



# Learning via variably scaled kernels

C. Campi<sup>1</sup> · F. Marchetti<sup>2</sup> · E. Perracchione<sup>1</sup> 

Received: 30 October 2019 / Accepted: 20 May 2021 / Published online: 26 June 2021

© The Author(s), under exclusive licence to Springer Science+Business Media, LLC, part of Springer Nature 2021

## Abstract

We investigate the use of the so-called variably scaled kernels (VSKs) for learning tasks, with a particular focus on support vector machine (SVM) classifiers and kernel regression networks (KRN). Concerning the kernels used to train the models, under appropriate assumptions, the VSKs turn out to be *more expressive* and *more stable* than the standard ones. Numerical experiments and applications to breast cancer and coronavirus disease 2019 (COVID-19) data support our claims. For the practical implementation of the VSK setting, we need to select a suitable *scaling function*. To this aim, we propose different choices, including for SVMs a probabilistic approach based on the naive Bayes (NB) classifier. For the classification task, we also numerically show that the VSKs inspire an alternative scheme to the sometimes computationally demanding feature extraction procedures.

**Keywords** Variably scaled kernels · Kernel ill-conditioning · Meshfree methods · Binary classification · Regression networks

**Mathematics Subject Classification (2010)** 65D15 · 41A05 · 68Q32

---

Communicated by: Robert Schaback

✉ E. Perracchione  
perracchione@dima.unige.it

C. Campi  
campi@dima.unige.it

F. Marchetti  
francesco.marchetti@unipd.it

<sup>1</sup> Dipartimento di Matematica DIMA, Università di Genova, Genoa, Italy

<sup>2</sup> Dipartimento di Matematica “Tullio Levi-Civita”, Università di Padova, Padua, Italy

## 1 Introduction

In the context of approximation theory, the variably scaled kernels (VSKs) were introduced in 2015 by [6]. The basic idea behind them is to map the initial set of examples via a scaling function and construct an augmented approximation space. Our main contribution consists in linking the VSKs to the field of machine learning, as the VSKs have a long-known equivalent in pattern analysis. Precisely, many methods based on feature augmentation, as, e.g., zero padding and feature replication [9, 21, 27], fall into the general VSK setting that we are going to investigate.

Focusing on kernel learning methods and specifically on KRN and SVMs (see, e.g., [16, 42]), we give a very general formulation of feature augmentation schemes via VSKs. In doing so, we drive our attention towards the Gaussian and linear kernels, being truly popular for learning issues. We provide theoretical results concerning their practical implementation and expressiveness [13] and we further analyze the spectrum of the kernel matrices constructed via VSKs. This study reveals the effectiveness of the proposed approach especially for the Gaussian kernel; indeed, the condition number of the VSK kernel matrix is less than or equal to the condition number of the matrix constructed via the standard kernel. This fact turns out to be relevant for KRN, where one may require to compute the inverse of the kernel matrix, which is usually affected by severe ill-conditioning. Moreover, for the selection of the scaling function of the KRN-VSK, one can refer to the available literature in the context of approximation theory [10, 36]. Indeed, the scaling function might be selected so that it *mimics* the samples and this might lead to an improvement in terms of accuracy and/or stability (see, e.g., [6, 10, 11]). Here, in particular, we propose to use a non-linear fitting of the function itself as augmented feature.

While for the KRN-VSK we can refer to some available literature for selecting the scaling function, for SVM-VSK, we consider a *probabilistic* solution. More precisely, focusing on binary classification problems, we first note that the VSK setting induces new feature maps and spaces that depend on the scaling function associated to the VSK. For being competitive with the accuracy of the classical SVMs, as well as with other common classifiers, we have to select a *suitable* scaling function for the VSKs. To this aim, we remark that the SVM is characterized by a *geometric* point of view. Nevertheless, methods based on probability distributions, as the NB classifiers, might outperform SVM. For that reason, many efforts are devoted to investigate which classifier performs better and under which conditions; for a general overview refer, e.g., to [7, 31, 47]. In this work, we thus fuse SVM and NB classifiers by means of VSKs, so that the mixed approach takes into account the probabilistic features of the NB algorithm and classifies geometrically with SVM. Moreover, we conclude the paper by presenting a feature extraction algorithm that is inspired by the VSK framework and that might be considered in place of other feature extraction schemes; refer, e.g., to [19, 45].

The paper is organized as follows. In Section 2, we briefly review the use of kernels in machine learning literature. In Section 3, we investigate the VSKs for two learning methods, specifically SVM and KRN. Then, in Sections 4 and 5, we drive our attention towards the Gaussian and linear VSKs as well as towards the problem of selecting the scaling function. Section 6 is devoted to numerical experiments with

both toy models and a real dataset. In Section 7, we present a feature extraction algorithm whose underlying idea is derived from the study of the variably scaled setting. The last section deals with conclusions and work in progress.

## 2 Preliminaries

We consider a learning problem with training examples

$$\Sigma = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\},$$

where  $\mathbf{x}_i \in \Omega \subseteq \mathbb{R}^n$  and  $y_i \in \mathbb{R}$ . For the particular case of the classification setting, we fix  $y_i \in \{-1, +1\}$ .

For both SVMs and KRNs, we drive our attention towards (strictly) positive definite kernels  $\kappa : \Omega \times \Omega \rightarrow \mathbb{R}$ , where  $\Omega$  is a bounded set, that can be decomposed via the Mercer’s Theorem as explained below (see, e.g., Theorem 2.2. [15] p. 107 or [26]).

**Theorem 1** *Let  $\kappa : \Omega \times \Omega \rightarrow \mathbb{R}$  be a continuous (strictly) positive definite kernel that satisfies*

$$\int_{\Omega} \kappa(\mathbf{x}, \mathbf{y})v(\mathbf{x})v(\mathbf{y})d\mathbf{x}d\mathbf{y} \geq 0, \quad \forall v \in L_2(\Omega), \mathbf{x}, \mathbf{y} \in \Omega,$$

then the kernel can be expressed as

$$\kappa(\mathbf{x}, \mathbf{y}) = \sum_{k \geq 0} \lambda_k \rho_k(\mathbf{x})\rho_k(\mathbf{y}), \quad \mathbf{x}, \mathbf{y} \in \Omega,$$

where  $\{\lambda_k\}_{k \geq 0}$  are the (non-negative) eigenvalues and  $\{\rho_k\}_{k \geq 0}$  are the ( $L_2$ -orthonormal) eigenfunctions of the operator  $T : L_2(\Omega) \rightarrow L_2(\Omega)$ , given by

$$T[v](\mathbf{x}) = \int_{\Omega} \kappa(\mathbf{x}, \mathbf{y})v(\mathbf{y})d\mathbf{y}.$$

Moreover, such expansion is absolutely and uniformly convergent.

We point out that many relevant kernels, e.g., cases where  $\Omega$  is unbounded or non-measurable, do not fall into the above Mercer decomposition. Thus, on one side, taking only Mercer kernels might be restrictive. On the other side, it provides the adequate background for our purposes and it offers an easy way to introduce feature maps and spaces. Indeed, for such kernels that admit a Mercer expansion (also called valid kernels according to the definition given by [42]), it is worth to note that we can interpret the series representation in terms of an inner product in the so-called *feature space*  $F$ , which is a Hilbert space. Indeed,

$$\kappa(\mathbf{x}, \mathbf{y}) = \langle \Phi(\mathbf{x}), \Phi(\mathbf{y}) \rangle_F, \quad \mathbf{x}, \mathbf{y} \in \Omega, \tag{2.1}$$

where  $\Phi : \Omega \rightarrow F$  is a *feature map*. For a given kernel, the feature map and space are not unique. A possible solution is the one of taking the map  $\Phi(\mathbf{x}) = \kappa(\cdot, \mathbf{x})$ , which is linked to the characterization of  $F$  as a reproducing kernel Hilbert space; see [16, 42] for further details.

In the classification context, many studies are devoted to investigate and measure the *complexity* of a chosen model, such as the so-called VC dimension [44] and the empirical Rademacher complexity [4]. The complexity of a method is usually referred to as *capacity* or *expressiveness*. Indeed, complex models have the capability to perform complex tasks, by determining elaborated decision functions, and thus to express sophisticated links between the data. In any case, the capacity of a method needs to be tailored to the considered task, in order to avoid overfitting; for a general overview, we refer, e.g., the reader to [39].

To better investigate the concept of expressiveness in the kernel setting, we introduce the kernel matrix  $K$  constructed via the dataset  $\mathcal{E} = \{x_1, \dots, x_N\} \subseteq \Omega$ , i.e., the matrix of entries

$$K_{ij} = \kappa(x_i, x_j), \quad i, j = 1, \dots, N, \tag{2.2}$$

where  $\kappa$  is a (strictly) positive definite kernel. Note that if  $\kappa$  is a strictly positive definite kernel then  $K$  is positive definite, while it is positive semi-definite if  $\kappa$  is a positive definite kernel.

*Remark 1* The expressiveness of a kernel-based model is related to the number of dichotomies achievable by a linear separator in the feature space. Moreover, concerning the rank of the kernel matrix, we have the following result [13, Theorem 2, p. 7].

**Theorem 2** *Let  $K$  be the kernel matrix as in (2.2) constructed via  $\mathcal{E} = \{x_1, \dots, x_N\} \subseteq \Omega$ , let us denote by  $\text{rank}(K)$  its rank. Then, there exists at least one subset of examples of size  $\text{rank}(K)$  that can be shattered by a linear function.*

As capacity measure dedicated to the kernel setting, we consider the *spectral ratio* that has been introduced in [13]. It is defined as

$$S(K) = \frac{\text{tr}(K)}{\|K\|_F} = \frac{\sum_{i=1}^N K_{ii}}{\sqrt{\sum_{i=1}^N \sum_{j=1}^N K_{ij}^2}}.$$

According to the following definition (see [13, Definition 1, p. 8]), such quantity is an expressiveness measure for kernels. As a remark, we also point out that it is connected to the empirical Rademacher complexity [13, Theorem 4, p. 9].

**Definition 1** Let  $\kappa_i, \kappa_j : \Omega \times \Omega \rightarrow \mathbb{R}$ , be two (strictly) positive definite kernels. We say that  $\kappa_j$  is more specific (or more expressive) than  $\kappa_i$  whenever for any dataset  $\mathcal{E} = \{x_1, \dots, x_N\} \subseteq \Omega$ , we have

$$S(K^i) \leq S(K^j), \tag{2.3}$$

where  $K^i$  and  $K^j$  are the kernel matrices on  $\mathcal{E}$  obtained via  $\kappa_i$  and  $\kappa_j$ , respectively.

*Remark 2* Technically the Definition 1, which is taken by [13], states that  $\kappa_j$  is more specific (or more expressive) than  $\kappa_i$  also when the equality in (2.3) holds true. In the latter case, we should use the term “equally or more specific than.” To take a common

notation with the native definition, we simply use the term “more expressive” also for the trivial case.

The spectral ratio being an expressiveness measure, it is related to the rank of the kernel matrix (see also Remark 1), indeed

$$1 \leq S(K) \leq \sqrt{\text{rank}(K)}.$$

We conclude this brief review on kernels for machine learning by pointing out that the kernel matrices introduced above might suffer from severe ill-conditioning. In order to partially overcome instability issues in the approximation framework, a possible solution comes from the use of VSKs (see below for their definition), which have been recently introduced in [6]; refer also to [10, 11].

**Definition 2** Let  $\Lambda \subseteq \mathbb{R}^m$ ,  $m > 0 \in \mathbb{N}$ . Let  $\kappa : \tilde{\Omega} \times \tilde{\Omega} \rightarrow \mathbb{R}$ ,  $\tilde{\Omega} = \Omega \times \Lambda \subseteq \mathbb{R}^{n+m}$ , be a continuous (strictly) positive definite kernel. Given a scaling function  $\psi : \Omega \rightarrow \Lambda$ , a variably scaled kernel  $\kappa^\Psi : \Omega \times \Omega \rightarrow \mathbb{R}$  is defined as

$$\kappa^\Psi(\mathbf{x}, \mathbf{y}) := \kappa((\mathbf{x}, \psi(\mathbf{x})), (\mathbf{y}, \psi(\mathbf{y}))), \tag{2.4}$$

for  $\mathbf{x}, \mathbf{y} \in \Omega$ .

We point out that in Definition 2, we present a multidimensional extension of the scaling function  $\psi$ , which has been introduced as a real-valued function in the previous literature [6].

When dealing with Mercer’s kernels, the construction of a VSK as in Definition 2 provides a valid kernel. We now extend this general setting to work with KRNs and SVMs.

### 3 Learning with VSKs

To have a clear theoretical framework, we investigate the use of VSKs as a feature augmentation algorithm, where *new* features are added to the original dataset in order to possibly increase the performances of learning schemes. According to Definition 2, we define a function  $\Psi : \Omega \rightarrow \tilde{\Omega}$  as

$$\Psi(\mathbf{x}) := (\mathbf{x}, \psi(\mathbf{x})).$$

The function  $\Psi$  extends the data vector  $\mathbf{x} \in \Omega$ , including  $m$  features that depend on the original ones. The VSK kernel defined in (2.4) is a valid kernel, as it corresponds to an inner product in the associated feature space  $F_\Psi$  (see [42, Proposition 3.22, p. 75]). Moreover, it induces a *new* feature map  $\Theta : \Omega \rightarrow F_\Psi$  so that

$$\kappa^\Psi(\mathbf{x}, \mathbf{y}) = \langle \Theta(\mathbf{x}), \Theta(\mathbf{y}) \rangle_{F_\Psi}. \tag{3.1}$$

Referring to (2.1), because of [6, Theorem 3.1], the spaces  $F_\Psi$  and the classical feature space  $F$ , associated to  $\kappa : \tilde{\Omega} \times \tilde{\Omega} \rightarrow \mathbb{R}$  and induced by the feature map  $\Upsilon : \tilde{\Omega} \rightarrow F$ , are isometric; see also [10, Proposition 2.3].

We now investigate the use of the VSKs for both SVMs and KRNs.

### 3.1 SVM-VSK

In this section, we present the VSK setting in the SVM algorithm. For this general overview, we also refer the reader to [16, 42].

We take  $\mathcal{E} = \{\mathbf{x}_i, i = 1, \dots, N\} \subseteq \Omega$ , where  $\Omega \subseteq \mathbb{R}^n$ . The associate function values are so that  $y_i \in \{-1, +1\}, i = 1, \dots, N$ . Indeed, for the binary classification problem via VSKs, we need to find a predictor, i.e., a decision function  $s^\Psi : \Omega \rightarrow \{-1, +1\}$ , that assigns appropriate labels, i.e.,  $\tilde{y}_i \in \{-1, +1\}$ , to other unknown samples  $\tilde{\mathbf{x}}_i, i = 1, \dots, t$ .

Given  $\mathbf{x} = (x_1, \dots, x_n)^\top \in \Omega$ , we define a non-linear SVM classifier that makes use of VSKs via the following decision function:

$$s^\Psi(\mathbf{x}) = \text{sign}(h^\Psi(\mathbf{x})) := \text{sign}(\langle \Theta(\mathbf{x}), \mathbf{w} \rangle_{F_\Psi} + b),$$

where  $\Theta : \Omega \rightarrow F_\Psi$  is the VSK feature map and

$$\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \Theta(\mathbf{x}_i) \in F_\Psi.$$

The coefficients  $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_N) \in \mathbb{R}^N$  are the solution of the following *soft margin* problem [16, Sect. 18, pp. 346–347]

$$\begin{cases} \min_{\boldsymbol{\alpha} \in \mathbb{R}^N} \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \kappa^\Psi(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i=1}^N \alpha_i, \\ \text{s.t. } \sum_{i=1}^N \alpha_i y_i = 0, \\ 0 \leq \alpha_i \leq \zeta, \quad i = 1, \dots, N, \end{cases}$$

where  $[0, \zeta]^N$  is known as *bounding box*, with  $\zeta \in \mathbb{R} = [0, +\infty)$ . Then,

$$b = y_i - \sum_{j=1}^N \alpha_j \kappa^\Psi(\mathbf{x}_i, \mathbf{x}_j). \tag{3.2}$$

Observe that in the computation of  $b$  any given index  $i$  so that  $0 < \alpha_i < \zeta$  can be used. However, to make  $b$  uniquely defined and for stability purposes, it is computed via an average over all such candidates.

The equation of the SVM decision function  $s^\Psi : \Omega \rightarrow \{-1, +1\}$ , i.e.,  $\mathbf{w}$  and  $b$  as in equation (3.1) and (3.2), is then found by imposing the Karush Kuhn Tucker conditions (see, e.g., [29]) and thanks to (3.1), for  $\mathbf{x} \in \Omega$ , it reads as follows:

$$s^\Psi(\mathbf{x}) = \text{sign}(h^\Psi(\mathbf{x})) = \text{sign}(\langle \Theta(\mathbf{x}), \mathbf{w} \rangle_{F_\Psi} + b) = \text{sign}\left(\sum_{i=1}^N y_i \alpha_i \kappa^\Psi(\mathbf{x}, \mathbf{x}_i) + b\right).$$

If one uses the standard kernel  $\kappa : \Omega \times \Omega \rightarrow \mathbb{R}$ , then we recover the classical SVM setting.

As a second test case for the use of VSKs in the machine learning context, we investigate regression networks.

### 3.2 KRN-VSK

Since here KRN are used for regression/interpolation tasks, given distinct data  $\mathcal{E} = \{\mathbf{x}_i, i = 1, \dots, N\} \subseteq \Omega$ , where  $\Omega \subseteq \mathbb{R}^n$ , we fix the output variables  $y_i \in \mathbb{R}$ ,  $i = 1, \dots, N$ .

Concerning supervised learning networks, the simplest strategy consists in learning the trend between inputs and outputs via a predictor  $s^\Psi : \Omega \rightarrow \mathbb{R}$  which is a linear combination of *some* basis functions, in this case VSKs. For a general overview on KRN, we refer the reader to [16, 30].

We keep the general framework of KRN and we adapt them to the use of VSKs. Here, we focus on kernels with centers at locations  $Z = \{\mathbf{z}_i, i = 1, \dots, M\} \subseteq \Omega$ ; and thus, our KRN-VSK predictor  $s^\Psi : \Omega \rightarrow \mathbb{R}$  is of the form

$$s^\Psi(\mathbf{x}) = \sum_{i=1}^M c_i \kappa^\Psi(\mathbf{x}, \mathbf{z}_i),$$

for (strictly) positive definite kernels  $\kappa^\Psi : \Omega \times \Omega \rightarrow \mathbb{R}$  and for some real coefficients  $c_1, \dots, c_M$ .

For KRN-VSK, we compute  $\mathbf{c} = (c_1, \dots, c_M)^\top \in \mathbb{R}^M$  via the following minimization problem [15]

$$\min_{\mathbf{c} \in \mathbb{R}^M} \left[ \sum_{i=1}^N \left( y_i - \sum_{j=1}^M c_j \kappa^\Psi(\mathbf{x}_i, \mathbf{z}_j) \right)^2 + \nu \sum_{j=1}^M c_j^2 \right],$$

where  $\nu \in \mathbb{R}_+$  is a regularization parameter.

In the following, we may take the set of kernel centers  $Z \equiv \mathcal{E}$ . In that case, the kernel matrix  $\mathbf{K}^\Psi$  of entries

$$\mathbf{K}_{ij}^\Psi = \kappa^\Psi(\mathbf{x}_i, \mathbf{x}_j), \quad i, j = 1, \dots, N,$$

is square. Furthermore, if a strictly positive definite kernel as the Gaussian function is used, then the matrix is non-singular. Therefore, we may look at the special setting for which  $\nu = 0$ . In that case, the solution can be found as  $\mathbf{c} = (\mathbf{K}^\Psi)^{-1} \mathbf{y}$ , where  $\mathbf{y} = (y_1, \dots, y_N)^\top$  and  $\mathbf{c} = (c_1, \dots, c_N)^\top$ .

In general, computing the inverse of the kernel matrix  $\mathbf{K}$  might lead to serious instability issues due to the typical ill-conditioning of the kernel matrix. This problem may be somehow overcome by selecting a *safe shape parameter*  $\gamma$ , formally introduced below, and/or by using stable bases; refer, e.g., to [20, 24, 34]. In the incoming sections, we will point out that the use of VSKs might reduce the usual ill-conditioning of the kernel matrices.

## 4 Gaussian and linear VSKs

In this section, we focus on specific kernels providing the practical implementation of the variably scaled setting. Furthermore, we also study the expressiveness and the conditioning induced by the VSKs.

### 4.1 Gaussian kernel

Radial kernels are truly common. They are kernels for whom there exists a radial basis function (RBF)  $\varphi : \mathbb{R}_+ \rightarrow \mathbb{R}$ , where  $\mathbb{R}_+ := [0, \infty)$ , and (possibly) a shape parameter  $\gamma > 0$  such that, for all  $\mathbf{x}, \mathbf{y} \in \Omega$ ,

$$\kappa(\mathbf{x}, \mathbf{y}) = \kappa_\gamma(\mathbf{x}, \mathbf{y}) = \varphi_\gamma(\|\mathbf{x} - \mathbf{y}\|_2) = \varphi(r),$$

where  $r := \|\mathbf{x} - \mathbf{y}\|_2$ .

Among all radial kernels, we remark that the Gaussian is given by

$$\kappa(\mathbf{x}, \mathbf{y}) = \kappa_\gamma(\mathbf{x}, \mathbf{y}) = e^{-\gamma\|\mathbf{x}-\mathbf{y}\|_2^2} = e^{-\gamma r^2}. \tag{4.1}$$

We now discuss its practical implementation in the variably scaled setting. We point out that the Gaussian kernel is strictly positive definite; and thus, its associated kernel matrix turns out to be positive definite, provided that the data are distinct; see, e.g., [16].

#### Practical implementation for the Gaussian VSK

Throughout this section, we take  $N$  data points  $\mathcal{E} = \{\mathbf{x}_i, i = 1, \dots, N\} \subseteq \Omega$ , where  $\Omega \subseteq \mathbb{R}^n$  and we consider a subset  $\Lambda \subseteq \mathbb{R}^m$ .

The Gaussian VSK matrix can be seen as a Hadamard product; indeed, we have the following result.

**Theorem 3** *Let  $\mathcal{E} = \{\mathbf{x}_i, i = 1, \dots, N\} \subseteq \Omega$  be a set of data points. Let  $\psi : \Omega \rightarrow \Lambda$  be the scaling function for the VSK setting. Let  $\kappa : \Omega \times \Omega \rightarrow \mathbb{R}$  be the Gaussian kernel. Then, the VSK matrix constructed on  $\mathcal{E}$  via  $\kappa^\Psi : \Omega \times \Omega \rightarrow \mathbb{R}$  is given by*

$$\mathbf{K}^\Psi = \mathbf{K} \circ \mathbf{K}^\psi,$$

where  $\mathbf{K}_{ij}^\psi = e^{-\|\psi(\mathbf{x}_i) - \psi(\mathbf{x}_j)\|_2^2}$ ,  $i, j = 1, \dots, N$ , and  $\circ$  denotes the Hadamard matrix product.

*Proof* For  $\mathbf{x}, \mathbf{y} \in \Omega$ , we have that

$$\kappa^\Psi(\mathbf{x}, \mathbf{y}) = e^{-(\|\mathbf{x}-\mathbf{y}\|_2^2 + \|\psi(\mathbf{x}) - \psi(\mathbf{y})\|_2^2)} = e^{-\|\mathbf{x}-\mathbf{y}\|_2^2} e^{-\|\psi(\mathbf{x}) - \psi(\mathbf{y})\|_2^2}.$$

Therefore, the entries of the VSK matrix built on  $\mathcal{E} = \{\mathbf{x}_i, i = 1, \dots, N\}$  are given by

$$\mathbf{K}_{ij}^\Psi = e^{-\|\mathbf{x}_i - \mathbf{x}_j\|_2^2} e^{-\|\psi(\mathbf{x}_i) - \psi(\mathbf{x}_j)\|_2^2}, \quad i, j = 1, \dots, N,$$

and thus

$$\mathbf{K}^\Psi = \mathbf{K} \circ \mathbf{K}^\psi.$$

□

About the Hadamard product, we report here a result that can be traced back to 1911 by Schur [40]. It will be helpful in what follows; refer also to [14, Lemma A.5] and [18, Lemma 2.1].



**Theorem 4** *If  $E$  and  $M \in \mathbb{R}^{N \times N}$  are positive definite matrices, denoting by  $\lambda_{\min}$  and  $\lambda_{\max}$  the smallest and largest eigenvalue of a matrix, we have that*

$$\lambda_{\min}(E) \min_{i=1, \dots, N} M_{ii} \leq \lambda_i(E \circ M) \leq \lambda_{\max}(E) \max_{i=1, \dots, N} M_{ii}.$$

This result allows us to infer about the spectrum of the kernel matrix (see [12]) and to show that with the Gaussian VSK we gain both in terms of stability and expressiveness of the kernel.

*Spectral ratio for the Gaussian VSK*

We now give upper and lower bounds for the Frobenius norm  $\|\cdot\|_F$  of the kernel matrix  $K$  in terms of its variably scaled setting. This turns out to be helpful when comparing the spectral ratio of the two matrices ( $K$  and  $K^\psi$ ).

**Theorem 5** *Let  $\mathcal{E} = \{x_i, i = 1, \dots, N\} \subseteq \Omega$  be a set of data points. Let  $\psi : \Omega \rightarrow \Lambda$  be the scaling function for the VSK setting. Let  $\kappa : \Omega \times \Omega \rightarrow \mathbb{R}$  be the Gaussian kernel. Given the VSK matrix  $K^\psi = K \circ K^\psi$  constructed on  $\mathcal{E}$  via  $\kappa^\psi : \Omega \times \Omega \rightarrow \mathbb{R}$ , we have that*

$$\|K^\psi\|_F \leq \|K\|_F \leq \|K\|_F \|K^\psi\|_F.$$

*Proof* Being the RBF  $\varphi : \mathbb{R}_+ \rightarrow \mathbb{R}$  associated to the Gaussian kernel  $\kappa$  non-increasing, for  $x, y \in \Omega$ , we obtain

$$\varphi(\|x - y\|_2^2) \geq \varphi(\|x - y\|_2^2 + \|\psi(x) - \psi(y)\|_2^2),$$

which in particular implies that

$$K_{ij} \geq K_{ij}^\psi \geq 0, \quad i, j = 1, \dots, N.$$

Thus, we get

$$\|K\|_F \geq \|K^\psi\|_F.$$

Moreover, since  $\varphi(0) = 1$ , i.e.,  $K_{ii}^\psi = 1, i = 1, \dots, N$ , we obtain

$$\|K^\psi\|_F \geq \sqrt{\sum_{i=1}^N (K_{ii}^\psi)^2} = \sqrt{N(\varphi(0))^2} \geq 1,$$

and therefore

$$\|K^\psi\|_F \leq \|K\|_F \leq \|K\|_F \|K^\psi\|_F.$$

□

From this theorem, we can easily infer on the spectral ratio in the VSK setting.

**Corollary 1** *Let  $\mathcal{E} = \{x_i, i = 1, \dots, N\} \subseteq \Omega$  be a set of data points. Let  $\psi : \Omega \rightarrow \Lambda$  be the scaling function for the VSK setting. Let  $\kappa : \Omega \times \Omega \rightarrow \mathbb{R}$  be the Gaussian kernel, then the VSK kernel  $\kappa^\psi : \Omega \times \Omega \rightarrow \mathbb{R}$  is more expressive than  $\kappa$ .*

*Proof* Let  $K^\Psi = K \circ K^\psi$  be the VSK matrix constructed on  $\mathcal{E}$  via  $\kappa^\psi : \Omega \times \Omega \rightarrow \mathbb{R}$ , we have that

$$\text{tr}(K^\Psi) = \text{tr}(K) = N\varphi(0) = N,$$

where  $\varphi : \mathbb{R}_+ \rightarrow \mathbb{R}$  is the RBF associated to the Gaussian kernel  $\kappa : \Omega \times \Omega \rightarrow \mathbb{R}$ . Taking into account Theorem 5, we obtain

$$S(K) = \frac{N}{\|K\|_F} \leq \frac{N}{\|K^\Psi\|_F} = S(K^\Psi). \tag{4.2}$$

□

*Remark 3* The concept of expressiveness when the equality in (4.2) holds true has already been clarified in Remark 2. We further point out that the equality is satisfied, for instance, in the trivial case for which  $\psi(x) \equiv \mathbf{0}$ , for all  $x \in \Omega$ .

On one side, the fact that the Gaussian VSK is more expressive than the standard one tells us that the VSK-based learning might be able to deal with more complex tasks. In the next subsection, we focus on the stability of the kernel matrix.

### Spectrum of the Gaussian VSK

The smallest eigenvalue of a positive definite kernel matrix is of course linked to the ill-conditioning. Moreover, given  $\mathcal{E} = \{x_i, i = 1, \dots, N\} \subseteq \Omega$ , the stability is also related to the separation distance

$$q_{\mathcal{E}} := \frac{1}{2} \min_{i \neq j} \|x_i - x_j\|_2,$$

which only depends on the data. As shown in, e.g., [6], we have that

$$q_{\mathcal{E}} \leq q_{\mathcal{E}}^\Psi,$$

where

$$q_{\mathcal{E}}^\Psi := \frac{1}{2} \min_{i \neq j} \|\Psi(x_i) - \Psi(x_j)\|_2,$$

is the separation distance in the VSK setting. This gives the intuition of the fact that the VSKs might lead to possible improvements in terms of stability [6]. Indeed, in general, it is well-known that the smallest eigenvalue of the kernel matrix is related to the separation distance, meaning that the ill-conditioning usually grows as the separation distance decreases; refer, e.g., to [28], where the authors make use of a result from [3] on the eigenvalues of distance matrices. These facts are the fruits on many studies on the so-called *trade-off* or *uncertainty principle* [37, 38], which could be summarized in a conflict between accuracy and stability.

As already mentioned, the VSKs are helpful for improving the stability, especially in view of the following property. We also refer the reader to [43, Corollary 3.1]. For a given matrix  $M$ , we focus on the 2-condition number defined as

$$\text{cond}(M) = \|M\|_2 \|M^{-1}\|_2.$$

**Proposition 1** *Let  $\mathcal{E} = \{x_i, i = 1, \dots, N\} \subseteq \Omega$  be a set of distinct data. Let  $\psi : \Omega \rightarrow \Lambda$  be the scaling function for the VSK setting. Let  $\kappa : \Omega \times \Omega \rightarrow \mathbb{R}$*

be the Gaussian kernel. Given the VSK matrix  $K^\Psi = K \circ K^\psi$  constructed on  $\mathcal{E}$  via  $\kappa^\psi : \Omega \times \Omega \rightarrow \mathbb{R}$ , we have that

$$\text{cond}(K^\Psi) \leq \text{cond}(K).$$

*Proof* First note that, since in this case the matrix is positive definite, the condition number can be computed as

$$\text{cond}(K^\Psi) = \frac{\lambda_{\max}(K^\Psi)}{\lambda_{\min}(K^\Psi)}.$$

Moreover, from Theorem 4 and since the RBF  $\varphi : \mathbb{R}_+ \rightarrow \mathbb{R}$  associated to the Gaussian kernel  $\kappa : \Omega \times \Omega \rightarrow \mathbb{R}$  is so that  $\varphi(0) = 1$ , i.e.,  $K_{ii}^\psi = 1, i = 1, \dots, N$ , we obtain

$$\text{cond}(K^\Psi) = \frac{\lambda_{\max}(K^\Psi)}{\lambda_{\min}(K^\Psi)} \leq \frac{\lambda_{\max}(K)}{\lambda_{\min}(K^\Psi)} \leq \frac{\lambda_{\max}(K)}{\lambda_{\min}(K)} = \text{cond}(K).$$

□

This result turns out to be meaningful especially for the KRN-VSK approach.

As a second case study, we now consider the linear kernel, which is truly popular for classification tasks.

### 4.2 The linear VSK

For  $\mathbf{x}, \mathbf{y} \in \Omega$ , the linear kernel  $\kappa : \Omega \times \Omega \rightarrow \mathbb{R}$  is given by

$$\kappa(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\top \mathbf{y}.$$

As for the Gaussian kernel, its implementation in the variably scaled setting turns out to be trivial. We remark that the linear kernel is positive definite; and thus, its associated kernel matrix turns out to be positive

*Practical implementation for the linear VSK*

The linear VSK can be written as sum of matrices; indeed, we have the following result.

**Theorem 6** Let  $\mathcal{E} = \{\mathbf{x}_i, i = 1, \dots, N\} \subseteq \Omega$  be a set of data points. Let  $\psi : \Omega \rightarrow \Lambda$  be the scaling function for the VSK setting. Let  $\kappa : \Omega \times \Omega \rightarrow \mathbb{R}$  be the linear kernel. Then, the VSK matrix constructed on  $\mathcal{E}$  via  $\kappa^\psi : \Omega \times \Omega \rightarrow \mathbb{R}$  is given by

$$K^\Psi = K + K^\psi,$$

where  $K_{ij}^\psi = \psi(\mathbf{x}_i)^\top \psi(\mathbf{x}_j), i, j = 1, \dots, N$ .

*Proof* For  $\mathbf{x}, \mathbf{y} \in \Omega$  we have that:

$$\kappa^\psi(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^\top, \psi(\mathbf{x})^\top) \begin{pmatrix} \mathbf{y} \\ \psi(\mathbf{y}) \end{pmatrix} = \mathbf{x}^\top \mathbf{y} + \psi(\mathbf{x})^\top \psi(\mathbf{y}),$$

and thus the kernel matrix is given by

$$K^\psi = K + K^\psi.$$

□

We now drive our attention towards the expressiveness of the linear VSK.

*Spectral ratio for the linear VSK*

Depending on the function  $\psi$ , we might have that the linear VSK is less expressive than the standard linear kernel; indeed, we have the following proposition.

**Proposition 2** *Let  $\mathcal{E} = \{x_i, i = 1, \dots, N\} \subseteq \Omega$  be a set of data points. Let  $\psi : \Omega \rightarrow \Lambda$  be the scaling function for the VSK setting. Let  $\kappa : \Omega \times \Omega \rightarrow \mathbb{R}$  be the linear kernel. Let us suppose that the associated kernel matrix  $K$  is non-negative, i.e., so that all the entries of  $K$  are non-negative. Given the VSK matrix  $K^\psi = K + K^\psi$  constructed on  $\mathcal{E}$  via  $\kappa^\psi : \Omega \times \Omega \rightarrow \mathbb{R}$ , if  $\psi$  is so that  $K^\psi$  is non-negative, then:*

$$\frac{\text{tr}(K)}{\text{tr}(K^\psi)} \leq \frac{S(K)}{S(K^\psi)} \leq \frac{\|K^\psi\|_F}{\|K\|_F}.$$

*Proof* Under our assumptions, if  $\psi : \Omega \rightarrow \Lambda$  is so that  $K^\psi$  is non-negative, we have that

$$\frac{\text{tr}(K)}{\text{tr}(K^\psi)} \leq 1.$$

Moreover, since we suppose  $K$  to be non-negative, we get

$$\frac{\|K^\psi\|_F}{\|K\|_F} \geq 1.$$

Finally, taking into account the definition of the spectral ratio, the statement follows. □

Note that the requirements of Proposition 2 are satisfied, e.g., if  $\Omega \subseteq \mathbb{R}_+^n$  and  $\Lambda \subseteq \mathbb{R}_+^m$ .

*Minimum eigenvalue of the linear VSK matrix*

Being Gramian matrices,  $K^\psi$  and  $K$  are positive semi-definite. Concerning the minimum eigenvalue of the VSK matrix  $K^\psi$ , by virtue of Weyl’s inequality (see, e.g., [5, Section III.2, p. 62]), we obtain that:

$$\lambda_{\min}(K) \leq \lambda_{\min}(K + K^\psi) = \lambda_{\min}(K^\psi).$$

As for the Gaussian kernel, one can make many different choices for the function  $\psi$ . Some of them are discussed in the next section.

*Remark 4* In this section, we provided some theoretical findings concerning the expressiveness of VSKs in terms of the spectral ratio, without taking into account the tuning of the shape parameter  $\gamma$  (see (4.1)), which is considered fixed. While such a theoretical investigation concerns a relevant topic in the theory of machine learning, as we pointed out in Section 2, the spectral ratio represents a *poor* choice for model

selection in practical applications. Indeed, in view of the maximization of a certain score (e.g., accuracy, AUC, f1-score), it is convenient to perform a classical tuning of the model parameters (e.g., SVM, KRN) instead of analyzing its capacity via the spectral ratio *as it is*.

## 5 Choices for the scaling function

In the framework of approximation theory, as well as for KRNs, the choice of the scaling function can be guided by some characteristics concerning the data distribution or the underlying function that needs to be reconstructed (see, e.g., [35, 36]). In the classification setting, the VSKs can be seen as feature augmentation methods. More precisely, our aim is to adopt this strategy to encode possible a priori information in the kernel. Let us take  $N$  data points  $\mathcal{E} = \{\mathbf{x}_i, i = 1, \dots, N\} \subseteq \Omega$ , where  $\Omega \subseteq \mathbb{R}^n$  and consider a subset  $\Lambda \subseteq \mathbb{R}^m$ , we now propose some techniques to define the scaling function of the VSK framework.

### 5.1 Scaling function for SVM-VSK

Depending on the task and on the available knowledge, different choices for the scaling function could be taken into account. Here, we construct the scaling function  $\psi : \Omega \rightarrow \Lambda$  as follows. Given the dataset

$$\Sigma = \{(\mathbf{x}_i, y_i), i = 0, \dots, N, \mathbf{x}_i \in \Omega, y_i \in \{-1, +1\}\},$$

we introduce the classes  $C_1$  and  $C_2$ , associated to the labels  $y = -1$  and  $y = +1$ , respectively. Let  $\tilde{\mathbf{x}} = (\tilde{x}_1, \dots, \tilde{x}_n)$  be a new example that we need to classify. Treating the features as mutually independent, the NB classifier (see, e.g., [1, 25]) computes

$$P_j(\tilde{\mathbf{x}}) := P(\tilde{\mathbf{x}} \in C_j | \tilde{\mathbf{x}}) = \frac{P(C_j) \prod_{i=1}^n P(\tilde{x}_i | C_j)}{P(\tilde{\mathbf{x}})},$$

classifying

$$C(\tilde{\mathbf{x}}) = \operatorname{argmax}_{j=1,2} P_j(\tilde{\mathbf{x}}).$$

The *likelihood*  $\prod_{i=1}^n P(\tilde{x}_i | C_j)$  and the *prior*  $P(C_j)$  are typically estimated from the dataset  $\Sigma$ . In other cases, especially when the dataset is not too large, they could be obtained as a priori knowledge, for example by consulting the literature.

In this view, for the SVM-VSK, we propose the scaling map  $\Psi : \Omega \rightarrow \tilde{\Omega}$  defined by

$$\Psi(\mathbf{x}) := (\mathbf{x}, P_1(\mathbf{x})),$$

and the kernel  $\kappa^\Psi : \Omega \times \Omega \rightarrow \mathbb{R}$

$$\kappa^\Psi(\mathbf{x}, \mathbf{y}) := \kappa(\Psi(\mathbf{x}), \Psi(\mathbf{y})).$$

For  $\mathbf{x} \in \Omega$ , since  $P_2(\mathbf{x}) = 1 - P_1(\mathbf{x})$  and  $P_1(\mathbf{x})$  are not independent, we observe that it is sufficient to consider one of the two probabilities.

Concerning the effectiveness of this scaling function  $\Psi : \Omega \rightarrow \tilde{\Omega}$  for the Gaussian VSK, we refer to the notation introduced in Theorem 3 and we point out that,

for  $\mathbf{x}_i, \mathbf{x}_j \in \mathcal{E}$ ,

$$\mathbf{K}_{ij}^\psi = e^{-(P_1(\mathbf{x}_i) - P_1(\mathbf{x}_j))^2}, \quad i, j = 1, \dots, N.$$

We observe that if  $P_1(\mathbf{x}_i) \approx P_1(\mathbf{x}_j)$ , then  $\mathbf{K}_{ij}^\psi \approx 1$  and so  $\mathbf{K}_{ij}^\psi \approx \mathbf{K}_{ij}$ . Considering instead the linear VSK  $\kappa^\psi : \Omega \times \Omega \rightarrow \mathbb{R}$  described in Section 4.2, we get

$$\mathbf{K}_{ij}^\psi = P_1(\mathbf{x}_i)P_1(\mathbf{x}_j), \quad i, j = 1, \dots, N.$$

We remark that, according to Proposition 2, with the linear VSK we construct kernels that might be less expressive than the standard ones.

For both kernels, this means that the matrices *change* according to our a priori knowledge on the dataset, leading to a different, possibly easier, learning task for SVM.

### 5.2 Scaling function for KRN-VSK

Here, we take again  $N$  distinct data  $\mathcal{E} = \{\mathbf{x}_i, i = 1, \dots, N\} \subseteq \Omega$ , where  $\Omega \subseteq \mathbb{R}^n$ , and the associated measurements  $y_i \in \mathbb{R}, i = 1, \dots, N$ , and consider a subset  $\Lambda \subseteq \mathbb{R}^m$ . We now investigate some ideas to define the scaling function for KRN.

Therefore, concerning the choice of the scaling function  $\psi : \Omega \rightarrow \Lambda$ , we suppose to know the trend of data, which can be modelled via a specific class of functions, i.e., a model  $\mathcal{M} : \Omega \times \mathbb{R}^l \rightarrow \mathbb{R}$  depending on  $\mathbf{x} \in \Omega$ , and on  $l$  parameters  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_l)$ . To determine  $\boldsymbol{\beta}$ , we compute:

$$\boldsymbol{\beta}^* = \min_{\boldsymbol{\beta} \in \mathbb{R}^l} \sum_{i=1}^N (y_i - \mathcal{M}(\mathbf{x}_i, \boldsymbol{\beta}))^2.$$

Then, one possible solution to define the function  $\psi : \Omega \rightarrow \Lambda$  is

$$\psi(\mathbf{x}) = \mathcal{M}(\mathbf{x}, \boldsymbol{\beta}^*). \tag{5.1}$$

Of course, this gives a recipe for the selection of the scaling function which is not unique.

## 6 Numerical tests for SVM-VSK and KRN-VSK

All the performed experiments have been carried out in PYTHON, using also the scientific module scikit-learn [32], on a Intel(R) Core(TM) i7 CPU 4712MQ 2.13 GHz processor.

### 6.1 Tests for SVM-VSK

In the following, we consider different toy datasets of various sizes, with precise probability information concerning the features' distributions, and we compare our SVM-VSK approach with standard SVM and NB classifiers. A freely available software can be downloaded at <https://github.com/emmaA89/SVM-VSK>.

The hyperparameters are validated by taking

$$\zeta \in \{2^{-6}, 2^{-5}, \dots, 2^6\},$$

$$\gamma \in \{10^{-6}, 10^{-5}, \dots, 10^2\}.$$

Moreover, in the validation and in the test steps, we evaluate the performance of the considered methods by means of the  $f_1$ -score, weighted with respect to the classes. We remind that the  $f_1$ -score is defined as the harmonic mean between precision and recall. More precisely, given the number of true positive (TP), false positive (FP), and false negative (FN) cases,

$$f_1\text{-score} = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}},$$

where  $\text{precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$  and  $\text{recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$ .

We proceed by constructing 12 toy datasets that differ in terms of number of features and examples. We now fix  $n = 64$ . Letting  $\Omega \subseteq \mathbb{R}^n$ , they are extracted from the dataset

$$\Gamma = \{(\mathbf{x}_i, y_i), i = 1, \dots, 5000, \mathbf{x}_i \in \Omega, y_i \in \{-1, +1\}\},$$

where the two classes  $C_1$  and  $C_2$ , associated to the labels  $y = -1$  and  $y = +1$  respectively, are exactly balanced. The construction of such a dataset is explained in the following steps.

1. Each class  $C_j, j = 1, 2$ , is characterized by two vectors

$$\boldsymbol{\mu}_j = (\mu_j^1, \dots, \mu_j^n), \quad \boldsymbol{\sigma}_j = (\sigma_j^1, \dots, \sigma_j^n).$$

More precisely, let us denote by  $\mathcal{U}(a, b)$  a univariate uniform random distribution on the interval  $(a, b) \subseteq \mathbb{R}$  and by  $p \sim \mathcal{U}(a, b)$  a sample from such distribution. Then,  $\mu_j^k$  and  $\sigma_j^k, k = 1, \dots, n, j = 1, 2$  are determined as follows:

$$\begin{aligned} \mu_1^k &\sim \mathcal{U}(0, 20), \\ \sigma_1^k &\sim \mathcal{U}(0, 2), \\ \mu_2^k &= \mu_1^k + u^k, \text{ with } u^k \sim \mathcal{U}(0, 2), \\ \sigma_2^k &\sim \mathcal{U}(0, 4.5). \end{aligned}$$

2. We denote by  $\mathcal{N}(\mu, \sigma)$  the univariate normal distribution with mean  $\mu$  and standard deviation  $\sigma$ . Let  $\mathbf{x}_i = (x_i^1, \dots, x_i^n)$  be an example in  $\Omega$  belonging to a class  $C_j, j = 1, 2$ . The elements  $x_i^k$  of  $\mathbf{x}_i \in \Omega, k = 1, \dots, n$ , are then randomly generated as samples of  $\mathcal{N}(\mu_j^k, \sigma_j^k + 1) = \mathcal{N}(\mu_j^k, \sigma_j^k) + \mathcal{N}(0, 1)$ , where  $\mathcal{N}(0, 1)$  is Gaussian white noise.
3. Finally, the data are normalized between  $[0, 1]$ .

From the so-constructed  $\Gamma$ , let  $n_k \in \{2, 4, 16, 64\}$  and let  $\mathbf{x}_{i,n_k}$  be the element  $\mathbf{x}_i \in \Gamma$  reduced in dimensions to  $n_k$  randomly pre-selected features (that are the same for all  $\mathbf{x}_i \in \Gamma$ ). We extract the datasets

$$\Gamma_\xi = \{(\mathbf{x}_{i,n_k}, y_i), i = 1, \dots, N_q, y_i \in \{-1, +1\}\},$$

with  $N_q \in \{100, 500, 2500\}$ ,  $\xi = (N_q, n_k)$ , and preserving the balance among the two classes. More precisely, fixed  $N_q$ , we randomly select  $N_q$  indices from the set  $\{1, \dots, 5000\}$ . Moreover, since all combinations examples-features are taken into account, we obtain 12 different datasets.

In the following description, we fix one of the extracted datasets  $\Gamma_\xi$  for some value of  $N_q$  and  $n_k$ . We divide such a dataset in a training set  $\Sigma_\xi$  and a test set  $T_\xi$ . These sets are so that  $\text{card}(\Sigma_\xi) \approx 2\text{card}(T_\xi)$ .

In this experiment, we suppose to have a priori information and to encode it in the SVM-VSK method by means of the NB algorithm. More precisely, the NB classifier is trained considering both  $\Sigma_\xi$  and  $\bar{\Gamma}_\xi$ , which is defined as the dataset containing the examples of  $\Gamma$ , whose number of features has been reduced to  $n_k$ , that are not in  $\Gamma_\xi$ , i.e.,

$$\bar{\Gamma}_\xi := \Gamma_{(5000, n_k)} \setminus \Gamma_\xi.$$

Therefore, in this test, we compare the performances on  $T_\xi$  of the three methods constructed as follows.

1. The NB classifier, which is trained on  $\bar{\Gamma}_\xi \cup \Sigma_\xi$ . Given  $\mathbf{x} = (x_1, \dots, x_{n_k})$ , we adopt the Gaussian likelihood [33]

$$P(x_i|C_j) = \frac{1}{\sqrt{2\pi(\sigma_j^i)^2}} e^{-\left(\frac{x_i - \mu_j^i}{\sqrt{2}\sigma_j^i}\right)^2}.$$

for  $i = 1, \dots, n_k, j = 1, 2$ .

2. The standard SVM method, which is trained on  $\Sigma_\xi$ .
3. The SVM-VSK classifier, which is trained on  $\Sigma_\xi$  and whose scaling map  $\psi : \Omega \rightarrow \Lambda$ , constructed as explained in Section 5, considers the probabilistic outcomes of the NB classifier.

In order to tune the SVM hyperparameters  $\zeta$  and  $\gamma$ , the latter in case of RBF kernel, we consider a 5-fold cross-validation on  $\Sigma_\xi$ .

We carry out the test for each dataset  $\Gamma_\xi$  and we show the obtained results in Fig. 1. The proposed SVM-VSK algorithm is competitive with the best among SVM and NB methods, slightly outperforming both in some cases.

For the Gaussian kernel, we numerically verify Corollary 1 by reporting in Table 1 the spectral ratios related to the matrices  $K$  and  $K^\psi$ , obtained from the training sets  $\Sigma_\xi$  with  $N_q = 100, 500, 2500$ , and  $n_k = 2$ . The results numerically confirm what was theoretically predicted, i.e., the Gaussian VSK is more expressive than the standard one for a fixed shape parameter.

Moreover, for the linear kernel, we are in the hypothesis of Proposition 2. The quantities involved in that proposition are reported in Table 2. The results support what was theoretically observed.

### 6.2 Tests for KRN-VSK

Here, we refer the reader to [16, Program 18.1, p. 340] for a detailed software that deals with KRNs.



	100;2	100;4	100;16	100;64	500;2	500;4	500;16	500;64	2500;2	2500;4	2500;16	2500;64
NB	0.718	0.812	1.000	1.000	0.795	0.861	0.988	1.000	0.766	0.829	0.982	1.000
SVM lin.	0.619	0.686	0.969	1.000	0.728	0.752	0.952	1.000	0.718	0.750	0.948	0.996
SVM-VSK lin.	0.750	0.812	1.000	1.000	0.800	0.842	0.976	1.000	0.782	0.825	0.982	1.000
SVM RBF	0.656	0.656	1.000	1.000	0.788	0.818	0.958	1.000	0.771	0.820	0.970	1.000
SVM-VSK RBF	0.812	0.781	1.000	1.000	0.777	0.849	0.982	1.000	0.762	0.823	0.982	1.000

**Fig. 1** The  $f_1$ -score of the experiments performed on various datasets using the linear (lin.) and Gaussian kernel (RBF). The considered number of examples and features are displayed on the top

As an example for KRN, we focus on the Italian data of the 2020 COVID-19 pandemic. The task we consider consists in learning the time series, i.e.,  $\Omega \subseteq \mathbb{R}$ , of people that in Italy were hospitalized as intensive care unit (ICU) patients from 24 February 2020 to 26 April 2020. The dataset, provided by the “Dipartimento della Protezione Civile,” is available at <https://github.com/pcm-dpc/COVID-19/tree/master/dati-andamento-nazionale>.

The dataset  $I$  consists of 63 samples and it is divided as follows. The first 58 days define the training set  $\Sigma$ , i.e., they are used to construct the regression model, which is then tested on the last  $t = 5$  days,  $\tilde{x}_i, i = 1, \dots, t$ . Referring to Section 3.2 we take the set of kernel centers  $Z$  as the set of available data in  $\mathcal{E}$  and we construct the model using the Gaussian kernel defined in (4.1). Moreover, the feature augmentation strategy outlined in (5.1) is carried out considering  $\mathcal{M} : \Omega \times \mathbb{R} \rightarrow \mathbb{R}$  given by

$$\mathcal{M}(x, \beta) = e^{-\beta|x-\bar{p}|},$$

where  $\bar{p} = 42$  is the peak of the considered time series. The model  $\mathcal{M}$  is constructed on  $\Sigma$ .

In addition, we encode into the kernel also other available data. Precisely, thinking of time series, one usually disposes of other available and dependent data sampled at the same locations which can be used as additional features (see, e.g., [8, 41]). In this direction, we take into account the total number of COVID19 infected (included death and recovered people), the daily number of new infected and the total number

**Table 1** The spectral ratios of the matrices  $K$  and  $K^\Psi$  related to the normalized training sets  $\Sigma_k$ , varying  $N_q = 100, 500, 2500$ . We set  $n_k = 2$  and we considered a Gaussian kernel with  $\gamma = 1$

N	$S(K)$	$S(K^\Psi)$
50	1.3782 E+00	1.5756E+00
500	1.3222 E+00	1.5197E+00
2500	1.2770 E+00	1.4954E+00

**Table 2** The ratios of the norms involved in Proposition 2 obtained via the linear kernel. The matrices  $K$  and  $K^\psi$  are related to the normalized training sets  $\Sigma_{\tilde{x}}$ , varying  $N_q = 100, 500, 2500$ , and with  $n_k = 2$

N	$\text{tr}(K)/\text{tr}(K^\psi)$	$S(K)/S(K^\psi)$	$\ K^\psi\ _F/\ K\ _F$
50	6.1373 E-01	8.8640E-01	1.4443E+00
500	5.2890 E-01	8.8021E-01	1.6642E+00
2500	5.5705 E-01	8.8461E-01	1.5880E+00

of infected (excluded death and recovered people). Of course, this selection of the scaling function means that we are adding a priori knowledge to the selected time series. Therefore, the scaling function  $\psi$  is so that  $\psi : \Omega \rightarrow \Lambda$ , where  $\Lambda \subseteq \mathbb{R}^4$ .

To analyze the performances of the variably scaled setting, we take the Gaussian kernel and we compute the condition number of the kernel matrix and the rounded mean error (RME). Let

$$\text{ME} = \frac{1}{t} \sum_{i=1}^t |y_i - A(\tilde{x}_i)|,$$

be the mean error, where  $A$  is a decision function as defined in Section 3.2 obtained via classical or variably scaled kernels. Since hospitalized patients are involved in the dynamics we consider, as accuracy indicator, the RME defined as

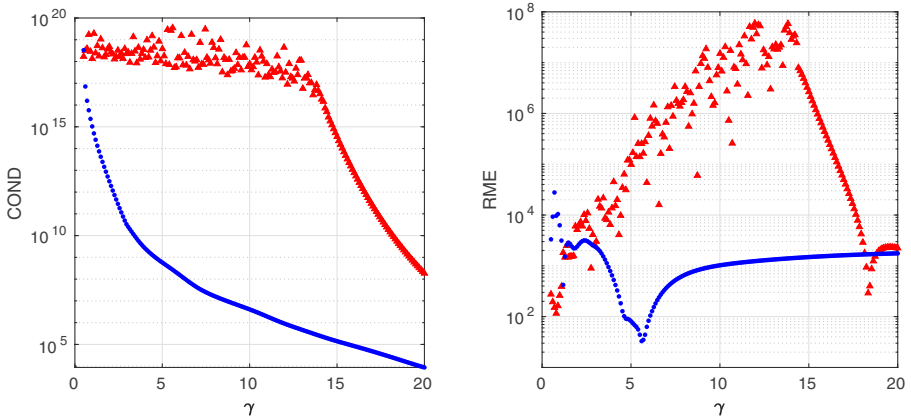
$$\text{RME} = \begin{cases} \lfloor \text{ME} \rfloor, & \text{if } \text{ME} - \lfloor \text{ME} \rfloor \leq 0.5, \\ \lceil \text{ME} \rceil, & \text{if } \text{ME} - \lfloor \text{ME} \rfloor > 0.5. \end{cases}$$

In the first experiments, we set the parameter  $\nu = 0$ . We remark that for regression networks the selection of the shape parameter plays a crucial role. Therefore, to make a fair comparison between classical and VSK regression networks, we report the condition numbers and the RME for 200 values of the shape parameter  $\gamma$  in the interval  $[0.5, 20]$ . The results are reported in Fig. 2. We observe that the computation carried out via VSKs is characterized by a lower condition number of the kernel matrix, as theoretically observed in Proposition 1. For such experiment, this directly reflects on the accuracy of the computation, meaning that the *safe* interval for the shape parameter  $\gamma$  is larger than for the classical method (see Fig. 2, right).

In Fig. 3, we report two graphical results corresponding to  $\nu = 0$  and  $\nu = 1e - 04$ , left and right respectively. In both cases we take the *optimal* shape parameter  $\gamma^*$ , meaning that it leads to the smallest RME, in the same framework of Fig. 2 (right). The associated RME is shown in Table 3. We note that the VSK setting outperforms the classical method for  $\nu = 0$ , while for  $\nu = 1e - 04$  the two approximations are comparable.

## 7 A VSK-like feature extraction algorithm

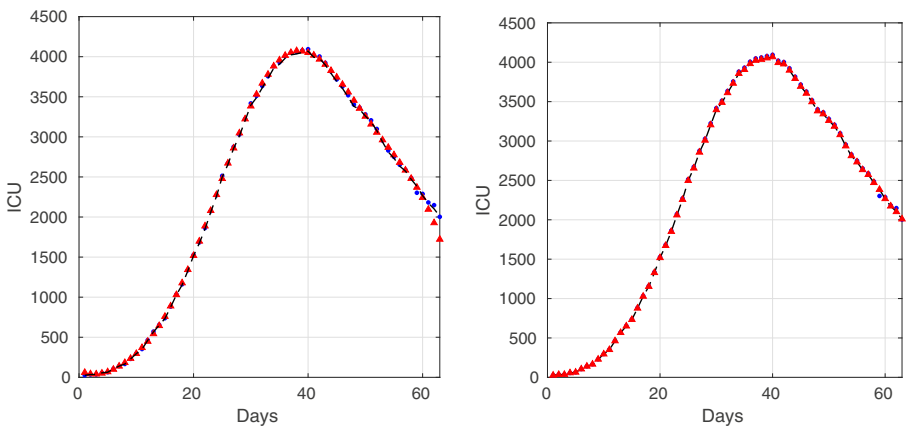
In this section, we propose a feature extraction method directly inspired by the presented variably scaled setting, which can be used as an alternative to other possible expensive feature extraction algorithms.



**Fig. 2** Left: The condition numbers for different values of the shape parameter of the classical kernel and VSK matrix denoted by red triangles and blue dots, respectively. Right: The RME for different values of the shape parameter of the classical KRN and KRN-VSK methods denoted by red triangles and blue dots, respectively. Both plots are in semi-logarithmic scale and obtained by considering the normalized dataset

To this aim, we consider the Wisconsin Breast Cancer Database [22, 23], which consists of 699 instances described by 9 features, extracted from a digitized image of a fine needle aspirate of a breast mass. The task consists in predicting if the mass is benign or malignant. From the original dataset, we exclude 16 instances that present missing values. The two classes are not equally distributed, presenting 444 benign instances and 239 malignant instances.

At first, we divide the dataset into a training set, consisting of 226 benign and 116 malignant cases, and a test set, which is composed of 218 benign and 123 malignant cases.



**Fig. 3** The ICU patients' curves reconstructed via KRN and KRN-VSK denoted by red triangles and blue dots, respectively. We fix  $\nu = 0$  and  $\nu = 1e - 04$ , left and right respectively. The true solution is plotted with a black solid line

**Table 3** The RME for the optimal shape parameter by using KRN and KRN-VSK in reconstructing the ICU curves

Method	$\nu$	
	$\nu = 0$	$\nu = 1e - 04$
KRN	116	13
KRN-VSK	33	9

Then, taking the hyperparameter grids adopted in Section 6.1, we compare the performances on the test set of the following four methods.

1. A NB classifier with Gaussian likelihood.
2. A standard SVM classifier, whose hyperparameters  $\zeta$  and  $\gamma$  (in the Gaussian case) are validated by means of 5-fold cross-validation on the training set.
3. A SVM classifier constructed after a feature selection process, as explained in what follows.

Analyzing the resulting weights of the SVM classifier (in the linear case), we can rank the features by their influence in the classification; see, e.g., [17]. Then, we choose the  $\bar{n}$  more relevant features, here we fix  $\bar{n} = 2$ , and we consequently reduce our training and test sets by restricting to the two most relevant features. Finally, we take both linear and Gaussian kernels, we train a SVM classifier via 5-fold cross-validation on the reduced training set and we evaluate the results on the reduced test set.

We denote this method with SVM-Selection (SVM-S).

4. A SVM classifier constructed after a VSK-like feature extraction process, as described in the following lines.

We randomly select  $\bar{n} - 1$  features (here  $\bar{n} = 2$ ). The training set restricted to the remaining 8 features is used to train a Gaussian NB classifier. Reduced training and test sets are obtained by juxtaposing the previously selected  $\bar{n} - 1$  features to the probabilistic output of the NB classifier. Then, we take both linear and Gaussian kernels, we train a SVM classifier via 5-fold cross-validation on the reduced training set and we evaluate the results on the reduced test set.

We denote this method with SVM-Extraction (SVM-E).

**Table 4** The  $f_1$ -score for the Wisconsin Breast Cancer Database via the SVM, NB, and SVM-S methods

	Linear		Gaussian	
	SVM	SVM-S	SVM	SVM-S
0.965	0.968	0.959	0.965	0.953

**Table 5** The  $f_1$ -score for the Wisconsin Breast Cancer Database via the SVM-E method

Random feature	Linear	Gaussian
1	0.965	0.965
2	0.962	0.968
3	0.959	0.977
4	0.962	0.965
5	0.965	0.965
6	0.965	0.962
7	0.965	0.962
8	0.959	0.956
9	0.968	0.965

We point out that both SVM-S and SVM-E consider reduced training and test sets that are characterized by the same number of features  $\bar{n}$ . Moreover, the SVM-E presents some advantages in terms of computational complexity with respect to SVM-S, since training an auxiliary NB classifier to perform feature extraction is cheaper than training a SVM classifier to carry out the feature selection.

In Table 4, we present the results obtained considering the SVM, NB, and SVM-S methods. In Table 5, we report the results concerning the SVM-E algorithm. For completeness, we vary the randomly selected feature, taking into account all the possibilities.

We observe that the best score is achieved by the SVM-E algorithm. Moreover for this dataset, we point out that such a method prefers the Gaussian kernel with respect to the linear one, while the standard SVM and SVM-S obtain better classification scores when the linear kernel is considered.

## 8 Conclusions and future work

We investigated the link between VSKs and feature augmentation strategies. In doing so, we tailored the VSKs for SVM and KRN. The proposed methods turn out to be flexible and easy to implement. For KRN, the use of VSKs takes advantage of being stable and for classification of merging the probabilistic features of NB and the geometric ones of SVM. This results in effective algorithms that can be used for many tasks. Applications to real datasets show the effectiveness of our approach.

Work in progress consists in extending this concept for support vector regression and as well as for greedy methods [2, 46].

**Acknowledgements** We sincerely thank the reviewers for helping us to significantly improve the manuscript. This research has been accomplished within Rete Italiana di Approssimazione (RITA) and partially funded by GNCS-INdAM, by the European Union's Horizon 2020 research and innovation programme ERA-PLANET, grant agreement no. 689443, via the GEOEssential project and by the ASI - INAF grant "Artificial Intelligence for the analysis of solar FLARES data (AI-FLARES)."

## References

1. Aggarwal, C.C.: *Data Classification: Algorithms and Applications*, Boca Raton, FL, USA CRC Press (2014)
2. Aminian Shahrokhbadi, M., Neisy, A., Perracchione, E., Polato, M.: Learning with subsampled kernel-based methods: Environmental and financial applications. *Dolomites Res. Notes Approx.* **12**, 17–27 (2019)
3. Ball, K.: Eigenvalues of Euclidean distance matrices. *J. Approx. Theory* **68**, 74–82 (1992)
4. Bartlett, P.L., Mendelson, S.: Rademacher and Gaussian complexities: risk bounds and structural results. *J. Mach. Learn. Res.* **3**, 463–482 (2002)
5. Bhatia, R.: *Matrix Analysis*. Springer-Verlag, New York (1997)
6. Bozzini, M., Lenarduzzi, L., Rossini, M., Schaback, R.: Interpolation with variably scaled kernels. *IMA J. Numer. Anal.* **35**, 199–219 (2015)
7. Bui, D.T., Pradhan, B., Lofman, O., Revhau, I.: Landslide susceptibility assessment in Vietnam using support vector machines, decision tree and Naïve Bayes models. *Math. Probl. Eng.* 1–26
8. Campagna, R., Conti, C., Cuomo, S.: Smoothing exponential-polynomial splines for multiexponential decay data. *Dolomites Res. Notes Approx.* **12**, 86–100 (2019)
9. Daumé, H.: Frustratingly easy domain adaptation. In: *Association for computational linguistics (ACL)* (2007)
10. De Marchi, S., Erb, W., Marchetti, F., Perracchione, E., Rossini, M.: Shape-Driven Interpolation with Discontinuous Kernels: Error Analysis, Edge Extraction and Applications in MPI. *SIAM J. Sci. Comput.* **42**, B472–B491 (2020)
11. De Marchi, S., Marchetti, F., Perracchione, E.: Jumping with variably scaled discontinuous kernels (VSDKs). *BIT Num. Math.* **60**, 441–463 (2020)
12. Diederichs, B., Iske, A.: Improved estimates for condition numbers of radial basis function interpolation matrices. *J. Approx. Theory* **238**, 38–51 (2019)
13. Donini, M., Aioli, F.: Learning deep kernels in the space of dot product polynomials. *Mach. Learn.* **106**, 1245–1269 (2017)
14. El Karoui, N.: The spectrum of kernel random matrices. *Ann. Statist.* **38**, 1–50 (2010)
15. Fasshauer, G.E.: *Meshfree Approximations Methods with Matlab*. World Scientific, Singapore (2007)
16. Fasshauer, G.E., McCourt, M.J.: *Kernel-based Approximation Methods Using Matlab*. World Scientific, Singapore (2015)
17. Hoffmann, H.: Kernel PCA for novelty detection. *Pattern Recogn.* **40**, 863–874 (2007)
18. Horn, R.A., Zhang, F.: Bounds on the spectral radius of a Hadamard product of nonnegative or positive semidefinite matrices. *Electron J. Linear Algebra* **20**, 90–94 (2010)
19. Kim, K.I., Jung, K., Kim, H.J.: Face recognition using kernel principal component analysis. *IEEE Signal Proc. Lett.* **9**, 40–42 (2002)
20. Larsson, E., Fornberg, B.: Theoretical and computational aspects of multivariate interpolation with increasingly flat radial basis functions. *Comput Math. Appl.* **49**, 103–130 (2005)
21. Li, W., Duan, L., Xu, D., Tsang, I.W.: Learning with augmented features for supervised and semi-supervised heterogeneous domain adaptation. *IEEE Trans. Pattern. Anal. Mach. Intell.* **36**, 1134–1148 (2014)
22. Mangasarian, O.L., Nick Street, W., Wolberg, W.H.: Wisconsin breast cancer database, UCI machine learning repository. <http://archive.ics.uci.edu/ml> University of Wisconsin (1991)
23. Mangasarian, O.L., Wolberg, W.H.: Cancer diagnosis via linear programming. *SIAM News* **106**, 1–18 (1990)
24. Marchetti, F.: The extension of Rippa’s algorithm beyond LOOCV. *Appl. Math. Lett.* **120**, 107262 (2021)
25. Maron, M.E.: Automatic indexing: An experimental inquiry. *J. ACM.* **8**, 404–417 (1961)
26. Mercer, J.: Functions of positive and negative type and their connection with the theory of integral equations. *Phil. Trans. Royal Society* **209**, 415–446 (1909)
27. Muquet, B., Wang, Z., Giannakis, G.B., De Courville, M., Duhamel, P.: Cyclic prefixing or zero padding for wireless multicarrier transmissions? *IEEE Trans. Commun.* **50**(12), 2136–2148 (2002)
28. Narcowich, F.J., Ward, J.F.: Norm estimates for the inverses of a general class of scattered-data radial-function interpolation matrices. *J. Approx. Theory* **69**, 84–109 (1992)
29. Nocedal, J., Wright, S.J.: *Numerical Optimization*. Springer-Verlag, New York (1999)

30. Orr, M.J.L.: Introduction to radial basis function networks, Tech. rep., University of Edinburgh Centre for Cognitive Sciences (1996)
31. Pang, B., Lee, B., Vaithyanathan, S.: Thumbs up? Sentiment Classification Using Machine Learning Techniques. In: Proc. of EMNLP, pp. 79–86 (2002)
32. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, E.: Scikit-learn: Machine learning in python. *J. Mach. Learn. Res.* **12**, 2825–2830 (2011)
33. Ren, J., Lee, S.D., Chen, X., Kao, B., Cheng, R., Cheung, D.: Naive Bayes classification of uncertain data. In: Proc. 9th IEEE Int. Conf. Data Mining (ICDM), pp. 944–949 (2009)
34. Rippa, S.: An algorithm for selecting a good value for the parameter  $c$  in radial basis function interpolation. *Adv. Comput. Math.* **11**, 193–210 (1999)
35. Romani, L., Rossini, M., Schenone, D.: Edge detection methods based on RBF interpolation. *J. Comput. Appl. Math.* **349**, 532–547 (2019)
36. Rossini, M.: Interpolating functions with gradient discontinuities via variably scaled kernels. *Dolom. Res. Notes Approx.* **11**, 3–14 (2018)
37. Schaback, R.: Error estimates and condition numbers for radial basis function interpolation. *Adv. Comput. Math.* **3**, 251–264 (1995)
38. Schaback, R. et al.: Multivariate interpolation and approximation by translates of a basis function. In: Chui, C. (ed.) *Approximation Theory VIII: Approximation and Interpolation*, pp. 491–514. World Scientific, Singapore (1995)
39. Schölkopf, B., Smola, A.J.: *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. MIT Press, Cambridge (2002)
40. Schur, J.: Bemerkungen zur Theorie der beschränkten Bilinearformen mit unendlich vielen Veränderlichen. *J. Reine Angew. Math.* **140**, 1–28 (1911)
41. Stura, I., Gabriele, D., Guiot, C.: A simple PSA-based computational approach predicts the timing of cancer relapse in prostatectomized patients. *Cancer Res.* **76**, 4941–4947 (2016)
42. Shawe-Taylor, J., Cristianini, N.: *Kernel Methods for Pattern Analysis*. Cambridge University Press, Cambridge (2004)
43. Styan, G.P.H.: Hadamard products and multivariate statistical analysis. *Linear Algebra Appl.* **6**, 217–240 (1973)
44. Vapnik, V.N., Chervonenkis, A.Y.: On the uniform convergence of relative frequencies of events to their probabilities. *Theory Probab. Appl.* **16**, 264–280 (1971)
45. Weston, J., Bengio, S., Usunier, N.: Large scale image annotation: Learning to rank with joint word-image embeddings. *Mach. Learn.* **81**, 21–35 (2010)
46. Wirtz, D., Haasdonk, B.: A vectorial kernel orthogonal greedy algorithm. *Dolomites Res. Notes Approx.* **6**, 83–100 (2013)
47. Zhang, D., Lee, W.S.: Question classification using support vector machines. In: *Proceedings of the 26th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, New York, NY, 26–32 (2003)

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.