

Learning via variably scaled kernels

*Original*

Learning via variably scaled kernels / Campi, C.; Marchetti, F.; Perracchione, E.. - In: ADVANCES IN COMPUTATIONAL MATHEMATICS. - ISSN 1019-7168. - 47:4(2021). [10.1007/s10444-021-09875-6]

*Availability:*

This version is available at: 11583/2987556 since: 2024-04-18T07:27:04Z

*Publisher:*

Springer

*Published*

DOI:10.1007/s10444-021-09875-6

*Terms of use:*

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

*Publisher copyright*

Springer postprint/Author's Accepted Manuscript

This version of the article has been accepted for publication, after peer review (when applicable) and is subject to Springer Nature's AM terms of use, but is not the Version of Record and does not reflect post-acceptance improvements, or any corrections. The Version of Record is available online at: <http://dx.doi.org/10.1007/s10444-021-09875-6>

(Article begins on next page)

# Learning via variably scaled kernels

C. Campi\*, F. Marchetti<sup>+</sup>, E. Perracchione<sup>°</sup>

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

<sup>\*</sup> Dipartimento di Salute della Donna e del Bambino,  
Università di Padova,

<sup>°</sup> Dipartimento di Matematica DIMA, Università di Genova,

`cristina.campi@unipd.it`  
`francesco.marchetti.1@phd.unipd.it`  
`perracchione@dim.unige.it`

## 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 COro-naVirus Disease 19 (COVID19) 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 can be seen as an alternative to the sometimes computationally demanding feature extraction procedures.

## 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. In this sense, they can be seen as a generalisation of *feature augmentation strategies*. Indeed, all methods based on feature augmentation, as e.g. zero padding and feature replication [9, 21, 26], 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, 41]), 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, 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 meaningful 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, 35]. 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 need to take into account also probabilistic solutions. 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, 30, 46]. 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. In view of this, our method can also be *substituted* to feature extraction schemes; refer e.g. to [19, 44]. Indeed, we numerically show that we are able to randomly reduce the dataset and then encode the missing information via SVM-VSK.

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 real datasets. 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 [25]).

**Theorem 2.1** *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.*

For such kernels that admit a Mercer expansion (also called valid kernels according to the definition given by [41]), 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,

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

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, 41] 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 [43] 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 [38].

To better investigate the concept of expressiveness in the kernel setting, we introduce the kernel matrix  $\mathbf{K}$  constructed via the dataset  $\Xi = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} \subseteq \Omega$ , i.e. the matrix of entries

$$K_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}_j), \quad i, j = 1, \dots, N, \quad (2.2)$$

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

**Remark 2.1** *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 2.2** *Let  $\mathbf{K}$  be the kernel matrix as in (2.2) constructed via  $\Xi = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} \subseteq \Omega$ , let us denote by  $\text{rank}(\mathbf{K})$  its rank. Then, there exists at least one subset of examples of size  $\text{rank}(\mathbf{K})$  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

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

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 2.1** *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  $\Xi = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} \subseteq \Omega$ , we have*

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

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

Being the spectral ratio an expressiveness measure, it is related to the rank of the kernel matrix (see also Remark 2.1), indeed

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

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 2.2** Let  $\Lambda \subseteq \mathbb{R}$  be a bounded set. Let  $\kappa : \tilde{\Omega} \times \tilde{\Omega} \rightarrow \mathbb{R}$ ,  $\tilde{\Omega} = \Omega \times \Lambda \subseteq \mathbb{R}^{n+1}$ , 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

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

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

When dealing with Mercer's kernels, the construction of a VSK as in Definition 2.2 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.

At first, we give a multidimensional extension of the scaling function. Thus, let  $\Lambda \subseteq \mathbb{R}^m$ ,  $m > 0 \in \mathbb{N}$  be a bounded set. Given  $\mathbf{x} \in \Omega$ , let  $\psi : \Omega \rightarrow \Lambda$  be a function that extracts  $m$  features from  $\mathbf{x}$ . Then, letting  $\tilde{\Omega} = \Omega \times \Lambda$ , we can 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.

Thus, in the VSK setting, we define the (strictly) positive definite kernel  $\kappa^\Psi : \Omega \times \Omega \rightarrow \mathbb{R}$ , given by

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

where  $\kappa : \tilde{\Omega} \times \tilde{\Omega} \rightarrow \mathbb{R}$  is a (strictly) positive definite kernel. The kernel  $\kappa^\Psi$  is a valid kernel, namely that it corresponds to an inner product in the associated feature space  $F_\psi$  (see [41, Proposition 3.22, page 75]). Moreover, it induces a *new* feature map  $\Theta : \Omega \rightarrow F_\psi$  so that

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

Referring to equation (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, 41].

We take  $\Xi = \{\mathbf{x}_i, i = 1, \dots, N\} \subseteq \Omega$ , where  $\Omega \subseteq \mathbb{R}^n$  is a bounded set. 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  $\mathbf{x} \in \Omega$ , we define a non-linear SVM classifier that makes use of VSKs via the following decision function:

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

where  $\Theta : \Omega \rightarrow F_\psi$  is the VSK feature map,  $\mathbf{w} = (w_1, \dots, w_n) \in \mathbb{R}^n$  and  $b \in \mathbb{R}$  are given by

$$\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \Theta(\mathbf{x}_i),$$

and

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

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

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

where  $\zeta \in \mathbb{R}_+ = [0, +\infty)$  is known as *bounding box*. 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. [28]) and thanks to (3.1), for  $\mathbf{x} \in \Omega$ , it reads as follows

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

Note that in equation (3.2)  $i$  denotes the index of an  $\alpha_i$  so that  $0 < \alpha_i < \zeta$ . For stability purposes,  $b$  is usually evaluated as an average among all candidates. 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 machine learning context, we investigate regression networks.

### 3.2 KRN-VSK

For the purposes of KRN, given  $\Xi = \{\mathbf{x}_i, i = 1, \dots, N\} \subseteq \Omega$ , where  $\Omega \subseteq \mathbb{R}^n$  is a bounded set, we fix the output variables  $y_i \in \mathbb{R}$ ,  $i = 1, \dots, N$ . Indeed, KRNs

are used for regression/interpolation purposes, hence, e.g. they turn out to be meaningful for studying the behaviour of longitudinal data.

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 KRNs, we refer the reader to [16, 29].

We keep the general framework of KRNs 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

$$s^\Psi(\mathbf{x}) = \sum_{i=1}^M c_i \kappa^\Psi(\mathbf{x}, \mathbf{z}_i), \quad (3.3)$$

for (strictly) positive definite kernels  $\kappa^\Psi : \Omega \times \Omega \rightarrow \mathbb{R}$  and for some real coefficients  $c_1, \dots, c_M$ . The learning function (3.3) is thus the simplest neural network which involves only one linear combination of basis functions, i.e. one layer.

For KRN-VSK, we compute  $\mathbf{c} = (c_1, \dots, c_M) \in \mathbb{R}^M$  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 centres  $Z \equiv \Xi$ . 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 to the special setting for which  $\nu = 0$ . In that case, the solution can be found as  $\mathbf{c}^\top = (\mathbf{K}^\Psi)^{-1} \mathbf{y}^\top$ , where  $\mathbf{y} = (y_1, \dots, y_N)$  and  $\mathbf{c} = (c_1, \dots, c_N)$ .

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, 33]. Moreover, we will point out both numerically and theoretically in the incoming sections that a performing alternative to reduce the ill-conditioning comes from the use of VSKs.

## 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)  $\varphi : [0, \infty) \rightarrow \mathbb{R}$  and (possibly) a shape parameter  $\gamma > 0$  such that, for all  $\mathbf{x}, \mathbf{y} \in \Omega$ ,

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

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

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

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; see e.g. [16].

*Practical implementation for the Gaussian VSK*

Throughout this section we take  $N$  distinct data  $\Xi = \{\mathbf{x}_i, i = 1, \dots, N\} \subseteq \Omega$ , where  $\Omega \subseteq \mathbb{R}^n$  is a bounded set. Moreover, let  $\Lambda \subseteq \mathbb{R}^m$  be a bounded set.

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

**Theorem 4.1** *Let  $\Xi = \{\mathbf{x}_i, i = 1, \dots, N\} \subseteq \Omega$  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. Then, the VSK matrix constructed on  $\Xi$  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  $\Xi = \{\mathbf{x}_i, i = 1, \dots, N\}$  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 [39]. It will be helpful in what follows; refer also to [14, Lemma A.5] and [18, Lemma 2.1].

**Theorem 4.2** *If  $\mathbf{E}$  and  $\mathbf{M} \in \mathbb{R}^{N \times N}$  are positive definite matrices, denoting by  $\lambda_{\min}$  and  $\lambda_{\max}$  the smallest and largest eigenvalue of a matrix, we have that*

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

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  $\mathbf{K}$  in terms of its variably scaled setting. This turns out to be helpful when comparing the spectral ratio of the two matrices ( $\mathbf{K}$  and  $\mathbf{K}^\Psi$ ).

**Theorem 4.3** *Let  $\Xi = \{\mathbf{x}_i, i = 1, \dots, N\} \subseteq \Omega$  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  $\mathbf{K}^\Psi = \mathbf{K} \circ \mathbf{K}^\psi$  constructed on  $\Xi$  via  $\kappa^\Psi : \Omega \times \Omega \rightarrow \mathbb{R}$ , we have that*

$$\|\mathbf{K}^\Psi\|_F \leq \|\mathbf{K}\|_F \leq \|\mathbf{K}\|_F \|\mathbf{K}^\psi\|_F.$$

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

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

which in particular implies that

$$\mathbf{K}_{ij} \geq \mathbf{K}_{ij}^\Psi \geq 0, \quad i, j = 1, \dots, N.$$

Thus, we get

$$\|\mathbf{K}\|_F \geq \|\mathbf{K}^\Psi\|_F.$$

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

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

and therefore

$$\|\mathbf{K}^\Psi\|_F \leq \|\mathbf{K}\|_F \leq \|\mathbf{K}\|_F \|\mathbf{K}^\psi\|_F.$$

□

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

**Corollary 4.1** *Let  $\Xi = \{\mathbf{x}_i, i = 1, \dots, N\} \subseteq \Omega$  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, then the VSK kernel  $\kappa^\Psi : \Omega \times \Omega \rightarrow \mathbb{R}$  is more expressive than  $\kappa$ .*

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

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

where  $\varphi : \Omega \longrightarrow \mathbb{R}$  is the RBF associated to the Gaussian kernel  $\kappa : \Omega \times \Omega \longrightarrow \mathbb{R}$ . Taking into account Theorem 4.3, we obtain

$$S(\mathbf{K}) = \frac{N}{\|\mathbf{K}\|_{\text{F}}} \leq \frac{N}{\|\mathbf{K}^\Psi\|_{\text{F}}} = S(\mathbf{K}^\Psi).$$

□

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  $\Xi = \{\mathbf{x}_i, i = 1, \dots, N\} \subseteq \Omega$ , the stability is also related to the separation distance

$$q_\Xi := \frac{1}{2} \min_{i \neq j} \|\mathbf{x}_i - \mathbf{x}_j\|_2,$$

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

$$q_\Xi \leq q_\Xi^\Psi,$$

where

$$q_\Xi^\Psi := \frac{1}{2} \min_{i \neq j} \|\Psi(\mathbf{x}_i) - \Psi(\mathbf{x}_j)\|_2,$$

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 [27], 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* [36, 37], 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 [42, Corollary 3.1]. For a given matrix  $\mathbf{M}$ , we focus on the 2-condition number defined as

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

**Proposition 4.1** *Let  $\Xi = \{\mathbf{x}_i, i = 1, \dots, N\} \subseteq \Omega$  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  $\mathbf{K}^\Psi = \mathbf{K} \circ \mathbf{K}^\psi$  constructed on  $\Xi$  via  $\kappa^\Psi : \Omega \times \Omega \rightarrow \mathbb{R}$ , we have that*

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

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

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

Moreover, from Theorem 4.2 and since the RBF  $\varphi : \Omega \rightarrow \mathbb{R}$  associated to the Gaussian kernel  $\kappa : \Omega \times \Omega \rightarrow \mathbb{R}$  is so that  $\varphi(0) = 1$ , i.e.  $\mathbf{K}_{ii}^\psi = 1, i = 1, \dots, N$ , we obtain

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

□

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

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

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 semi-definite; see e.g. [16].

*Practical implementation for the linear VSK*

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

**Theorem 4.4** *Let  $\Xi = \{\mathbf{x}_i, i = 1, \dots, N\} \subseteq \Omega$  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 linear kernel. Then, the VSK matrix constructed on  $\Xi$  via  $\kappa^\Psi : \Omega \times \Omega \rightarrow \mathbb{R}$  is given by*

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

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

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

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

and thus the kernel matrix is given by

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

□

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 4.2** *Let  $\Xi = \{\mathbf{x}_i, i = 1, \dots, N\} \subseteq \Omega$  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 linear kernel. Let us suppose that the associated kernel matrix  $\mathbf{K}$  is non-negative, i.e. so that all the entries of  $\mathbf{K}$  are non-negative. Given the VSK matrix  $\mathbf{K}^\Psi = \mathbf{K} + \mathbf{K}^\psi$  constructed on  $\Xi$  via  $\kappa^\Psi : \Omega \times \Omega \rightarrow \mathbb{R}$ , if  $\psi$  is so that  $\mathbf{K}^\psi$  is non-negative, then:*

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

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

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

Moreover, since we suppose  $\mathbf{K}$  non-negative, we get

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

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

Note that the requirements of Proposition 4.2 are satisfied e.g. if  $\Omega \subseteq \mathbb{R}_+^n$  and  $\Lambda \subseteq \mathbb{R}_+^m$ .

*Spectrum of the linear VSK*

Being Gramiam matrices,  $\mathbf{K}^\Psi$  and  $\mathbf{K}^\psi$  are positive semi-definite. Concerning the spectrum of the VSK matrix  $\mathbf{K}^\Psi$ , by virtue of Weyl's inequality (see e.g. [5, Section III.2, p. 62]), we obtain that:

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

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

## 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. [34, 35]). 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  $N$  distinct data  $\Xi = \{\mathbf{x}_i, i = 1, \dots, N\} \subseteq \Omega$ , where  $\Omega \subseteq \mathbb{R}^n$  is a bounded set. Moreover, let  $\Lambda \subseteq \mathbb{R}^m$  be a bounded set, 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  $\psi : \Omega \longrightarrow \Lambda$  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, 24]) 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  $\Psi : \Omega \longrightarrow \tilde{\Omega}$  defined by

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

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

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

For  $\mathbf{x} \in \Omega$ , since  $P_2(\mathbf{x}) = 1 - P_1(\mathbf{x})$  is correlated to  $P_1(\mathbf{x})$ , we observe that it is sufficient to consider one of the two probabilities.

Concerning the effectiveness of this scaling function  $\Psi : \Omega \longrightarrow \tilde{\Omega}$  for the Gaussian VSK, we refer to the notation introduced in Theorem 4.1 and we point out that, for  $\mathbf{x}_i, \mathbf{x}_j \in \Xi$ ,

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

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

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

We remark that, according to Proposition 4.2, 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  $\Xi = \{\mathbf{x}_i, i = 1, \dots, N\} \subseteq \Omega$ , where  $\Omega \subseteq \mathbb{R}^n$  is a bounded set, and the associated measurements  $y_i \in \mathbb{R}$ ,  $i = 1, \dots, N$ . Moreover, let  $\Lambda \subseteq \mathbb{R}^m$  be a bounded set, we now propose some ideas to define the scaling function for KRN.

Regression networks are also used to learn longitudinal data, i.e. time series, see e.g. [8]. This allows to extrapolate the evolution of samples and consequently to give short time predictions on the dynamics of the considered process. Therefore, concerning the choice of the scaling function  $\psi : \Omega \longrightarrow \Lambda$ , we suppose to know the trend of data, which can be modelled via a specific class of functions, i.e. a model  $\mathcal{M} : \Omega \times \mathbb{R}^l \longrightarrow \mathbb{R}$  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  $\psi : \Omega \longrightarrow \Lambda$  is

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

Of course, this gives a recipe for the selection of the scaling function which is not unique. Moreover, thinking of time series, one usually disposes of other available and correlated data sampled at the same locations which can be used as additional features. For instance, when studying the evolution of a tumor mass in time, usually also other bio-markers are available and can be taken into account; see e.g. [40].

## 6 Numerical tests

Experiments have been carried out in PYTHON using also the scientific module scikit-learn [31] on a Intel(R) Core(TM) i7 CPU 4712MQ 2.13 GHz processor. For the classification setting, we provide a freely available software that can be downloaded at <https://github.com/emmaA89/SVM-VSK>. For KRN-VSK, we refer the reader to [16, Program 18.1, p. 340].

## 6.1 Tests for SVM-VSK

We present two examples of binary classification. In the first one, 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. For the second experiment, we take a real dataset and we show a feature extraction strategy inspired by our framework.

In both the examples, 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,

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

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

### *Test 1: Simulated datasets and a priori knowledge*

Here we construct 12 toy datasets that differ in terms of number of features and examples. We now fix  $n = 64$ . Letting  $\Omega \subseteq \mathbb{R}^n$ , where  $\Omega$  is a bounded set, 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  $C_2$  are exactly balanced. The construction of such 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 real measurement sampled from such distribution. Then,  $\mu_j^i$  and  $\sigma_j^i$ ,  $i = 1, \dots, n$ ,  $j = 1, 2$  are determined as follows:

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



2. We denote by  $\mathcal{N}(\mu, \sigma)$  the univariate normal distribution with mean  $\mu$  and standard deviation  $\sigma$ . Let  $\mathbf{x}_k = (x_k^1, \dots, x_k^n)$  be an example in  $\Omega$  belonging to a class  $C_j$ ,  $j = 1, 2$ . The elements  $x_k^i$  of  $\mathbf{x}_k \in \Omega$ ,  $i = 1, \dots, n$ , are then randomly generated as samples of  $\mathcal{N}(\mu_j^i, \sigma_j^i)$ .
3. Finally, Gaussian white noise, distributed according to  $\mathcal{N}(0, 1)$ , is added to each feature and example.

From the so-constructed  $\Gamma$ , letting  $\Omega_k$  be a bounded set of  $\mathbb{R}^{n_k}$ ,  $n_k \in \{2, 4, 16, 64\}$ , we extract the datasets

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

with  $N_q \in \{100, 500, 2500\}$ ,  $\boldsymbol{\xi} = (N_q, n_k)$ , and preserving the balance among the two classes. Fixed  $N_q$  and  $n_k$ , we point out that the examples and the features are randomly selected. 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  $\Gamma_{\boldsymbol{\xi}}$  for some value of  $N_q$  and  $n_k$ . We divide such a dataset in a training set  $\Sigma_{\boldsymbol{\xi}}$  and a test set  $T_{\boldsymbol{\xi}}$ . These sets are so that  $\text{card}(\Sigma_{\boldsymbol{\xi}}) \approx 2\text{card}(T_{\boldsymbol{\xi}})$ .

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  $\Sigma_{\boldsymbol{\xi}}$  and  $\bar{\Gamma}_{\boldsymbol{\xi}}$ , which is defined as the dataset containing the examples of  $\Gamma$  that are not in  $\Gamma_{\boldsymbol{\xi}}$ , i.e.

$$\bar{\Gamma}_{\boldsymbol{\xi}} := \Gamma \setminus \Gamma_{\boldsymbol{\xi}}.$$

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

1. The NB classifier, which is trained on  $\bar{\Gamma}_{\boldsymbol{\xi}} \cup \Sigma_{\boldsymbol{\xi}}$ . Given  $\mathbf{x} = (x_1, \dots, x_{n_k})$ , we adopt the Gaussian likelihood [32]

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

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

2. The standard SVM method, which is trained on  $\Sigma_{\boldsymbol{\xi}}$ .
3. The SVM-VSK classifier, which is trained on  $\Sigma_{\boldsymbol{\xi}}$  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  $\Sigma_{\boldsymbol{\xi}}$ .

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

	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

Figure 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.

For the Gaussian kernel, we numerically verify Corollary 4.1 by reporting in Table 1 the spectral ratios related to the matrices  $\mathbf{K}$  and  $\mathbf{K}^\Psi$ , obtained from the training sets  $\Sigma_{\mathbf{x}}$  with  $N_q = 100, 500, 2500$  and  $n_k = 2$ . The results numerically confirm what theoretically observed, i.e. the Gaussian VSK is more expressive than the standard one.

Moreover, for the linear kernel, we are in the hypothesis of Proposition 4.2 since in our experiments the data are normalized in  $[0, 1]$ . In order to numerically verify its consistency, we also evaluate the quantities involved in such proposition. They are reported in Table 2. The results show once more what theoretically observed.

N	$S(\mathbf{K})$	$S(\mathbf{K}^\Psi)$
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 1: The spectral ratios of the matrices  $\mathbf{K}$  and  $\mathbf{K}^\Psi$  related to the normalized training sets  $\Sigma_{\mathbf{x}}$ , varying  $N_q = 100, 500, 2500$ . We set  $n_k = 2$  and we considered a Gaussian kernel with  $\gamma = 1$ .

N	$\text{tr}(\mathbf{K})/\text{tr}(\mathbf{K}^\Psi)$	$S(\mathbf{K})/S(\mathbf{K}^\Psi)$	$\ \mathbf{K}^\Psi\ _F/\ \mathbf{K}\ _F$
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

Table 2: The ratios of the norms involved in Proposition 4.2 obtained via the linear kernel. The matrices  $\mathbf{K}$  and  $\mathbf{K}^\Psi$  are related to the normalized training sets  $\Sigma_{\xi}$ , varying  $N_q = 100, 500, 2500$  and with  $n_k = 2$ .

*Test 2: Feature selection with a real dataset*

The Wisconsin Breast Cancer Database [22, 23] 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.

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. To this aim, 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 by 218 benign and 123 malignant cases.

Then, 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).

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

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

Table 3: The  $f_1$ -score for the Wisconsin Breast Cancer Database via the SVM, NB and SVM-S methods.

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

Table 4: The  $f_1$ -score for the Wisconsin Breast Cancer Database via the SVM-E method.

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.

## 6.2 Tests for KRN-VSK

As an example for KRN, we focus on the Italian data of the 2020 COVID19 pandemic. The task we consider consists in learning the longitudinal data, i.e.  $\Omega \subseteq \mathbb{R}$ , of people that in Italy were hospitalized as Intensive Care Unit (ICU) patients from 24/02/2020 to 26/04/2020. 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  $\Gamma$  consists of 63 samples and it is divided as follows. The first 58 days are used as training set  $\Sigma$  and we test the model on the last  $t = 5$  days,  $\tilde{x}_i$ ,  $i = 1, \dots, t$ . Referring to Subsection 3.2 we take the set of kernel centres  $Z$  as the set of available data in  $\Xi$  and we focus on the Gaussian kernel. The feature augmentation strategy outlined in (5.1) is carried out considering  $\mathcal{M} : \Omega \times \mathbb{R} \rightarrow \mathbb{R}$  defined as

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

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

Furthermore, we encode into the kernel also other available data. Specifically, we take the total number of COVID19 infected (included death and recovered people), the daily number of new infected and the total number 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). Precisely, since hospitalized patients are involved in the dynamics, letting the Mean Error

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

where  $A$  is a decision function as defined in Subsection 3.2 obtained via classical of variably scaled kernels, the RME is 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 Figure 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 4.1. For such experiment, this directly reflects on the accuracy of

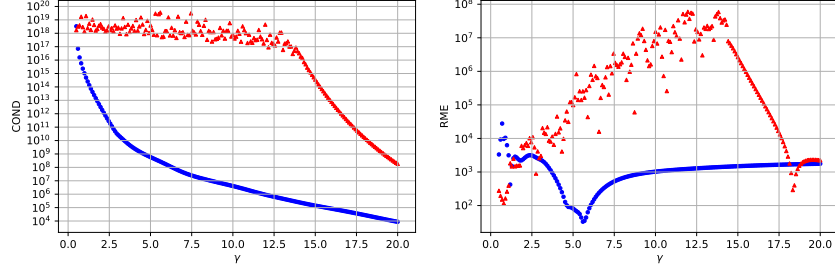


Figure 2: Left: the condition numbers for different values of the shape parameter of the classical kernel and VSK matrix denoted by red triangles and blue dots, respectively. Right: the RME for different values of the shape parameter of the classical KRN and KRN-VSK methods denoted by red triangles and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

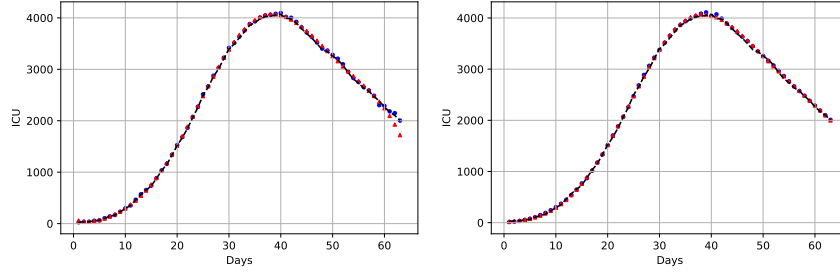


Figure 3: The ICU patients curves reconstructed via KRN and KRN-VSK denoted by red triangles and blue dots, respectively. We fix  $\nu = 0$  and  $\nu = 1e-04$ , left and right respectively. The true solution is plotted with black solid line.

the computation, meaning that the *safe* interval for the shape parameter  $\gamma$  is larger than for the classical method (see Figure 2, right).

In Figure 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 Figure 2 (right). The associated RME is shown in Table 5. We note that, the VSK setting outperforms the classical method for  $\nu = 0$ , while for  $\nu = 1e-04$  the two approximations are comparable.

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

Table 5: The RME for the optimal shape parameter by using KRN and KRN-VSK in reconstructing the ICU curves.

## 7 Conclusions and future work

We presented an original approach for learning issues via VSKs. 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, 45].

## 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] C.C. AGGARWAL, *Data Classification: Algorithms and Applications*, Boca Raton, FL, USA, CRC Press, 2014.
- [2] M. AMINIAN SHAHROKHABADI, A. NEISY, E. PERRACCHIONE, M. POLATO, *Learning with subsampled kernel-based methods: Environmental and financial applications*, Dolomites Res. Notes Approx. **12** (2019), 17–27.
- [3] K. BALL, *Eigenvalues of Euclidean distance matrices*, J. Approx. Theory **68** (1992), 74–82.
- [4] P.L. BARTLETT, S. MENDELSON *Rademacher and Gaussian complexities: risk bounds and structural results*, Journal of Machine Learning Research **3** (2002), 463–482.

- [5] R. BHATIA, *Matrix Analysis*, Springer-Verlag, New York, 1997.
- [6] M. BOZZINI, L. LENARDUZZI, M. ROSSINI, R. SCHABACK, *Interpolation with variably scaled kernels*, IMA J. Numer. Anal. **35** (2015), 199–219.
- [7] D.T. BUI, B. PRADHAN, O. LOFMAN, I. REVHAUG, *Landslide susceptibility assessment in Vietnam using support vector machines, decision tree and Naïve Bayes models*, Math. Probl. Eng. 1–26.
- [8] R. CAMPAGNA, C. CONTI, S. CUOMO, *Smoothing exponential-polynomial splines for multiexponential decay data*, Dolomites Res. Notes Approx. **12** (2019), 86–100.
- [9] H. DAUMÉ, *Frustratingly easy domain adaptation*. In: Association for computational linguistics (ACL), 2007.
- [10] S. DE MARCHI, W. ERB, F. MARCHETTI, E. PERRACCHIONE, M. ROSSINI, *Shape-Driven Interpolation with Discontinuous Kernels: Error Analysis, Edge Extraction and Applications in MPI*, SIAM J. Sci. Comput. **42** (2020), B472–B491.
- [11] S. DE MARCHI, F. MARCHETTI, E. PERRACCHIONE, *Jumping with Variably Scaled Discontinuous Kernels (VSDKs)*, to appear on BIT Numerical Mathematics, 2020, <https://doi.org/10.1007/s10543-020-00800-9>.
- [12] B. DIEDERICHS, A. ISKE, *Improved estimates for condition numbers of radial basis function interpolation matrices*, J. Approx. Theory **238** (2019), 38–51.
- [13] M. DONINI, F. AIOLLI, *Learning deep kernels in the space of dot product polynomials*, Machine Learning, **106** (2017), 1245–1269.
- [14] N. EL KAROUI, *The spectrum of kernel random matrices*, Ann. Statist. **38** (2010), 1–50.
- [15] G.E. FASSHAUER, *Meshfree Approximations Methods with MATLAB*, World Scientific, Singapore, 2007.
- [16] G.E. FASSHAUER, M.J. MCCOURT, *Kernel-based Approximation Methods Using MATLAB*, World Scientific, Singapore, 2015.
- [17] H. HOFFMANN, *Kernel PCA for novelty detection*, Pattern Recognition, **40** (2007), 863–874.
- [18] R.A. HORN, F. ZHANG, *Bounds on the spectral radius of a Hadamard product of nonnegative or positive semidefinite matrices*, Electron. J. Linear Algebra **20** (2010), 90–94.
- [19] K.I. KIM, K. JUNG, H.J. KIM, *Face recognition using kernel principal component analysis*, IEEE Signal Processing Letters, **9** (2002), 40–42.



- [20] E. LARSSON, B. FORNBERG, *Theoretical and computational aspects of multivariate interpolation with increasingly flat radial basis functions*, Comput. Math. Appl. **49** (2005), 103–130.
- [21] W. LI, L. DUAN, D. XU, I.W. TSANG, *Learning with augmented features for supervised and semi-supervised heterogeneous domain adaptation*, IEEE Trans. Pattern. Anal. Mach. Intell. **36** (2014), 1134–1148.
- [22] O.L. MANGASARIAN, W. NICK STREET, W.H. WOLBERG, *Wisconsin Breast Cancer Database*, UCI Machine Learning Repository, <http://archive.ics.uci.edu/ml>, University of Wisconsin, 1991.
- [23] O. L. MANGASARIAN, W. H. WOLBERG, *Cancer diagnosis via linear programming*, SIAM News, **106** (1990), 1–18.
- [24] M.E. MARON, *Automatic Indexing: An Experimental Inquiry*, Journal of the ACM. **8** (1961), 404–417.
- [25] J. MERCER, *Functions of positive and negative type and their connection with the theory of integral equations*, Phil. Trans. Royal Society **209** (1909), 415–446.
- [26] B. MUQUET, Z. WANG, G. B. GIANNAKIS, M. DE COURVILLE, P. DUHAMEL, *Cyclic prefixing or zero padding for wireless multicarrier transmissions?* IEEE Transactions on Communications, vol. 50, no. 12, 2136–2148, 2002.
- [27] F.J. NARCOWICH, J.F. WARD, *Norm estimates for the inverses of a general class of scattered-data radial-function interpolation matrices*, J. Approx. Theory **69** (1992), 84–109.
- [28] J. NOCEDAL, S.J. WRIGHT, *Numerical Optimization*, Springer-Verlag, New York, 1999.
- [29] M.J.L. ORR, *Introduction to radial basis function networks*, Tech. rep., University of Edinburgh, Centre for Cognitive Sciences, 1996.
- [30] B. PANG, B. LEE, S. VAITHYANATHAN, *Thumbs up? Sentiment Classification Using Machine Learning Techniques*, In Proc. of EMNLP 2002, 79–86.
- [31] F. PEDREGOSA, G. VAROQUAUX, A. GRAMFORT, V. MICHEL, B. THIRION, O. GRISEL, M. BLONDEL, P. PRETTENHOFER, R. WEISS, V. DUBOURG, J. VANDERPLAS, A. PASSOS, D. COURNAPEAU, M. BRUCHER, M. PERROT, E. DUCHESNAY, *Scikit-learn: Machine Learning in Python*, Journal of Machine Learning Research, **12** (2011) 2825–2830.
- [32] J. REN, S. D. LEE, X. CHEN, B. KAO, R. CHENG, D. CHEUNG, *Naive Bayes classification of uncertain data*, in Proc. 9th IEEE Int. Conf. Data Mining (ICDM), 2009, 944–949.