

# Explicit Universal and Approximate-Universal Kernels on Compact Metric Spaces

Eloi Tanguy<sup>1</sup>

<sup>1</sup>Université Paris Cité, CNRS, MAP5, F-75006 Paris, France

June 13, 2025

## Abstract

Universal kernels, whose Reproducing Kernel Hilbert Space is dense in the space of continuous functions are of great practical and theoretical interest. In this paper, we introduce an explicit construction of universal kernels on compact metric spaces. We also introduce a notion of approximate universality, and construct tractable kernels that are approximately universal.

## Table of Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Kernels in Practice and Related Works . . . . .	2
1.2	Elements of RKHS Theory . . . . .	2
1.3	Paper Outline and Contributions . . . . .	4
<b>2</b>	<b>Explicit Universal Taylor and Radial Kernels on a Compact Metric Space</b>	<b>5</b>
2.1	Injection of $\mathcal{X}$ into $\ell^2$ . . . . .	5
2.2	Universal Kernels on $\mathcal{X}$ . . . . .	7
<b>3</b>	<b>Approximate Universal Kernels</b>	<b>8</b>
3.1	Constructing a Smaller RKHS $\hat{H}$ . . . . .	8
3.2	Showing that $\hat{H}$ is Approximately Universal . . . . .	13
3.3	An Approximate Universal Truncated Kernel . . . . .	16

# 1 Introduction

## 1.1 Kernels in Practice and Related Works

Kernels Methods at large are a ubiquitous tool in statistics, starting with Kernel Density Estimation [Ros56; Par62] and Kernel Regression [Nad64; Wat64]. For an overview of the use of kernel methods in statistics and probability, we refer to the monograph [BT11]. In Machine Learning, the first uses of kernels hinged on the “kernel trick” [Aiz64; SSM98], which allows high expressivity of models without the need of an explicit feature map into the underlying infinite-dimensional space. A cornerstone model is the Support Vector Machine [CV95], whose statistical properties have garnered extensive attention, see for example the monograph [CS08]. A useful tool is the Kernel Mean embedding (we refer to the review [Mua+17]) which maps a measure  $\mu$  to a point  $M(\mu)$  in a Hilbert space of features, and can be used to compare measures with the Maximum Mean Discrepancy defined as  $\text{MMD}(\mu, \nu) = \|M(\mu) - M(\nu)\|_H$  which fostered numerous applications [Góm+09; Zha+11; Mua+12; FSG13; Gre+12; Dor+14; Li+17]. Theoretical guarantees for the MMD depending on properties of the kernel have been reviewed in [SFL11].

From a theoretical standpoint, Reproducing Kernel Hilbert Spaces (RKHS) introduced by Aronszajn [Aro50] have been the object of several monographs [SS02; CS08; SS16]. Some questions remain open, in particular constructing suitable kernels on non-euclidean metric spaces is a challenging problem that is the subject of ongoing research. For compact metric spaces, [CS10] show the existence of universal kernels (i.e. such that the associated RKHS is dense in the space of continuous functions) when the space is continuously embedded into a separable Hilbert space, and [SZ21] relate the notions of universality and strictly proper kernel scores. On complete Riemannian manifolds, [Jay+15] observe (Theorem 6.2) that the natural Gaussian kernel  $k(x, y) = \exp(-s d(x, y)^2)$  is indeed a kernel only in the very restricted case where the manifold is isometric to  $\mathbb{R}^d$ . On Hilbert and Banach spaces, [ZGD24] introduce radial kernels and show universality-adjacent properties. Regarding universality, [MXZ06] study conditions on the feature maps that ensure universality.

Our contribution first consists in an *explicit* construction of universal kernels on a compact metric space  $(\mathcal{X}, d_{\mathcal{X}})$ , in some sense extending [CS10] whose construction is not explicit and relied on the existence of an embedding. The constructed kernels use known kernels known as *Taylor* and *radial* kernels, which are defined on compact subsets of separable Hilbert spaces. Noticing that our kernels are not tractable in practice, we introduce a notion of *approximate universality* and construct other explicit kernels that are approximately universal and tractable.

## 1.2 Elements of RKHS Theory

For a set  $\mathcal{X}$ , a *kernel*  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is a *positive-definite symmetric* function, which is to say a function that verifies  $k(x, y) = k(y, x)$  and:

$$\forall n \in \mathbb{N}^*, \forall (x_1, \dots, x_n) \in \mathcal{X}^n, \forall a \in \mathbb{R}^n, \sum_{i=1}^n \sum_{j=1}^n a_i k(x_i, x_j) a_j \geq 0.$$

By the Moore-Aronszajn theorem [Aro50], there exists a unique Hilbert space  $(H, \langle \cdot, \cdot \rangle_H)$  of functions  $\mathcal{X} \rightarrow \mathbb{R}$ , such that  $H$  contains all basic functions  $k(\cdot, x)$ , and its inner product is characterised by the “reproducing property”  $\langle k(\cdot, x), k(\cdot, y) \rangle_H = k(x, y)$ . Denoting by  $\overline{\text{Span}}$  the Hilbertian completion of the linear span of a set, it follows that  $H = \overline{\text{Span}}\{k(\cdot, x), x \in \mathcal{X}\}$ . The space  $H$  is referred to as the Reproducing Kernel Hilbert Space (RKHS) associated to

the kernel  $k$ . The reproducing property of the kernel implies that for any  $h \in H$  and  $x \in \mathcal{X}$ , we have  $\langle h, k(\cdot, x) \rangle_H = h(x)$ .

If  $k$  is continuous (w.r.t. a metric on  $\mathcal{X}$ ), the RKHS  $H$  is contained in the space of continuous functions from  $\mathcal{X}$  to  $\mathbb{R}$ , denoted  $\mathcal{C}(\mathcal{X})$ . In this work, we will always consider continuous kernels. Some continuous kernels have an additional property called *universality*:

**Definition 1.** A continuous kernel  $k$  on a compact metric space  $(\mathcal{X}, d_{\mathcal{X}})$  is said to be *universal* if the RKHS  $H$  is dense in  $(\mathcal{C}(\mathcal{X}), \|\cdot\|_{\infty})$ , the space of continuous functions from  $\mathcal{X}$  to  $\mathbb{R}$  equipped with the supremum norm. In other words, for any  $\varepsilon > 0$  and  $f \in \mathcal{C}(\mathcal{X})$ , there exists  $h \in H$  such that  $\|f - h\|_{\infty} \leq \varepsilon$ .

Another equivalent definition of kernels uses the notion of feature map / feature space pairs: through these lens, a kernel is any map  $\mathcal{X}^2 \rightarrow \mathbb{R}$  such that there exists a Hilbert space  $H_0$  and a map  $\Phi_0 : \mathcal{X} \rightarrow H_0$  such that Eq. (1) holds.

$$\forall x, y \in \mathcal{X}, k(x, y) = \langle \Phi_0(x), \Phi_0(y) \rangle_{H_0}. \quad (1)$$

The pair  $(\Phi_0, H_0)$  is called a *feature map / feature space pair* (or simply *feature pair*) for  $k$ , and any kernel can be written in this form ([CS08], Theorem 4.16). The associated RKHS is then defined as:

$$H = \{x \mapsto \langle h_0, \Phi_0(x) \rangle_{H_0}, h_0 \in H_0\}. \quad (2)$$

The RKHS  $H$  in Eq. (2) is unique ([CS08], Theorem 4.21), and equal to  $\overline{\text{Span}} \{k(\cdot, x), x \in \mathcal{X}\}$  as stated above. The *canonical feature map* is defined as  $\Phi(x) = k(\cdot, x)$ , and the pair  $(\Phi, H)$  is called the *canonical feature pair* for  $k$ .

From the space viewpoint, an RKHS can equivalently be defined as a Hilbert space of functions  $\mathcal{X} \rightarrow \mathbb{R}$  in which the evaluation  $\delta_x : h \mapsto h(x)$  is continuous for all  $x \in \mathcal{X}$ , as is done in [CS08], Section 4.2. The kernel is then defined as  $k(x, y) = \langle L\delta_x, L\delta_y \rangle_H$ , where  $L\delta_x \in H$  is the Riesz representation of  $\delta_x \in H'$ . In this paper, we stick to the (equivalent) kernel viewpoint.

For a compact metric space  $(E, d_E)$ , we will denote by  $\text{diam}(E)$  its diameter, which is defined by  $\text{diam}(E) := \max_{(x,y) \in E^2} d_E(x, y)$ . Throughout this work,  $\mathcal{X}$  will be assumed to be a compact metric space, and we denote  $D_{\mathcal{X}} := \text{diam}(\mathcal{X})$ .

The first type of universal kernels of interest in this work are Taylor kernels (see [CS08] Lemma 4.8 and Corollary 4.57 for their study on compact subsets of  $\mathbb{R}^d$ ).

**Definition 2.** Let  $W \subset \ell^2$  be a non-empty compact set and  $D_W^2 := \text{diam}(W)^2 > 0$  the square of its diameter. Take a sequence  $(a_n)_{n \in \mathbb{N}} \in (0, +\infty)^{\mathbb{N}}$  such that  $K(t) := \sum_n a_n t^n$  converges absolutely on  $[-D_W^2, D_W^2]$ . The Taylor kernel associated to  $K$  is the map

$$k_W := \begin{cases} W^2 & \longrightarrow \mathbb{R} \\ (u, v) & \longmapsto K(\langle u, v \rangle_{\ell^2}) \end{cases} \quad (3)$$

Taylor kernels are shown to be universal on compact subsets of  $\ell^2$  in [CS10] Theorem 2.1. The second type of universal kernels we will consider are radial kernels<sup>1</sup>.

<sup>1</sup>Radial kernels can be defined (and shown to be universal) on separable Hilbert spaces and more [ZGD24], but we will use compactness for other reasons, and thus restrict to compact subsets of  $\ell^2$  for our purposes.

**Definition 3.** Let  $W \subset \ell^2$  be a non-empty compact set and  $\mu \in \mathcal{M}([0, +\infty))$  a finite Borel measure on  $[0, +\infty)$  with  $\text{supp}(\mu) \neq \{0\}$ . The associated radial function  $K$  and the radial kernel  $k_W$  are defined as follows:

$$K := \left\{ \begin{array}{ccc} \mathbb{R}_+ & \longrightarrow & \mathbb{R} \\ t & \longmapsto & \int_0^{+\infty} e^{-st} d\mu(s) \end{array} \right., \quad k_W := \left\{ \begin{array}{ccc} W^2 & \longrightarrow & \mathbb{R} \\ (u, v) & \longmapsto & K(\|u - v\|_{\ell^2}^2) \end{array} \right. . \quad (4)$$

The universality of radial kernels on  $W$  is a consequence of [ZGD24] Proposition 5.2 combined with [SZ21] Theorem 3.13. Note that the well-known Gaussian (or RBF) kernel  $\exp(-\|\cdot - \cdot\|_{\ell^2}^2 / (2\sigma^2))$  is a particular radial kernel with  $\mu := \delta_{1/(2\sigma^2)}$ .

### 1.3 Paper Outline and Contributions

The objective of this paper is to construct kernels  $k$  on a compact metric space  $(\mathcal{X}, d_{\mathcal{X}})$  that are *universal* (see Definition 1). We also introduce a notion of approximate universality (Definition 7), and introduce other (tractable) explicit kernels  $\hat{k}$  and  $k_t$  that verify this property.

**Construction of universal kernels in Section 2** To construct universal kernels on  $\mathcal{X}$ , we first introduce an explicit continuous injection  $\varphi : \mathcal{X} \rightarrow \ell^2$  in Proposition 5. Given any universal kernel  $k_V$  on  $V := \varphi(\mathcal{X}) \subset \ell^2$  we show in Theorem 6 that  $k(x, y) := k_V(\varphi(x), \varphi(y))$ , is universal on  $\mathcal{X}$ .

The construction of  $\varphi$  in Section 2 is based on a countable basis of  $\mathcal{X}$ , and the associated kernel requires inner products in  $\ell^2$ . In Section 3, we explain how we can use instead a (finite)  $\eta$ -covering of  $\mathcal{X}$ , yielding a finite-dimensional approximation of the embedding  $\varphi$ , with theoretical guarantees. We also investigate the natural idea of truncating the sequence  $\varphi(x)$ .

**Approximate universal kernels in Section 3** We introduce a notion of *approximate universal kernels* on  $\mathcal{X}$ , which are kernels  $\hat{k}$  of RKHS  $\hat{H}$  such that for all  $\varepsilon > 0$  and  $f \in \mathcal{C}(\mathcal{X})$ , there exists  $\hat{h} \in \hat{H}$  such that  $\|f - \hat{h}\|_{\infty} \leq \varepsilon + \rho(f)$ , where  $\rho(f) > 0$  is an error term depending on  $\hat{k}$  and  $f$ . We construct a simpler map  $\hat{\varphi} : \mathcal{X} \rightarrow \mathbb{R}^J$  as a surrogate for the embedding  $\varphi : \mathcal{X} \rightarrow \ell^2$ , and embed  $\mathbb{R}^J$  into  $\ell^2$  appropriately to compare  $\varphi$  and  $B \circ \hat{\varphi}$ . This allows us to introduce the kernel  $\hat{k}(x, y) := k_W(B \circ \hat{\varphi}(x), B \circ \hat{\varphi}(y))$  for a compact set  $W \supset \varphi(\mathcal{X}) \cup B(\hat{\varphi}(\mathcal{X}))$  and a Taylor or radial kernel  $k_W$  on  $W$ . In Corollary 11, we provide a tractable (as in numerically computable) expression for  $\hat{k}$ . Finally, we show in Theorem 15 that  $\hat{k}$  is an approximate universal kernel on  $\mathcal{X}$  with an explicit error term  $\rho$  depending on discretisation parameters and the “complexity” of the function  $f$ . In Section 3.3, we introduce a simple truncation of  $\varphi$  which leads to another approximate universal kernel  $k_t$  on  $\mathcal{X}$ .

## 2 Explicit Universal Taylor and Radial Kernels on a Compact Metric Space

As a preliminary to our main constructions, we begin with two elementary general properties of RKHS which will be useful throughout this section. We remind that a *homeomorphism* is a continuous bijection with a continuous inverse.

**Lemma 4.** i) Let  $\mathcal{X}, \mathcal{Y}$  be two sets and  $k_{\mathcal{Y}} : \mathcal{Y}^2 \rightarrow \mathbb{R}$  a kernel on  $\mathcal{Y}$ , and  $\varphi : \mathcal{X} \rightarrow \mathcal{Y}$  a function. The map  $k_{\mathcal{X}}$  defined in Eq. (5) is a kernel on  $\mathcal{X}$ :

$$k_{\mathcal{X}} := \begin{cases} \mathcal{X}^2 & \longrightarrow \mathbb{R} \\ (x, x') & \longmapsto k_{\mathcal{Y}}(\varphi(x), \varphi(x')) \end{cases} . \quad (5)$$

ii) If additionally  $\mathcal{X}$  and  $\mathcal{Y}$  are compact metric spaces,  $\varphi$  is a homeomorphism and  $k_{\mathcal{Y}} : \mathcal{Y}^2 \rightarrow \mathbb{R}$  is universal, the kernel  $k_{\mathcal{X}}$  is universal.

*Proof. Proof of i):* We verify immediately using the definition that  $k_{\mathcal{X}}$  is a positive-definite symmetric function on  $\mathcal{X}$  using the fact that  $k_{\mathcal{Y}}$  is such on  $\mathcal{Y}$ .

**Proof of ii):** Let  $H_{\mathcal{Y}}$  be the unique RKHS associated to the kernel  $k_{\mathcal{Y}}$  on  $\mathcal{Y}$ , and  $\Phi_{\mathcal{Y}} : \mathcal{Y} \rightarrow H_{\mathcal{Y}}$  its canonical feature map (i.e.  $\forall y \in \mathcal{Y}, \Phi_{\mathcal{Y}}(y) = k_{\mathcal{Y}}(\cdot, y)$ ). Since  $k_{\mathcal{X}}(x, x') = \langle \Phi_{\mathcal{Y}} \circ \varphi(x), \Phi_{\mathcal{Y}} \circ \varphi(x') \rangle_{H_{\mathcal{Y}}}$ , the map  $\Phi_{\mathcal{Y}} \circ \varphi$  and the space  $H_{\mathcal{Y}}$  are a feature pair for  $k_{\mathcal{X}}$ . In the following, given a set  $\mathcal{F}$  of functions and  $g$  a function, we write  $\mathcal{F} \circ g := \{f \circ g, f \in \mathcal{F}\}$ . By uniqueness ([CS08] Theorem 4.21), it follows that the RKHS  $H_{\mathcal{X}}$  associated to  $k_{\mathcal{X}}$  can be written

$$H_{\mathcal{X}} = \{x \mapsto \langle h_{\mathcal{Y}}, \Phi_{\mathcal{Y}} \circ \varphi(x) \rangle_{H_{\mathcal{Y}}}, h_{\mathcal{Y}} \in H_{\mathcal{Y}}\} = H_{\mathcal{Y}} \circ \varphi,$$

where the second equality comes from the reproducing property: for any  $x \in \mathcal{X}$  and  $h_{\mathcal{Y}} \in H_{\mathcal{Y}}$ , we have  $h_{\mathcal{Y}} \circ \varphi(x) = \langle h_{\mathcal{Y}}, \Phi_{\mathcal{Y}} \circ \varphi(x) \rangle_{H_{\mathcal{Y}}}$ . Since  $\varphi$  is a homeomorphism, we also have  $\mathcal{C}(\mathcal{X}) = \mathcal{C}(\mathcal{Y}) \circ \varphi$ . Now for  $\varepsilon > 0$  and  $f_{\mathcal{X}} \in \mathcal{C}(\mathcal{X})$ , take  $f_{\mathcal{Y}} := f_{\mathcal{X}} \circ \varphi^{-1} \in \mathcal{C}(\mathcal{Y})$ . By universality of  $k_{\mathcal{Y}}$ , there exists  $h_{\mathcal{Y}} \in H_{\mathcal{Y}}$  such that  $\|f_{\mathcal{Y}} - h_{\mathcal{Y}}\|_{\infty} \leq \varepsilon$ . Taking  $h_{\mathcal{X}} := h_{\mathcal{Y}} \circ \varphi \in H_{\mathcal{X}}$  yields  $\|f_{\mathcal{X}} - h_{\mathcal{X}}\|_{\infty} \leq \varepsilon$ , and as a result  $k_{\mathcal{X}}$  is universal.  $\square$

Lemma 4 is useful for the construction of universal kernels on compact metric spaces  $\mathcal{X}$  using universal kernels on another space. In the following, we will consider a space  $\mathcal{Y}$  which is a compact subspace of the Hilbert space  $\ell^2$  of square-summable sequences.

### 2.1 Injection of $\mathcal{X}$ into $\ell^2$

Let  $(\mathcal{X}, d_{\mathcal{X}})$  be a non-empty compact metric space, and let  $D_{\mathcal{X}} > 0$  be its diameter. We take a *basis* of  $\mathcal{X}$ , i.e. a countable sequence  $(x_n)_{n \in \mathbb{N}}$  such that for any  $x \in \mathcal{X}$  and  $\varepsilon > 0$ , there exists  $n \in \mathbb{N}$  such that  $d_{\mathcal{X}}(x, x_n) \leq \varepsilon$ . Using a basis, we construct an implicit continuous injection  $\varphi$  from  $\mathcal{X}$  into  $\ell^2$  (Proposition 5), then use universal kernels on  $V := \varphi(\mathcal{X})$  to build a universal kernel  $k$  on  $\mathcal{X}$  in Theorem 6. In Fig. 1, we illustrate the injection.

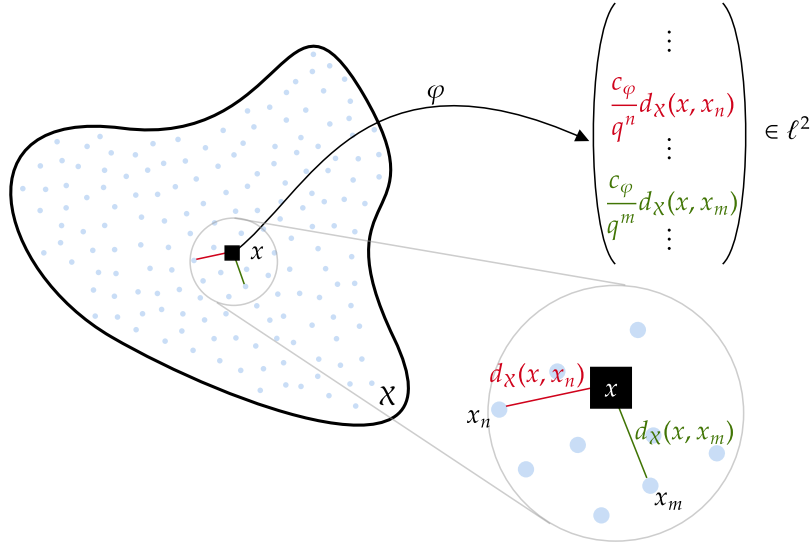


Figure 1: Given a basis  $(x_n)_{n \in \mathbb{N}}$  of  $\mathcal{X}$ , the mapping  $\varphi : \mathcal{X} \rightarrow \ell^2$  maps a point  $x \in \mathcal{X}$  to the sequence of its distances to the points of the basis.

**Proposition 5.** Let  $(x_n)$  a basis of  $\mathcal{X}$  and  $q > 1$ . The map

$$\varphi := \begin{cases} \mathcal{X} & \rightarrow \ell^2 \\ x & \mapsto \left( \frac{c_\varphi d_{\mathcal{X}}(x, x_n)}{q^n} \right)_{n \in \mathbb{N}} \end{cases}, \quad c_\varphi := \frac{\sqrt{q^2 - 1}}{q}$$

is 1-Lipschitz and injective.

*Proof.* The fact that  $\varphi(\mathcal{X}) \subset \ell^2$  comes from the compactness of  $\mathcal{X}$ . Take now  $x, y \in \mathcal{X}$ :

$$\|\varphi(x) - \varphi(y)\|_{\ell^2}^2 = c_\varphi^2 \sum_{n=0}^{+\infty} \frac{|d_{\mathcal{X}}(x, x_n) - d_{\mathcal{X}}(y, x_n)|^2}{q^{2n}} \leq c_\varphi^2 \sum_{n=0}^{+\infty} \frac{d_{\mathcal{X}}(x, y)^2}{q^{2n}} = \frac{c_\varphi^2 q^2}{q^2 - 1} d_{\mathcal{X}}(x, y)^2,$$

showing 1-Lipschitzness. As for injectivity, consider  $x \neq y \in \mathcal{X}^2$  and  $\varepsilon := d_{\mathcal{X}}(x, y)/3 > 0$ . Since  $(x_n)$  is a basis of  $\mathcal{X}$ , there exists  $n \in \mathbb{N}$  such that  $d_{\mathcal{X}}(x, x_n) \leq \varepsilon$ . The triangle inequality then shows

$$d_{\mathcal{X}}(y, x_n) \geq \underbrace{d_{\mathcal{X}}(y, x)}_{=3\varepsilon} - \underbrace{d_{\mathcal{X}}(x, x_n)}_{\in [0, \varepsilon]} \geq 2\varepsilon,$$

and thus  $\underbrace{|d_{\mathcal{X}}(y, x_n) - d_{\mathcal{X}}(x, x_n)|}_{\geq 2\varepsilon} \geq \varepsilon$ , allowing us to conclude

$$\|\varphi(y) - \varphi(x)\|_{\ell^2}^2 \geq c_\varphi^2 \frac{|d_{\mathcal{X}}(y, x_n) - d_{\mathcal{X}}(x, x_n)|^2}{q^{2n}} \geq c_\varphi^2 \frac{\varepsilon^2}{q^{2n}} > 0.$$

□

## 2.2 Universal Kernels on $\mathcal{X}$

We can now build universal kernels  $k : \mathcal{X}^2 \rightarrow \mathbb{R}$  using  $\varphi$  and a universal kernel  $k_W : W^2 \rightarrow \mathbb{R}$  (for example Taylor or radial) on  $W$  a compact subset of  $\ell^2$  containing  $V := \varphi(\mathcal{X})$ . The technique follows closely that of [CS10] Theorem 2.2. Note that thanks to the 1-Lipschitzness of  $\varphi$ , we have  $\text{diam}(\varphi(\mathcal{X})) \leq \text{diam}(\mathcal{X}) =: D_{\mathcal{X}}$ .

**Theorem 6.** Let  $V := \varphi(\mathcal{X}) \subset \ell^2$  and  $W$  be a compact subset of  $\ell^2$  containing  $V$ . Consider  $k_W : W^2 \rightarrow \mathbb{R}$  a universal kernel on  $W$  (e.g. Taylor as in Definition 2 or radial as in Definition 3). The kernel

$$k := \begin{cases} \mathcal{X}^2 & \longrightarrow \mathbb{R} \\ (x, y) & \longmapsto k_W(\varphi(x), \varphi(y)) \end{cases}$$

is universal on  $\mathcal{X}$ .

*Proof.* We introduce  $k_V : V^2 \rightarrow \mathbb{R}$  the restriction of  $k_W$  to  $V^2$ . By [CS08] Lemma 4.55 item iii),  $k_V$  remains universal. Since  $\mathcal{X}$  is a compact metric space and  $\ell^2$  is Hausdorff, the co-restriction of  $\varphi$  to  $V$  denoted  $\varphi_V : \mathcal{X} \rightarrow V$  is a homeomorphism. Noticing that  $k = (x, y) \in \mathcal{X}^2 \mapsto k_V(\varphi_V(x), \varphi_V(y))$ , by Lemma 4,  $k$  is thus a universal kernel on  $\mathcal{X}$ .  $\square$

**A strictly convex functional on  $\mathcal{P}(\mathcal{X})$**  We consider the set  $\mathcal{P}(\mathcal{X})$  of probability measures on  $\mathcal{X}$ . As a universal kernel,  $k$  is also *characteristic* (see [SFL11] and use the compactness of  $\mathcal{X}$ ), which is to say that the map

$$M := \begin{cases} \mathcal{P}(\mathcal{X}) & \longrightarrow H \\ \mu & \longmapsto \int_{\mathcal{X}} k(\cdot, x) d\mu(x) \end{cases},$$

known as the *kernel mean embedding* [Sri+10], is injective. One can show that the map

$$F := \begin{cases} \mathcal{P}(\mathcal{X}) & \longrightarrow \mathbb{R}_+ \\ \mu & \longmapsto \|M(\mu)\|_H^2 \end{cases}$$

is continuous with respect to the weak convergence of measures (apply [HC11] Theorem A.1 using that  $x \mapsto k(\cdot, x)$  is continuous and bounded). Furthermore, by linearity of  $M$  and strict convexity of  $\|\cdot\|_H^2$ , the function  $F$  is strictly convex with respect to the “vertical” convex combination of probability measures:

$$\forall \mu, \nu \in \mathcal{P}(\mathcal{X}), \forall t \in (0, 1), F((1-t)\mu + t\nu) < (1-t)F(\mu) + tF(\nu).$$

Note that the fact that  $M$  is injective is required to prove *strict* convexity.



### 3 Approximate Universal Kernels

In practice, the function  $\varphi$  introduced in [Proposition 5](#) is not tractable (in the sense that computing and storing a full sequence  $\varphi(x) \in \ell^2$  is numerically impossible), limiting the use of the kernels proposed in [Theorem 6](#) in their exact formulation. The natural idea of truncating the sequence  $\varphi(x)$  to the first  $N$  terms is tackled later in [Section 3.3](#), we begin with the main approach of the paper, which relies on a discretisation of the space  $\mathcal{X}$ .

We will now introduce a family of tractable kernels which are approximately universal on  $\mathcal{X}$ . Throughout this section, the kernels  $k_W$  on a compact subset  $W$  of  $\ell^2$  that we will consider are Taylor or radial (see [Definitions 2](#) and [3](#)). Our objective is to construct another kernel  $\hat{k}$  with a simpler explicit mapping  $\hat{\varphi} : \mathcal{X} \rightarrow \mathbb{R}^J$ , yielding an RKHS  $\hat{H}$  which we will show to be approximately universal in the sense of [Definition 7](#).

**Definition 7.** Let  $\hat{k} : \mathcal{X}^2 \rightarrow \mathbb{R}$  a kernel on  $\mathcal{X}$  of RKHS  $\hat{H}$  and  $\rho : \mathcal{C}(\mathcal{X}) \rightarrow \mathbb{R}_+$  an error function. We say that  $\hat{k}$  is an *approximate universal kernel* on  $\mathcal{X}$  if for all  $\varepsilon > 0$  and  $f \in \mathcal{C}(\mathcal{X})$ , there exists  $\hat{h} \in \hat{H}$  such that  $\|f - \hat{h}\|_\infty \leq \varepsilon + \rho(f)$ .

#### 3.1 Constructing a Smaller RKHS $\hat{H}$

In this section, we provide a principled method to “sub-sample” the sequence  $\varphi(x)$ : we will begin with a well-chosen finite family  $(y_j) \in \mathcal{X}^J$  and construct a well-suited basis  $(x_n) \in \mathcal{X}^\mathbb{N}$  such that a distance sequence  $(d_{\mathcal{X}}(x, x_n))_{n \in \mathbb{N}}$  is adequately approximable by the finite number of distances  $(d_{\mathcal{X}}(x, y_j))_{j \in \llbracket 1, J \rrbracket}$ . We illustrate this discretisation concept in [Fig. 2](#).

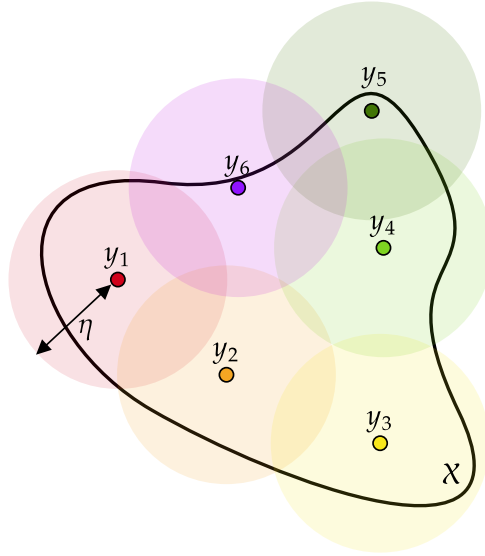


Figure 2: Discretisation of the space  $\mathcal{X}$  into a cover of  $J$  balls of radius  $\eta > 0$  centred at each  $(y_j)_{j \in \llbracket 1, J \rrbracket}$ .

Instead of a basis of  $\mathcal{X}$ , we will now fix  $\eta \in (0, D_{\mathcal{X}}]$  and consider  $(y_j)_{j \in \llbracket 1, J \rrbracket}$  a family of distinct points of  $\mathcal{X}$  such that the family of (closed) balls  $B_{d_{\mathcal{X}}}(y_j, \eta)$  covers  $\mathcal{X}$ . In [Eq. \(6\)](#), we introduce a map  $\hat{\varphi} : \mathcal{X} \rightarrow \mathbb{R}^J$  in the spirit of  $\varphi$  defined in [Proposition 5](#), which we visualise in [Fig. 3](#).

$$\hat{\varphi} := \begin{cases} \mathcal{X} & \rightarrow \mathbb{R}^J \\ x & \mapsto \left( \frac{d_{\mathcal{X}}(x, y_j)}{\sqrt{J}} \right)_{j \in \llbracket 1, J \rrbracket} \end{cases} . \quad (6)$$



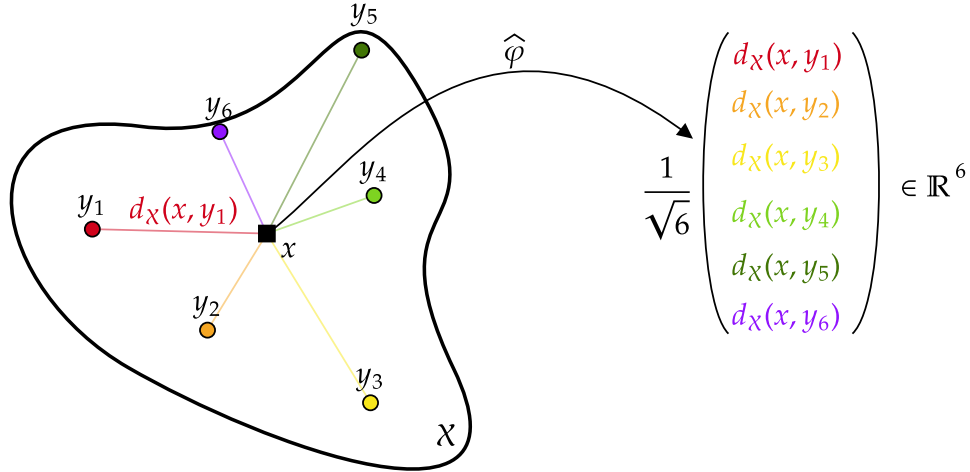


Figure 3: The mapping  $\hat{\varphi} : \mathcal{X} \rightarrow \mathbb{R}^J$  maps a point  $x \in \mathcal{X}$  to the vector of normalised distances between  $x$  and the centres  $y_j$  of the covering.

It is immediate to verify that  $\hat{\varphi} : (\mathcal{X}, d_{\mathcal{X}}) \rightarrow (\mathbb{R}^J, \|\cdot\|_2)$  is 1-Lipschitz, thanks to the  $J^{-1/2}$  normalisation. In [Proposition 10](#), we show how to embed  $\mathbb{R}^J \supset \hat{\varphi}(\mathcal{X})$  into  $\ell^2$  with a mapping  $B$ , which will allow us to compare the RKHS induced by  $\hat{\varphi}$  and a particular  $\varphi : \mathcal{X} \rightarrow \ell^2$ . To construct  $B$ , we first begin with a geometric series separation lemma, which will be convenient to deal with the factor  $\frac{1}{q^n}$  in  $\varphi$ .

**Lemma 8.** Let  $J \geq 2$ ,  $q \in (1, 1 + \frac{1}{J-1})$  and coefficients  $(\lambda_1, \dots, \lambda_J) \in (0, 1)^J$  such that  $\sum_j \lambda_j = 1$ , there exists  $\alpha : \mathbb{N} \rightarrow \llbracket 1, J \rrbracket$  with for all  $j \in \llbracket 1, J \rrbracket$ ,  $\alpha^{-1}(\{j\})$  infinite such that:

$$\forall j \in \llbracket 1, J \rrbracket, \quad \sum_{n \in \alpha^{-1}(\{j\})} \frac{1}{q^n} = \lambda_j \frac{q}{q-1}. \quad (7)$$

*Proof.* Set  $S := \frac{q}{q-1}$ . We will construct a sequence  $(\alpha(N))_{N \in \mathbb{N}}$  by induction over  $N$ , verifying the property

$$P_N : \forall j \in \llbracket 1, J \rrbracket, \quad \sum_{n \in \llbracket 0, N \rrbracket : \alpha(n)=j} \frac{1}{q^n} < \lambda_j S.$$

*Initialisation:* set  $\alpha(0)$  the first  $j \in \llbracket 1, J \rrbracket$  such that  $1 < \lambda_j S$ . Note that such a  $j$  exists, otherwise summing over  $j \in \llbracket 1, J \rrbracket$  yields

$$J \geq \frac{q}{q-1} > \frac{1 + \frac{1}{J-1}}{1 + \frac{1}{J-1} - 1} = J,$$

which is a contradiction. We have defined  $\alpha(0) := j$  verifying  $P_0$ .

*Induction step:* let  $N \in \mathbb{N}$ , suppose  $P_N$  true. We show that there exists  $j \in \llbracket 1, J \rrbracket$  such that

$$\sum_{n \in \llbracket 0, N \rrbracket : \alpha(n)=j} \frac{1}{q^n} + \frac{1}{q^{N+1}} < \lambda_j S \quad (8)$$

by contradiction. If that were not the case, we would have by summing Eq. (8) over  $j \in \llbracket 1, J \rrbracket$ :

$$\sum_{n=0}^N \frac{1}{q^n} + \frac{J}{q^{N+1}} \geq S,$$

which by computation is equivalent to  $q \geq 1 + \frac{1}{J-1}$ , obtaining a contradiction. Selecting  $j \in \llbracket 1, J \rrbracket$  such that Eq. (8) holds, we can set  $\alpha(N+1) := j$  which satisfies  $P_{N+1}$ .

Now that  $\alpha : \mathbb{N} \rightarrow \llbracket 1, J \rrbracket$  verifying  $(P_N)$  is constructed, we introduce the convergent series

$$\forall j \in \llbracket 1, J \rrbracket, \forall N \in \mathbb{N}, S_N^{(j)} := \sum_{n \in \llbracket 0, N \rrbracket : \alpha(n)=j} \frac{1}{q^n}, S_\infty^{(j)} := \lim_{N \rightarrow +\infty} S_N^{(j)}.$$

Thanks to  $(P_N)$ , for all  $j \in \llbracket 1, J \rrbracket$  taking the limit yields  $S_\infty^{(j)} \leq \lambda_j S$ , and summing over  $j \in \llbracket 1, J \rrbracket$  gives  $\sum_j S_\infty^{(j)} = S$ , hence necessarily for all  $j \in \llbracket 1, J \rrbracket$ ,  $S_\infty^{(j)} = \lambda_j S$ .

Finally, observing the strict inequality in  $P_N$  at each  $N \in \mathbb{N}$  shows that  $\alpha^{-1}(\{j\})$  has to be infinite, concluding the proof.  $\square$

We now turn to constructing an embedding  $B : \mathbb{R}^J \rightarrow \ell^2$ , which will allow us to compare  $\hat{\varphi}$  and  $\varphi$ . An important property of  $B$  will be the correspondence between the inner products in  $\mathbb{R}^J$  and  $\ell^2$  (i.e.  $B$  will be an isometry). The construction of this embedding revolves around the construction of an adapted basis  $(x_n)_{n \in \mathbb{N}}$  of  $\mathcal{X}$  which is balanced with respect to the covering by the balls  $B(y_j, \eta)$ , as illustrated in Fig. 4.

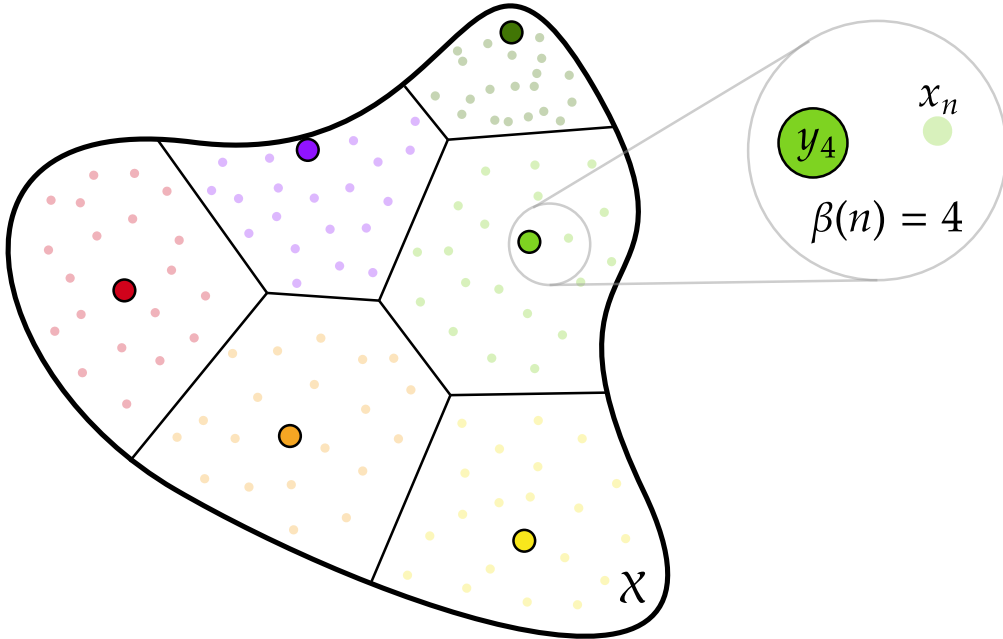


Figure 4: The basis  $(x_n)_{n \in \mathbb{N}}$  is such that there are equally as many  $(x_n)$  in each region  $\mathcal{X}_j$  of points closest to  $y_j$ . In the figure, we observe a zoom on the region  $\mathcal{X}_4$ , where the example point  $x_n$  is closest to  $y_4$ . In mathematical terms, we write this property as  $\beta(n) = y_j$ , and in Proposition 9 we will construct  $(x_n)$  such that the sum  $\sum_n q^{-2n}$  is split evenly between the sets  $\beta^{-1}(\{j\})$ .

**Proposition 9.** Take  $q \in (1, \sqrt{1 + \frac{1}{J-1}})$ . There exists a basis  $(x_n)_{n \in \mathbb{N}}$  of  $\mathcal{X}$  and a map-

ping  $\beta : \mathbb{N} \longrightarrow \llbracket 1, J \rrbracket$  with infinite pre-images which verifies  $\forall n \in \mathbb{N}, d_{\mathcal{X}}(x_n, y_{\beta(n)}) = \min_j d_{\mathcal{X}}(x_n, y_j)$ , and with the following property:

$$\forall j \in \llbracket 1, J \rrbracket, \sum_{n \in \beta^{-1}(\{j\})} \frac{1}{q^{2n}} = \frac{1}{J} \frac{q^2}{q^2 - 1}. \quad (9)$$

*Proof.* Consider for  $j \in \llbracket 1, J \rrbracket$  the set  $\mathcal{X}_j := \{x \in \mathcal{X} : \operatorname{argmin}_m d_{\mathcal{X}}(x, y_m) = j\}$  (with disambiguation by taking the smallest minimiser if multiple exist). By definition, the sets  $\mathcal{X}_j$  are disjoint and cover  $\mathcal{X}$ . Since  $(\mathcal{X}, d_{\mathcal{X}})$  is a compact metric space, each subset  $\mathcal{X}_j$  is separable, allowing us to choose a basis  $(z_n^{(j)})_{n \in \mathbb{N}}$  of  $\mathcal{X}_j$  for each  $j \in \llbracket 1, J \rrbracket$ . By Lemma 8, we can choose  $\beta : \mathbb{N} \longrightarrow \llbracket 1, J \rrbracket$  with infinite pre-images which verifies Eq. (9). Since for each  $j \in \llbracket 1, J \rrbracket$ , the set  $\beta^{-1}(\{j\}) \subset \mathbb{N}$  is infinite, we can choose  $\omega_j : \beta^{-1}(\{j\}) \longrightarrow \mathbb{N}$  a bijection. We can now define  $\forall n \in \mathbb{N}, x_n := z_{\omega_{\beta(n)}(n)}^{(\beta(n))}$ , which is a basis of  $\mathcal{X}$  since  $\cup_j \mathcal{X}_j = \mathcal{X}$  and

$$\{x_n\}_{n \in \mathbb{N}} = \bigcup_j \{z_{\omega_j(m)}^{(j)}\}_{m \in \beta^{-1}(\{j\})} = \bigcup_j \{z_n^{(j)}\}_{n \in \mathbb{N}},$$

by construction. Furthermore, by definition, we have  $\forall n \in \mathbb{N}, \operatorname{argmin}_j d_{\mathcal{X}}(x_n, y_j) = \beta(n)$ , which shows that the mapping  $\beta$  satisfies the desired properties.  $\square$

Using the adapted basis from Proposition 9, we can finally construct an isometry  $B : \mathbb{R}^J \longrightarrow \ell^2$ :

**Proposition 10.** Take a basis  $(x_n)$  of  $\mathcal{X}$  and  $\beta : \mathbb{N} \longrightarrow \llbracket 1, J \rrbracket$  as in Proposition 9. The mapping  $B$  defined below is an isometry:

$$B := \left\{ \begin{array}{ccc} \mathbb{R}^J & \longrightarrow & \ell^2 \\ (u_j)_{j=1}^J & \longmapsto & \left( c_B \frac{u_{\beta(n)}}{q^n} \right)_{n \in \mathbb{N}} \end{array} \right., \quad c_B := \frac{\sqrt{J(q^2 - 1)}}{q}. \quad (10)$$

*Proof.* The mapping  $B$  is clearly linear, and for  $u, v \in \mathbb{R}^J$  we compute using Eq. (9):

$$\langle B(u), B(v) \rangle_{\ell^2} = \sum_{n=0}^{+\infty} \frac{c_B^2}{q^{2n}} u_{\beta(n)} v_{\beta(n)} = c_B^2 \sum_{j=1}^J u_j v_j \sum_{n \in \beta^{-1}(\{j\})} \frac{1}{q^{2n}} = c_B^2 \frac{1}{J} \frac{q^2}{q^2 - 1} \langle u, v \rangle_{\mathbb{R}^J} = \langle u, v \rangle_{\mathbb{R}^J},$$

which shows that  $B$  is an isometry.  $\square$

In the following, we draw a correspondence between a RKHS  $\hat{H}$  built with  $\hat{\varphi}$  from Eq. (6) and another RKHS  $H$  built using  $\varphi$  from Proposition 5. Let  $U := \hat{\varphi}(\mathcal{X})$ , which is a compact subset of  $\mathbb{R}^J$ , then let  $\hat{V} := B(U)$ , it is a compact subset of  $\ell^2$ . Consider the injection  $\varphi$  introduced in Proposition 5 with basis  $(x_n)$  and scale  $q$  as in Proposition 10. Define  $V := \varphi(\mathcal{X})$ ,  $W := V \cup \hat{V}$ , which are also compact subsets of  $\ell^2$ . We now summarise our objects in the following diagram:

$$\begin{array}{ccc} \mathcal{X} & \xrightarrow{\varphi} & V \subset W \subset \ell^2 \\ \downarrow \hat{\varphi} & & \\ U \subset \mathbb{R}^J & \xrightarrow{B} & \hat{V} \subset W \subset \ell^2 \end{array} \quad (11)$$

We fix a kernel  $k_W : W^2 \longrightarrow \mathbb{R}$  which is of Taylor type or radial (see Definitions 2 and 3) and thus in particular universal on  $W$ , and introduce its canonical feature map:

$$\Phi_W := \left\{ \begin{array}{ccc} W & \longrightarrow & H_W \\ u & \longmapsto & k_W(\cdot, u) \end{array} \right., \quad (12)$$

where  $H_W = \overline{\text{Span}} \{k_W(\cdot, u), u \in W\} \subset \mathcal{C}(W)$  is the unique RKHS associated to the kernel  $k_W$  ([CS08] Theorem 4.21). Consider the kernels  $k, \hat{k}$  on  $\mathcal{X}$  defined respectively as:

$$k := \left\{ \begin{array}{ccc} \mathcal{X}^2 & \longrightarrow & \mathbb{R} \\ (x, y) & \longmapsto & k_W(\varphi(x), \varphi(y)) \end{array} \right., \quad \hat{k} := \left\{ \begin{array}{ccc} \mathcal{X}^2 & \longrightarrow & \mathbb{R} \\ (x, y) & \longmapsto & k_W(B \circ \hat{\varphi}(x), B \circ \hat{\varphi}(y)) \end{array} \right. . \quad (13)$$

By definition of the feature pair  $(H_W, \Phi_W)$  for  $k_W$ , we observe that for  $x, y \in \mathcal{X}$ :

$$k(x, y) = \langle \Phi_W \circ \varphi(x), \Phi_W \circ \varphi(y) \rangle_{H_W}, \quad \hat{k}(x, y) = \langle \Phi_W \circ B \circ \hat{\varphi}(x), \Phi_W \circ B \circ \hat{\varphi}(y) \rangle_{H_W}. \quad (14)$$

The RKHS spaces  $H, \hat{H}$  associated to  $k, \hat{k}$  are both subspaces of  $\mathcal{C}(\mathcal{X})$  and can be written with the following respective feature pairs  $(H_W, \Phi), (H_W, \hat{\Phi})$  (use Eq. (14) with [CS08] Theorem 4.21):

$$H = \{x \mapsto \langle h_W, \Phi_W \circ \varphi(x) \rangle_{H_W}, h_W \in H_W\}, \quad \Phi := \left\{ \begin{array}{ccc} \mathcal{X} & \longrightarrow & H_W \\ x & \longmapsto & \Phi_W \circ \varphi(x) \end{array} \right. \quad (15)$$

$$\hat{H} = \{x \mapsto \langle h_W, \Phi_W \circ B \circ \hat{\varphi}(x) \rangle_{H_W}, h_W \in H_W\}, \quad \hat{\Phi} := \left\{ \begin{array}{ccc} \mathcal{X} & \longrightarrow & H_W \\ x & \longmapsto & \Phi_W \circ B \circ \hat{\varphi}(x) \end{array} \right. . \quad (16)$$

Notice that the feature space  $H_W$  is shared. To finish the diagram, we introduce the “feature-to-map” functionals:

$$\Psi := \left\{ \begin{array}{ccc} H_W & \longrightarrow & H \\ h_W & \longmapsto & x \mapsto \langle h_W, \Phi(x) \rangle_{H_W} \end{array} \right., \quad \hat{\Psi} := \left\{ \begin{array}{ccc} H_W & \longrightarrow & \hat{H} \\ h_W & \longmapsto & x \mapsto \langle h_W, \hat{\Phi}(x) \rangle_{H_W} \end{array} \right. . \quad (17)$$

By Eqs. (15) and (16),  $\Psi$  and  $\hat{\Psi}$  are surjective. Extending the diagram in Eq. (11), we obtain:

Using the inner product correspondence induced by the isometry  $B$  from Proposition 10, a tractable formula for  $\hat{k}$  is obtained immediately for Taylor and radial kernels.

**Corollary 11.** The kernel  $\hat{k}$  on  $\mathcal{X}$  is given by, for all  $x, y \in \mathcal{X}$ :

- if  $k_W$  is a Taylor kernel (Definition 2):

$$\hat{k}(x, y) = K(\langle \hat{\varphi}(x), \hat{\varphi}(y) \rangle_{\mathbb{R}^J}) = \sum_{n=0}^{+\infty} a_n \left( \frac{1}{J} \sum_{j=1}^J d_{\mathcal{X}}(x, y_j) d_{\mathcal{X}}(y, y_j) \right)^n ; \quad (19)$$

- if  $k_W$  is a radial kernel (Definition 3):

$$\hat{k}(x, y) = K(\|\hat{\varphi}(x) - \hat{\varphi}(y)\|_{\mathbb{R}^J}^2) = \int_0^{+\infty} \exp\left(-\frac{s}{J} \sum_{j=1}^J (d_{\mathcal{X}}(x, y_j) - d_{\mathcal{X}}(y, y_j))^2\right) d\mu(s). \quad (20)$$

*Proof.* Let  $x, y \in \mathcal{X}$ , we remind that from Eq. (13) that  $\hat{k}(x, y) := k_W(B \circ \hat{\varphi}(x), B \circ \hat{\varphi}(y))$ . Now by Proposition 10,  $B$  is an isometry, yielding:

$$\langle B \circ \hat{\varphi}(x), B \circ \hat{\varphi}(y) \rangle_{\ell^2} = \langle \hat{\varphi}(x), \hat{\varphi}(y) \rangle_{\mathbb{R}^J}; \quad \|B \circ \hat{\varphi}(x) - B \circ \hat{\varphi}(y)\|_{\ell^2}^2 = \|\hat{\varphi}(x) - \hat{\varphi}(y)\|_{\mathbb{R}^J}^2.$$

Eqs. (19) and (20) are then obtained by replacing  $k_W$  and  $K$  by their definitions in the Taylor and radial cases.  $\square$

We refer to the expressions in Eqs. (19) and (20) as “tractable” since they can be computed explicitly on a computer or approximated efficiently to numerical precision (note that the measure  $\mu$  in the radial kernel can be discrete with finite support). The numerical computation of  $\hat{\varphi}$  is explicit and tractable in the sense that it can be done with a finite amount of closed-form expressions. For the Taylor kernel, the infinite series can be approximated to numerical precision, which we also refer to as “tractable”, as would be said of the exponential function for instance.

### 3.2 Showing that $\hat{H}$ is Approximately Universal

In this section, we show that the RKHS  $\hat{H}$  introduced in Section 3.1 is approximately universal on  $\mathcal{X}$ . We use the notation and objects constructed in Section 3.1 extensively, in particular, the mapping  $\varphi : \mathcal{X} \rightarrow \ell^2$  is defined using Proposition 5 with a suitable basis  $(x_n)$  and scale  $q$  from Proposition 9. The first approximation result we will show concerns a comparison in  $\ell^2$  between  $\varphi(x)$  and  $B \circ \hat{\varphi}(x)$ :

**Proposition 12.** For  $x \in \mathcal{X}$ , we have  $\|\varphi(x) - B \circ \hat{\varphi}(x)\|_{\ell^2} \leq \eta$ .  
The diameter of  $W$  verifies  $D_W := \text{diam}(W) \leq 2D_{\mathcal{X}}$ .

*Proof.* Let  $n \in \mathbb{N}$ , we look at the terms of the sequences  $\varphi(x), B \circ \hat{\varphi}(x) \in W \subset \ell^2$ :

$$|[\varphi(x)]_n - [B \circ \hat{\varphi}(x)]_n| = \left| \frac{c_{\varphi} d_{\mathcal{X}}(x, x_n)}{q^n} - \frac{c_B d_{\mathcal{X}}(x, y_{\beta(n)})}{\sqrt{J} q^n} \right| \leq \frac{1}{q^n} c_{\varphi} d_{\mathcal{X}}(x_n, y_{\beta(n)}).$$

By construction of the covering  $(B_{d_{\mathcal{X}}}(y_j, \eta))_j$  and of  $\beta$  (see Proposition 9),  $d_{\mathcal{X}}(x_n, y_{\beta(n)}) \leq \eta$ . Summing the squares over  $n \in \mathbb{N}$  and replacing  $c_{\varphi} = \frac{c_B}{\sqrt{J}}$  yields:

$$\|\varphi(x) - B \circ \hat{\varphi}(x)\|_{\ell^2}^2 \leq \eta^2 \sum_{n=0}^{+\infty} \frac{c_{\varphi}^2}{q^{2n}} = \eta^2.$$

For the diameter of  $W := \varphi(\mathcal{X}) \cup (B \circ \hat{\varphi}(\mathcal{X}))$ , we have by 1-Lipschitzness of  $\varphi, \hat{\varphi}$  and  $B$ :  $\text{diam}(\varphi(\mathcal{X})) \leq D_{\mathcal{X}}$  and  $\text{diam}(B \circ \hat{\varphi}(\mathcal{X})) \leq D_{\mathcal{X}}$ . Using the inequality in the above display and the fact that  $\eta \leq D_{\mathcal{X}}$ , we conclude:

$$D_W = \max(\text{diam}(\varphi(\mathcal{X})), \sup_{x, y \in \mathcal{X}} \|\varphi(x) - B \circ \hat{\varphi}(y)\|_{\ell^2}) \leq \max(D_{\mathcal{X}}, 2\eta) \leq 2D_{\mathcal{X}}. \quad \square$$

Using regularity properties of Taylor and radial kernels, we will show that the kernel  $\hat{k}$  is approximately universal on  $\mathcal{X}$  by relating it to  $k$  which is universal by [Theorem 6](#). First, we see in [Lemma 13](#) that the canonical feature map  $\Phi_W$  is Hölder-continuous for Taylor kernels, and Lipschitz for radial kernels. We introduce the radius of  $W$ :  $R_W := \max_{w \in W} \|w\|_{\ell^2}$ . Using the definition of  $W$  and of  $\varphi, \hat{\varphi}$  and  $B$  with their well-chosen normalisations, it is easy to see that  $R_W \leq D_{\mathcal{X}}$ .

**Lemma 13.** The feature map  $\Phi_W : (W, \|\cdot\|_{\ell^2}) \longrightarrow (H_W, \|\cdot\|_{H_W})$  has the following regularity:

- If  $k_W$  is a Taylor kernel, then  $\Phi_W$  is  $\frac{1}{2}$ -Hölder continuous:

$$\forall u, v \in W, \|\Phi_W(u) - \Phi_W(v)\|_{H_W} \leq \sqrt{2D_{\mathcal{X}}C_{K'}} \|u - v\|_{\ell^2}^{\frac{1}{2}},$$

where  $C_{K'} := \max_{t \in [-4D_{\mathcal{X}}^2, 4D_{\mathcal{X}}^2]} |K'(t)|$ .

- If  $k_W$  is a radial kernel, then  $\Phi_W$  is  $\sqrt{2C_{K'}}$ -Lipschitz:

$$\forall u, v \in W, \|\Phi_W(u) - \Phi_W(v)\|_{H_W} \leq \sqrt{2C_{K'}} \|u - v\|_{\ell^2},$$

where  $C_{K'} := \max_{t \in [0, 4D_{\mathcal{X}}^2]} |K'(t)|$ .

*Proof.* First, we remind that by [Proposition 12](#), we have  $\text{diam}(W) \leq 2D_{\mathcal{X}}$ . For the proof, we take inspiration from [\[Fie23\]](#) Section 4.2. Using the reproducing property, we begin computations for both kernel types, letting  $u, v \in W$ :

$$\begin{aligned} \|\Phi_W(u) - \Phi_W(v)\|_{H_W}^2 &= k_W(u, u) - 2k_W(u, v) + k_W(v, v) \\ &\leq |k_W(u, u) - k_W(u, v)| + |k_W(v, v) - k_W(u, v)|. \end{aligned}$$

For Taylor kernels, we use the fact that  $K$  is  $C_{K'}$ -Lipschitz on  $[-4D_{\mathcal{X}}^2, 4D_{\mathcal{X}}^2]$  and the Cauchy-Schwarz inequality for  $\langle \cdot, \cdot \rangle_{\ell^2}$ :

$$\begin{aligned} \|\Phi_W(u) - \Phi_W(v)\|_{H_W}^2 &\leq C_{K'} (|\langle u, u \rangle_{\ell^2} - \langle u, v \rangle_{\ell^2}| + |\langle v, v \rangle_{\ell^2} - \langle u, v \rangle_{\ell^2}|) \\ &\leq C_{K'} (\|u\|_{\ell^2} + \|v\|_{\ell^2}) \|u - v\|_{\ell^2} \\ &\leq 2R_W C_{K'} \|u - v\|_{\ell^2}, \end{aligned}$$

and we conclude using  $R_W \leq D_{\mathcal{X}}$ . For radial kernels, we use the fact that  $K$  is  $C_{K'}$ -Lipschitz on  $[0, D_W^2]$  (we remind that  $K$  is non-increasing on  $[0, +\infty)$ ):

$$\|\Phi_W(u) - \Phi_W(v)\|_{H_W}^2 = 2(K(0) - K(\|u - v\|_{\ell^2}^2)) \leq 2C_{K'} \|u - v\|_{\ell^2}^2.$$

□

We now use [Lemma 13](#) to approximate any  $h \in H$  with a  $\hat{h} \in \hat{H}$  with a certain error, which we approach by comparing the feature-to-map functionals  $\Psi$  and  $\hat{\Psi}$  from [Eqs. \(15\) and \(16\)](#).

**Proposition 14.** For  $h \in H$ , take  $h_W \in H_W$  such that  $h = x \longmapsto \langle h_W, \Phi(x) \rangle_{H_W} = \Psi(h_W)$ . Then let  $\hat{h} := x \longmapsto \langle h_W, \hat{\Phi}(x) \rangle_{H_W} = \hat{\Psi}(h_W)$ . Denoting  $\|\cdot\|_{\infty}$  the supremum norm on  $\mathcal{X}$ , we have:

$$\|h - \hat{h}\|_{\infty} \leq \rho_0 \|h_W\|_{H_W}, \quad (21)$$

where  $\rho_0 = \eta^{\frac{1}{2}} \sqrt{2D_{\mathcal{X}}C_{K'}}$  for a Taylor kernel and  $\rho_0 = \eta \sqrt{2C_{K'}}$  for a radial kernel..

*Proof.* First, we use the regularity of  $\Phi_W$  from [Lemma 13](#): we have for  $x \in X$ ,

$$\begin{aligned} |h(x) - \hat{h}(x)| &= \langle h_W, \Phi(x) - \hat{\Phi}(x) \rangle_{H_W} \leq \|h_W\|_{H_W} \|\Phi(x) - \hat{\Phi}(x)\|_{H_W} \\ &= \|h_W\|_{H_W} \|\Phi_W \circ \varphi(x) - \Phi_W \circ B \circ \hat{\varphi}(x)\|_{H_W} \\ &\leq \tilde{c} \|h_W\|_{H_W} \|\varphi(x) - B \circ \hat{\varphi}(x)\|_{\ell^2}^s, \end{aligned}$$

where  $(\tilde{c}, s) = (\sqrt{2D_{\mathcal{X}}C_{K'}}, \frac{1}{2})$  for a Taylor kernel and  $(\tilde{c}, s) = (\sqrt{2C_{K'}}, 1)$  for a radial kernel. Combining with [Proposition 12](#), we obtain [Eq. \(21\)](#).  $\square$

Using the universality of the kernel  $k$  (thanks to [Theorem 6](#)), we can frame the result of [Proposition 14](#) as an approximate universality property of  $\hat{k}$ . Again, the approximation error functions depend on the type of kernel  $k_W$ .

**Theorem 15.** Let  $\varepsilon > 0$  and  $f \in \mathcal{C}(\mathcal{X})$ , the element  $h[\varepsilon, f] \in H$  defined by

$$h[\varepsilon, f] := \operatorname{argmin}_{h \in H: \|h-f\|_{\infty} \leq \varepsilon} \|h\|_H^2 \quad (22)$$

is well-defined, and there exists  $\hat{h} \in \hat{H}$  such that:

$$\|f - \hat{h}\|_{\infty} \leq \varepsilon + \rho_0 \|h[\varepsilon, f]\|_H, \quad (23)$$

where  $\rho_0 = \eta^{\frac{1}{2}} \sqrt{2D_{\mathcal{X}}C_{K'}}$  for a Taylor kernel and  $\rho_0 = \eta \sqrt{2C_{K'}}$  for a radial kernel..

*Proof.* First, we introduce:

$$h_W[\varepsilon, f] := \operatorname{argmin}_{h_W \in H_W: \|\Psi(h_W) - f\|_{\infty} \leq \varepsilon} \|h_W\|_{H_W}^2$$

We show that  $h_W[\varepsilon, f]$  and  $h[\varepsilon, f]$  are well-defined. The triangle inequality ensures that the sets  $\mathcal{B}_{H_W} := \{h_W \in H_W : \|\Psi(h_W) - f\|_{\infty} \leq \varepsilon\}$  and  $\mathcal{B}_H := \{h \in H : \|h - f\|_{\infty} \leq \varepsilon\}$  are convex.

We now show the continuity of  $\Psi$  as a mapping  $(H_W, \|\cdot\|_{H_W}) \rightarrow (\mathcal{C}(\mathcal{X}), \|\cdot\|_{\infty})$ . Fixing  $x \in \mathcal{X}$ , we upper-bound by the Cauchy-Schwarz inequality:

$$|\Psi(h_W)[x]| = |\langle h_W, \Phi(x) \rangle_{H_W}| \leq \|h_W\|_{H_W} \|\Phi(x)\|_{H_W}, \quad (24)$$

then we use the definition  $\Phi = \Phi_W \circ \varphi$  to show the continuity of  $\Phi : (\mathcal{X}, d_{\mathcal{X}}) \rightarrow (H_W, \|\cdot\|_{H_W})$ : by [Proposition 5](#),  $\varphi : (\mathcal{X}, d_{\mathcal{X}}) \rightarrow (W, \|\cdot\|_{\ell^2})$  is continuous, and by [Lemma 13](#),  $\Phi_W : (W, \|\cdot\|_{\ell^2}) \rightarrow (H_W, \|\cdot\|_{H_W})$  is also continuous. Combining [Eq. \(24\)](#) with the continuity of  $\Phi$  and the compactness of  $\mathcal{X}$  ensures that there exists  $C > 0$  independent of  $x \in \mathcal{X}$  and  $h_W \in H_W$  such that  $|\Psi(h_W)[x]| \leq C \|h_W\|_{H_W}$ , thus  $\Psi$  as a mapping  $(H_W, \|\cdot\|_{H_W}) \rightarrow (\mathcal{C}(\mathcal{X}), \|\cdot\|_{\infty})$  is continuous.

Thanks to the continuity of  $\Psi : (H_W, \|\cdot\|_{H_W}) \rightarrow (\mathcal{C}(\mathcal{X}), \|\cdot\|_{\infty})$ , we conclude that  $\mathcal{B}_{H_W} := \{h_W \in H_W : \|\Psi(h_W) - f\|_{\infty} \leq \varepsilon\}$  is closed in  $(H_W, \|\cdot\|_{H_W})$ . Regarding  $\mathcal{B}_H$ , by [\[CS08\]](#) Lemma 4.23, since the kernel  $k$  is bounded on  $\mathcal{X}$ , the inclusion  $\iota : (H, \|\cdot\|_H) \rightarrow (\mathcal{C}(\mathcal{X}), \|\cdot\|_{\infty})$  is continuous. We deduce the closedness of  $\mathcal{B}_H = \iota^{-1}(\mathcal{B}_{\mathcal{C}(\mathcal{X})}(f, \varepsilon))$  in  $(\mathcal{C}(\mathcal{X}), \|\cdot\|_{\infty})$ .

Finally, the sets  $\mathcal{B}_{H_W}$  and  $\mathcal{B}_H$  are non-empty since  $H = \Psi(H_W)$  is dense in  $(\mathcal{C}(X), \|\cdot\|_{\infty})$ .

We conclude that  $\mathcal{B}_{H_W}$ , resp.  $\mathcal{B}_H$  is a non-empty closed convex set in the Hilbert space  $(H_W, \|\cdot\|_{H_W})$ , resp.  $(H, \|\cdot\|_H)$ , and the Hilbert projection theorem (or directly [Theorem 4.10](#) in [\[Rud87\]](#)) ensures that  $h_W[\varepsilon, f]$ , resp.  $h[\varepsilon, f]$  is uniquely defined.



Now, we show that  $\|h[\varepsilon, f]\|_H = \|h_W[\varepsilon, f]\|_{H_W}$ . By [CS08] Theorem 4.21, we have for all  $h \in H$ :

$$\|h\|_H = \inf\{\|h_W\|_{H_W}, h = \Psi(h_W)\}.$$

By the same argument as before (using Theorem 4.10 in [Rud87]), we show that the infimum is attained. The equality between norms is then straightforward by separating both inequalities and using  $H = \Psi(H_W)$ .

To obtain Eq. (23), we take  $h := \Psi(h_W[\varepsilon, f])$  in Eq. (21) and  $\hat{h} := \hat{\Psi}(h_W[\varepsilon, f]) \in \hat{H}$ , and apply the triangle inequality for  $\|\cdot\|_\infty$ , using  $\|h - f\|_\infty \leq \varepsilon$  and  $\|h[\varepsilon, f]\|_H = \|h_W[\varepsilon, f]\|_{H_W}$ .  $\square$

The approximation result in Eq. (23) shows that  $\hat{k}$  is  $\rho$ -approximately universal (Definition 7) for  $\rho(f) := \rho_0 \|h[\varepsilon, f]\|_H$ . In the case where  $\mathcal{X}$  is of dimension  $d$  (or has intrinsic dimension  $d$ ), the number of covering balls scales as  $J = \mathcal{O}(\eta^{-d})$ , which does not impact the approximation rate, as is commonly the case in kernel methods which do not suffer from the curse of dimensionality (see for example [Gre+12] Section 4.1). However, as is typically the case for discretisation methods, the rate  $J = \mathcal{O}(\eta^{-d})$  is computationally prohibitive for small discretisation step  $\eta$  in high dimension  $d$ .

From a functional standpoint, a larger oscillation (a large value for  $C_{K'} = \max_{t \in [-D_{\mathcal{X}}^2, D_{\mathcal{X}}^2]} |K'(t)|$  e.g. for the Taylor case), of the function  $K$  worsens the error, which could be understood as excessive locality or over-fitting. Finally, the error term  $\rho(f)$  is relative in the sense that it depends on  $\|h[\varepsilon, f]\|_H$ , which is the smallest possible norm of an  $\varepsilon$ -approximation of  $f$  within  $H$ , and can be seen as a measure of complexity of  $f$  (in loose terms). This term depends on  $q$ , and while the exact dependence is unclear, we expect it to grow as  $q$  increases.

### 3.3 An Approximate Universal Truncated Kernel

In this section, we consider another approximate universal kernel which is obtained by truncation of  $\varphi$ . We will undergo a similar process as the construction of  $\hat{H}$  in Section 3.1 and follow closely the proof methods of Section 3.2.

A natural idea is to simply consider a truncation of the mapping  $\varphi$  from Proposition 5: fixing a basis  $(x_n)_{n \in \mathbb{N}}$  of  $\mathcal{X}$ , a discretisation size  $N \geq 2$  and scale  $q > 1$ , consider the mapping:

$$\varphi_t := \begin{cases} \mathcal{X} & \longrightarrow \mathbb{R}^N \\ x & \longmapsto \left( \frac{c_\varphi d_{\mathcal{X}}(x, x_n)}{q^j} \right)_{n \in \llbracket 0, N-1 \rrbracket} \end{cases}.$$

Straightforward computation shows that  $\varphi_t : (\mathcal{X}, d_{\mathcal{X}}) \longrightarrow (\ell^2, \|\cdot\|_{\ell^2})$  is  $\sqrt{1 - q^{-2N}}$ -Lipschitz. We introduce the “padding” isometry:

$$B_t := \begin{cases} \mathbb{R}^N & \longrightarrow \ell^2 \\ (u_n)_{n=0}^{N-1} & \longmapsto (u_0, \dots, u_{N-1}, 0, \dots) \end{cases},$$

Similarly to Section 3.1, we take  $V := \varphi(\mathcal{X})$ ,  $U_t := \varphi_t(\mathcal{X})$  and  $V_t := B_t(U_t)$ , allowing us to introduce the compact set  $W := V \cup V_t \subset \ell^2$  (we use the same notation as in Section 3.1 to alleviate notation). Take  $k_W$  a Taylor or radial kernel on  $W$ , and introduce the kernel:

$$k_t := \begin{cases} \mathcal{X}^2 & \longrightarrow \mathbb{R} \\ (x, y) & \longmapsto k_W(B_t \circ \varphi_t(x), B_t \circ \varphi_t(y)) \end{cases}.$$

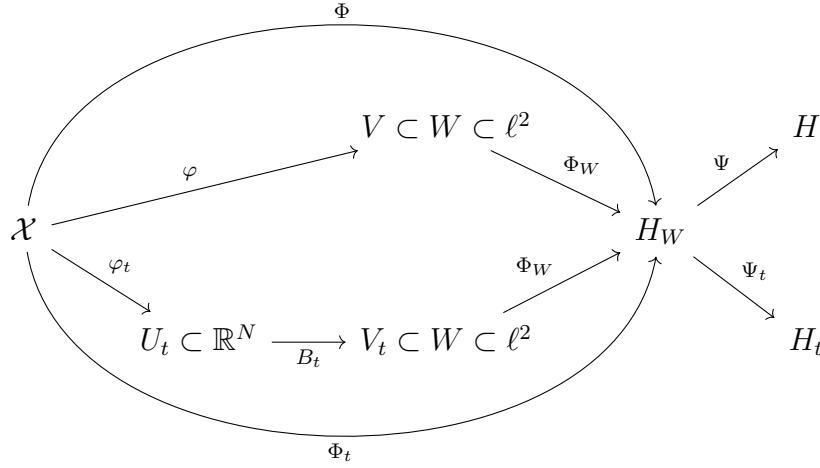
We continue with the feature pair  $(H_W, \Phi_t)$  for the RKHS  $H_t$  associated to  $k_t$ , where:

$$\Phi_t := \begin{cases} \mathcal{X} & \longrightarrow & H_W \\ x & \longmapsto & \Phi_W \circ B_t \circ \varphi_t(x) \end{cases} .$$

As in [Eq. \(17\)](#) we introduce the “feature-to-map” functionals:

$$\Psi := \begin{cases} H_W & \longrightarrow & H \\ h_W & \longmapsto & x \mapsto \langle h_W, \Phi(x) \rangle_{H_W} \end{cases}, \quad \Psi_t := \begin{cases} H_W & \longrightarrow & H_t \\ h_W & \longmapsto & x \mapsto \langle h_W, \Phi_t(x) \rangle_{H_W} \end{cases}, \quad (25)$$

and finish the diagram:



The computation in the proof of [Corollary 11](#) stands, but the coefficients in the expression of  $\varphi_t$  lead to a different expression for  $k_t$ , which is a truncated version of  $k$ : if  $k_W$  is a Taylor kernel, we have:

$$k_t(x, y) = K(\langle \varphi_t(x), \varphi_t(y) \rangle_{\mathbb{R}^N}) = \sum_{n=0}^{+\infty} a_n \left( \sum_{m=0}^{N-1} \frac{c_\varphi^2 d_{\mathcal{X}}(x, x_m) d_{\mathcal{X}}(y, x_m)}{q^{2m}} \right)^n,$$

and likewise for radial kernels:

$$k_t(x, y) = K(\|\varphi_t(x) - \varphi_t(y)\|_{\mathbb{R}^N}^2) = \int_0^{+\infty} \exp\left(-s \sum_{n=0}^{N-1} \frac{c_\varphi^2 (d_{\mathcal{X}}(x, x_n) - d_{\mathcal{X}}(y, x_n))^2}{q^{2n}}\right) d\mu(s).$$

We now adapt [Proposition 12](#) to  $k_t$ :

**Proposition 16.** For  $x \in \mathcal{X}$ , we have  $\|\varphi(x) - B_t \circ \varphi_t(x)\|_{\ell^2} \leq \frac{D_{\mathcal{X}}}{q^N}$ .  
The diameter of  $W$  verifies  $D_W \leq 2D_{\mathcal{X}}$ .

*Proof.* For  $n \in \llbracket 0, N-1 \rrbracket$ , by construction  $[\varphi(x)]_n = [B_t \circ \varphi_t(x)]_n$ . For  $n \geq N$ , we have:

$$|[\varphi(x)]_n - [B_t \circ \varphi_t(x)]_n| = \frac{c_\varphi d_{\mathcal{X}}(x, y_n)}{q^n},$$

and by bounding the distance term by  $D_{\mathcal{X}}$ , and summing the squares, we obtain:

$$\|\varphi(x) - B_t \circ \varphi_t(x)\|_{\ell^2}^2 \leq \sum_{n=N}^{+\infty} \frac{c_\varphi^2 D_{\mathcal{X}}^2}{q^{2n}} = \frac{c_\varphi^2 D_{\mathcal{X}}^2 q^2}{q^{2N}(q^2 - 1)} = \frac{D_{\mathcal{X}}^2}{q^{2N}}.$$

As for the result on  $D_W$ , it follows from 1-Lipschitzness as done in [Proposition 12](#).  $\square$

As in [Section 3.2](#), it is easy to verify that  $R_W := \max_{w \in W} \|w\|_{\ell^2} \leq D_{\mathcal{X}}$ . Following the same steps as in [Theorem 15](#), we show a similar result for  $k_t$ , replacing  $\eta$  with  $D_{\mathcal{X}} q^{-N}$ :

**Theorem 17.** Let  $\varepsilon > 0$  and  $f \in \mathcal{C}(\mathcal{X})$ , there exists  $h_t \in H_t$  such that:

$$\|f - h_t\|_\infty \leq \varepsilon + \rho_t \|h[\varepsilon, f]\|_H, \quad (26)$$

where  $\rho_t = q^{-N/2} D_{\mathcal{X}} \sqrt{2C_{K'}}$  for a Taylor kernel and  $\rho_t = q^{-N} D_{\mathcal{X}} \sqrt{2C_{K'}}$  for a radial kernel, with the constants  $C_{K'}$  as in Lemma 13.

*Proof.* We follow the same progression as in the proof of Theorem 15. First, we follow the proof of Proposition 14, applying Lemma 13, then upper-bounding the term  $\|\varphi(x) - B_t \circ \varphi_t(x)\|_{\ell_2}^2$  using Proposition 16, and the only difference is the replacement of the term  $\eta$  with  $D_{\mathcal{X}} q^{-N}$ . Having adapted Proposition 14, the proof of Theorem 17 follows as in Theorem 15.  $\square$

To compare with the rate from Theorem 15, we see that the term  $\eta$  is replaced by  $D_{\mathcal{X}} q^{-N}$ . While  $q^{-N}$  becomes rapidly smaller as  $q$  increases, we suspect the term  $\|h[\varepsilon, f]\|_H$  to grow quickly as  $q$  increases, which would favour the kernel  $\hat{k}$  from Section 3.1. From an intuitive standpoint, the quality of the truncation approximation depends on how well the truncated basis  $(x_n)_{n=0}^{N-1}$  represents  $\mathcal{X}$ . For example, if  $\mathcal{X}$  is a manifold of  $\mathbb{R}^d$  with two connected components, and the first  $N$  elements are all in the first component, it can be expected that a substantial part of the information about the space is lost, hindering function approximation. This issue can arise because the “basis” property of  $(x_n)$  relates to the full sequence, whereas truncation focuses on the first  $N$  terms, which are not assumed to satisfy particular conditions. We saw in Sections 3.1 and 3.2 a more principled discretisation approach, which is in some sense a refinement of the truncation principle.

## Acknowledgements

We thank Joan Glaunès for carefully proofreading this work and for his valuable insight.

This research was funded in part by the Agence nationale de la recherche (ANR), Grant ANR-23-CE40-0017 and by the France 2030 program, with the reference ANR-23-PEIA-0004.

## References

- [Aiz64] A Aizerman. “Theoretical foundations of the potential function method in pattern recognition learning”. In: *Automation and remote control* 25 (1964), pp. 821–837.
- [Aro50] Nachman Aronszajn. “Theory of reproducing kernels”. In: *Transactions of the American mathematical society* 68.3 (1950), pp. 337–404.
- [BT11] Alain Berlinet and Christine Thomas-Agnan. *Reproducing kernel Hilbert spaces in probability and statistics*. Springer Science & Business Media, 2011.
- [CS08] Andreas Christmann and Ingo Steinwart. *Support vector machines*. Springer, 2008.
- [CS10] Andreas Christmann and Ingo Steinwart. “Universal kernels on non-standard input spaces”. In: *Advances in neural information processing systems* 23 (2010).
- [CV95] Corinna Cortes and Vladimir Vapnik. “Support-vector networks”. In: *Machine learning* 20 (1995), pp. 273–297.
- [Dor+14] G Doran, K Muandet, K Zhang, and B Schölkopf. “A Permutation-Based Kernel Conditional Independence Test”. In: *30th Conference on Uncertainty in Artificial Intelligence (UAI 2014)*. AUAI Press. 2014, pp. 132–141.
- [Fie23] Christian Fiedler. “Lipschitz and Hölder Continuity in Reproducing Kernel Hilbert Spaces”. In: *arXiv preprint arXiv:2310.18078* (2023).

- [FSG13] Kenji Fukumizu, Le Song, and Arthur Gretton. “Kernel Bayes’ rule: Bayesian inference with positive definite kernels”. In: *The Journal of Machine Learning Research* 14.1 (2013), pp. 3753–3783.
- [Góm+09] Luis Gómez-Chova, Gustavo Camps-Valls, Lorenzo Bruzzone, and Javier Calpe-Maravilla. “Mean map kernel methods for semisupervised cloud classification”. In: *IEEE Transactions on Geoscience and Remote Sensing* 48.1 (2009), pp. 207–220.
- [Gre+12] Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola. “A kernel two-sample test”. In: *The Journal of Machine Learning Research* 13.1 (2012), pp. 723–773.
- [HC11] Robert Hable and Andreas Christmann. “On qualitative robustness of support vector machines”. In: *Journal of Multivariate Analysis* 102.6 (2011), pp. 993–1007.
- [Jay+15] Sadeep Jayasumana, Richard Hartley, Mathieu Salzmann, Hongdong Li, and Mehrtash Harandi. “Kernel methods on Riemannian manifolds with Gaussian RBF kernels”. In: *IEEE transactions on pattern analysis and machine intelligence* 37.12 (2015), pp. 2464–2477.
- [Li+17] Chun-Liang Li, Wei-Cheng Chang, Yu Cheng, Yiming Yang, and Barnabás Póczos. “Mmd gan: Towards deeper understanding of moment matching network”. In: *Advances in neural information processing systems* 30 (2017).
- [MXZ06] Charles A Micchelli, Yuesheng Xu, and Haizhang Zhang. “Universal Kernels.” In: *Journal of Machine Learning Research* 7.12 (2006).
- [Mua+12] Krikamol Muandet, Kenji Fukumizu, Francesco Dinuzzo, and Bernhard Schölkopf. “Learning from distributions via support measure machines”. In: *Advances in neural information processing systems* 25 (2012).
- [Mua+17] Krikamol Muandet, Kenji Fukumizu, Bharath Sriperumbudur, Bernhard Schölkopf, et al. “Kernel mean embedding of distributions: A review and beyond”. In: *Foundations and Trends® in Machine Learning* 10.1-2 (2017), pp. 1–141.
- [Nad64] E. A. Nadaraya. “On Estimating Regression”. In: *Theory of Probability & Its Applications* 9.1 (1964), pp. 141–142. eprint: <https://doi.org/10.1137/1109020>.
- [Par62] Emanuel Parzen. “On estimation of a probability density function and mode”. In: *The annals of mathematical statistics* 33.3 (1962), pp. 1065–1076.
- [Ros56] Murray Rosenblatt. “Remarks on some nonparametric estimates of a density function”. In: *Ann. Math. Stat* 27 (1956), pp. 832–837.
- [Rud87] Walter Rudin. *Real and complex analysis*. McGraw-Hill, Inc., 1987.
- [SS16] Saburo Saitoh and Yoshihiro Sawano. *Theory of reproducing kernels and applications*. Vol. 44. Springer, 2016.
- [SSM98] Bernhard Schölkopf, Alexander Smola, and Klaus-Robert Müller. “Nonlinear component analysis as a kernel eigenvalue problem”. In: *Neural computation* 10.5 (1998), pp. 1299–1319.
- [SS02] Bernhard Schölkopf and Alexander J Smola. *Learning with kernels: support vector machines, regularization, optimization, and beyond*. MIT press, 2002.
- [SFL11] Bharath K Sriperumbudur, Kenji Fukumizu, and Gert RG Lanckriet. “Universality, Characteristic Kernels and RKHS Embedding of Measures.” In: *Journal of Machine Learning Research* 12.7 (2011).
- [Sri+10] Bharath K Sriperumbudur, Arthur Gretton, Kenji Fukumizu, Bernhard Schölkopf, and Gert RG Lanckriet. “Hilbert space embeddings and metrics on probability measures”. In: *The Journal of Machine Learning Research* 11 (2010), pp. 1517–1561.

- [SZ21] Ingo Steinwart and Johanna F. Ziegel. “Strictly proper kernel scores and characteristic kernels on compact spaces”. In: *Applied and Computational Harmonic Analysis* 51 (2021), pp. 510–542.
- [Wat64] Geoffrey S. Watson. “Smooth Regression Analysis”. In: *Sankhyā: The Indian Journal of Statistics, Series A (1961-2002)* 26.4 (1964), pp. 359–372.
- [Zha+11] K Zhang, J Peters, D Janzing, and B Schölkopf. “Kernel-based Conditional Independence Test and Application in Causal Discovery”. In: *27th Conference on Uncertainty in Artificial Intelligence (UAI 2011)*. AUAI Press. 2011, pp. 804–813.
- [ZGD24] Johanna Ziegel, David Ginsbourger, and Lutz Dümbgen. “Characteristic kernels on Hilbert spaces, Banach spaces, and on sets of measures”. In: *Bernoulli* 30.2 (2024), pp. 1441–1457.