_Q[f(Y)] \Big) \qquad (9.32)$$

If the function space is rich enough, a null MMD implies equality between $P$ and $Q$. As previously, it can be shown that if $\mathcal{F}$ is restricted to

be the unit ball of a RKHS associated to a universal kernel, this result is true, $MMD(P, Q; \mathcal{F}) = 0 \Leftrightarrow P = Q$. This can be used to test for independence by measuring the maximum mean discrepancy between a joint probability and the product of its marginals. Furthermore, and like HSIC, there is a nice way of estimating MMD in a RKHS. In fact, thanks to Schwartz inequality, MMD can be expressed as

$$
\begin{aligned}
MMD(P, Q; V)^2 &= \sup_{f \in V, \|f\| \leq 1} \left( E_P[f(X)] - E_Q[f(Y)] \right)^2 \\
&= \sup_{f \in V, \|f\| \leq 1} \left( \langle \mu_P - \mu_Q, f \rangle \right)^2 \\
&= \|\mu_P - \mu_Q\|_V^2 \quad\quad\quad (9.33)
\end{aligned}
$$

where $\mu_P$, $\mu_Q$ are the mean elements of respectivelyy $P$ and $Q$ in $V$. When data $(x_i, y_i), i = 1, \ldots, N$ are to be tested for independence, empirical estimators of MMD are very easy to develop and implement. Furthermore, asymptotic results for their efficiency exist that allows a complete development of independence testing [26].

## 9.4   Conditional independence measures

Remarkably, the measure of independence presented in the preceding section has an extension which allows to quantify conditional independence. Conditional independence is a fundamental concept in different problems such as graphical modeling or dimension reduction. For example, graphical Markov properties of Gaussian random vectors are revealed by conditional independence relations between these components [39, 74]. Since measuring independence was found to be elegantly and efficiently done by embedding measures into RKHS, it was natural to work on the extension to conditional independence. It turns out that conditional independence can also be assessed in RKHS.

**Some recalls [39, 74].** Let $X, Y, Z$ be three real random vectors of arbitrary finite dimensions, and $\hat{X}(Z)$ and $\hat{Y}(Z)$ the best linear MMSE (minimum mean square error) estimates of $X$ and $Y$ based on $Z$. It is well-known that these are given by $\hat{X}(Z) = \Sigma_{XZ}\Sigma_{ZZ}^{-1}Z$ and $\hat{Y}(Z) =$

$\Sigma_{YZ}\Sigma_{ZZ}^{-1}Z$, where $\Sigma_{AB} := \text{Cov } [A, B]$ stands for the covariance matrix of vectors $A$ and $B$. The errors $X - \hat{X}(Z)$ and $Y - \hat{Y}(Z)$ are orthogonal to the linear subspace generated by $Z$, and this can be used to show the well-known relations

$$
\begin{aligned}
\Sigma_{XX|Z} &:= \text{Cov } \left[X - \hat{X}(Z), X - \hat{X}(Z)\right] \\
&= \Sigma_{XX} - \Sigma_{XZ}\Sigma_{ZZ}^{-1}\Sigma_{ZX} \quad (9.34) \\
\Sigma_{XY|Z} &:= \text{Cov } \left[X - \hat{X}(Z), Y - \hat{Y}(Z)\right] \\
&= \Sigma_{XY} - \Sigma_{XZ}\Sigma_{ZZ}^{-1}\Sigma_{ZY} \quad (9.35)
\end{aligned}
$$

$\Sigma_{XX|Z}$ is the covariance of the error in estimating $X$ linearly from $Z$. It is also called the partial covariance and it is equal to the conditional covariance in the Gaussian case. The second term measures the correlation remaining between $X$ and $Y$ once the effect of their possibly common observed cause $Z$ has been linearly removed from them. $\Sigma_{XY|Z}$ is called the partial cross-covariance matrix and is equal to the conditional cross-covariance in the Gaussian case (*i.e.* $X$, $Y$ and $Z$ are jointly Gaussian.)

Therefore, in the Gaussian case, conditional independence can be assessed using linear prediction and the partial cross-covariance matrix. This has led to extensive development in the field of graphical modeling.

**Using kernels.** The approach above can be extended to assess conditional independence for nonGaussian variables by using embeddings in RKHS. The extension relies on the notion of conditional cross-covariance operators, a natural extension of the covariance operators. Having in mind that cross-covariance operators suffices to assess independence (as cross-covariance does in the finite dimensional Gaussian case), the idea is consider

$$\Sigma_{XY|Z} := \Sigma_{XY} - \Sigma_{XZ}\Sigma_{ZZ}^{-1}\Sigma_{ZY} \quad (9.36)$$

as a potential candidate to assess conditional independence the operator. The first remark concerns the existence of this operator.

$\Sigma_{ZZ}$ is an operator from $V_z$ to $V_z$. Let $\mathcal{N}(\Sigma_{ZZ})$ and $\mathcal{R}(\Sigma_{ZZ})$ be respectively its null space and its range. The operator is supposed to

be invertible on its range and the inverse is abusively denoted as $\Sigma_{ZZ}^{-1}$. The inverse exits in full generality if and only if $\mathcal{N}(\Sigma_{ZZ}) = \{0\}$ and $\mathcal{R}(\Sigma_{ZZ}) = V_z$, corresponding to injectivity and surjectivity. Thus in the sequel, when dealing with ensemble operators, covariance operator will be supposed invertible. To avoid working with inverses, normalized covariance operators $V_{XY}$ should be considered. They are defined using $\Sigma_{XY} = \Sigma_{XX}^{1/2} V_{XY} \Sigma_{YY}^{1/2}$ [6]. The conditional covariance operator then reads

$$\Sigma_{XY|Z} := \Sigma_{XY} - \Sigma_{XX}^{1/2} V_{XZ} V_{ZY} \Sigma_{YY}^{1/2} \tag{9.37}$$

and the normalized version is given by

$$V_{XY|Z} := V_{XY} - V_{XZ} V_{ZY} \tag{9.38}$$

This last definition is the only theoretically well grounded, since the $V$ operators are shown to exist in [6] under some assumptions, but without relying on the invertibility of the covariance operators. However for practical purposes, we will use the other form, knowing that the existence of the inverse is subject to caution.

Several theorems show the meaning of the conditional covariance operators, and how we can assess conditional independence with them. They are all mainly due to Fukumizu, Bach and Jordan [22, 23]. The first result links conditional expectation to covariance and cross-covariance operators.

**Theorem 9.1.** For all $g \in V_y$,

$$\langle g, \Sigma_{YY|X} g \rangle = \inf_{f \in V_x} E\left[ \left( (g(Y) - E[g(Y)]) - (f(X) - E[f(X)]) \right)^2 \right] \tag{9.39}$$

If furthermore the direct sum $V_x + \mathbb{R}$ is dense in $L^2(P_X)$, then

$$\langle g, \Sigma_{YY|X} g \rangle = E_X \left[ \mathrm{Var}[g(Y)|X] \right] \tag{9.40}$$

The density assumption means than any random variable of $L^2(P_X)$ can be approximated as closely as desired by a function of $V_x$ plus a real. Adding the real is necessary since very often, constants do not belong to the RKHS under study (Remind that $L^2(P_X)$ is the space of square integrable functions with respect to $P_X$, or otherwise stated, the space of

functions of $X$ with finite expected squared norm $E[\|f(X)\|^2] < +\infty$.) The result of the theorem is an extension of what was recalled above, but stated in RKHS. The operator $\Sigma_{YY|X}$ measures the power of the error made in approximating a function of a random variable embedded in a RKHS by a function of another random variable embedded in its own RKHS. The second result generalizes the Gaussian case since under the assumption of density the operator evaluates a conditional variance. An informal proof is given, needing hypothesis not present in the statement. A full proof may be found in [23, Prop. 2 and 3].

*Proof.* Let $\mathcal{E}_g(f) = E\left[\left((g(Y) - E[g(Y)]) - (f(X) - E[f(X)])\right)^2\right]$. Then $f_0$ provides the infimum if $\mathcal{E}_g(f_0 + f) - \mathcal{E}_g(f_0) \geq 0$ for all $f \in V_x$. But we have

$$\mathcal{E}_g(f_0 + f) - \mathcal{E}_g(f_0) \quad = \quad \langle \Sigma_{XX} f, f \rangle + 2\langle \Sigma_{XX} f_0 - \Sigma_{XY} g, f \rangle \quad (9.41)$$

Obviously, $\Sigma_{XX} f_0 - \Sigma_{XY} g = 0$ satisfies the condition. It is also necessary. Indeed, suppose $\Sigma_{XX} f_0 - \Sigma_{XY} g \neq 0$. $\Sigma_{XX}$ is auto-ajoint and thus only has positive or null eigen values. Thus $\Sigma_{XX} f = -f$ has no solution and the null space of $\Sigma_{XX} + I$ is reduced to 0. Thus $\Sigma_{XX} + I$ is invertible. Therefore there is a non zero $f$ such that $\Sigma_{XX} f + f = -2(\Sigma_{XX} f_0 - \Sigma_{XY} g)$, and this $f$ satisfies $\mathcal{E}_g(f_0 + f) - \mathcal{E}_g(f_0) = -\langle f, f \rangle < 0$, giving a contradiction. Thus, this gives the result. Note we use the fact that $\Sigma_{XX}$ is invertible, at least on its range. The fact that $\langle g, \Sigma_{YY|X} g \rangle = E_X \left[ \text{Var}[g(Y)|X] \right]$ is shown hereafter as a particular case of conditional crosscovariance operator. $\square$

Since the conditional operator is linked to optimal estimation (in the mean square sense) of a function $g(Y)$ from a transformation of $X$, $\Sigma_{XX} E[g(.)|X] = \Sigma_{XY} g(.)$ should be a solution. However, this requires that the conditional expectation $E[g(.)|X]$ lies in $V_x$, a fact that is absolutely not guaranteed. If it is supposed so, the statement and results are more direct. In that case, for any $g \in V_y$,

$$\Sigma_{XX} E[g(.)|X] = \Sigma_{XY} g(.) \tag{9.42}$$

this provides a means of calculating the conditional mean in a RKHS if the covariance is invertible.

The following set of relations highlights the effect of the conditional covariance operators on function.

$$
\begin{aligned}
\langle f, \Sigma_{XY|Z} g \rangle &= \text{Cov}\,[f(X), g(Y)] - \langle f, \Sigma_{XZ} \Sigma_{ZZ}^{-1} \Sigma_{ZY} g \rangle \\
&= \text{Cov}\,[f(X), g(Y)] - \langle \Sigma_{ZX} f, \Sigma_{ZZ}^{-1} \Sigma_{ZY} g \rangle \\
&= \text{Cov}\,[f(X), g(Y)] - \langle \Sigma_{ZX} f, E[g(.)|Z] \rangle \\
&= \text{Cov}\,[f(X), g(Y)] - \langle \Sigma_{ZZ} E[f(.)|Z], E[g(.)|Z] \rangle \\
&= \text{Cov}\,[f(X), g(Y)] - \text{Cov}\,_Z [E[f(.)|Z], E[g(.)|Z]] \\
&= E[f(X)g(Y)] - E_Z [E[f(.)|Z] E[g(.)|Z]] \\
&= E_Z \Big[ \text{Cov}\,[f(X), g(Y)|Z] \Big]
\end{aligned}
\tag{9.43}
$$

The next question concerns whether conditional independence can be measured using the conditional covariance operator or not? The previous result and the first one in the following theorem show that a zero conditional covariance operator is not equivalent to conditional independence, but equivalent to a weaker form. The second result in the theorem below shows how to slightly modify the covariance operator to obtain the equivalence. This theorem is also from Fukumizu and his colleagues [22]. We suppose in the following that all the kernels used are characteristic, and that the conditional mean involved belongs to the proper RKHS.

**Theorem 9.2.** Let $X, Y, Z$ be three random vectors embedded in corresponding RKHS. Then we have

1. $\Sigma_{XY|Z} = 0 \iff P_{XY} = E_Z[P_{X|Z} \otimes P_{Y|Z}]$

2. $\Sigma_{(XZ)Y|Z} = 0 \iff X \perp Y | Z.$

*Proof.* **First assertion.** We have seen that

$$
\langle f, \Sigma_{XY|Z} g \rangle = E[f(X)g(Y)] - E_Z [E[f(.)|Z] E[g(.)|Z]] \tag{9.44}
$$

which can be written as

$$
\langle f, \Sigma_{XY|Z} g \rangle = \tag{9.45}
$$
$$
\int f(x)g(y) \Big( P_{XY}(dx, dy) - \int P_Z(dz) P_{X|Z}(dx, z) P_{Y|Z}(dy, z) \Big)
$$

Thus obviously, if for all $A$ and $B$ in the adequate sigma algebra

$$P_{XY}(A, B) = \int P_Z(dz) P_{X|Z}(A, z) P_{Y|Z}(B, z) \qquad (9.46)$$

we have $\langle f, \Sigma_{XY|Z} g \rangle = 0$ for all $f$ and $g$ leading necessarily to $\Sigma_{XY|Z} = 0$. Now if the covariance operator is zero then we have for all $f$ and $g$ $E_{P_{XY}}[f(X)g(Y)] = E_Q[f(X)g(Y)]$ where $Q = E_Z[P_{X|Z} \otimes P_{Y|Z}]$. Working in the tensorial product $V_x \otimes V_y$ where we have assumed $K_x K_y$ as a characteristic kernel allows to conclude that $Q = P_{XY}$.

**Second assertion.** Let $A, B, C$ be elements of the sigma algebra related to $X, Y$ and $Z$ respectively. Let $\mathbf{1}_A$ the characteristic function of set $A$. Then we have

$$
\begin{aligned}
&P_{XZY}(A, C, B) - E_Z[P_{XZ|Z}(A, C) P_{Y|Z}(B)] \\
&= E[\mathbf{1}_{A \times C}(X, Z) \mathbf{1}_B(Y)] - E_Z\big[E[\mathbf{1}_{A \times C}(X, Z)|Z] E[\mathbf{1}_B(Y)|Z]\big] \\
&= E_Z\Big[\mathbf{1}_C(Z)\big(E[\mathbf{1}_A(X)\mathbf{1}_B(Y)|Z] - E[\mathbf{1}_A(X)|Z] E[\mathbf{1}_B(Y)|Z]\big)\Big] \\
&= \int_C P_Z(dz)\big(P_{X,Y|Z}(A, B, z) - P_{X|Z}(A, z) P_{X|Z}(B, z)\big) \qquad (9.47)
\end{aligned}
$$

If $\Sigma_{(XZ)Y|Z} = 0$ then the first assertion implies $P_{XZY} = E_Z[P_{XZ|Z} \otimes P_{Y|Z}]$ and the previous integral is equal to zero for any $C$, which in turn implies that $P_{X,Y|Z}(A, B, z) - P_{X|Z}(A, z) P_{X|Z}(B, z)$ almost everywhere $(P_Z)$ for any $A, B$. But this is precisely the defintion of conditional independence. The converse is evident. $\qquad\square$

The conclusion of this theorem is that the variables $X$ and $Y$ has to be extended using *Z prior* conditioning. Assessing conditional independence relies on the conditional covariance operators (extended as above.) However, as done for independence testing, a measure is needed. The Hilbert-Schmidt norm $\|\Sigma_{(XZ)(YZ)|Z}\|^2$ is used for that purpose.

**Estimation.** The estimators of the conditional measures have representations in terms of Gram matrices. In the following, the indication of the RKHS in the inner product is suppressed for the sake of readability. For $N$ identically distributed observations $(x_i, y_i, z_i)$ the application of

the empirical covariance estimator to a function is given by

$$\widehat{\Sigma}_{XY} f = \frac{1}{N} \sum_j \tilde{K}_x(., x_j) \langle \tilde{K}_y(., y_j), f \rangle \tag{9.48}$$

where the tildas mean the kernel are centered. The effect of the regularized inverse of this operator on a function $\sum_i \beta_i K_x(., x_i)$ is to produce the function $\sum_i \alpha_i K_x(., x_i)$ with $\boldsymbol{\alpha} = N(\widetilde{\boldsymbol{K}}_x + N\lambda I)^{-1}\boldsymbol{\beta} := N\widetilde{\boldsymbol{K}}_{r,x}^{-1}\boldsymbol{\beta}$. The inner product $\langle f, \Sigma_{XY|Z} g \rangle$ for $f(.) = \sum_i \alpha_i \tilde{K}_x(., x_i)$ and $g(.) = \sum_i \beta_i \tilde{K}_y(., y_i)$ is evaluated for $N$ triple $x_i, y_i, z_i$ identically distributed. The result for the covariance operator is known from the first section

$$\langle f, \widehat{\Sigma}_{XY} g \rangle = \frac{1}{N} \boldsymbol{\beta}^\top \widetilde{\boldsymbol{K}}_y \widetilde{\boldsymbol{K}}_x \boldsymbol{\alpha} \tag{9.49}$$

Then the remaining term in $\widehat{\Sigma}_{XY|Z} = \widehat{\Sigma}_{XY} - \widehat{\Sigma}_{XZ} \widehat{\Sigma}_{ZZ}^{-1} \widehat{\Sigma}_{ZY}$ leads to

$$\begin{aligned}
\langle f, \widehat{\Sigma}_{XZ} \widehat{\Sigma}_{ZZ}^{-1} \widehat{\Sigma}_{ZY} g \rangle &= \sum_{i,j} \alpha_i \beta_j \langle \tilde{K}_x(., x_i), \widehat{\Sigma}_{XZ} \widehat{\Sigma}_{ZZ}^{-1} \widehat{\Sigma}_{ZY} \tilde{K}_y(., y_j) \rangle \\
&= \frac{1}{N^2} \sum_{i,j,k,l} \alpha_i \beta_j (\widetilde{\boldsymbol{K}}_y)_{kj} (\widetilde{\boldsymbol{K}}_x)_{li} \langle \tilde{K}_z(., z_l), \widehat{\Sigma}_{ZZ}^{-1} \tilde{K}_z(., z_k) \rangle \\
&= \frac{1}{N} \boldsymbol{\beta}^\top \widetilde{\boldsymbol{K}}_y \widetilde{\boldsymbol{K}}_{r,z}^{-1} \widetilde{\boldsymbol{K}}_z \widetilde{\boldsymbol{K}}_x \boldsymbol{\alpha} \tag{9.50}
\end{aligned}$$

Thus the final result is then

$$\langle f, \widehat{\Sigma}_{XY|Z} g \rangle = \frac{1}{N} \boldsymbol{\beta}^\top \left( \widetilde{\boldsymbol{K}}_y \widetilde{\boldsymbol{K}}_x - \widetilde{\boldsymbol{K}}_y \widetilde{\boldsymbol{K}}_{r,z}^{-1} \widetilde{\boldsymbol{K}}_z \widetilde{\boldsymbol{K}}_x \right) \boldsymbol{\alpha} \tag{9.51}$$

**Hilbert-Schmidt norms.** Practically, a measure using the cross-covariance is preferable. Like for independence testing, a nice measure is provided by measuring the norm of the operator. For the Hilbert-Schmidt norm, we have

$$\left\| \widehat{\Sigma}_{XY|Z} \right\|_{HS}^2 = \sum_i \langle \widehat{\Sigma}_{XY|Z} \varphi_i, \widehat{\Sigma}_{XY|Z} \varphi_i \rangle \tag{9.52}$$

$$= \left\| \widehat{\Sigma}_{XY} \right\|_{HS}^2 + \left\| \widehat{\Sigma}_{XZ} \widehat{\Sigma}_{ZZ}^{-1} \widehat{\Sigma}_{ZY} \right\|_{HS}^2 - 2 \sum_i \langle \widehat{\Sigma}_{XY} \varphi_i, \widehat{\Sigma}_{XZ} \widehat{\Sigma}_{ZZ}^{-1} \widehat{\Sigma}_{ZY} \varphi_i \rangle$$

where we recall that $\{\varphi_i\}_{i\in\mathbb{N}}$ is an orthonormal basis of $V_y$. The double

product term is denoted as $P$. It reads

$$
\begin{aligned}
P \; &:= \; \sum_i \langle \widehat{\Sigma}_{XY}\varphi_i, \widehat{\Sigma_{XZ}\Sigma_{ZZ}^{-1}\Sigma_{ZY}}\varphi_i \rangle \\
&= \; \frac{1}{N}\sum_{i,k} \langle \tilde{K}_y(.,y_k), \varphi_i \rangle \langle \tilde{K}_x(.,x_k), \widehat{\Sigma_{XZ}\Sigma_{ZZ}^{-1}\Sigma_{ZY}}\varphi_i \rangle \\
&= \; \frac{1}{N^2}\sum_{i,k,l} \langle \tilde{K}_y(.,y_k), \varphi_i \rangle \langle \tilde{K}_y(.,y_l), \varphi_i \rangle \langle \tilde{K}_x(.,x_k), \widehat{\Sigma_{XZ}\Sigma_{ZZ}^{-1}}\tilde{z}_y(.,z_l) \rangle \\
&= \; \frac{1}{N^2}\sum_{k,l,m,n} (\widetilde{\boldsymbol{K}}_y)_{kl} (\widetilde{\boldsymbol{K}}_{r,z}^{-1})_{lm} (\widetilde{\boldsymbol{K}}_x)_{kn} (\widetilde{\boldsymbol{K}}_z)_{mn} \\
&= \; \frac{1}{N^2}\mathrm{Tr}\left( \widetilde{\boldsymbol{K}}_y \widetilde{\boldsymbol{K}}_{r,z}^{-1} \widetilde{\boldsymbol{K}}_z \widetilde{\boldsymbol{K}}_x \right) \quad\quad (9.53)
\end{aligned}
$$

Carrying the same calculation for the last term allows to obtain

$$
\begin{aligned}
\left\| \widehat{\Sigma}_{XY|Z} \right\|_{HS}^2 \; &= \; \frac{1}{N^2}\mathrm{Tr}\Big( \widetilde{\boldsymbol{K}}_x \widetilde{\boldsymbol{K}}_y - 2\widetilde{\boldsymbol{K}}_y \widetilde{\boldsymbol{K}}_{r,z}^{-1} \widetilde{\boldsymbol{K}}_z \widetilde{\boldsymbol{K}}_x \quad (9.54) \\
&\quad + \; \widetilde{\boldsymbol{K}}_y \widetilde{\boldsymbol{K}}_z \widetilde{\boldsymbol{K}}_{r,z}^{-1} \widetilde{\boldsymbol{K}}_x \widetilde{\boldsymbol{K}}_{r,z}^{-1} \widetilde{\boldsymbol{K}}_z \Big) \quad\quad (9.55)
\end{aligned}
$$

If the normalized version $V_{XY|Z} = V_{XY} - V_{XZ}V_{ZY}$ is used, the estimator of the Hilbert-Schmidt norm $\|V_{XY|Z}\|_{HS}^2$ evaluated using the empirical estimate $\widehat{V}_{XY|Z} = \widehat{\Sigma}_{r,XX}^{-1/2}\widehat{\Sigma}_{XY|Z}\widehat{\Sigma}_{r,YY}^{-1/2}$ is given by

$$
\left\| \widehat{V}_{XY|Z} \right\|_{HS}^2 \; = \; \mathrm{Tr}\Big( \boldsymbol{N}_x\boldsymbol{N}_y - 2\boldsymbol{N}_y\boldsymbol{N}_z\boldsymbol{N}_x + \boldsymbol{N}_y\boldsymbol{N}_z\boldsymbol{N}_x\,\boldsymbol{N}_z \Big) (9.56)
$$

where $\boldsymbol{N}_u$ is the normalized centered Gram matrix for variable $u = x, y$ or $z$, and reads $\boldsymbol{N}_u = \widetilde{\boldsymbol{K}}_u \widetilde{\boldsymbol{K}}_{r,u}^{-1}$. The proof follows the same line as before. This estimator has been shown to converge to $\left\| \Sigma_{XY|Z} \right\|_{HS}^2$ in probability in [24]. To obtain the result, the regularization parameter $\lambda$ must of course depend on $N$ and goes to zero at an appropriate rate. Refer to [22, 23, 24] for results concerning the consistency of all the estimates seen so far.

Following theorem 9.2, the use of these measures to assess conditional independence is not enough. The extension of the random variables including $Z$ must be considered. The measure to be tested is thus $\left\| \widehat{V}_{(XZ)Y|Z} \right\|_{HS}^2$.

**A simple illustration.** The simplest toy example to illustrate conditional HSIC is to test that three random variables $X, Y, Z$ constitute a Markov chain, *e.g.* $X$ and $Z$ are independent conditionally to $Y$. Consider the simple generative model

$$\begin{aligned} X &= U_1 \\ Y &= a(X^2 - 1) + U_2 \\ Z &= Y + U_3 \end{aligned}$$

where $U_1, U_2, U_3$ are three independent zero mean, unit variance Gaussian random variables, and $a \in \mathbb{R}$ is a coupling parameter. $N = 512$ independent and identically distributed samples $(x_i, y_i, z_i)$ were generated and used to estimate $\left\| \widehat{V}_{(XY)Z|Y} \right\|_{HS}^2$ using the equations above.

Practically with finite length data, the distribution of the measure is not known under both hypothesis $H_0$ : Markov and $H_1$: non Markov. To simulate the null hypothesis $X - Y - Z$ is a Markov chain, equivalently, $X, Z$ are independent conditionally to $Y$ or $\left\| \widehat{V}_{(XY)Z|Y} \right\|_{HS}^2 = 0$ , we use random permutations of the realizations of one variable, say $X$. This is usually done for independence test, since randomly permuting preserve empirical distributions. However here, some care must be taken because the distributions that must be preserved under permutations are the conditional distributions. Thus, to create the permuted data, the range of the conditioning variable $Y$ is partitioned into $L$ domains $Y_1, \ldots, Y_L$ such that each domain contains the same number of observations. The permuted observations $\tilde{x}_i$ are obtained *per* domain, *i.e.* , $\forall l = 1, \ldots, L$, $\tilde{x}_j = x_{\sigma(j)}$ for those $j$s such that $y_j \in Y_l$, where $\sigma(.)$ is a random permutation of these $j$s. For this toy problem, 100 permuted data sets were created. This allows to find the threshold corresponding to a 5% false alarm probability (level of the test).

Figure 9.2 displays the result of the simulation for this simple example. We test the three possibilities of having the Markov property among $X, Y, Z$. The coupling parameter $a$ is varied from 0 to 1. The plot displays the conditional HSIC measure, as well as the threshold that ensures at most 5% of false alarms (evaluated as mentioned above with $L = 8$ domains.) In the left plot, the squared norm of $\widehat{V}_{(XZ)Y|Z}$ is plotted to test the Markov property for $X - Z - Y$. In the right
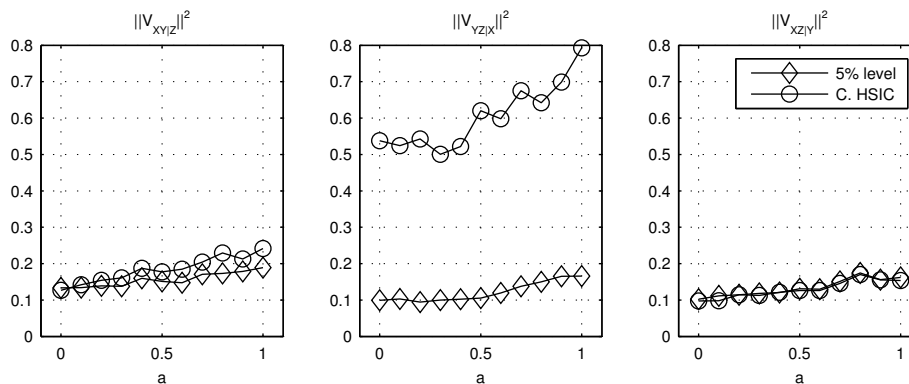
**Figure 9.2:** Hilbert-Schmidt norms of conditional covariance operator to assess conditional independence. Diamond represent the threshold insuring a 5% level test. Left plot: Testing for the chain $X \to Z \to Y$. Middle plot: Testing for the chain $Y \to X \to Z$. Right plot: Testing for the chain $X \to Y \to Z$. For the two leftmost plots, the measures is above the threshold meaning that the alternative $H_1$ is chosen by the test, except for $a = 0$ in the case $X \to Z \to Y$. These results are perfectly consistent with the simple model used. In the right plot, the measure is always below the diamond curves, meaning that $H_1$ is rejected, and that $X - Y-$ is Markov, in agreement with the model.

plot, the squared norm of $\widehat{V}_{(XY)Z|Y}$ is displayed to test if $X - Y - Z$ is a Markov chain. Finally, in the middle plot, the squared norm of $\widehat{V}_{(YX)Z|X}$ is plotted to test the Markov property for $Y - X - Z$. In the left plot for $a = 0$, $X$ and $Y$ are independent and $X - Z - Y$ is a particular Markov chain. However, as soon as $a > 0$, $X - Z - Y$ is not a Markov chain, and this is correctly inferred by the measure for $a$ as small as 0.1. The fact that $Y - X - Z$ is not a Markov chain is clearly assessed also, as illustrated in the middle plot. Finally, the Markov case $X - Y - Z$ is also correctly assessed in the right plot, since the squared norm of $\widehat{V}_{(XY)Z|Y}$ is always below the threshold insuring a 5% level test.

## 9.5   Kernel Bayesian filtering

The idea developed mainly by Song, Fukumizu and Gretton is to transfer into reproducing kernel Hilbert spaces the manipulation of proba-

bility measures used in inference, namely the sum rule, the chain rule and their application in Bayesian inference [25, 63]. Since conditional distribution may be embedded into RKHS, it seems natural to seek for the generalization of Bayesian inference in RKHS.

The starting point is theorem 9.1 which states that the conditional covariance operator $\Sigma_{YY|X} : V_y \to V_y$ is deeply linked to optimal estimation of random variables of $V_y$ from transformation of random variables in $V_x$. Precisely, $\Sigma_{XX}E[g(Y)|X = .] = \Sigma_{XY}g$ provided that the function $E[g(Y)|X = .]$ belongs to $V_x$, a fact that is not guaranteed for any kernel, but a fact assumed to be true in the sequel (see [22] for a sufficient condition.) Recall that the conditional covariance operator is defined as

$$\Sigma_{YY|X} = \Sigma_{YY} - \Sigma_{YX}\Sigma_{XX}^{-1}\Sigma_{XY} \tag{9.57}$$

a formula which assumes the invertibility of the covariance operator in $V_x$, and which exactly match the conditional covariance formula for Gaussian finite dimensional vectors.

The expression of the conditional mean allows to obtain an explicit form for the conditional kernel mean, that is, for the embedding of the conditional distribution. Let $\mu_X = E[K_x(., X)]$ and $\mu_Y = E[K_y(., Y)]$ be the embeddings of $X$ and $Y$ respectively in $V_x$ and $V_y$. The two embeddings are linked by

$$\mu_Y = \Sigma_{YX}\Sigma_{XX}^{-1}\mu_X \tag{9.58}$$

To show this relation, the conditional expectation of $g(Y)$ given $X$ which satisfies $E[g(Y)|X = .] = (\Sigma_{XX}^{-1}\Sigma_{XY}g)(.)$ is used in the following set of equations, valid for all $g \in V_y$,

$$
\begin{aligned}
\langle \Sigma_{YX}\Sigma_{XX}^{-1}\mu_X, g \rangle_{V_y} &= \langle \mu_X, \Sigma_{XX}^{-1}\Sigma_{XY}g \rangle_{V_x} \\
&= \langle \mu_X, E[g(Y)|X] \rangle_{V_x} \\
&= E_X[E[g(Y)|X]] \\
&= E[g(Y)] \\
&= \langle \mu_Y, g \rangle_{V_y} \tag{9.59}
\end{aligned}
$$

In particular, since $\mu_Y = E_Y[K_y(., Y)] = E_X E_{Y|X}[K_y(., Y)] =$

$E_X[\mu_{Y|X}]$, setting $P(dX) = \delta_x(dX)$ in (9.58) allows to obtain the conditional kernel mean as

$$\mu_{Y|x}(.) = E_{Y|X=x}[K_y(.,Y)] \quad = \quad \Sigma_{YX}\Sigma_{XX}^{-1}K_x(.,x)$$
$$:= \quad \Sigma_{Y|X}K_x(.,x) \qquad (9.60)$$

Note that the embedding $\mu_{Y|x}(.)$ is a function which belongs to $V_y$.

In these formula, the couple $X, Y$ is distributed according to the joint probability $P(X, Y)$. The covariance operator $\Sigma_{YX}$ is defined as the expected value of the tensor product $K_y(.,Y) \otimes K_x(.,X)$ over the joint probability. Recall that it has three interpretations. The first considers the covariance as a tensor, it is to say a bilinear functional over the product $V_y \times V_x$ (precisely their duals). The second well-know fact is based on the very definition of the tensor product $(K_y(.,Y) \otimes K_x(.,X))(g,f) = \langle g(.) \,|\, K_y(.,Y) \rangle \langle f(.) \,|\, K_x(.,X) \rangle$, which allows to write $\Sigma_{YX}f = E[\langle f(.) \,|\, K_x(.,X) \rangle K_y(.,Y)]$ and to consider $\Sigma_{YX}$ as a linear operator from $V_x$ to $V_y$. The third interpretation considers $\Sigma_{YX}$ as the embedding of the joint probability into the tensor product space $V_y \otimes V_x$. Since under this interpretation the covariance $\Sigma_{XX}$ can be seen as the embedding of $P(X)$ into the tensor product $V_x \otimes V_x$, this point of view allows to consider $\Sigma_{YX}\Sigma_{XX}^{-1}$ as a representation of

$$P(Y|X) = \frac{P(X,Y)}{P(X)} = \frac{P(X|Y)\int P(dX,Y)}{\int P(X,dY)} \qquad (9.61)$$

Writing the conditional kernel embedding as $\mu_{Y|x}(.) = \Sigma_{YX}\Sigma_{XX}^{-1}K_x(.,x)$ is at the root of the embedding of Bayesian inference into RKHS. It can be seen as the embedding of Bayes law when the likelihood is $P(X|Y)$ and the *prior* probability is $P(Y) = \int P(dX,Y)$. However, in Bayesian inference, if the likelihood is generally given, the *prior* is not given and in general not equal to the marginal of a given joint probability distribution.

Thus for Bayesian inference, the previous formula for the conditional embedding can be used, but for

$$P(Y|X) = \frac{Q(X,Y)}{\int Q(X,dY)} = \frac{P(X|Y)\pi(Y)}{\int Q(X,dY)} \qquad (9.62)$$

The joint probability to be considered is no longer $P(X,Y) = P(X|Y)P(Y)$ but instead $Q(X,Y) = P(X|Y)\pi(Y)$. The embedding of the *a posteriori* probability is then given by

$$\mu_{Y|x}^{\pi} = E^{\pi}[K_y(.,Y)|x] = \Sigma_{YX}^{\pi}\Sigma_{XX}^{\pi-1}K_x(.,x) \qquad (9.63)$$

where the superscript $\pi$ reminds that the *a priori* probability is $\pi$ instead of $\int P(dX,Y)$.

It is possible to relate $\Sigma_{YX}^{\pi}$ (or its adjoint $\Sigma_{XY}^{\pi}$) and $\Sigma_{XX}^{\pi-1}$ to embeddings evaluated on the joint probability $P$, using the following set of equations

$$
\begin{aligned}
\Sigma_{XY}^{\pi} &= E_Q[K_x(.,X) \otimes K_y(.,Y)] \\
&= E_{\pi}\left[E[K_x(.,X) \otimes K_y(.,Y)|Y]\right] \\
&= E_{\pi}\left[E[K_x(.,X)|Y] \otimes K_y(.,Y)\right] \\
&= E_{\pi}\left[\mu_{X|Y} \otimes K_y(.,Y)\right] \\
&= E_{\pi}\left[\Sigma_{X|Y}K_y(.,Y) \otimes K_y(.,Y)\right] \\
&= \Sigma_{X|Y}\Sigma_{YY}^{\pi} \qquad (9.64)
\end{aligned}
$$

The second line in this equation can also be interpreted as the average of the embedding of $P(Y,X|Y)$: this interpretation offers an alternative expression as

$$
\begin{aligned}
\Sigma_{XY}^{\pi} &= E_{\pi}\left[E[K_x(.,X) \otimes K_y(.,Y)|Y]\right] \\
&= E_{\pi}\left[\mu_{XY|Y}\right] \\
&= E_{\pi}\left[\Sigma_{XY|Y}K_y(.,Y)\right] \\
&= \Sigma_{XY|Y}\mu_Y^{\pi} \qquad (9.65)
\end{aligned}
$$

Likewise, the covariance operator reads

$$\Sigma_{XX}^{\pi} = \Sigma_{XX|Y}\mu_Y^{\pi} \qquad (9.66)$$

**Interpretations.** The operators $\Sigma_{XY}^{\pi}$ and $\Sigma_{XX}^{\pi}$ have simple interpretations when considered as embeddings. $\Sigma_{XX}^{\pi}$ corresponds to the embedding of the law $Q(X) = \int P(X|Y)\pi(dY)$ into the tensorial product $V_x \otimes V_x$. $\Sigma_{XY}^{\pi}$ is the embedding into $V_x \otimes V_y$ of

$Q(X, Y) = P(X|Y)\pi(Y)$. Thus $\Sigma^\pi_{XX}$ can be seen as the embedding of the sum rule, and is thus called **kernel sum rule**, whereas $\Sigma^\pi_{XY}$ is the embedding of the chain rule, and is thus called **kernel chain rule**. Obviously, Bayesian manipulation are a succession of applications of these rules.

To sum up, the embedding of the *a posteriori* probability reads

$$\mu_{Y|x} \quad = \quad \Sigma^\pi_{YX}\Sigma^{\pi-1}_{XX}K_x(.,x) \quad = \quad (\Sigma^\pi_{XY})^\top\Sigma^{\pi-1}_{XX}K_x(.,x) \quad (9.67)$$

where

$$\textbf{Chain rule} \quad : \quad \begin{cases} \Sigma^\pi_{XY} & = \; \Sigma_{X|Y}\Sigma^\pi_{YY} \quad = \quad \Sigma_{XY}\Sigma^{-1}_{YY}\Sigma^\pi_{YY} \\ \quad \text{or} & = \; \Sigma_{XY|Y}\mu^\pi_Y \quad = \quad \Sigma_{(XY)Y}\Sigma^{-1}_{YY}\mu^\pi_Y \end{cases} \quad (9.68)$$

$$\textbf{Sum rule} \quad : \; \Sigma^\pi_{XX} = \Sigma_{XX|Y}\mu^\pi_Y \quad = \quad \Sigma_{(XX)Y}\Sigma^{-1}_{YY}\mu^\pi_Y \quad (9.69)$$

**Estimators.** Estimators for

$$\mu_{Y|x} \quad = \quad \Sigma^\pi_{YX}\Sigma^{\pi-1}_{XX}K_x(.,x) \quad (9.70)$$

$$\Sigma^\pi_{XY} \quad = \quad \Sigma_{(XY)Y}\Sigma^{-1}_{YY}\mu^\pi_Y \quad (9.71)$$

$$\Sigma^\pi_{XX} \quad = \quad \Sigma_{(XX)Y}\Sigma^{-1}_{YY}\mu^\pi_Y \quad (9.72)$$

are obtained using empirical estimators of the different covariance operators. The last two operators are seen as linear operator from $V_y$ into respectively $V_x \otimes V_y$ and $V_x \otimes V_x$. Let us find an estimator for $\Sigma^\pi_{XY}$. The other will be obtained immediately by replacing $\Sigma_{(XY)Y}$ with $\Sigma_{(XX)Y}$.

The estimators are based on the observation of $N$ i.i.d. samples $(x_i, y_i)$ of the couple $(X, Y)$. We denote by $\boldsymbol{K}_x$ and $\boldsymbol{K}_y$ the Gram matrices evaluated on this sample. Furthermore, since information about the *prior* is needed, a number $N_\pi$ of i.i.d. samples $Y^\pi_i$ from the *prior* $\pi$ are assumed to be observed. This seems a strange assumption, but in many situations, these samples are at hand. For example, in recursive nonlinear filtering, the *posterior* probability at a particular time step serves as *prior* probability for the following time step. This will be detailed in the next section. The estimator for the function $\mu^\pi_Y$ is written as

$$\mu^\pi_Y(.) = \sum_{i=1}^{N_\pi} \gamma_i K_y(., Y^\pi_i) \quad (9.73)$$

Let $\boldsymbol{\mu}_Y^\pi$ the vector containing the $\mu_Y^\pi(y_k)$. Using the results of §9.1, we know that $\Sigma_{YY}^{-1}\mu_Y^\pi(.) = \sum_i \beta_i K_y(., y_i)$ where

$$\boldsymbol{\mu}_Y^\pi = \frac{1}{N}(\boldsymbol{K}_y + N\lambda I)\boldsymbol{K}_y\boldsymbol{\beta} \tag{9.74}$$

Applying $\Sigma_{(XY)Y}$ to $\Sigma_{YY}^{-1}\mu_Y^\pi$ leads to

$$\begin{aligned}
\Sigma_{XY}^\pi &= \frac{1}{N}\sum_{ij} \beta_i \boldsymbol{K}_{y,ij} K_x(., x_j) \otimes K_y(., y_j) \\
&= \sum_j \mu_j K_x(., x_j) \otimes K_y(., y_j) \text{ where} \tag{9.75} \\
\boldsymbol{\mu} &= (\boldsymbol{K}_y + \lambda I)^{-1}\boldsymbol{\mu}_Y^\pi \tag{9.76}
\end{aligned}$$

Likewise

$$\begin{aligned}
\Sigma_{XX}^\pi &= \sum_j \mu_j K_x(., x_j) \otimes K_y(., x_j) \text{ where} \tag{9.77} \\
\boldsymbol{\mu} &= (\boldsymbol{K}_y + \lambda I)^{-1}\boldsymbol{\mu}_Y^\pi \tag{9.78}
\end{aligned}$$

To get an estimate for $\mu_{Y|x}$ note that $\Sigma_{YX}^\pi = \Sigma_{XY}^{\pi,\top} = \sum_j \mu_j K(., y_j) \otimes K(., x_j)$. Since $\Sigma_{XX}^\pi$ is not insured to be positive definite, the regularization $(\Sigma^2 + \varepsilon I)^{-1}\Sigma$ of the inverse is used. Doing as above, searching for $\mu_{Y|x}(.) = \sum_j \zeta_j(x)K_y(., y_j)$, the vector

$$\begin{aligned}
\boldsymbol{\zeta}(x) &= \Lambda\left((\boldsymbol{K}_x\Lambda)^2 + \varepsilon I\right)^{-1} \boldsymbol{K}_x \Lambda \boldsymbol{k}_X(x) \\
&= \Lambda\boldsymbol{K}_x\left((\boldsymbol{K}_x\Lambda)^2 + \varepsilon I\right)^{-1} \Lambda \boldsymbol{k}_X(x) \tag{9.79}
\end{aligned}$$

is finally obtained, where $\boldsymbol{k}_X(x) = (K_x(x_1, x), \ldots, K_x(x_N, x))^\top$ and $\Lambda = \text{Diag}(\boldsymbol{\mu})$ is a diagonal matrix, the diagonal elements of which are the entries of $\boldsymbol{\mu}$.

Note that the matrix representation presented here has been shown to converge to the true embedding when the number of data goes to infinity, and when the regularization parameters goes to zero at correct speed (see [25].)

Some application of kernel Bayes rule were presented in [25, 63], among which a kernel belief propagation for inference on graphs, Bayes inference problems with unknown likelihood (an alternative solution

to Approximate Bayesian Calculation), and to filtering. In the sequel, kernel Bayes filtering is developed and applied to the prediction of time series.

**Application in filtering.** The problem of filtering is to estimate an hidden state $x_k$ from past observations $y_{1:k} := (y_1, \ldots, y_k)$. Assuming the state is Markovian and the observation conditionally white, the solution of the problem is given by the well-known recursion for the *a posteriori* probability

$$p(x_k|y_{1:k}) = \frac{p(y_k|x_k)p(x_k|y_{1:k-1})}{\int p(y_k|x_k)p(x_k|y_{1:k-1})dx_k} \qquad (9.80)$$

which is nothing but Bayes law where the *prior* probability is $p(x_k|y_{1:k-1})$. Therefore, kernel Bayes rules can realize this recursion in a RKHS. Let $m_{z,k|l}$ be the embedding of $p(z_k|y_{1:l})$ in $V_z$ where $z$ is either $x$ or $y$, and $V_z$ is the RKHS associated to kernel $K_z(.,.)$.

Embedding the previous recursion in a RKHS amounts to applying kernel Bayes rule (9.67) with *prior* probability $p(x_k|y_{1:k-1})$ and likelihood $p(y_k|x_k)$, to obtain the embedding $m_{x,k|k}$ of the *posterior* probability.

Firstly, the link between the embedding $m_{x,k|k-1}$ of the *prior* probability, and $m_{x,k-1|k-1}$ is obtained by applying (9.58), or

$$m_{x,k|k-1} = \Sigma_{x_k x_{k-1}} \Sigma_{x_{k-1} x_{k-1}}^{-1} m_{x,k-1|k-1} \qquad (9.81)$$

Then the kernel sum rule for $p(y_k|y_{1:k-1}) = \int p(y_k|x_k)p(x_k|y_{1:k-1})dx_k$ and the kernel chain rule for $p(y_k|x_k)p(x_k|y_{1:k-1})$ have to be used. The application of (9.68) and (9.69) will respectively give the operators $c_{x,k|k}$ and $c_{yy,k|k-1}$ needed for kernel Bayes rule (9.67),

$$m_{x,k|k} = c_{x,k|k} c_{yy,k|k-1}^{-1} K_y(.,y_k) \qquad (9.82)$$

The operator $c_{yy,k|k-1}$ satisfies, according to the sum rule (9.69)

$$c_{yy,k|k-1} = \Sigma_{(y_k y_k) x_k} \Sigma_{x_k x_k}^{-1} m_{x,k|k-1} \qquad (9.83)$$

whereas the operator $c_{x,k|k}$ is provided by the chain rule (9.68), or

$$c_{x,k|k} = \Sigma_{y_k x_k} \Sigma_{x_k x_k}^{-1} m_{x,k|k-1} \qquad (9.84)$$

These last four equations provide the embedding of the optimal filtering solution into the RKHS $V_x$.

To obtain a matrix representation for all these rules, $N+1$ samples of the couple $(x_k, y_k)$ are supposed to be observed. At time $k-1$ the kernel conditional mean is given by

$$m_{x,k-1|k-1}(.) = \sum_{i=1}^{N} \alpha_i^{k-1} K_x(., x_i) = \boldsymbol{k}_X(.)\boldsymbol{\alpha}^{k-1} \tag{9.85}$$

and therefore, the matrix representation of (9.81) is given by

$$m_{x,k|k-1}(.) = \boldsymbol{k}_{X+}(.)(\boldsymbol{K}_x + \lambda I)^{-1}\boldsymbol{K}_x\boldsymbol{\alpha}^{k-1} \tag{9.86}$$

where $\boldsymbol{k}_{X+}(.) = (K_x(., x_2), \dots, K_x(., x_{N+1}))$ and $\boldsymbol{K}_x$ is the Gram matrix built on $x_1, \dots, x_N$.

Then, the operator $c_{x,k|k}$ given by equation (9.84) has the representation

$$
\begin{aligned}
c_{x,k|k} &= \boldsymbol{k}_Y(.)(\boldsymbol{K}_x + \lambda I)^{-1}\big(m_{x,k|k-1}(x_1), \dots m_{x,k|k-1}(x_N)\big)^\top \\
&= \boldsymbol{k}_Y(.)(\boldsymbol{K}_x + \lambda I)^{-1}\boldsymbol{K}_{XX+}(\boldsymbol{K}_x + \lambda I)^{-1}\boldsymbol{K}_x\boldsymbol{\alpha}^{k-1} \\
&= \sum_{i=1}^{N} \boldsymbol{\mu}_i^k K_y(., y_i) \quad \text{where} \tag{9.87}
\end{aligned}
$$

$$\boldsymbol{\mu}^k = (\boldsymbol{K}_x + \lambda I)^{-1}\boldsymbol{K}_{XX+}(\boldsymbol{K}_x + \lambda I)^{-1}\boldsymbol{K}_x\boldsymbol{\alpha}^{k-1} \tag{9.88}$$

Likewise the operator $c_{yy,k|k-1}$ in (9.83) has the representation

$$c_{yy,k|k-1} = \sum_i \boldsymbol{\mu}_i^k K_y(., y_i) \otimes K_y(., y_i) \tag{9.89}$$

Finally, $m_{x,k|k}(.) = \boldsymbol{k}_X(.)\boldsymbol{\alpha}^{k-1}$ where parameter $\boldsymbol{\alpha}^k$ reads

$$\boldsymbol{\alpha}^k = \Lambda^k \boldsymbol{K}_y \left((\boldsymbol{K}_y\Lambda^k)^2 + \varepsilon I\right)^{-1}\Lambda^k\boldsymbol{k}_Y(y_k) \tag{9.90}$$

where $\Lambda^k = \text{Diag}(\boldsymbol{\mu}^k)$.

To synthesize: the kernel conditional mean is given by

$$m_{x,k|k}(.) = \sum_{i=1}^{N} \alpha_i^k K_x(., x_i) = \boldsymbol{k}_X(.)\boldsymbol{\alpha}^k \tag{9.91}$$

where the vector $\boldsymbol{\alpha}^k$ satisfies the recursion

$$\boldsymbol{\mu}^k = (\boldsymbol{K}_x + \lambda I)^{-1}\boldsymbol{K}_{XX+}(\boldsymbol{K}_x + \lambda I)^{-1}\boldsymbol{K}_x\boldsymbol{\alpha}^{k-1} \tag{9.92}$$

$$\Lambda^k = \text{Diag}(\boldsymbol{\mu}^k) \tag{9.93}$$

$$\boldsymbol{\alpha}^k = \Lambda^k\boldsymbol{K}_y\left((\boldsymbol{K}_y\Lambda^k)^2 + \varepsilon I\right)^{-1}\Lambda^k\boldsymbol{k}_Y(y_k) \tag{9.94}$$

The matrix $(\boldsymbol{K}_x + \lambda I)^{-1}\boldsymbol{K}_{XX+}(\boldsymbol{K}_x + \lambda I)^{-1}\boldsymbol{K}_x$ can obviously be pre-computed. With the notation taken, the first useful time for a real estimation is $k = N + 2$, since the first $N + 1$ dates are used for learning. To initialize, $\widehat{\pi}(x_{N+1}) = E[K_x(., x_{N+1})] = \Sigma_{xy}\Sigma_{yy}^{-1}K(., y_{N+1})$ can be used, and thus $\boldsymbol{\alpha}^{N+1} = (\boldsymbol{K}_y + \lambda I)^{-1}\boldsymbol{k}_Y(y_{N+1})$.

The outcome of the algorithm is an estimation of the embedding of the *a posteriori* measure. If an estimate of $E[f(x_k)|y_1, \ldots, y_k]$ where $f \in V_x$ is seeked for, the definition $E[f(x_k)|y_1, \ldots, y_k] = \langle f, m_{k|k}\rangle$ is applied. However, if $f$ does not belong to the RKHS, this can not be applied. A possibility is to find the pre-image $x^k$ whose image $K_x(., x^k)$ is the closest to the embedding of the *posterior* probability. For radial kernel $K_x(., x) = f(\|. - x\|^2)$ this can be solved efficiently if closeness is measured using the RKHS norm [60]. Indeed, searching for $\min_x \|K_x(., x) - \sum_i K_x(., x_i)\alpha_i^t\|$ leads to the fixed point condition $x = \sum_i x_i f'(\|x - x_i\|^2)\alpha_i^t / \sum_i f'(\|x - x_i\|^2)\alpha_i^t$. A solution can be obtained sequentially as

$$x_n^t = \frac{\sum_i x_i f'(\|x_{n-1}^t - x_i\|^2)\alpha_i^t}{\sum_i f'(\|x_{n-1}^t - x_i\|^2)\alpha_i^t} \tag{9.95}$$

No convergence guarantees exist for this procedure, and *ad-hoc* stategies are usually called for, such as running the algorithms with several different initial conditions, ...

**Illustration in prediction.** An example of kernel Bayes filtering on a prediction problem is presented. In the example taken from [58], a signal $z_n$ is generated using the following nonlinear autoregressive model

$$z_n = \frac{(8 - 5e^{-z_{n-1}^2})z_{n-1}}{10} - \frac{(3 + 9e^{-z_{n-1}^2})z_{n-2}}{10} + \frac{\sin(\pi z_{n-1})}{10} + \varepsilon_n \tag{9.96}$$

where $\varepsilon_n$ is chosen from an i.i.d. sequence of zero mean Gaussian random variables of standard deviation set to 0.1 in the illustration. A

plot of the phase space $(z_{n-1}, z_n)$ is depicted in the left plot in figure 9.3. The learning set is composed of the first 512 samples of the signal $z_n$. The data set is composed by the state $x_n = (z_n, z_{n+1})^\top$ and the observation is simply $y_n = z_n$. The state $x_n$ is estimated and $\widehat{z_{n+1}}$ is defined as the second coordinate of the estimated state. The parameters chosen for the simulation are not optimized at all. The kernels are Gaussian kernels with variance parameters set to 0.5. The regularisation parameters were set to $10^{-4}$.
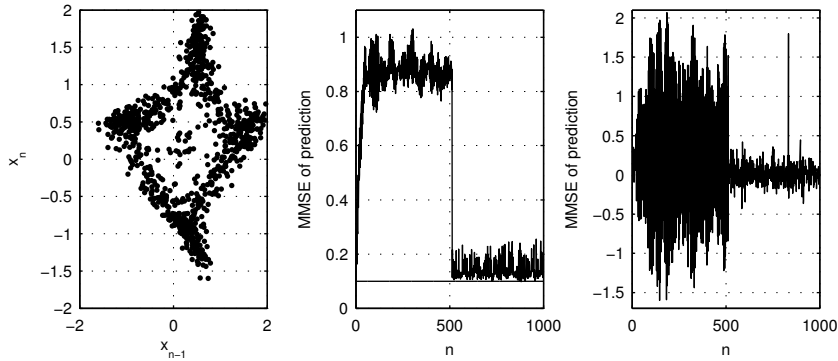


**Figure 9.3:** Left plot: phase diagram of the time series used to illustrate kernel Bayes filtering. Middle plot: Square root of the mean square error of the predictor obtained bby averagin gover 200 snapshots. The first 512 sample represent the power of the signal to be estimated since this interval correspond to learning, and no estimationis performed: the estimator is set to zero during this period of time. The horizontal line marks the amplitude 0.1 which corresponds to the standard deviation of the dynamical noise, and is therefore the optimal power to reach. Right plot: a particular snaphot. The outlier corresponds to a bad estimate by the pre-image algortihm.

As shown in the right plot of figure 9.3, the error $e_{n+1|n} = z_{n+1} - \widehat{z_{n+1}}$ is almost equal to to the dynamical noise $\varepsilon_{n+1}$. The standard deviation of $e_{n+1|n}$ is estimated by averaging the above procedure over 100 snaphots; its value after convergence is 0.13, close to 0.1, the standard deviation of the dynamical noise. The square root of the mean square error of prediction is depicted in the middle plot. As recalled earlier learning is performed on the first 512 samples and estimation begins at time 514. Therefore, the first 512 samples presented are respectively the signal $z_{n+1}$ in the right plot, and its power in the middle plot.

Finally, the presence of an outlier at time around 850 in the right plot is interesting. It comes from a bad estimation of the pre-image. The algorithm to go back from the RKHS to the input space is very important to produce a good estimation. In the implementation used for this illustration, the convergence of the pre-image algorithm is controlled. If divergence occurs, or if no convergence (oscillations) occurs, the algorithm is re-initialized and run again. The initial condition is randomly chosen. The outlier in the figure corresponds to a point which is just below the threshold used to decide of divergence.

This example illustrates the effectiveness of kernel Bayes rules here applied in prediction. Note that as developped, kernel Bayes filtering is not an on-line procedure *per se*, since learning is done on a fixed set of data. Thus the approach here should be compared with gaussian processes regression or kernel regression. All these methods will have the same level of quality. On-line learning with kernel Bayes rule remains to be developped, especially for cases in which the *a priori* and/or the observation are nonstationary. For these cases, on-line filtering approaches have been developed in a non Bayesian framework, as evoked in the following.

## 9.6   On-line filtering using RKHS

Filtering may be seen is an estimation problem. Before discussing on-line filtering in RKHS it is worth presenting some consideration of estimation in RKHS.

**On estimation in RKHS.** Consider for example the following problem which is at the root of filtering. Let $y$ be a real square integrable random variable, and consider using this random variable to estimate another random variable, say $X$ with values in $\mathbb{X}$. Minimizing the mean square error $E[(Y-h(X))^2]$ is often used. With the only constraint that $h$ is measurable, the optimal function $h$ is given by $h(X) = E[Y|X]$, the conditional mean.

Usually the conditional mean is very difficult to evaluate so that the function $h$ is searched for in some more constrained classes of functions. If $h$ is restricted to be a linear functional, the solution is

Wiener filter. For example, if $\mathbb{X} = \mathbb{R}^n$, the vector $\boldsymbol{h}$ which minimizes $E[(Y - \boldsymbol{h}^\top X)^2]$ is seeked for. The solution satisfies the well known equation $E[YX] = E[XX^\top]\boldsymbol{h}$. Here, the problem is formalized and solved when functions $h$ is searched for in a RKHS whose kernel is a kernel on $\mathbb{X}$. This corresponds to the important situation where the optimal filter $h$ is as a linear functional of the transformed observation (by the kernel.)

Thus let $K$ be a kernel on $\mathbb{X}$ and $\mathcal{H}$ its associated reproducing kernel Hilbert space of functions from $\mathbb{X}$ to $\mathbb{R}$. The best estimator satisfies

$$f_0 = \arg\min_{f \in \mathcal{H}} E\left[(Y - f(X))^2\right] \tag{9.97}$$

Let $R(f) = E\left[(Y - f(X))^2\right]$. To find the minimun, a necessary condition is to set to zero the Gateaux derivative in every possible direction [14], *i.e.* since $R$ is a functional, to solve

$$\frac{dR(f + \varepsilon\varphi)}{d\varepsilon}\bigg|_{\varepsilon=0} = 0, \quad \forall\varphi \in V \tag{9.98}$$

This is equivalent to

$$\Sigma_{YX}\varphi = \langle\varphi, \Sigma_{XX}f\rangle, \quad \forall\varphi \in V \tag{9.99}$$

where

$$\Sigma_{XX} : V \longrightarrow V$$
$$f \longmapsto \Sigma_{XX}f := E[\langle f, K(.,X)\rangle K(.,X)] \tag{9.100}$$
$$\Sigma_{YX} : V \longrightarrow \mathbb{R}$$
$$f \longmapsto \Sigma_{YX}f := E[\langle f, K(.,X)\rangle Y] \tag{9.101}$$

are the covariance and the correlation operators. The correlation operator is in this particular case a bounded linear functional and thanks to the Riesz representation theorem can be written as $\Sigma_{YX}f := \langle f, E[YK(.,X)]\rangle$. Then the optimal function is found if we solve

$$\langle\varphi, E[YK(.,X)]\rangle = \langle\varphi, \Sigma_{XX}f\rangle \quad \forall\varphi \in V \tag{9.102}$$

Since this is valid for any $\varphi$ the solution is given by the function that solves $\Sigma_{XX}f = E[YK(.,X)]$. Practically of course, the covariance and

correlation operators have to be estimated from a finite set of data. Let $(y_i, x_i)$ these data. From the representer theorem we know that the functions are to be searched for in the subspace of $\mathcal{H}$ generated by the $K(., x_i)$. Then the preceding equation has an empirical counterpart which reads

$$\langle \sum_i \alpha_i K(., x_i), \sum_j y_j K(., x_j)] \rangle \;=\; \boldsymbol{\alpha}^T K_X K_X \boldsymbol{\beta}, \quad \forall \boldsymbol{\alpha} \quad (9.103)$$

$$\boldsymbol{\alpha}^T K_X \boldsymbol{y} \;=\; \boldsymbol{\alpha}^T K_X K_X \beta, \quad \forall \boldsymbol{\alpha} \quad (9.104)$$

where $K_X$ is the Gram matrix, and $f(.) := \sum_j \beta_j K(., x_j)$. This is the same solution as the regression case in machine learning, as developed in the following paragraph.

**Filtering in a RKHS.** Filtering in a RKHS is a particular problem of regression. The data $(y_i, x_i)$ considered above can be seen as observations that can be explained using a regression $Y = f(X)$ where $f$ is searched for in a RKHS. From a signal processing perspective, $y_i$ may represent the value of signal $y$ at time $i$, whereas vector $x_i$ contains past samples of another signal. $y = f(X)$ is thus considered as a black box model, and the aim of the optimal estimation is to identify $f(.)$.

Thus, let $y_n$ and $z_n$ be two real valued random signals, and suppose we look for a function which minimizes the power of the error $y_n - f(x_n)$, where $x_n = (z_n, z_{n-1}, \ldots, z_{n-d+1})^\top \in \mathbb{R}^d$. $\mathbb{R}^d$ is embedded into a RKHS using a kernel $K$ on $\mathbb{R}^d$, and the optimizing function is seeked for in the RKHS. The optimization is realized on a set of $N$ observed data $(y_i, x_i)$. Furthermore, since most of the interesting RKHS are of either high or infinite dimension, overfitting is likely to occur, and some regularization must be included. The norm in the RKHS is used to constrain the function. Thus, the optimization problem can be written as

$$f_0(.) = \arg \min_{f \in V} \sum_i \left| y_i - f(x_i) \right|^2 + \lambda \left\| f \right\|_V^2 \qquad (9.105)$$

The representer theorem [37, 60] states that the optimizing function $f_0$ belongs to the the subspace of $V$ generated by $\{K(., x_i)\}, i = 1, \ldots, N$. If $f_0 = \sum_i \alpha_{0,i} K(., x_i)$, the preceding program is equivalent to

$$\boldsymbol{\alpha}_0 = \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^N} (\boldsymbol{y}_i - \boldsymbol{K}_x \boldsymbol{\alpha})^\top (\boldsymbol{y}_i - \boldsymbol{K}_x \boldsymbol{\alpha}) + \lambda \boldsymbol{\alpha}^\top \boldsymbol{K}_x \boldsymbol{\alpha} \qquad (9.106)$$

the solution of which is given by

$$\boldsymbol{\alpha}_0 = (\boldsymbol{K}_x + \lambda\boldsymbol{I})^{-1}\boldsymbol{y} \tag{9.107}$$

Then when a new datum $x_k$ is observed, $k > N$, the corresponding estimate $y_k$ writes

$$y_k = \sum_{i=1}^{N} \alpha_{0,i}K(x_k, x_i) = \boldsymbol{K}_x(x_k)^\top(\boldsymbol{K}_x + \lambda\boldsymbol{I})^{-1}\boldsymbol{y} \tag{9.108}$$

where $\boldsymbol{K}_x(x_k)^\top = \big(K(x_k, x_1), \ldots, K(x_k, x_N)\big)$. This regularized solution is also known as the ridge regression. The result may appear strange since in the linear case it reads $y_k = x_k^\top\boldsymbol{X}(\boldsymbol{X}\boldsymbol{X}^\top + \lambda\boldsymbol{I})^{-1}\boldsymbol{y}$ where $\boldsymbol{X}$ is the design matrix $(x_1, \ldots, x_N)^\top$. But the well known trick $\boldsymbol{X}(\boldsymbol{X}\boldsymbol{X}^\top + \lambda\boldsymbol{I})^{-1} = (\boldsymbol{X}^\top\boldsymbol{X} + \lambda\boldsymbol{I})^{-1}\boldsymbol{X}^T$ allows to recover the usual Wiener filtering form $y_k = x_k^\top(\boldsymbol{X}^\top\boldsymbol{X} + \lambda\boldsymbol{I})^{-1}\boldsymbol{X}^\top\boldsymbol{y}$.

**On-line filtering.** As an application the problem of on-line filtering is considered. The aim here is to use the previous filter in situations where data arrives in streams or on-line, and can not be processed in batch mode. Thus a recursive structure is needed to refresh the filter when a new datum is acquired, so as to produce a new estimate without the need of re-calculating. This problem has seen some important developments in the last decade, with the presentation of the kernel Recursive Least Square (kRLS) algorithm [19], followed by many works such as those reviewed recently in [61].

The biggest difficulty for on-line kernel algorithms lies in the fact that they inherently manipulates Gram matrices whose sizes grow linearly with the size of the data. Thus, on-line kernel algorithms necessarily include a sparsification procedure which allow to approximate the Gram matrices needed as well as possible. For a recent thorough review of on-line kernel algorithms for signal processing problems, we refer to [61]. Here, the kRLS in which the sparsification is the Approximate Linear Dependence criterion is presented. The kLMS which uses the coherence criterion as sparsification criterion is proposed as a simplified alternative. We already used the coherence criterion in §9.3.

The sparsification procedure is based on a recursive building of a dictionary using the ALD. Initializing the dictionary $\mathcal{D}_1 = 1$ with the

first datum acquired $x_1$, the rule to create the dictionary is

$$\mathcal{D}_n = \begin{cases} \mathcal{D}_{n-1} \cup \{n\} & \text{if } K(.,x_n) \text{ is not ALD of } \{K(.,x_\alpha)\}_{\alpha \in \mathcal{D}_{n-1}} \\ \mathcal{D}_{n-1} & \text{otherwise} \end{cases} \tag{9.109}$$

Approximate linear dependence is measured according to the mean square error in the prediction of $K(.,x_n)$ from the members $K(.,x_\alpha)$ of the dictionary at time $n-1$. If this error is greater than a user specified threshold, the ALD is rejected, and $K(.,x_n)$ carries enough new information to be incorporated into the dictionary.

Thus at each time $n$, the approximation is $\hat{k}(.,x_n) = \sum_{i=1}^{d_{n-1}} a_{n,i} K(.,x_{\alpha_i})$, where $d_n = |\mathcal{D}_n|$, and where the coefficients are obtained as in the regression framework presented above, that is

$$\boldsymbol{a}_n = \widehat{\boldsymbol{K}}_{n-1}^{-1} \boldsymbol{k}_{n-1}(x_n) \tag{9.110}$$

and the minimum mean square error reached is given by

$$\begin{aligned} e_n &= K(x_n,x_n) - \boldsymbol{k}_{n-1}(x_n)^\top \boldsymbol{a}_n \\ &= K(x_n,x_n) - \boldsymbol{k}_{n-1}(x_n)^\top \widehat{\boldsymbol{K}}_{n-1}^{-1} \boldsymbol{k}_{n-1}(x_n) \end{aligned} \tag{9.111}$$

In these equations, $(\widehat{\boldsymbol{K}}_{n-1})_{ij} = K(x_{\alpha_i}, x_{\alpha_j})$ is the Gram matrix evaluated on the dictionary at time $n-1$ and $\boldsymbol{k}_{n-1}(x_n)^\top = \left( K(x_n,x_{\alpha_1}), \ldots, K(x_n,x_{\alpha_{d_{n-1}}}) \right)$. The test for ALD consists in comparing $e_n$ to a given value $e_0$. If the ALD is accepted or $e_n < e_0$, the dictionary is not modified and $\mathcal{D}_n = \mathcal{D}_{n-1}$. If the ALD is rejected, $\mathcal{D}_n = \mathcal{D}_{n-1} \cup \{n\}$ and obviously $\hat{k}(.,x_n) = K(.,x_n)$.

Interestingly, this allows to obtain an approximation of the Gram matrix calculated over all the observation measured up to time $n$. Indeed, since for all $n$, $\hat{k}(.,x_n) = \sum_{i=1}^{d_{n-1}} a_{n,i} K(.,x_{\alpha_i})$, the $(l,m)$ entry of the full Gram matrix reads approximately

$$\begin{aligned} K(x_l, x_m) &\approx \langle \hat{k}(.,x_l), \hat{k}(.,x_m) \rangle \tag{9.112} \\ &= \sum_{i,j=1}^{d_{m-1}} a_{l,i} a_{m,j} K(x_{\alpha_i}, x_{\alpha_j}) \tag{9.113} \end{aligned}$$

This supposes without loss of generality that $l \leq m$, and implicitely set $a_{l,i} = 0$ as soon as $i > d_{l-1}$. This allows to store the coefficient $a_{l,i}$

into a matrix $\boldsymbol{A}_l$ of appropriate dimension and to write

$$\boldsymbol{K}_n \approx \boldsymbol{A}_n \widehat{\boldsymbol{K}}_n \boldsymbol{A}_n^\top \tag{9.114}$$

The matrix $\boldsymbol{A}_n$ is updated as $(\boldsymbol{A}_{n-1}^\top \boldsymbol{a}_n^\top)^\top$ if ALD is accepted, and as

$$\left( \begin{array}{c|c} \boldsymbol{A}_{n-1} & 0 \\ \hline 0 & 1 \end{array} \right) \tag{9.115}$$

if ALD is rejected.

For the problem of on-line regression, the following problem has to be solved

$$\boldsymbol{\alpha}_n = \arg\min_{\boldsymbol{\alpha}} (\boldsymbol{y}_n - \boldsymbol{K}_n \boldsymbol{\alpha})^\top (\boldsymbol{y}_n - \boldsymbol{K}_n \boldsymbol{\alpha}) \tag{9.116}$$

a program replaced by

$$\boldsymbol{\alpha}_n = \arg\min_{\boldsymbol{\alpha}} (\boldsymbol{y}_n - \boldsymbol{A}_n \widehat{\boldsymbol{K}}_n \boldsymbol{\alpha})^\top (\boldsymbol{y}_n - \boldsymbol{A}_n \widehat{\boldsymbol{K}}_n \boldsymbol{\alpha}) \tag{9.117}$$

where the substitution $\boldsymbol{\alpha} \leftrightarrow \boldsymbol{A}_n^\top \boldsymbol{\alpha}$ has been made. The optimal parameter $\boldsymbol{\alpha}_n$ at time $n$ is thus given by the pseudo-inverse of $\boldsymbol{A}_n \widehat{\boldsymbol{K}}_n$ applied to $\boldsymbol{y}_n$, or after elementary manipulations

$$\boldsymbol{\alpha}_n = \widehat{\boldsymbol{K}}_n^{-1} (\boldsymbol{A}_n^\top \boldsymbol{A}_n)^{-1} \boldsymbol{A}_n^\top \boldsymbol{y}_n \tag{9.118}$$

This form allows an easy transformation into a recursive form. Basically, at time $n$, if ALD is accepted and the dictionary does not change $\mathcal{D}_n = \mathcal{D}_{n-1}$, then $\widehat{\boldsymbol{K}}_n = \widehat{\boldsymbol{K}}_{n-1}$, $\boldsymbol{A}_n$ is updated as $(\boldsymbol{A}_{n-1}^\top \boldsymbol{a}_n^\top)^\top$ and $\boldsymbol{P}_n = (\boldsymbol{A}_n^\top \boldsymbol{A}_n)^{-1}$ is updated as usual as

$$\boldsymbol{P}_n = \boldsymbol{P}_{n-1} - \frac{\boldsymbol{P}_{n-1} \boldsymbol{a}_n \boldsymbol{a}_n^\top \boldsymbol{P}_{n-1}}{1 + \boldsymbol{a}_n^\top \boldsymbol{P}_{n-1} \boldsymbol{a}_n} \tag{9.119}$$

However, if at time $n$ ALD is rejected, the dictionary is increased $\mathcal{D}_n = \mathcal{D}_{n-1} \cup \{n\}$, $\boldsymbol{K}_n^{-1}$ is updated as

$$\begin{aligned} \widehat{\boldsymbol{K}}_n^{-1} &= \left( \begin{array}{c|c} \widehat{\boldsymbol{K}}_{n-1} & \boldsymbol{k}_{n-1}(x_n) \\ \hline \boldsymbol{k}_{n-1}(x_n)^\top & K(x_n, x_n) \end{array} \right)^{-1} \\ &= \frac{1}{e_n} \left( \begin{array}{c|c} e_n \widehat{\boldsymbol{K}}_{n-1}^{-1} + \boldsymbol{a}_n \boldsymbol{a}_n^\top & -\boldsymbol{a}_n \\ \hline -\boldsymbol{a}_n^\top & 1 \end{array} \right) \end{aligned} \tag{9.120}$$

$\boldsymbol{A}_{n-1}$ updated according to (9.115), and $\boldsymbol{P}_n$ follows

$$\boldsymbol{P}_n = \left( \begin{array}{c|c} \boldsymbol{P}_{n-1} & 0 \\ \hline 0 & 1 \end{array} \right) \tag{9.121}$$

Using some algebra then leads to the following kRLS algorithm. Initialize $\mathcal{D}_1 = \{1\}, \boldsymbol{\alpha}_1 = y_1/K(x_1, x_1)$ and then for $n \geq 2$

1. $\boldsymbol{a}_n = \widehat{\boldsymbol{K}}_{n-1}^{-1} \boldsymbol{k}_{n-1}(x_n)$ and $e_n = K(x_n, x_n) - \boldsymbol{k}_{n-1}(x_n)^{\top} \boldsymbol{a}_n$

2. **if** $e_n \leq e_0$, $\mathcal{D}_n = \mathcal{D}_{n-1}$, update $\boldsymbol{P}_n$ according to (9.119), $\widehat{\boldsymbol{K}}_n^{-1} = \widehat{\boldsymbol{K}}_{n-1}^{-1}$, $\boldsymbol{A}_n = (\boldsymbol{A}_{n-1}^{\top} \boldsymbol{a}_n^{\top})^{\top}$ and set

$$\boldsymbol{\alpha}_n = \boldsymbol{\alpha}_{n-1} + \frac{\widehat{\boldsymbol{K}}_{n-1}^{-1} \boldsymbol{P}_{n-1} \boldsymbol{a}_n}{1 + \boldsymbol{a}_n^{\top} \boldsymbol{P}_{n-1} \boldsymbol{a}_n} (y_n - \boldsymbol{k}_{n-1}(x_n)^{\top} \boldsymbol{\alpha}_{n-1}) \tag{9.122}$$

**else** $\mathcal{D}_n = \mathcal{D}_{n-1} \cup \{n\}$, update $\boldsymbol{P}_n$ according to (9.121), $\widehat{\boldsymbol{K}}_n^{-1}$ according to (9.120), $\boldsymbol{A}_n$ according to (9.115) and set

$$\boldsymbol{\alpha}_n = \begin{cases} \boldsymbol{\alpha}_{n-1} - \frac{\boldsymbol{a}_n}{e_n}(y_n - \boldsymbol{k}_{n-1}(x_n)^{\top} \boldsymbol{\alpha}_{n-1}) \\ \frac{1}{e_n}(y_n - \boldsymbol{k}_{n-1}(x_n)^{\top} \boldsymbol{\alpha}_{n-1}) \end{cases} \tag{9.123}$$

The complexity of the kRLS is dominated by inverting the matrices. If a lower complexity is required, simpler approaches may be used, such as the kernel Least Mean Square (kLMS) algorithm. Its expression is just written down, referring to [40, 61] for more details and variants. It uses a dictionary as well. However, since the ALD requires the propagation of the inverse of $\widehat{\boldsymbol{K}}_n$, the criterion can not be used here because a motivation for the kLMS is to reduce complexity. Thus another criterion must be used. The simpler one up to today is the coherence criterion proposed in this context in [58], and that we detailed previously. If the new datum is not coherent enough with the dictionary it is included. Coherence is measured in the RKHS and therefore simply uses the kernel evaluated at the new datum and at the members of the dictionary. The normalized kLMS recursively and adaptively computes the coefficient vector according to the following steps:

Initialize $\mathcal{D}_1 = \{1\}, \boldsymbol{\alpha}_1 = y_1/K(x_1, x_1)$ and then for $n \geq 2$

1. $e_n = \max_{\alpha \in \mathcal{D}_{n-1}} |K(x_n, x_\alpha)|$

2. **if** $e_n \geq e_0$, $\mathcal{D}_n = \mathcal{D}_{n-1}$, set

$$\boldsymbol{\alpha}_n = \boldsymbol{\alpha}_{n-1} + \frac{\lambda\big(y_n - \boldsymbol{k}_{n-1}(x_n)^\top \boldsymbol{\alpha}_{n-1}\big)}{\varepsilon + \boldsymbol{k}_{n-1}(x_n)^\top \boldsymbol{k}_{n-1}(x_n)} \boldsymbol{k}_{n-1}(x_n) \qquad (9.124)$$

**else** $\mathcal{D}_n = \mathcal{D}_{n-1} \cup \{n\}$, set

$$\boldsymbol{\alpha}_n = \begin{pmatrix} \boldsymbol{\alpha}_{n-1} \\ 0 \end{pmatrix} + \frac{\lambda\big(y_n - \boldsymbol{k}_{n-1}(x_n)^\top \boldsymbol{\alpha}_{n-1}\big)}{\varepsilon + \boldsymbol{k}_n(x_n)^\top \boldsymbol{k}_n(x_n)} \begin{pmatrix} \boldsymbol{k}_{n-1}(x_n) \\ K(x_n, x_n) \end{pmatrix} (9.125)$$

where recall that $\boldsymbol{k}_n(x_n)^\top = (K(x_n, x_{\alpha_1}), \ldots, K(x_n, x_{\alpha_{d_n}})$

$\lambda$ is the step size parameter that controls as usual the variance-bias trade-off. A study of convergence of this algorithm as been proposed recently in the Gaussian case for Gaussian kernels [47].

**An illustration.** The two algorithms are illustrated on the example developped for illustrating kernel Bayes filtering. In figure 9.4 the left plot depicts the square root of the mean square error as a function of time. The kernel used is the Gaussian $\exp(-\|x - y\|^2/\sigma^2)$. Parameters chosen are $\sigma^2 = 1/3.73$, $\lambda = 0.09$, $\varepsilon = 0.03$. These values were used in [58] and are taken here for the sake of simplicity. For the left plot, $e_0$ has been set to 0.55 for the kRLS and the kLMS as well. This choice was made since it ensures the dictionaries are of the same size on average for the two algorithms (27 for this value of $e_0$). The convergence curves were obtained by averaging 500 snapshots. The middle plot displays the logarithm of the size of the dictionary after convergence as a function of $e_0$. The same parameter is used here but has a different meaning for the two algorithms. For the kLMS, the larger the parameter, the easier the dictionary is increased. For the kRLS, this goes the other way around. $e_0$ is varied from 0.01 to 1 in a non uniform way. The right plot shows the asymptotic mean square error (evaluated by averaging the last 500 samples of the convergence curves, these curves being obtained by averaging 500 snapshots). Interestingly with these two plots is the fact that the loss in MMSE is small for a high compression. For example, in the kRLS, the optimal performance of 0.1 is nearly obtained as soon
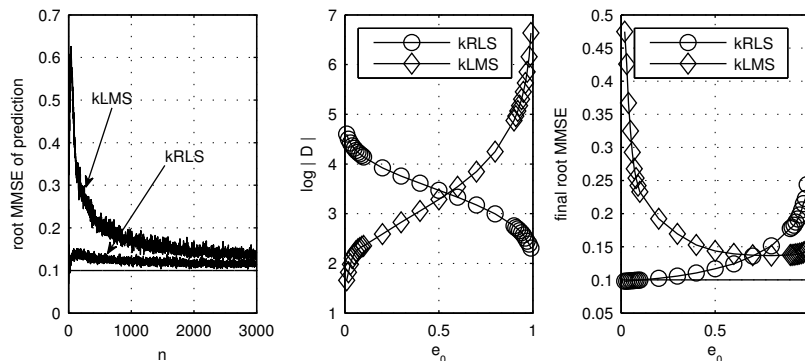
**Figure 9.4:** Left plot: Square root of the mean square error obtained by the kLMS and kRLS for the prediction problem studied in the text. The horizontal line marks the amplitude 0.1 which corresponds to the standard deviation of the dynamical noise, and is therefore the optimal power to reach. Middle plot: Size of the dictionary at convergence as a function of parameter $e_0$. Right plot: Square root of the mean square error obtained by the kLMS and kRLS at convergence as a function of parameter $e_0$.

as $e_0 = 0.1$ , for which the size of the dictionary is only 55! Likewise, the best performance of the kLMS (about 0.13) is obtained for $e_0 = 0.7$ for which the size of the dictionary is 47. This simple example show the effectiveness of the algorithms coupled with simple yet powerful sparsification techniques. These sparsification procedures open up the use of kernel methods to very large data sets as well as on-line efficient algorithms.

# Bibliography

[1] P. O. Amblard. Versions récursives et adaptatives d'une mesure d'indépendance. In *GRETSI 2013, Brest, France*, 2013.

[2] P. O. Amblard and J. H. Manton. On data sparsification and a recursive algorithm for estimating a kernel-based measure of independence. In *proc. ICASSP, Vancouver, Canada*, 2013.

[3] N. Aronszajn. Theory of Reproducing Kernels. *Transactions of the American Mathematical Society*, 68:337–404, 1950.

[4] F. R. Bach. Sharp analysis of low-rank kernel matrix approximations. *JMLR: Workshop and Conference Proceedings*, 30:1–25, 2013.

[5] F. R. Bach and M. I. Jordan. Kernel independent component analysis. *Journal of Machine Learning Research*, 3(1–48), 2002.

[6] C. R. Baker. Joint measures and cross-covariance operators. *Trans. Am. Math. Soc.*, 186:273–289, 1973.

[7] R. J. Barton and H. V. Poor. An RKHS approach to robust $L^2$ estimation and signal detection. *Information Theory, IEEE Transactions on*, 36(3):485–501, 1990.

[8] S. Bergman. *The kernel function and conformal mapping*, volume 5. American Mathematical Soc., 1970.

[9] A. Berlinet and C. Thomas-Agnan. *Reproducing kernel Hilbert spaces in probability and statistics.* Kluwer Academic Publishers, 2004.

[10] A. Berlinet and C. Thomas-Agnan. *Reproducing Kernel Hilbert Spaces in Probability and Statistics.* Springer, Dec. 2012.

[11] R. F. Brown. *A topological introduction to nonlinear analysis.* Birkhäuser Boston, Inc., Boston, MA, second edition, 2004.

[12] R. F. Curtain and H. Zwart. *An introduction to infinite-dimensional linear systems theory*, volume 21 of *Texts in Applied Mathematics.* Springer-Verlag, New York, 1995.

[13] G. Da Prato. *An Introduction to Infinite-Dimensional Analysis.* Universitext. Springer, 2006.

[14] L. Debnath and P. Mikusiński. *Introduction to Hilbert spaces with applications, 3rd ed.* Elsevier Academic Press, 2005.

[15] M. F. Driscoll. The reproducing kernel Hilbert space structure of the sample paths of a Gaussian process. *Z. Wahrscheinlichkeitstheorie und Verw. Gebiete*, 26:309–316, 1973.

[16] D. L. Duttweiler and T. Kailath. RKHS approach to detection and estimation problems. IV. Non-Gaussian detection. *Information Theory, IEEE Transactions on*, IT-19(1):19–28, 1973.

[17] D. L. Duttweiler and T. Kailath. RKHS approach to detection and estimation problems. V. Parameter estimation. *Information Theory, IEEE Transactions on*, IT-19(1):29–36, 1973.

[18] R. Elliott, L. Aggoun, and J. Moore. *Hidden Markov Models: Estimation and Control.* Applications of mathematics. Springer, 1995.

[19] Y. Engel, S. Mannor, and R. Meir. The kernel recursive least-squares algorithm. *IEEE Trans. on Signal Processing*, 52(8):2275–2284, 2004.

[20] S. Fine and K. Scheinberg. Efficient svm training using low-rank kernel representations. *Journal of Machine Learning Research*, 2:243–264, 2001.

[21] R. Fortet. *Vecteurs, fonctions et distributions aléatoires dans les espaces de Hilbert (Random Vectors, functions and distributions in Hilbert spaces)*. (in French), Hermès, 1995.

[22] K. Fukumizu, F. R. Bach, and M. I. Jordan. Dimensionality reduction for supervised learning with reproducing kernel hilbert spaces. *Journal of Machine Learning Research*, 4:73–99, 2004.

[23] K. Fukumizu, F. R. Bach, and M. I. Jordan. Kernel dimension reduction in regression. *The Annals of Statistics*, 37(4):1871–1905, 2009.

[24] K. Fukumizu, A. Gretton, X. Sun, and B. Scholkopf. Kernel measures of conditional dependence. In *NIPS*, 2007.

[25] K. Fukumizu, L. Song, and A. Gretton. Kernel bayes'rules: Bayesian inference with positive definite kernels. *Journal of Machine Learning Research*, 14:3753–3783, 2013.

[26] A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, and A. Smola. A kernel two sample test. *Journal of Machine Learning Research*, 13:723–773, 2012.

[27] A. Gretton, R. Herbrich, A. Smola, O. Bousquet, and B. Schölkopf. Kernel methods for measuring independence. *Journal of Machine Learning Research*, 6:2075–2129, 2005.

[28] D. Han. *Frames for Undergraduates*. Student mathematical library. American Mathematical Society, 2007.

[29] T. Hida and N. Ikeda. Analysis on Hilbert space with reproducing kernel arising from multiple Wiener integral. In *Proc. Fifth*

*Berkeley Sympos. Math. Statist. and Probability (Berkeley, Calif.,
1965/66). Vol. II: Contributions to Probability Theory, Part 1*,
pages 117–143. Univ. California Press, Berkeley, Calif., 1967.

[30] E. Hille. Introduction to general theory of reproducing kernels.
*The Rocky Mountain Journal of Mathematics*, 2(3):321–368, 1972.

[31] J. Hoffmann-Jørgensen and G. Pisier. The Law of Large Numbers
and the Central Limit Theorem in Banach Spaces. *The Annals of
Probability*, 4(4):587–599, Aug. 1976.

[32] T. Kailath. RKHS approach to detection and estimation problems.
I. Deterministic signals in Gaussian noise. *Information Theory,
IEEE Transactions on*, IT-17(5):530–549, 1971.

[33] T. Kailath and D. Duttweiler. An RKHS approach to detection
and estimation problems. III. Generalized innovations representa-
tions and a likelihood-ratio formula. *Information Theory, IEEE
Transactions on*, IT-18(6):730–745, 1972.

[34] T. Kailath, R. Geesey, and H. Weinert. Some relations among
RKHS norms, Fredholm equations, and innovations representa-
tions. *Information Theory, IEEE Transactions on*, 18(3):341–348,
May 1972.

[35] T. Kailath and H. L. Weinert. An RKHS approach to detection
and estimation problems. II. Gaussian signal detection. *Informa-
tion Theory, IEEE Transactions on*, IT-21(1):15–23, 1975.

[36] G. Kallianpur. The role of reproducing kernel Hilbert spaces in
the study of Gaussian processes. In *Advances in Probability and
Related Topics, Vol. 2*, pages 49–83. Dekker, New York, 1970.

[37] G. S. 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.

[38] F. M. Larkin. Optimal Approximation in Hilbert Spaces with Reproducing Kernel Functions. *Mathematics of Computation*, 24(112):911–921, Oct. 1970.

[39] S. Lauritzen. *Graphical models.* Oxford University Press, 1996.

[40] W. Liu, P. Pokharel, and J. Principe. The kernel least mean squares algorithm. *IEEE Trans. on Signal Processing*, 56(2):543–554, 2008.

[41] D. G. Luenberger. *Optimization by vector space methods.* John Wiley & Sons Inc., New York, 1969.

[42] M. Lukic and J. Beder. Stochastic processes with sample paths in reproducing kernel Hilbert spaces. *Transactions of the American Mathematical Society*, 353(10):3945–3969, 2001.

[43] J. H. Manton. A Primer on Stochastic Differential Geometry for Signal Processing. *Selected Topics in Signal Processing, IEEE Journal of*, 7(4):681–699, 2013.

[44] J. Matoušek. *Thirty-three Miniatures: Mathematical and Algorithmic Applications of Linear Algebra.* Student mathematical library. American Mathematical Society, 2010.

[45] J. 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*, 209:415–446, Jan. 1909.

[46] H. J. Newton. A Conversation with Emanuel Parzen. *Statistical Science*, 17(3):357–378, Aug. 2002.

[47] W. D. Parreira, J. C. M. Bermudez, C. Richard, and J.-Y. Tourneret. Stochastic behavior analysis of the gaussian kernel least mean square algorithm. *IEEE Trans. on Sig. Proc.*, 60(5):2208–2222, 2012.

[48] E. Parzen. Statistical inference on time series by Hilbert space methods. Technical Report 23, Applied Mathematics and Statistics Laboratory, Stanford University, Jan. 1959.

[49] E. Parzen. Extraction and detection problems and reproducing kernel Hilbert spaces. *Journal of the Society for Industrial & Applied Mathematics, Series A: Control*, 1(1):35–62, 1962.

[50] E. Parzen. Probability density functionals and reproducing kernel hilbert spaces. In *Proceedings of the Symposium on Time Series Analysis*, volume 196, pages 155–169. Wiley, New York, 1963.

[51] E. Parzen. Statistical inference on time series by RKHS methods. In *Proc. Twelfth Biennial Sem. Canad. Math. Congr. on Time Series and Stochastic Processes; Convexity and Combinatorics (Vancouver, B.C., 1969)*, pages 1–37. Canad. Math. Congr., Montreal, Que., 1970.

[52] V. I. Paulsen. An introduction to the theory of reproducing kernel Hilbert spaces. `http://www.math.uh.edu/~vern/rkhs.pdf`, 2009.

[53] B. Picinbono and P. Duvaut. Optimal linear-quadratic systems for detection and estimation. *IEEE Trans. on Info. Theory*, 34(2):304–311, 1988.

[54] B. Picinbono and P. Duvaut. Geometrical properties of optimal volterra filters for signal detection. *IEEE Trans. on Info. Theory*, 36(5):1061–1068, 1990.

[55] T. S. Pitcher. Likelihood ratios of Gaussian processes. *Arkiv för Matematik*, 4:35–44 (1960), 1960.

[56] H. Poor. *An Introduction to Signal Detection and Estimation*. A Dowden & Culver book. Springer, 1994.

[57] A. Rényi. On measures of dependence. *Acta Math. Acad. Sci. Hungar.*, 10(2):441–451, september 1959.

[58] C. Richard, J. C. M. Bermudez, and P. Honeine. Online prediction of time series data with kernels. *IEEE Trans. on Signal Processing*, 57(3):1058–1067, Mar 2009.

[59] S. Saitoh. *Theory of reproducing kernels and its applications*, volume 189 of *Pitman Research Notes in Mathematics Series*. Longman Scientific & Technical, Harlow; copublished in the United States with John Wiley & Sons, Inc., New York, 1988.

[60] B. Schölkopf and A. J. Smola. *Learning with kernels*. MIT Press, Cambridge, Ma, USA, 2002.

[61] K. Slavakis, P. Bouboulis, and S. Theodoridis. *Academic Press Library in Signal Processing: Volume 1, Signal Processing Theory and Machine Learning*, chapter ch. 17, Online learning in reproducing kernel Hilbert spaces, pages 883–987. Elsevier, 2014.

[62] C. G. Small and D. L. McLeish. *Hilbert Space Methods in Probability and Statistical Inference*. John Wiley & Sons, Sept. 2011.

[63] L. Song, K. Fukumizu, and A. Gretton. Kernel embeddings of conditional distributions. *IEEE Signal Processing Magazine*, 30(4):98–111, 2013.

[64] L. Song, A. Smola, A. Gretton, J. Bedo, and K. Borgwardt. Feature selection via dependence maximization. *Journal of Machine Learning Research*, 13:1393–1434, 2012.

[65] B. K. Sriperumbudur, K. Fukumizu, and G. R. G. Lanckriet. Universality, characteristic kernels and rkhs embeddings of measures. *Journal of Machine Learning Research*, 12:2389–2410, 2011.

[66] F. Steinke and B. Schölkopf. Kernels, regularization and differential equations. *Pattern Recognition*, 41(11):3271–3286, 2008.

[67] I. Steinwart. On the influence of the kernel on the consistency of support vector machines. *Journal of Machine Learning Research*, 2:67–93, 2001.

[68] F. Stulajter. RKHS approach to nonlinear estimation of random variables. In *Transactions of the Eighth Prague Conference on Information Theory, Statistical Decision Functions, Random Processes (Prague, 1978), Vol. B*, pages 239–246. Reidel, Dordrecht, 1978.

[69] Y. Sun, F. Gomez, and J. Schmidhuber. On the size of the online kernel sparsification dictionary. In *proceedings of ICML, Edinburgh, Scotland*, 2012.

[70] H. J. Sussmann and J. C. Willems. 300 years of optimal control: from the brachystochrone to the maximum principle. *Control Systems, IEEE*, 17(3):32–44, 1997.

[71] A. N. Tikhonov and V. Y. Arsenin. *Solutions of ill-posed problems.* John Wiley&Sons, 1977.

[72] G. Wahba. Spline interpolation and smoothing on the sphere. *SIAM Journal on Scientific and Statistical Computing*, 2(1):5–16, 1981.

[73] G. G. Wahba. *Interpolating Spline Methods for Density Estimation.* Department of Statistics, University of Wisconsin, 1973.

[74] J. Whittaker. *Graphical models in applied multivariate statistics.* Wiley&Sons, 1989.

[75] D. Williams. *Probability with Martingales.* Cambridge mathematical textbooks. Cambridge University Press, 1991.

[76] E. Wong and B. Hajek. *Stochastic Processes in Engineering Systems.* Springer Texts in Electrical Engineering. Springer New York, 2011.

[77] K. Yao. Applications of reproducing kernel Hilbert spaces — bandlimited signal models. *Information and Control*, 11(4):429–444, 1967.