

# Input Space vs. Feature Space in Kernel-Based Methods

Bernhard Schölkopf, Sebastian Mika, Chris J. C. Burges,  
Philipp Knirsch, Klaus-Robert Müller, Gunnar Rätsch, Alex J. Smola

*Abstract*— This paper collects some ideas targeted at advancing our understanding of the feature spaces associated with Support Vector (SV) kernel functions. We first discuss the geometry of feature space. In particular, we review what is known about the shape of the image of input space under the feature space map, and how this influences the capacity of SV methods. Following this, we describe how the metric governing the intrinsic geometry of the mapped surface can be computed in terms of the kernel, using the example of the class of inhomogeneous polynomial kernels, which are often used in SV pattern recognition.

We then discuss the connection between feature space and input space by dealing with the question of how one can, given some vector in feature space, find a pre-image (exact or approximate) in input space. We describe algorithms to tackle this issue, and show their utility in two applications of kernel methods. First, we use it to reduce the computational complexity of SV decision functions; second, we combine it with the Kernel PCA algorithm, thereby constructing a nonlinear statistical denoising technique which is shown to perform well on real-world data.

*Keywords*— Kernel methods, support vector machines, PCA, sparse representation, denoising, reduced set method

## I. INTRODUCTION

REPRODUCING KERNELS are functions  $k : \mathcal{X}^2 \rightarrow \mathbb{R}$  which for all pattern sets

$$\{\mathbf{x}_1, \dots, \mathbf{x}_\ell\} \subset \mathcal{X} \quad (1)$$

give rise to positive matrices  $K_{ij} := k(\mathbf{x}_i, \mathbf{x}_j)$  [1]. Here,  $\mathcal{X}$  is some compact set in which the data lives, typically (but not necessarily) a subset of  $\mathbb{R}^N$ . In the SV community, reproducing kernels are often referred to as *Mercer kernels* (section II-B will show why). They provide an elegant way of dealing with nonlinear algorithms by reducing

B. Schölkopf is with GMD FIRST, Rudower Chausee 5, 12489 Berlin, Germany, and Microsoft Research, 1 Guildhall Street, Cambridge CB2 3NH, UK. E-mail: bs@first.gmd.de.

S. Mika is with GMD FIRST, Rudower Chausee 5, 12489 Berlin, Germany. E-mail: mika@first.gmd.de.

C. Burges is with Bell Laboratories, 101 Crawfords Corner Rd., Holmdel NJ, USA. E-mail: burges@lucent.com

P. Knirsch is with the Max-Planck-Institut für biologische Kybernetik, Spemannstr. 38, 72076 Tübingen, Germany. E-mail: phil@kyb.tuebingen.mpg.de

K. Müller is with GMD FIRST, Rudower Chausee 5, 12489 Berlin, Germany. E-mail: klaus@first.gmd.de.

G. Rätsch is with GMD FIRST, Rudower Chausee 5, 12489 Berlin, Germany. E-mail: raetsch@first.gmd.de.

A. Smola is with GMD FIRST, Rudower Chausee 5, 12489 Berlin, Germany. E-mail: smola@first.gmd.de.

Part of this work was done while PK was with Bell Labs and BS and AS were in the Department of Engineering, Australian National University, Canberra. The work was supported by the ARC and the DFG (# Ja 379/52,71,91). Thanks to André Elisseeff for discussions.

them to linear ones in some feature space  $F$  nonlinearly related to input space: Using  $k$  instead of a dot product in  $\mathbb{R}^N$  corresponds to mapping the data into a possibly high-dimensional dot product space  $F$  by a (usually nonlinear) map  $\Phi : \mathbb{R}^N \rightarrow F$ , and taking the dot product there, i.e. [2]

$$k(\mathbf{x}, \mathbf{y}) = (\Phi(\mathbf{x}) \cdot \Phi(\mathbf{y})). \quad (2)$$

By virtue of this property, we shall call  $\Phi$  a *feature map associated with  $k$* . Any linear algorithm which can be carried out in terms of dot products can be made nonlinear by substituting an a priori chosen kernel. Examples of such algorithms include the potential function method, SV Machines, and kernel PCA [3], [4], [5]. The price that one has to pay for this elegance, however, is that the solutions are only obtained as expansions in terms of input patterns mapped into feature space. For instance, the normal vector of an SV hyperplane is expanded in terms of Support Vectors, just as the kernel PCA feature extractors are expressed in terms of training examples,

$$\Psi = \sum_{i=1}^{\ell} \alpha_i \Phi(\mathbf{x}_i). \quad (3)$$

When evaluating an SV decision function or a kernel PCA feature extractor, this is normally not a problem: due to (2), multiplying  $\Psi$  with some mapped test point  $\Phi(\mathbf{x})$  transforms (3) into a kernel expansion which can be evaluated even if  $\Psi$  lives in an infinite-dimensional space. In some cases, however, there are reasons mandating a more comprehensive understanding of what exactly is the connection between patterns in input space and elements of feature space, given as expansions such as (3). This field being far from understood, the current paper attempts to gather some ideas elucidating the problem, and simultaneously proposes some algorithms for situations where the above connection is important. These are the problem of denoising by kernel PCA, and the problem of speeding up Support Vector decision functions.

The remainder of this article is organized as follows. Section II discusses different ways of understanding the mapping from input space into feature space. Following that, we briefly review two feature space algorithms, SV machines and kernel PCA (section III). The focus of interest of the paper, the way back from feature space to input space, is described in section IV, and, in the more general form of constructing sparse approximations of feature space expansions, in section V. The algorithms proposed in these two sections are experimentally evaluated in section VI and discussed in section VII.

## II. FROM INPUT SPACE TO FEATURE SPACE

In this section we show how the feature spaces in question are defined by choice of a suitable kernel function. Insight into the structure of feature space can then be gained by considering their relation to reproducing kernel Hilbert spaces, by how they can be approximated by an empirical map, by their extrinsic geometry, which leads to useful new capacity results, and by their intrinsic geometry, which can be computed solely in terms of the kernel.<sup>1</sup>

### A. The Mercer Kernel Map

We start by stating the version of Mercer's theorem given in [6]. We assume  $(\mathcal{X}, \mu)$  to be a finite measure space.<sup>2</sup> By "almost all" we mean except for sets of measure zero.

**Theorem 1 (Mercer)** *Suppose  $k \in L_\infty(\mathcal{X}^2)$  is a symmetric real-valued kernel<sup>3</sup> such that the integral operator  $T_k : L_2(\mathcal{X}) \rightarrow L_2(\mathcal{X})$ ,*

$$(T_k f)(\mathbf{x}) := \int_{\mathcal{X}} k(\mathbf{x}, \mathbf{y}) f(\mathbf{y}) d\mu(\mathbf{y}) \quad (4)$$

is positive, i.e. for all  $f \in L_2(\mathcal{X})$ , we have  $\int_{\mathcal{X}} k(\mathbf{x}, \mathbf{y}) f(\mathbf{x}) f(\mathbf{y}) d\mu(\mathbf{x}) d\mu(\mathbf{y}) \geq 0$ .

Let  $\psi_j \in L_2(\mathcal{X})$  be the normalized eigenfunctions of  $T_k$  associated with the eigenvalues  $\lambda_j > 0$ , sorted in nonincreasing order. Then

1.  $(\lambda_j)_j \in l_1$ ,
2.  $\psi_j \in L_\infty(\mathcal{X})$  and  $\sup_j \|\psi_j\|_{L_\infty} < \infty$ ,
3.  $k(\mathbf{x}, \mathbf{y}) = \sum_{j=1}^{N_F} \lambda_j \psi_j(\mathbf{x}) \psi_j(\mathbf{y})$  holds for almost all  $(\mathbf{x}, \mathbf{y})$ . Either  $N_F \in \mathbb{N}$ , or  $N_F = \infty$ ; in the latter case, the series converges absolutely and uniformly for almost all  $(\mathbf{x}, \mathbf{y})$ .

From statement 3 it follows that  $k(\mathbf{x}, \mathbf{y})$  corresponds to a dot product in  $l_2^{N_F}$ , i.e.  $k(\mathbf{x}, \mathbf{y}) = (\Phi(\mathbf{x}) \cdot \Phi(\mathbf{y}))$  with

$$\begin{aligned} \Phi : \mathcal{X} &\rightarrow l_2^{N_F} \\ \mathbf{x} &\mapsto (\sqrt{\lambda_j} \psi_j(\mathbf{x}))_j, \quad j = 1, \dots, N_F, \end{aligned} \quad (5)$$

for almost all  $\mathbf{x} \in \mathcal{X}$ .

In fact, the uniform convergence of the series implies that given  $\epsilon > 0$ , there exists an  $n \in \mathbb{N}$  such that even if the range of  $\Phi$  is infinite-dimensional,  $k$  can be approximated within accuracy  $\epsilon$  as a dot product in  $\mathbb{R}^n$ , between images of

$$\tilde{\Phi}^n : \mathbf{x} \mapsto (\sqrt{\lambda_1} \psi_1(\mathbf{x}), \dots, \sqrt{\lambda_n} \psi_n(\mathbf{x})).$$

<sup>1</sup>Those readers who are chiefly interested in applications and algorithms might want to consider skipping over the present section.

<sup>2</sup>A finite measure space is a set  $\mathcal{X}$  with a  $\sigma$ -algebra defined on it, and a measure defined on the latter, satisfying  $\mu(\mathcal{X}) < \infty$  (i.e. up to a scaling factor,  $\mu$  is a probability measure). A  $\sigma$ -algebra  $\Sigma$  on  $\mathcal{X}$  is a family of subsets of  $\mathcal{X}$ , which is closed under elementary set-theoretic operations (countable unions, intersections, and complements), and which contains  $\mathcal{X}$  as a member. A measure is a function  $\mu : \Sigma \rightarrow \mathbb{R}_+ \cup \{\infty\}$  which assigns 0 to the empty set, and moreover is  $\sigma$ -additive, i.e. the measure of a set which is the disjoint unions of some other sets equals the sum of the latter's measures.

<sup>3</sup>i.e. a function of two variables which gives rise to an integral operator

### B. The Reproducing Kernel Map

We can also think of the feature space as a *reproducing kernel Hilbert space* (RKHS). To see this [7], [1], [8], [9], [10], recall that a RKHS is a Hilbert space of functions  $f$  on some set  $\mathcal{X}$  such that all evaluation functionals, i.e. the maps  $f \mapsto f(\mathbf{y})$  ( $\mathbf{y} \in \mathcal{X}$ ), are continuous. In that case, by the Riesz representation theorem, for each  $\mathbf{y} \in \mathcal{X}$  there exists a unique function of  $\mathbf{x}$ , call it  $k(\mathbf{x}, \mathbf{y})$ , such that

$$f(\mathbf{y}) = \langle f, k(\cdot, \mathbf{y}) \rangle \quad (6)$$

(here,  $k(\cdot, \mathbf{y})$  is the function on  $\mathcal{X}$  obtained by fixing the second argument of  $k$  to  $\mathbf{y}$ , and  $\langle \cdot, \cdot \rangle$  is the dot product of the RKHS. In contrast, we use  $(\cdot, \cdot)$  to denote the canonical (Euclidean) dot product). In view of this property,  $k$  is called a *reproducing kernel*.

Note that by (6),  $\langle f, k(\cdot, \mathbf{y}) \rangle = 0$  for all  $\mathbf{y}$  implies that  $f$  is identically zero. Hence the set of functions  $\{k(\cdot, \mathbf{y}) : \mathbf{y} \in \mathcal{X}\}$  spans the whole RKHS. The dot product on the RKHS thus only needs to be defined on  $\{k(\cdot, \mathbf{y}) : \mathbf{y} \in \mathcal{X}\}$  and can then be extended to the whole RKHS by linearity and continuity. From (6), it follows that in particular

$$\langle k(\cdot, \mathbf{x}), k(\cdot, \mathbf{y}) \rangle = k(\mathbf{y}, \mathbf{x}) \quad (7)$$

for all  $\mathbf{x}, \mathbf{y} \in \mathcal{X}$  (this implies that  $k$  is symmetric). Note that this means that any reproducing kernel  $k$  corresponds to a dot product in another space.

Let us now consider a Mercer kernel, i.e. one which satisfies the condition of Theorem 1, and construct a dot product such that  $k$  becomes a reproducing kernel for the Hilbert space  $H$  containing the functions

$$f(\mathbf{x}) = \sum_{i=1}^{\infty} \alpha_i k(\mathbf{x}, \mathbf{x}_i) = \sum_{i=1}^{\infty} \alpha_i \sum_{j=1}^{N_F} \lambda_j \psi_j(\mathbf{x}) \psi_j(\mathbf{x}_i). \quad (8)$$

By linearity, we have

$$\langle f, k(\cdot, \mathbf{y}) \rangle = \sum_{i=1}^{\infty} \alpha_i \sum_{j,n=1}^{N_F} \lambda_j \psi_j(\mathbf{x}_i) \langle \psi_j, \psi_n \rangle \lambda_n \psi_n(\mathbf{y}). \quad (9)$$

Since  $k$  is a Mercer kernel, the  $\psi_i$  ( $i = 1, \dots, N_F$ ) can be chosen to be orthogonal with respect to the dot product in  $L^2(\mathcal{X})$ . Hence it is straightforward to construct a dot product  $\langle \cdot, \cdot \rangle$  such that

$$\langle \psi_j, \psi_n \rangle = \delta_{jn} / \lambda_j \quad (10)$$

(using the Kronecker symbol  $\delta_{jn}$ ), in which case (9) reduces to the reproducing kernel property (6) (using (8)).

Therefore, (7) provides us with a feature map  $\tilde{\Phi}$  associated with  $k$ :

**Proposition 2** *For any Mercer kernel  $k$ , there exists a RKHS  $H$  such that for*

$$\tilde{\Phi} : \mathbb{R}^N \rightarrow H, \quad \mathbf{x} \mapsto k(\cdot, \mathbf{x}), \quad (11)$$

we have

$$\langle \tilde{\Phi}(\mathbf{x}), \tilde{\Phi}(\mathbf{y}) \rangle = k(\mathbf{x}, \mathbf{y}). \quad (12)$$

### C. The Empirical Kernel Map

While giving an interesting alternative theoretical viewpoint, the map  $\tilde{\Phi}$  does not appear all that useful at first sight. In practice, it would seem pointless to first map the inputs into functions, i.e. into infinite-dimensional objects. However, for a given dataset, it is possible to approximate  $\tilde{\Phi}$  from (11) by only evaluating it on these points (cf. [11], [12], [13], [14]):

**Definition 3** For a given set  $\{\mathbf{z}_1, \dots, \mathbf{z}_m\}$ , we call

$$\begin{aligned} \Phi_m : \mathbb{R}^N &\rightarrow \mathbb{R}^m \\ \mathbf{x} &\mapsto k(\cdot, \mathbf{x})|_{\{\mathbf{z}_1, \dots, \mathbf{z}_m\}} \\ &= (k(\mathbf{z}_1, \mathbf{x}), \dots, k(\mathbf{z}_m, \mathbf{x})) \end{aligned} \quad (13)$$

the empirical kernel map w.r.t.  $\{\mathbf{z}_1, \dots, \mathbf{z}_m\}$ .

As an example, consider first the case where  $k$  is a Mercer kernel, and  $\{\mathbf{z}_1, \dots, \mathbf{z}_m\} = \{\mathbf{x}_1, \dots, \mathbf{x}_\ell\}$ , i.e. we evaluate  $k(\cdot, \mathbf{x})$  on the training patterns. If we carry out a linear algorithm in feature space, then everything will take place in the linear span of the mapped training patterns. Therefore, we can represent the  $k(\cdot, \mathbf{x})$  of (11) as  $\Phi_\ell(\mathbf{x})$  without losing information. However, the dot product to use in that representation is not simply the canonical dot product in  $\mathbb{R}^\ell$ , since the  $\Phi(\mathbf{x}_i)$  will usually not form an orthonormal system. To turn  $\Phi_\ell$  into a feature map associated with  $k$ , we need to endow  $\mathbb{R}^\ell$  with a dot product  $\langle \cdot, \cdot \rangle$  such that

$$k(\mathbf{x}, \mathbf{y}) = \langle \Phi_\ell(\mathbf{x}), \Phi_\ell(\mathbf{y}) \rangle. \quad (14)$$

To this end, we use the ansatz  $\langle \cdot, \cdot \rangle = (\cdot, A \cdot)$ , with  $A$  being a positive matrix.<sup>4</sup> Enforcing (14) on the training patterns, this yields the self-consistency condition (cf. [15], [16])

$$K = KAK. \quad (15)$$

Here, we have used  $K$  to denote the  $\ell \times \ell$  kernel Gram matrix

$$K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j). \quad (16)$$

The condition (15) can be satisfied for instance by the pseudo-inverse  $A = K^{-1}$ .

Equivalently, we could have incorporated this rescaling operation, which corresponds to a kernel PCA whitening [5], [15], [14], directly into the map, by modifying (13) to

$$\Phi_\ell^w : \mathbf{x} \mapsto K^{-1/2}(k(\mathbf{x}_1, \mathbf{x}), \dots, k(\mathbf{x}_\ell, \mathbf{x})). \quad (17)$$

This simply amounts to dividing the eigenvector basis vectors of  $K$  by  $\sqrt{\lambda_i}$ , where the  $\lambda_i$  are the eigenvalues of  $K$ .<sup>5</sup> It parallels the rescaling of the eigenfunctions of the integral operator belonging to the kernel, given by (10).

For data sets where the number of examples is smaller than their dimensionality, it can actually be computationally attractive to carry out  $\Phi_\ell^w$  explicitly rather than using

<sup>4</sup>Note that every dot product can be written in this form. Moreover, we do not require definiteness of  $A$ , as the null space can be projected out, leading to a lower-dimensional feature space.

<sup>5</sup>It is understood that if  $K$  is singular, we use the pseudo-inverse of  $K^{1/2}$ .

kernels in whatever subsequent algorithm (SVMs, or kernel PCA, say) one wants to use. As an aside, note that in the case of kernel PCA (to be described in section III), one does not even need to worry about the whitening step: Using the canonical dot product in  $\mathbb{R}^\ell$  (rather than  $\langle \cdot, \cdot \rangle$ ) will simply lead to diagonalizing  $K^2$  instead of  $K$ , which yields the same eigenvectors with squared eigenvalues. This was pointed out by [12] and [13].

We end this section with two notes which illustrate why the use of (13) need not be restricted to the special case we just discussed.

- *More general kernels.* When using non-symmetric kernels  $k$  in (13), together with the canonical dot product, one will effectively work with a matrix  $K^\top K$ , with general  $K$ . Note that each positive semidefinite matrix can be written as  $K^\top K$ .

If we wanted to carry out the whitening step, it would have to be using  $(K^\top K)^{-1/4}$  (cf. footnote 5 concerning potential singularities).

- *Different evaluation sets.* [13] has performed experiments to speed up kernel PCA by choosing  $\{\mathbf{z}_1, \dots, \mathbf{z}_m\}$  as a proper subset of  $\{\mathbf{x}_1, \dots, \mathbf{x}_\ell\}$ .

Now that we have described the kernel map in some detail, including variations on it, we shall next look at its properties. Specifically, we study its effect on the capacity of kernel methods (section II-D) and the induced geometry in feature space (section II-E).

### D. The Capacity of the Kernel Map

[4], [17] gives a bound on the capacity, measured by the VC-dimension  $h$ , of optimal margin classifiers. It takes the form

$$h \leq R^2 \Lambda^2 + 1, \quad (18)$$

where  $\Lambda$  is an upper bound constraining the length of the weight vector of the hyperplane in canonical form, and  $R$  is the radius of the smallest sphere containing the data in the space where the hyperplane is constructed. The smaller this sphere is, the smaller is also the capacity, with beneficial effects on the generalization error bounds.

If the data is distributed in a reasonably isotropic way, which is often the case in input space, then (18) can be fairly precise.<sup>6</sup> If, however, the distribution of the data is such that it does not fill the sphere, then (18) is wasteful. The argument in the remainder of the section, which is summarized from [19], shows that using a kernel typically entails that the data in fact lies in some box with rapidly decaying sidelengths, which can be much smaller than the above sphere.

From statement 2 of Theorem 1, there exists some constant  $C_k > 0$  depending on the kernel  $k$  such that

$$|\psi_j(\mathbf{x})| \leq C_k \text{ for all } j \in \mathbb{N} \text{ and almost all } \mathbf{x} \in \mathcal{X}. \quad (19)$$

Therefore,  $\Phi(\mathcal{X})$  is essentially contained in an axis parallel parallelepiped in  $l_2$  with side lengths  $2C_k \sqrt{\lambda_j}$  (cf. (5)).

<sup>6</sup>In terms of entropy numbers, this is due to the tightness of the rate given by a famous theorem of Maurey (e.g. [18]).

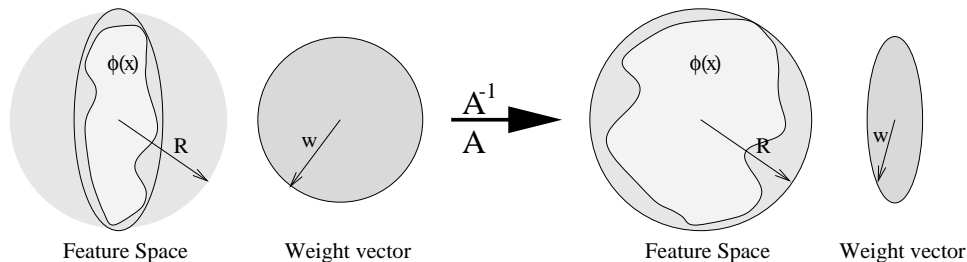


Fig. 1. Since everything is done in terms of dot products, scaling up the data by  $A$  can be compensated by scaling the weight vectors with  $A^{-1}$ . By choosing  $A$  such that the data is still contained in a ball of the same radius  $R$ , we effectively reduce our function class (parameterized by the weight vector), which leads to better generalization bounds which depend on the kernel inducing the map  $\Phi$ .

To see how this effectively restricts the class of functions we are using, we first note that everything in  $F$  is done in terms of dot products. Therefore, we can compensate any invertible linear transformation of the data in  $F$  by the corresponding inverse adjoint transformation on the set of admissible weight vectors in  $F$ , i.e. for any invertible operator  $A$  on  $l_2$ , we have

$$\langle A^* \mathbf{w}, A^{-1} \Phi(\mathbf{x}) \rangle = \langle \mathbf{w}, \Phi(\mathbf{x}) \rangle \quad \text{for all } \mathbf{x} \in \mathcal{X}. \quad (20)$$

Hence, we may construct a diagonal scaling operator  $A^{-1}$  which inflates the sides of the above parallelepiped as much as possible, while ensuring that it still lives in a sphere of the original radius  $R$  in  $l_2$  (figure 1). This will not change the  $R^2$  factor on the right hand side of (18), but it buys us something regarding the second factor: one can show that the function class essentially behaves as if it was finite-dimensional, with a cut-off determined by the rate of decay of  $T_k$ 's eigenvalues.

The reasoning that leads to the improved bounds is somewhat intricate and cannot presently be explained in full detail. In a nutshell, the idea is to compute the capacity (measured in terms of covering numbers of the SV function class, evaluated on an  $\ell$ -sample) via the entropy numbers of a suitable linear operator.<sup>7</sup> Using  $\mathbf{X} := (\mathbf{x}_1, \dots, \mathbf{x}_\ell)$ , consider first the operator

$$\begin{aligned} S_{\Phi(\mathbf{X})} : l_2 &\rightarrow l_\infty^\ell \\ \mathbf{w} &\mapsto \left( \langle \mathbf{w}, \Phi(\mathbf{x}_1) \rangle, \dots, \langle \mathbf{w}, \Phi(\mathbf{x}_\ell) \rangle \right). \end{aligned} \quad (21)$$

For our purpose, the entropy numbers of  $S_{\Phi(\mathbf{X})}$  applied to a sphere of radius  $\Lambda$  are crucial. These can be computed as the entropy numbers of the operator  $\Lambda A S_{A^{-1} \Phi(\mathbf{X})}$ . Using factorization properties of entropy numbers, these can

<sup>7</sup>Consider two normed spaces  $E$  and  $F$ . For  $n \in \mathbb{N}$ , the  $n$ th entropy number of a set  $M \subset E$  is defined as  $\epsilon_n(M) := \inf\{\epsilon > 0: \text{there exists an } \epsilon\text{-cover for } M \text{ in } E \text{ containing } n \text{ or fewer points}\}$ . (Recall that the covering number  $\mathcal{N}(\epsilon, M)$ , being essentially its functional inverse, measures how many balls of radius  $\epsilon$  one needs to cover  $M$ .) The entropy numbers  $\epsilon_n(T)$  of an operator  $T: E \rightarrow F$  are defined as the entropy numbers of the image of the unit ball under  $T$ . Note that  $\epsilon_1(T) = \|T\|$ ; intuitively, the higher entropy number allow a finer characterization of the complexity of the image of  $T$  (e.g. [18], [6]). Note that entropy numbers have some nice properties that covering numbers are lacking. For instance, scaling a subset of a normed vector space by some factor simply scales its entropy numbers by the same factor.

be upper bounded taking into account the above scaling operator  $A$  in a rather precise way. The faster the eigenvalues of the integral operator  $T_k$  associated with a given kernel decay, the smaller the entropy numbers (and hence the capacity) of the corresponding feature space algorithm function class, and the stronger the generalization error bounds that one can prove.

As an example, we now consider how the entropy numbers  $\epsilon_n(A: l_2 \rightarrow l_2)$  depend asymptotically on the eigenvalues of  $T_k$ .

**Proposition 4 ([6])** *Suppose  $k$  is a Mercer kernel with  $\ln(\lambda_j) \sim -j^p$  for some  $p > 0$ . Then  $\ln \epsilon_n^{-1}(A: l_2 \rightarrow l_2) = O(\ln^{\frac{p}{p+1}} n)$*

An example of such a kernel (for  $p = 2$ ) is the Gaussian  $k(x, y) = e^{-|x-y|^2}$ .

This proposition allows the formulation of a priori generalization error bounds depending on the eigenvalues of the kernel. Using similar entropy number methods, it is also possible to give rather precise data-dependent bounds in terms of the eigenvalues of the kernel Gram matrix [20].

Entropy numbers are a promising tool for studying the capacity of feature space methods. This is due to the fact that in the linear case, which is what we are interested in for feature space algorithms, they can be studied using powerful methods of functional analysis (e.g. [18]).

#### E. The Metric of the Kernel Map

Another way to gain insight into the structure of feature space is to consider the intrinsic shape of the manifold to which one's data is mapped. It is important here to distinguish between the feature space  $F$  and the surface in that space to which points in input space  $\mathbb{R}^N$  actually map, which we will call  $S$ . In general  $S$  will be an  $N$  dimensional submanifold embedded in  $F$ . For simplicity, we assume here that  $S$  is sufficiently smooth that structures such as a Riemannian metric can be defined on it. Here we will follow the analysis of [21], to which the reader is referred for more details, although the application to the class of inhomogeneous polynomial kernels is new.

We first note that all intrinsic geometrical properties of  $S$  can be derived once we know the Riemannian metric induced by the embedding of  $S$  in  $F$ . The Riemannian metric can be defined by a symmetric metric tensor  $g_{\mu\nu}$ .

Interestingly, we do not need to know the explicit mapping  $\Phi$  to construct  $g_{\mu\nu}$ ; it can be written solely in terms of the kernel. To see this consider the line element:

$$ds^2 = g_{ab}d\Phi^a(\mathbf{x})d\Phi^b(\mathbf{x}) = g_{\mu\nu}dx^\mu dx^\nu, \quad (22)$$

where indices  $a, b$  correspond to the vector space  $F$  and  $\mu, \nu$  to input space  $\mathbb{R}^N$ . Now letting  $d\mathbf{x}$  represent a small but finite displacement in  $\mathbb{R}^N$ , we have

$$\begin{aligned} ds^2 &= \|\Phi(\mathbf{x} + d\mathbf{x}) - \Phi(\mathbf{x})\|^2 \\ &= k(\mathbf{x} + d\mathbf{x}, \mathbf{x} + d\mathbf{x}) - 2k(\mathbf{x}, \mathbf{x} + d\mathbf{x}) + k(\mathbf{x}, \mathbf{x}) \\ &\approx \left( (1/2)\partial_{x_\mu}\partial_{x_\nu}k(\mathbf{x}, \mathbf{x}) - \partial_{y_\mu}\partial_{y_\nu}k(\mathbf{x}, \mathbf{y}) \right)_{\mathbf{y}=\mathbf{x}} dx^\mu dx^\nu \\ &\equiv g_{\mu\nu}dx^\mu dx^\nu \end{aligned}$$

Thus we can read off the components of the metric tensor:

$$g_{\mu\nu} = (1/2)\partial_{x_\mu}\partial_{x_\nu}k(\mathbf{x}, \mathbf{x}) - \{\partial_{y_\mu}\partial_{y_\nu}k(\mathbf{x}, \mathbf{y})\}_{\mathbf{y}=\mathbf{x}} \quad (23)$$

For the class of kernels which are functions only of dot products between points in  $\mathbb{R}^N$ , both covariant and contravariant components of the metric take a simple form:

$$g_{\mu\nu} = \delta_{\mu\nu}k'(\|\mathbf{x}\|^2) + x_\mu x_\nu k''(\|\mathbf{x}\|^2) \quad (24)$$

$$g^{\mu\nu} = \frac{\delta_{\mu\nu}}{k'} - x_\mu x_\nu \frac{k''/k'}{k' + \|\mathbf{x}\|^2 k''} \quad (25)$$

where the prime denotes derivative with respect to the argument  $\|\mathbf{x}\|^2$ . To illustrate, let us compute the intrinsic geometrical properties of the surface  $S$  corresponding to the class of inhomogeneous polynomial Mercer kernels:

$$k(\mathbf{x}, \mathbf{y}) = (\mathbf{x} \cdot \mathbf{y} + c)^p, \quad (26)$$

where  $c$  is a constant. Properties relating to the intrinsic curvature of a surface are completely captured by the Riemann curvature tensor

$$R_{\nu\alpha\beta}{}^\mu = \partial_\alpha\Gamma_{\beta\nu}^\mu - \partial_\beta\Gamma_{\alpha\nu}^\mu + \Gamma_{\alpha\nu}^\rho\Gamma_{\beta\rho}^\mu - \Gamma_{\beta\nu}^\rho\Gamma_{\alpha\rho}^\mu \quad (27)$$

where the Christoffel symbols of the second kind are defined by:

$$\Gamma_{\beta\gamma}^\alpha \equiv g^{\alpha\mu}\Gamma_{\beta\gamma\mu} = \frac{1}{2}g^{\alpha\mu}(\partial_\beta g_{\gamma\mu} - \partial_\mu g_{\beta\gamma} + \partial_\gamma g_{\mu\beta}) \quad (28)$$

Thus we find that for this class of kernels, the Riemann curvature, for arbitrary input space dimension  $N$  and polynomial order  $p$ , is given by

$$\begin{aligned} R_{\nu\alpha\beta}{}^\mu &= \frac{-2(p-1)}{(\|\mathbf{x}\|^2 + c)^2}(x_\alpha x_\nu \delta_{\mu\beta} - x_\beta x_\nu \delta_{\mu\alpha}) \\ &+ \frac{(p-1)}{(\|\mathbf{x}\|^2 + c)}(\delta_{\nu\alpha}\delta_{\mu\beta} - \delta_{\nu\beta}\delta_{\mu\alpha}) \\ &- \frac{(p-1)}{(\|\mathbf{x}\|^2 + c)^2(\|\mathbf{x}\|^2 + c/p)} \cdot \\ &\quad (x_\beta x_\nu \delta_{\mu\alpha} - x_\alpha x_\nu \delta_{\mu\beta} + x_\beta x_\mu \delta_{\nu\alpha} - x_\alpha x_\mu \delta_{\nu\beta}) \\ &+ \frac{(p-1)^2}{(\|\mathbf{x}\|^2 + c)^2} \cdot \frac{c/p}{\|\mathbf{x}\|^2 + c/p}(x_\alpha x_\nu \delta_{\beta\mu} - x_\beta x_\nu \delta_{\alpha\mu}) \end{aligned}$$

It is interesting to compare this result with that for the homogeneous case ( $c = 0$ ) [21]:

$$\begin{aligned} R_{\mu\nu\rho}{}^\sigma(c=0) &= \frac{(p-1)}{\|\mathbf{x}\|^2}(\delta_{\rho\nu}\delta_{\sigma\mu} - \delta_{\rho\mu}\delta_{\sigma\nu}) \\ &- \frac{(p-1)}{\|\mathbf{x}\|^4}(x_\nu x_\rho \delta_{\mu\sigma} - x_\mu x_\rho \delta_{\sigma\nu} + x_\sigma x_\mu \delta_{\rho\nu} - x_\sigma x_\nu \delta_{\rho\mu}) \end{aligned} \quad (29)$$

This analysis shows that adding the constant  $c$  to the kernel results in striking differences in the geometries. Both curvatures vanish for  $p = 1$ , as expected. However for  $N = 2$ ,  $R_{\mu\nu\rho}{}^\sigma(c=0)$  vanishes for all powers  $p$ , whereas  $R_{\mu\nu\rho}{}^\sigma(c \neq 0)$  does not vanish for any  $p > 1$ . Furthermore, all surfaces for homogeneous kernels with non-vanishing curvature (i.e. those with  $N > 2$  and  $p > 1$ ) have a singularity at  $\|\mathbf{x}\|^2 = 0$ , whereas none of the corresponding surfaces for inhomogeneous kernels do, provided  $c > 0$ .<sup>8</sup>

Thus beyond providing insight into the geometrical structure of the surfaces in  $F$  generated by choice of kernel, the geometrical analysis also gives more concrete results. For example, one can expect that the density on the surface  $S$  will become ill behaved for data  $\mathbf{x}$  whose norm is small, for homogeneous polynomial kernels, but not for the inhomogeneous case. Similar problems may arise in the pattern recognition case if one's data lies near the singularity. These considerations can be extended to compute the intrinsic volume element on  $S$ , which may be used to compute the density  $p(\mathbf{x})$  on  $S$ , and to give some simple necessary tests that must be satisfied by any kernel if it is to be a Mercer kernel: the reader is referred to [21] for details.

### III. FEATURE SPACE ALGORITHMS

We next describe the two algorithms used in this paper: SV machines, and kernel PCA. The SV algorithm shall not be described in detail, let us only briefly fix the notation.

*Support Vector (SV) classifiers* [4] construct a maximum margin hyperplane in  $F$ . In input space, this corresponds to a nonlinear decision boundary of the form

$$f(\mathbf{x}) = \text{sgn} \left( \sum_{i=1}^{\ell} \alpha_i k(\mathbf{x}, \mathbf{x}_i) + b \right), \quad (30)$$

where the  $\mathbf{x}_i$  are the training examples. Those with  $\alpha_i \neq 0$  are called *Support Vectors*; in many applications, most of the  $\alpha_i$ , which are found by solving a quadratic program, turn out to be 0. Excellent classification accuracies in both OCR and object recognition have been obtained using SV machines [9], [22]. A generalization to the case of regression estimation, leading to similar function expansion, exists [4], [23].

*Kernel Principal Component Analysis* [5] carries out a linear PCA in the feature space  $F$ . The extracted features take the nonlinear form

$$f_k(\mathbf{x}) = \sum_{i=1}^{\ell} \alpha_i^k k(\mathbf{x}_i, \mathbf{x}), \quad (31)$$

<sup>8</sup>Note that if  $c < 0$  the kernel is not a Mercer kernel and the above analysis does not apply.

where, up to a normalization, the  $\alpha_i^k$  are the components of the  $k$ -th eigenvector of the matrix  $(k(\mathbf{x}_i, \mathbf{x}_j))_{ij}$ .

This can be understood as follows. We wish to find eigenvectors  $\mathbf{V}$  and eigenvalues  $\lambda$  of the covariance matrix  $C$  in the feature space,<sup>9</sup> where

$$C := \frac{1}{\ell} \sum_{i=1}^{\ell} \Phi(\mathbf{x}_i) \Phi(\mathbf{x}_i)^\top. \quad (32)$$

In the case when  $F$  is very high dimensional this will be impossible to compute directly. To be still able to solve this problem, one uses Mercer kernels. To this end, we need to derive a formulation which only uses  $\Phi$  in dot products. We then replace any occurrence of  $(\Phi(\mathbf{x}) \cdot \Phi(\mathbf{y}))$  by  $k(\mathbf{x}, \mathbf{y})$ . This way we avoid dealing with the mapped data explicitly, which may be intractable in terms of memory and computational cost.

To find a formulation for PCA which uses only dot products, we first substitute the covariance matrix (32) into the eigenequation  $C\mathbf{V} = \lambda\mathbf{V}$ . Note that all solutions to this equation with  $\lambda \neq 0$  must lie in the span of  $\Phi$ -images of the training data. Thus, we may consider the equivalent system

$$\lambda(\Phi(\mathbf{x}_k) \cdot \mathbf{V}) = (\Phi(\mathbf{x}_k) \cdot C\mathbf{V}) \text{ for all } k = 1, \dots, \ell, \quad (33)$$

and expand the solution  $\mathbf{V}$  as

$$\mathbf{V} = \sum_{i=1}^{\ell} \alpha_i \Phi(\mathbf{x}_i). \quad (34)$$

Substituting (32) and (34) into (33), and defining a  $\ell \times \ell$  matrix  $K$  by  $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$  we arrive at a problem which is cast in terms of dot products: solve

$$\ell\lambda\boldsymbol{\alpha} = K\boldsymbol{\alpha}, \quad (35)$$

where  $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_\ell)^\top$  (for details on the last step see [5]). Normalizing the solution  $\mathbf{V}_k$  in  $F$  translates into  $\lambda_k(\boldsymbol{\alpha}^k \cdot \boldsymbol{\alpha}^k) = 1$ . To extract features, we compute the projection of the  $\Phi$ -image of a test point  $\mathbf{x}$  onto the  $k$ -th eigenvector in the feature space by

$$\beta_k := (\mathbf{V}^k \cdot \Phi(\mathbf{x})) = \sum_{i=1}^{\ell} \alpha_i^k k(\mathbf{x}_i, \mathbf{x}). \quad (36)$$

Usually, this will be much cheaper than taking the dot product in the feature space explicitly.

To conclude the brief summary of kernel PCA, we state a characterization which involves the same regularizer (the length of the weight vector) as the one used in SV machines.

**Proposition 5** *For all  $k$ , the  $k$ -th kernel PCA feature extractor, scaled by  $1/\lambda_k$ , is optimal among all feature extractors of the form*

$$f(\mathbf{x}) = \sum_i \alpha_i k(\mathbf{x}_i, \mathbf{x})$$

<sup>9</sup>Here we assume that the mapped data is centered, too. In general this will not be true, but all computations can easily be reformulated to perform an explicit centering in  $F$  [5].

*in the sense that it has the shortest weight vector*

$$\|\mathbf{V}\|^2 = \sum_i \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j)$$

*subject to the conditions that*

(1) *it is orthogonal to the first  $k - 1$  kernel PCA feature extractors (in feature space)*

(2) *applied to the training set  $\mathbf{x}_1, \dots, \mathbf{x}_\ell$ , it leads to a unit variance set of outputs.*

Both SV machines and kernel PCA utilize Mercer kernels to generalize a linear algorithm to a nonlinear setting; moreover, both use the same regularizer, albeit in different domains of learning — supervised vs. unsupervised. Nevertheless, feature extraction experiments on handwritten digit images using kernel PCA have shown that a linear hyperplane classifier trained on the extracted features performs as well as a nonlinear SV machine trained directly on the inputs [5].

#### IV. FROM FEATURE SPACE TO INPUT SPACE

Unlike section II, which described how to get from input space into feature space, we now study the way back. There has been a fair amount of work on aspects of this problem in the context of developing so-called reduced set methods (e.g. [24], [25], [26], [27], [28]). For pedagogical reasons, we shall postpone reduced set methods to section V, as they focus on a problem that is already more complex than the one we would like to start with.

##### A. The Pre-Image Problem

As stated in the introduction, feature space algorithms express their solutions as expansions in terms of mapped input points (3). However, since the map  $\Phi$  into the feature space  $F$  is nonlinear, we cannot generally assert that each such expansion will have a pre-image under  $\Phi$ , i.e. a point  $\mathbf{z} \in \mathbb{R}^N$  such that  $\Phi(\mathbf{z}) = \Psi$  (figure 2). If the pre-image existed, it would be easy to compute, as shown by the following result [9]:

**Proposition 6** *Consider a feature space expansion  $\Psi = \sum_{j=1}^{\ell} \alpha_j \Phi(\mathbf{x}_j)$ . If there exists a  $\mathbf{z} \in \mathbb{R}^N$  such that*

$$\Phi(\mathbf{z}) = \Psi,$$

*and if  $k$  is an invertible function  $f_k$  of  $(\mathbf{x} \cdot \mathbf{y})$ , then we can compute  $\mathbf{z}$  as*

$$\mathbf{z} = \sum_{i=1}^N f_k^{-1} \left( \sum_{j=1}^{\ell} \alpha_j k(\mathbf{x}_j, \mathbf{e}_i) \right) \mathbf{e}_i,$$

*where  $\{\mathbf{e}_1, \dots, \mathbf{e}_N\}$  is any orthonormal basis of input space.*

*Proof:* We expand  $\mathbf{z}$  as

$$\mathbf{z} = \sum_{i=1}^N (\mathbf{z} \cdot \mathbf{e}_i) \mathbf{e}_i$$

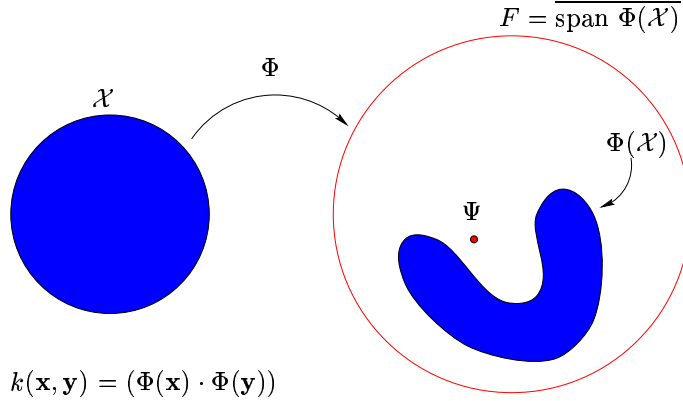


Fig. 2. The pre-image problem. Not each point in the span of the mapped input data is necessarily the image of some input pattern. Therefore, not each point that can be written as an expansion in terms of mapped input patterns (e.g. a kernel PCA eigenvector, or a SVM hyperplane normal vector), can necessarily be expressed as the image of a single input pattern.

$$\begin{aligned} &= \sum_{i=1}^N f_k^{-1}(k(\mathbf{z}, \mathbf{e}_i)) \mathbf{e}_i \\ &= \sum_{i=1}^N f_k^{-1} \left( \sum_{j=1}^{\ell} \alpha_j k(\mathbf{x}_j, \mathbf{e}_i) \right) \mathbf{e}_i. \end{aligned} \quad (37)$$

■

Several remarks on this proposition should be given. First, examples of kernels which are invertible functions of  $(\mathbf{x} \cdot \mathbf{y})$  are polynomial kernels

$$k(\mathbf{x}, \mathbf{y}) = (\mathbf{x} \cdot \mathbf{y} + c)^d, \quad c \geq 0, \quad d \text{ odd}, \quad (38)$$

and sigmoid kernels

$$k(\mathbf{x}, \mathbf{y}) = \sigma(\kappa \cdot (\mathbf{x} \cdot \mathbf{y}) + \Theta), \quad \kappa, \Theta \in \mathbb{R}. \quad (39)$$

A similar result holds for RBF kernels (using the polarization identity) — all we need is a kernel which allows the reconstruction of  $(\mathbf{x} \cdot \mathbf{y})$  from  $k$ , evaluated on some input points which we are allowed to choose (for details, cf. [9]).

The crucial assumption, clearly, is the *existence* of the pre-image. Unfortunately, there are many situations where there are no pre-images. To illustrate this, we consider the feature map  $\tilde{\Phi}$  (11). Clearly, only points in feature space which can be written as  $k(\cdot, \mathbf{x})$  do have a pre-image under this map. To characterize this set of points in a specific example, consider Gaussian kernels

$$k(\mathbf{x}, \mathbf{y}) = e^{-\frac{\|\mathbf{x}-\mathbf{y}\|^2}{2\sigma^2}}. \quad (40)$$

In this case,  $\tilde{\Phi}$  maps each input into a Gaussian sitting on that point. However, it is known [29] that no Gaussian can be written as a linear combination of Gaussians centered at other points. Therefore, in the Gaussian case, none of the expansions (3), excluding trivial cases with only one term, has an exact pre-image.

The problem that we had initially set out to solve has turned out to be insolvable in the general case. Let us try to ask for less. Rather than trying to find exact pre-images, we now consider approximate ones. We call  $\mathbf{z} \in \mathbb{R}^N$  an *approximate pre-image* for  $\Psi$  if

$$\rho(\mathbf{z}) = \|\Psi - \Phi(\mathbf{z})\|^2 \quad (41)$$

is small.<sup>10</sup>

Are there vectors  $\Psi$  for which good approximate pre-images exist? As described in section III, kernel PCA is nothing but PCA in  $F$ . Therefore, for  $n = 1, 2, \dots, p$ , it provides projections

$$P_n \Phi(\mathbf{x}) = \sum_{k=1}^n (\Phi(\mathbf{x}) \cdot \mathbf{V}^k) \mathbf{V}^k \quad (42)$$

with the following optimal approximation property (e.g. [9]) Here, we assume that the  $\mathbf{V}^k$  are sorted according to nonincreasing eigenvalues  $\lambda_k$ , with  $\lambda_p$  being the smallest nonzero eigenvalue.

**Proposition 7**  $P_n$  is the  $n$ -dimensional projection minimizing

$$\sum_{i=1}^{\ell} \|P_n \Phi(\mathbf{x}_i) - \Phi(\mathbf{x}_i)\|^2. \quad (43)$$

Therefore,  $P_n \Phi(\mathbf{x})$  can be expected to *have* a good approximate pre-image: trivially, already  $\mathbf{x}$  is a good approximate pre-image. As we shall see in experiments, however, even better pre-images can be found, which makes some interesting applications possible [30], [31]:

*Denoising.* Given a noisy  $\mathbf{x}$ , map it into  $\Phi(\mathbf{x})$ , discard higher components to obtain  $P_n \Phi(\mathbf{x})$ , and then compute a pre-image  $\mathbf{z}$ . Here, the hope is that the main structure in the data set is captured in the first  $n$  directions, and the remaining components mainly pick up the noise — in this sense,  $\mathbf{z}$  can be thought of as a denoised version of  $\mathbf{x}$ .

*Compression.* Given the eigenvectors  $\alpha^k$  and a small number of features  $\beta_k$  (cf. (36)) of  $\Phi(\mathbf{x})$ , but not  $\mathbf{x}$ , compute a pre-image as an approximate reconstruction of  $\mathbf{x}$ . This is useful if  $n$  is smaller than the dimensionality of the input data.

*Interpretation.* Visualize a nonlinear feature extractor  $\mathbf{V}^k$  by computing a pre-image.

<sup>10</sup>Just how small it needs to be in order to form a satisfactory approximation will depend on the problem at hand. Therefore, we have refrained from giving a formal definition.

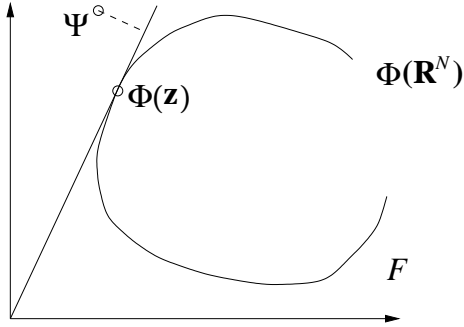


Fig. 3. Given a vector  $\Psi \in F$ , we try to approximate it by a multiple of a vector  $\Phi(\mathbf{z})$  in the image of input space  $(\mathbb{R}^N)$  under the nonlinear map  $\Phi$  by finding  $\mathbf{z}$  such that the projection distance of  $\Psi$  onto  $\text{span}(\Phi(\mathbf{z}))$  is minimized.

In this paper, we focus on the first point. In the next section, we shall develop a method for minimizing (41), which we will later, in the experimental section, apply to the case where  $\Psi = P_n \Phi(\mathbf{x})$ .

### B. An Algorithm for Approximate Pre-Images

The present section [26] gives an analysis for the case of the Gaussian kernel, which has proven to perform very well in applications [32], and proposes an iteration procedure for computing pre-images of kernel expansions.

We start by considering a problem slightly more general than the pre-image problem: we are seeking to approximate  $\Psi = \sum_{i=1}^{N_x} \alpha_i \Phi(\mathbf{x}_i)$  by  $\Psi' = \beta \Phi(\mathbf{z})$ . First observe that rather than minimizing

$$\|\Psi - \Psi'\|^2 = \sum_{i,j=1}^{N_x} \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) + \beta^2 k(\mathbf{z}, \mathbf{z}) - 2 \sum_{i=1}^{N_x} \alpha_i \beta k(\mathbf{x}_i, \mathbf{z}), \quad (44)$$

we can minimize the distance between  $\Psi$  and the orthogonal projection of  $\Psi$  onto  $\text{span}(\Phi(\mathbf{z}))$  (figure 3),

$$\left\| \frac{(\Psi \cdot \Phi(\mathbf{z}))}{(\Phi(\mathbf{z}) \cdot \Phi(\mathbf{z}))} \Phi(\mathbf{z}) - \Psi \right\|^2 = \|\Psi\|^2 - \frac{(\Psi \cdot \Phi(\mathbf{z}))^2}{(\Phi(\mathbf{z}) \cdot \Phi(\mathbf{z}))}. \quad (45)$$

To this end, we maximize

$$\frac{(\Psi \cdot \Phi(\mathbf{z}))^2}{(\Phi(\mathbf{z}) \cdot \Phi(\mathbf{z}))}, \quad (46)$$

which can be expressed in terms of the kernel. The maximization of (46) over  $\mathbf{z}$  is preferable to the one of (44) over  $\mathbf{z}$  and  $\beta$ , since it comprises a lower-dimensional problem, and since  $\mathbf{z}$  and  $\beta$  have different scaling behavior. Once the maximum of (46) is found, it is extended to the minimum of (44) by setting (cf. (45))  $\beta = (\Psi \cdot \Phi(\mathbf{z})) / (\Phi(\mathbf{z}) \cdot \Phi(\mathbf{z}))$ . The function (46) can either be minimized using standard techniques (as [24]), or, for particular choices of kernels, using fixed-point iteration methods, as shown presently.

For kernels which satisfy  $k(\mathbf{z}, \mathbf{z}) = 1$  for all  $\mathbf{z} \in \mathbb{R}^N$  (e.g. Gaussian kernels), (46) reduces to

$$(\Psi \cdot \Phi(\mathbf{z}))^2. \quad (47)$$

For the extremum, we have  $0 = \nabla_{\mathbf{z}}(\Psi \cdot \Phi(\mathbf{z}))^2 = 2(\Psi \cdot \Phi(\mathbf{z})) \nabla_{\mathbf{z}}(\Psi \cdot \Phi(\mathbf{z}))$ . To evaluate the gradient in terms of  $k$ , we substitute (53) to get the sufficient condition  $0 = \sum_{i=1}^{N_x} \alpha_i \nabla_{\mathbf{z}} k(\mathbf{x}_i, \mathbf{z})$ . For  $k(\mathbf{x}_i, \mathbf{z}) = k(\|\mathbf{x}_i - \mathbf{z}\|^2)$  (e.g. Gaussians, or  $(\|\mathbf{x}_i - \mathbf{z}\|^2 + 1)^c$  for  $c = -1, -1/2$ ), we obtain  $0 = \sum_{i=1}^{N_x} \alpha_i k'(\|\mathbf{x}_i - \mathbf{z}\|^2)(\mathbf{x}_i - \mathbf{z})$ , leading to

$$\mathbf{z} = \frac{\sum_{i=1}^{N_x} \alpha_i k'(\|\mathbf{x}_i - \mathbf{z}\|^2) \mathbf{x}_i}{\sum_{i=1}^{N_x} \alpha_i k'(\|\mathbf{x}_i - \mathbf{z}\|^2)}. \quad (48)$$

For the Gaussian kernel  $k(\mathbf{x}_i, \mathbf{z}) = \exp(-\|\mathbf{x}_i - \mathbf{z}\|^2 / (2\sigma^2))$  we thus arrive at

$$\mathbf{z} = \frac{\sum_{i=1}^{N_x} \alpha_i \exp(-\|\mathbf{x}_i - \mathbf{z}\|^2 / (2\sigma^2)) \mathbf{x}_i}{\sum_{i=1}^{N_x} \alpha_i \exp(-\|\mathbf{x}_i - \mathbf{z}\|^2 / (2\sigma^2))}, \quad (49)$$

and devise an iteration

$$\mathbf{z}_{n+1} = \frac{\sum_{i=1}^{N_x} \alpha_i \exp(-\|\mathbf{x}_i - \mathbf{z}_n\|^2 / (2\sigma^2)) \mathbf{x}_i}{\sum_{i=1}^{N_x} \alpha_i \exp(-\|\mathbf{x}_i - \mathbf{z}_n\|^2 / (2\sigma^2))}. \quad (50)$$

The denominator equals  $(\Psi \cdot \Phi(\mathbf{z}_n))$  and thus is nonzero in a neighborhood of the extremum of (47), unless the extremum itself is zero. The latter only occurs if the projection of  $\Psi$  on the linear span of  $\Phi(\mathbb{R}^N)$  is zero, in which case it is pointless to try to approximate  $\Psi$ . Numerical instabilities related to  $(\Psi \cdot \Phi(\mathbf{z}))$  being small can thus be approached by restarting the iteration with different starting values.

Interestingly, (50) can be interpreted in the context of clustering (e.g. [33]). It determines the center of a single Gaussian cluster, trying to capture as many of the  $\mathbf{x}_i$  with positive  $\alpha_i$  as possible, and simultaneously avoids those  $\mathbf{x}_i$  with negative  $\alpha_i$ . For SV classifiers, the sign of the  $\alpha_i$  equals the label of the pattern  $\mathbf{x}_i$ . It is this sign which distinguishes (50) from plain clustering or parametric density estimation. The occurrence of negative signs is related to the fact that we are not trying to estimate a parametric density but the *difference* between two densities (modulo normalization constants).

To see this, we define the sets  $pos = \{i : \alpha_i > 0\}$  and  $neg = \{i : \alpha_i < 0\}$ , and the shorthands  $p_{pos}(\mathbf{z}) = \sum_{pos} \alpha_i \exp(-\|\mathbf{x}_i - \mathbf{z}\|^2 / (2\sigma^2))$  and  $p_{neg}(\mathbf{z}) = \sum_{neg} |\alpha_i| \exp(-\|\mathbf{x}_i - \mathbf{z}\|^2 / (2\sigma^2))$ . The target (47) then reads  $(p_{pos}(\mathbf{z}) - p_{neg}(\mathbf{z}))^2$ , i.e. we are trying to find a point  $\mathbf{z}$  where the difference between the (unnormalized) “probabilities” for the two classes is maximal, and estimate the approximation to (53) by a Gaussian centered at  $\mathbf{z}$ . Moreover, note that we can rewrite (49) as

$$\mathbf{z} = \frac{p_{pos}(\mathbf{z})}{p_{pos}(\mathbf{z}) - p_{neg}(\mathbf{z})} \mathbf{x}_{pos} + \frac{p_{neg}(\mathbf{z})}{p_{neg}(\mathbf{z}) - p_{pos}(\mathbf{z})} \mathbf{x}_{neg}, \quad (51)$$

where

$$\mathbf{x}_{pos/neg} = \frac{\sum_{pos/neg} \alpha_i \exp(-\|\mathbf{x}_i - \mathbf{z}\|^2 / (2\sigma^2)) \mathbf{x}_i}{\sum_{pos/neg} \alpha_i \exp(-\|\mathbf{x}_i - \mathbf{z}\|^2 / (2\sigma^2))}. \quad (52)$$



## V. REDUCED SET (RS) METHODS

### A. The Problem

We now move on to a slightly more general problem, first studied by [24], where we are no longer only looking for single pre-images, but expansions. It will turn out that one can also use the method developed in the last section to design an algorithm for the more general case.

Assume we are given a vector  $\Psi \in F$ , expanded in images of input patterns  $\mathbf{x}_i \in \mathbb{R}^N$ ,

$$\Psi = \sum_{i=1}^{N_x} \alpha_i \Phi(\mathbf{x}_i), \quad (53)$$

with  $\alpha_i \in \mathbb{R}, \mathbf{x}_i \in \mathbb{R}^N$ . Rather than looking for a single pre-image, we now try to approximate it by a *reduced set* expansion [24]

$$\Psi' = \sum_{i=1}^{N_z} \beta_i \Phi(\mathbf{z}_i), \quad (54)$$

with  $N_z < N_x, \beta_i \in \mathbb{R}, \mathbf{z}_i \in \mathbb{R}^N$ . To this end, one can minimize [24]

$$\begin{aligned} \|\Psi - \Psi'\|^2 &= \sum_{i,j=1}^{N_x} \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) + \sum_{i,j=1}^{N_z} \beta_i \beta_j k(\mathbf{z}_i, \mathbf{z}_j) \\ &\quad - 2 \sum_{i=1}^{N_x} \sum_{j=1}^{N_z} \alpha_i \beta_j k(\mathbf{x}_i, \mathbf{z}_j). \end{aligned} \quad (55)$$

The crucial point is that even if  $\Phi$  is not given explicitly, (55) can be computed (and minimized) in terms of the kernel.

In the NIST benchmark of 60000 handwritten digits, SV machines are more accurate than any other single classifier [9]; however, they are inferior to neural nets in run-time classification speed [25]. In applications where the latter is an issue, it is thus desirable to come up with methods to speed up things by making the SV expansion more sparse, i.e. replacing (53) by (54).

### B. Finding the Coefficients

Evidently, the RS problem consists of two parts. One has to determine the RS vectors  $\mathbf{z}_i$ , and one has to compute the expansion coefficients  $\beta_i$ . We start with the latter; partly, as it is easier, partly, as it is common to different RS methods.

**Proposition 8 ([26])** *The optimal coefficients  $\beta = (\beta_1, \dots, \beta_m)$  for approximating  $\Psi = \sum_{i=1}^{\ell} \alpha_i \Phi(\mathbf{x}_i)$  by  $\sum_{i=1}^m \beta_i \Phi(\mathbf{z}_i)$  (for linearly independent  $\Phi(\mathbf{z}_1), \dots, \Phi(\mathbf{z}_m)$ ) in the 2-norm are given by*

$$\beta = (K^z)^{-1} K^{zx} \alpha. \quad (56)$$

Here,  $K_{ij}^z := (\Phi(\mathbf{z}_i) \cdot \Phi(\mathbf{z}_j))$  and  $K_{ij}^{zx} := (\Phi(\mathbf{z}_i) \cdot \Phi(\mathbf{x}_j))$ .

Note that if the  $\Phi(\mathbf{z}_i)$  are linearly independent, as they should be if we want to use them for approximation, then

$K^z$  has full rank. Otherwise, one can use the pseudo-inverse, or select the solution which has the largest number of zero components.

*Proof:* We evaluate the derivative of the distance in  $F$ ,

$$\frac{\partial}{\partial \beta_k} \left\| \Psi - \sum_{i=1}^m \beta_i \Phi(\mathbf{z}_i) \right\|^2 = -2 \Phi(\mathbf{z}_k) \left( \Psi - \sum_{i=1}^m \beta_i \Phi(\mathbf{z}_i) \right), \quad (57)$$

and set it to 0. Substituting  $\Psi = \sum_{i=1}^{\ell} \alpha_i \Phi(\mathbf{x}_i)$ , we obtain (using  $\alpha = (\alpha_1, \dots, \alpha_{\ell})$ )

$$K^{zx} \alpha = K^z \beta, \quad (58)$$

hence

$$\beta = (K^z)^{-1} K^{zx} \alpha. \quad (59)$$

No RS algorithm using the 2-norm optimality criterion can circumvent this result. For instance, suppose we are given an algorithm that computes the  $\beta_i$  and  $\mathbf{z}_i$  simultaneously and comes up with a solution. Then we can always use the proposition to recompute the optimal coefficients to get a solution which is at least as good as the original one. Different algorithms can, however, differ in the way they determine the vectors  $\mathbf{z}_i$  in the first place. The one dealt with in the next section simply *selects* subsets of the  $\mathbf{x}_i$ , while the one in section V-D uses vectors different from the original  $\mathbf{x}_i$ . ■

### C. Reduced Set Selection

#### C.1 Selection via Kernel PCA

The idea for the first algorithm arises from the observation that the null space of the Gram matrix  $K_{ij} = (\Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j))$  precisely tells us how many vectors can be removed from an expansion while committing zero approximation error (assuming we correctly adjust the coefficients), i.e. how sparse we can make an SV expansion, say, without changing it the least [9], [27]. (Here, the Gram matrix can either be computed only from those examples which have a nonzero  $\alpha_j$ , or from a larger set which we want to use for the expansion.) Interestingly, it will turn out that this problem is closely related to kernel PCA.

Let us start with the simplest case. Assume there exists an eigenvector  $\alpha \neq 0$  of  $K$  with eigenvalue 0, i.e.  $K\alpha = 0$ . Using (2), this reads

$$\sum_{j=1}^{\ell} (\Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)) \alpha_j = 0 \quad \text{for all } i = 1, \dots, \ell, \quad (60)$$

hence

$$\sum_{j=1}^{\ell} \alpha_j \Phi(\mathbf{x}_j) = 0. \quad (61)$$

This means that the  $\Phi(\mathbf{x}_j)$  are linearly dependent, and therefore any of the  $\Phi(\mathbf{x}_j)$  which comes with a nonzero  $\alpha_j$  can be expressed in terms of the others. Hence, we may use the eigenvectors with eigenvalue 0 to eliminate certain terms from any expansion in the  $\Phi(\mathbf{x}_j)$ .

What happens if we do not have nonzero eigenvalues, such as in the case of Gaussian kernels [29]? Intuitively, we would still believe that even though the above is no longer precisely true, it should give a good approximation. However, the crucial difference is that in order to get the best possible approximation, we now need to take into account the coefficients of the expansion  $\Psi = \sum_{j=1}^{\ell} \alpha_j \Phi(\mathbf{x}_j)$ : if we commit an error by removing  $\Phi(\mathbf{x}_n)$ , say, then this error will depend on  $\alpha_n$ , too. How do we then select the optimal  $n$ ?

Clearly, we would like to find coefficients  $\beta_j$  minimizing the error we commit by replacing  $\alpha_n \Phi(\mathbf{x}_n)$  with  $\sum_{j \neq n} \beta_j \Phi(\mathbf{x}_j)$ ,

$$\rho(\beta, n) = \left\| \alpha_n \Phi(\mathbf{x}_n) - \sum_{j \neq n} \beta_j \Phi(\mathbf{x}_j) \right\|^2. \quad (62)$$

To establish a connection to kernel PCA, we make a change of variables. First, define  $\eta_j = 1$  for  $j = n$ ,  $\eta_j = -\beta_j/\alpha_n$  for  $j \neq n$ . Hence (62) equals  $|\alpha_n|^2 \|\sum_{j=1}^{\ell} \eta_j \Phi(\mathbf{x}_j)\|^2$ . Normalizing  $\boldsymbol{\eta}$  to obtain  $\boldsymbol{\gamma} := \boldsymbol{\eta}/\|\boldsymbol{\eta}\|$ , hence  $\gamma_n = 1/\|\boldsymbol{\eta}\|$ , this leads to the problem of minimizing

$$\rho(\boldsymbol{\gamma}, n) = \left| \frac{\alpha_n}{\gamma_n} \right|^2 (\boldsymbol{\gamma} \cdot K \boldsymbol{\gamma}), \quad (63)$$

over  $\|\boldsymbol{\gamma}\| = 1$  (note that  $\rho$  is invariant when rescaling  $\boldsymbol{\gamma}$ ).<sup>11</sup> A straightforward calculation shows that we can recover the approximation coefficients for  $\alpha_n \Phi(\mathbf{x}_n)$ , i.e. the values to add to the  $\alpha_j$  ( $j \neq n$ ) for leaving out  $\alpha_n \Phi(\mathbf{x}_n)$ , as  $\beta_j = -\frac{\alpha_n \gamma_j}{\gamma_n}$ ,  $j \neq n$ .

Rather than minimizing the nonlinear function (63), we now devise a computationally attractive approximate solution. It is motivated by the observation that  $(\boldsymbol{\gamma} \cdot K \boldsymbol{\gamma})$  alone is minimized for the eigenvector with minimal eigenvalue, consistent with the special case discussed above (cf. (61)). In that case,  $(\boldsymbol{\gamma} \cdot K \boldsymbol{\gamma}) = \lambda_{\min}$ . More generally, if  $\boldsymbol{\gamma}^i$  is any normalized eigenvector of  $K$ , with eigenvalue  $\lambda_i$ , then

$$\rho(i, n) = \left| \frac{\alpha_n}{\gamma_n^i} \right|^2 \lambda_i. \quad (64)$$

This can be minimized in  $\mathcal{O}(\ell^3)$  operations by performing kernel PCA and scanning through the matrix  $(\rho(i, n))_{i,n}$ . The complexity can be reduced to  $\mathcal{O}(\ell^2)$  by only considering the smallest  $m$  eigenvalues, with  $m < \ell$  chosen a priori. Hence, we can eliminate  $\Phi(\mathbf{x}_n)$ , with  $n$  chosen in a principled yet efficient way.

Setting all computational considerations aside, the *optimal* greedy solution to the above selection problem, equivalent to (63), can also be obtained by using Proposition 8: compute the optimal solution for all possible patterns that one could leave out (i.e. use subsets of  $\{\mathbf{x}_1, \dots, \mathbf{x}_\ell\}$ , of size  $\ell - 1$ , as  $\{\mathbf{z}_1, \dots, \mathbf{z}_m\}$ , and evaluate (57) in each case.

<sup>11</sup>The idea of approximating the support vector expansion by optimally removing individual support vectors, and then adjusting the coefficients of those that remain to minimize the resulting error, was arrived at independently by Olivier Chapelle, who also derived the expression to be minimized, (63) (Private Communication).

The same applies to subsets of any size. If we have the resources to exhaustively scan through all subsets of size  $m$  ( $1 \leq m \leq \ell - 1$ ), then Proposition 8 provides the *optimal* way of selecting the best expansion of a given size. Better expansions can only be obtained if we drop the restriction that the approximation is written in terms of the original patterns, as done in section V-D.

No matter how we end up choosing  $n$ , we approximate  $\Psi$  by

$$\begin{aligned} \Psi &= \sum_{j \neq n} \alpha_j \Phi(\mathbf{x}_j) + \alpha_n \Phi(\mathbf{x}_n) \\ &\approx \sum_{j \neq n} \left( \alpha_j - \frac{\alpha_n \gamma_j}{\gamma_n} \right) \Phi(\mathbf{x}_j). \end{aligned} \quad (65)$$

The whole scheme can be iterated until the expansion of  $\Psi$  is sufficiently sparse. If one wants to avoid having to find the smallest eigenvalue of  $K$  at each step anew, then approximate schemes using heuristics can be conceived.

In our experiments to be reported below, we did compute all eigenvectors at each step, using the Gram matrix computed from the SVs and then selected  $n$  according to (64).

## C.2 Selection via $L_1$ Penalization

We next consider a method for enforcing sparseness which is inspired by  $L_1$  shrinkage penalizers (cf. [10]).

Given some expansion  $\sum_i \alpha_i \Phi(\mathbf{x}_i)$ , we approximate it by  $\sum_i \beta_i \Phi(\mathbf{x}_i)$  by minimizing

$$\left\| \sum_{i=1}^{\ell} \alpha_i \Phi(\mathbf{x}_i) - \sum_{i=1}^{\ell} \beta_i \Phi(\mathbf{x}_i) \right\|^2 + \lambda \sum_{i=1}^{\ell} c_i |\beta_i| \quad (66)$$

over all  $\beta_i$ . Here,  $\lambda > 0$  is a constant determining the trade-off between sparseness and quality of approximation. The constants  $c_i$  can be set to 1 or  $\alpha/|\alpha_i|$  (where  $\alpha$  is the mean of all  $|\alpha_i|$ ), say. In the latter case, we are hoping for a sparser decomposition, since more emphasis is put on shrinking terms which are already small. This reflects the intuition that it is less promising to try to shrink very large terms. Ideally, we would like to *count* the number of nonzero coefficients, rather than sum their moduli; however, the former does not lead to an efficiently solvable optimization problem.

To dispose of the modulus, we rewrite  $\beta_i$  as

$$\beta_i := \beta_i^+ - \beta_i^-, \quad (67)$$

where  $\beta_i^\pm \geq 0$ . In terms of the new variables, we end up with the quadratic programming problem

$$\min_{\beta^+, \beta^-} \sum_{ij} (\beta_i^+ - \beta_i^-) (\beta_j^+ - \beta_j^-) K_{ij} + \quad (68)$$

$$\sum_j \left( \beta_j^+ (\lambda c_j - 2 \sum_i K_{ij} \alpha_i) + \beta_j^- (\lambda c_j + 2 \sum_i K_{ij} \alpha_i) \right)$$

$$\text{subject to } \beta_j^+, \beta_j^- \geq 0. \quad (69)$$

This problem can be solved with standard quadratic programming tools. The solution (67) could be used directly as expansion coefficients. For optimal precision, however, we merely use it to select which patterns to use for the expansion (those with nonzero coefficients), and re-compute the optimal coefficients using Proposition 8.

### C.3 The Multi-Class Case

In many applications, we face the problem of *simultaneously* approximating a set of  $k$  feature space expansions. For instance, in digit classification, a common approach is to train  $k = 10$  binary recognizers, one for each digit. To this end, the quadratic programming formulation of section V-C.2 can be modified to

$$\min_{\beta} \sum_j \left\| \sum_i \alpha_i^j \Phi(\mathbf{x}_i) - \sum_i \beta_i^j \Phi(\mathbf{x}_i) \right\|^2 + \lambda \sum_i c_i \xi_i \quad (70)$$

$$\text{subject to } |\beta_i^j| \leq \xi_i, \quad \xi_i \geq 0. \quad (71)$$

The indices  $i, j$  are understood to range over all SVs (i.e. expansion vectors which have a nonzero coefficient for at least one of the recognizers), and over the  $k$  classes, respectively. For  $c_i$ , we use either 1 or, with the same rationale as above,  $\alpha / \max_j |\alpha_i^j|$  (here,  $\alpha$  is the mean of all  $\max_j |\alpha_i^j|$  over  $i$ ).

The term  $\sum_i c_i \xi_i$  together with the constraint (71) ensures that we only penalize the largest of the coefficients pertaining to each individual SV. Therefore, as soon as one coefficient is being used by one the expansions, it does not cost anything to also use it for another one. This is precisely what we want for speed-up purposes: if we have to compute a certain dot product anyway, we might just as well re-use it for the other expansions.

Using  $\beta_i^j = \beta_i^{+j} - \beta_i^{-j}$ , with  $\beta_i^{\pm j} \geq 0$ , we arrive at the following problem:

$$\min_{\beta^{\pm}} \sum_{j,i,q} (\beta_i^{+j} - \beta_i^{-j})(\beta_q^{+j} - \beta_q^{-j}) K_{iq} + \quad (72)$$

$$2 \sum_{j,i} (\beta_i^{-j} - \beta_i^{+j}) \sum_q K_{qi} \alpha_q^j + \lambda \sum_i c_i \xi_i$$

$$\text{subject to } \beta_i^{+j} + \beta_i^{-j} \leq \xi_i, \quad \beta_i^{\pm j} \geq 0. \quad (73)$$

Here,  $q$  ranges over all SVs.

### C.4 A Note on the Utility of Reduced Set Selection

Why should we expect these procedures to be useful? First, for SV machines, where  $\alpha_j \in [-C, C]$  for some positive value of the regularization constant  $C$ , there is reason to believe that the SV expansion can be made sparser by removing the constraint on  $\alpha_j$  (note that (65) does not care about the constraint).<sup>12</sup> Second, the number of eigenvalues which are zero (and thus the number of patterns

<sup>12</sup>Imagine a case where a certain pattern appears twice in the training set, and the SV expansion has to utilize both of the copies only because the upper bound constraint limits the coefficient of each of them to  $C$ .

that can be removed without loss), or at least small, depends on the problem at hand and the kernel used. For instance, Gaussian kernel Gram matrices do not have zero eigenvalues unless some of the patterns are duplicates [29]. Nevertheless, good approximations are possible, since the eigenvalues of Gaussian kernels decay rapidly (e.g. [6]).

### D. Reduced Set Construction

So far, we have dealt with the problem of how to *select* a reduced set of expansion vectors from the original one. We now return to the originally posed problem, which includes the *construction* of new vectors to reach high reduction rates. To this end, suppose we want to approximate a vector

$$\Psi_1 = \sum_{i=1}^{\ell} \alpha_i \Phi(\mathbf{x}_i) \quad (74)$$

by an expansion of the type (54) with  $N_z > 1$ . To this end, we iterate the procedure of section IV-B by  $\Psi_{m+1} := \Psi_m - \beta_m \Phi(\mathbf{z}_m)$ . Here,  $\mathbf{z}_m$  denotes the  $\mathbf{z}$  found for  $\Psi_m$ , obtained by iterating (50). To apply (50),  $\Psi_m$  needs to be utilized in its representation in terms of mapped input images,

$$\Psi_m = \sum_{i=1}^{\ell} \alpha_i \Phi(\mathbf{x}_i) - \sum_{i=1}^{m-1} \beta_i \Phi(\mathbf{z}_i), \quad (75)$$

i.e. we need to set  $N_x = \ell + m - 1$ , and

$$(\alpha_1, \dots, \alpha_{N_x}) = (\alpha_1, \dots, \alpha_{\ell}, -\beta_1, \dots, -\beta_{m-1}), \quad (76)$$

$$(\mathbf{x}_1, \dots, \mathbf{x}_{N_x}) = (\mathbf{x}_1, \dots, \mathbf{x}_{\ell}, \mathbf{z}_1, \dots, \mathbf{z}_{m-1}). \quad (77)$$

The coefficient  $\beta_m$  could be computed as  $(\Psi \cdot \Phi(\mathbf{z})) / (\Phi(\mathbf{z}) \cdot \Phi(\mathbf{z}))$ . However, if the vectors  $\Phi(\mathbf{z}_1), \dots, \Phi(\mathbf{z}_m)$  are not orthogonal in  $F$ , then the best approximation of  $\Psi_1$  in their span is *not* obtained by computing orthogonal projections onto each direction. Instead, we need to compute the optimal coefficients  $\beta = (\beta_1, \dots, \beta_m)$  anew in each step, using Proposition 8 (if the discrepancy  $\Psi_{m+1}$  has not yet reached zero, then  $K^z$  will be invertible).

The iteration is stopped after  $N_z$  steps, either specified in advance, or by monitoring when  $\|\Psi_{m+1}\|$  (i.e.  $\|\Psi_1 - \sum_{i=1}^m \beta_i \Phi(\mathbf{z}_i)\|$ ) falls below a specified threshold. The solution vector takes the form (54).

We conclude this section by noting that in many cases, such as multiclass SV machines, or multiple kernel PCA feature extractors, we may actually want to approximate several vectors simultaneously. This leads to more complex equations, given in [26].

## VI. EXPERIMENTS

To see how the proposed methods work in practice we ran several toy and real-world experiments. In section VI-A, we give denoising results for the approach of finding approximate pre-images presented in section IV-A. In section VI-B, we present some experiments for the reduced set methods described in section V-C and section V-D.

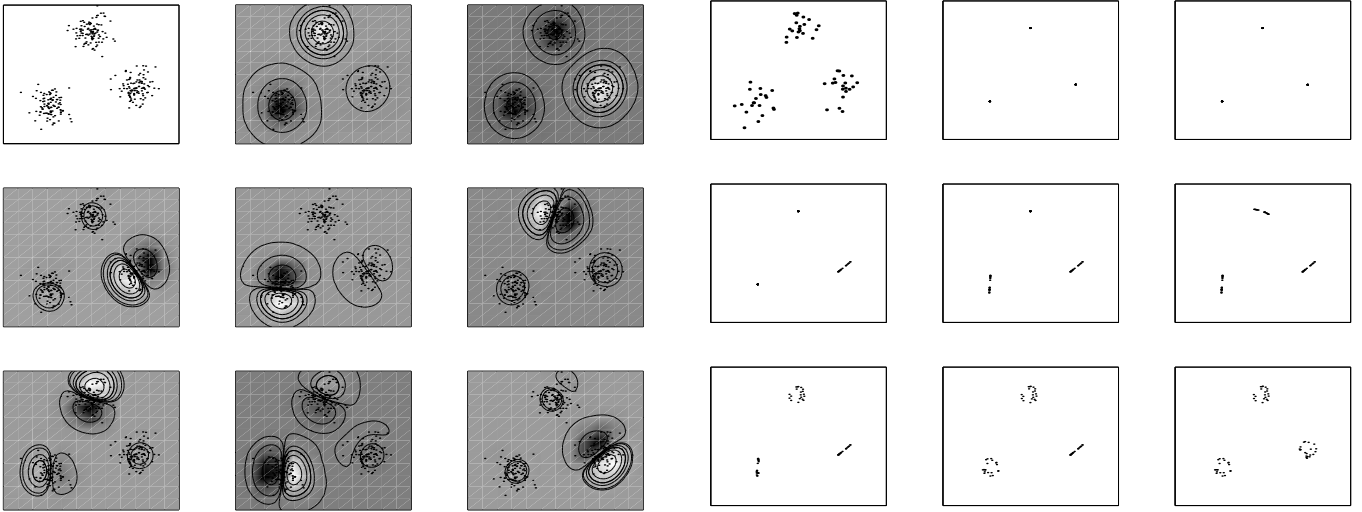


Fig. 4. Kernel PCA toy example (see text): lines of constant feature value for the first 8 nonlinear principal components extracted with  $k(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|^2/0.1)$ . The first 2 principal components (top middle/right) separate the three clusters. Components 3–5 split the clusters. Components 6–8 split them again, orthogonal to the above splits.

## A. Kernel PCA Denoising

### A.1 Toy Examples

All experiments reported were carried out with Gaussian kernels (40), minimizing (41) with the iteration scheme given by (50). However, similar results were obtained with polynomial kernels. Matlab code for computing kernel PCA is available on the web from <http://svm.first.gmd.de>.

We generated an artificial data set from three point sources at  $(-0.5, -0.1)$ ,  $(0, 0.7)$ ,  $(0.5, 0.1)$  (100 points each) with Gaussian noise ( $\sigma = 0.1$ ), and performed kernel PCA on it (figure 4). Using the resulting eigenvectors, we extracted nonlinear principal components from a set of test points generated from the same model, and reconstructed the points from varying numbers of principal components. Figure 5 shows that discarding higher-order components leads to removal of the noise — the points move towards their respective sources.

In a second experiment (table I), we generated a data set from eleven Gaussians in  $\mathbb{R}^{10}$  with zero mean and variance  $\sigma^2$  in each component, by selecting from each source 100 points as a training set and 33 points for a test set (centers of the Gaussians randomly chosen in  $[-1, 1]^{10}$ ). Then we applied kernel PCA to the training set and computed the projections  $\beta_k$  of the points in the test set. With these, we carried out denoising, yielding an approximate pre-image in  $\mathbb{R}^{10}$  for each test point. This procedure was repeated for different numbers of components in reconstruction, and for different values of  $\sigma$  (also used in the kernel). We compared the results provided by our algorithm to those of linear PCA via the mean squared distance of all denoised test points to their corresponding center. Table I shows the *ratio* of these values; here and below, ratios larger than one indicate that kernel PCA performed better than linear PCA. For almost every choice of  $n$  and  $\sigma$ , kernel PCA did

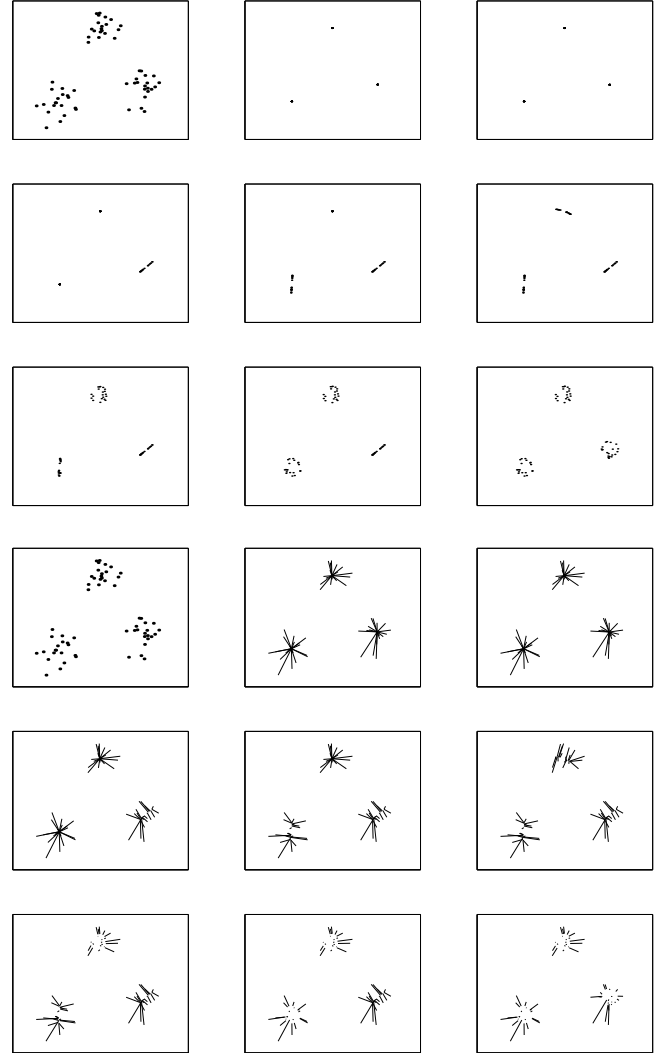


Fig. 5. Kernel PCA denoising by reconstruction from projections onto the eigenvectors of figure 4. We generated 20 new points from each Gaussian, represented them in feature space by their first  $n = 1, 2, \dots, 8$  nonlinear principal components, and computed approximate pre-images, shown in the upper 9 pictures (top left: original data, top middle:  $n = 1$ , top right:  $n = 2$ , etc.). Note that by discarding higher order principal components (i.e. using a small  $n$ ), we remove the noise inherent in the nonzero variance  $\sigma^2$  of the Gaussians. The lower 9 pictures show how the original points “move” in the denoising. Unlike the corresponding case in linear PCA, where where we obtain lines (see figure 6), in kernel PCA clusters shrink to points.

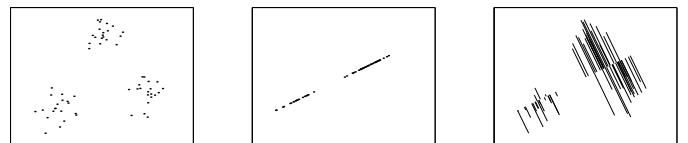


Fig. 6. Reconstructions and point movements (cf. figure 5) for linear PCA, based on the first principal component.

better. Note that using all 10 components, linear PCA is just a basis transformation and hence cannot denoise. The extreme superiority of kernel PCA for small  $\sigma$  is due to the fact that all test points are in this case located close to the eleven spots in input space, and linear PCA has to cover

TABLE I

DENOISING GAUSSIANS IN  $\mathbb{R}^{10}$  (SEE TEXT). PERFORMANCE RATIOS LARGER THAN ONE INDICATE HOW MUCH BETTER KERNEL PCA DID, COMPARED TO LINEAR PCA, FOR DIFFERENT CHOICES OF THE GAUSSIANS' STANDARD DEVIATION  $\sigma$ , AND DIFFERENT NUMBERS OF COMPONENTS USED IN RECONSTRUCTION.

$\sigma$	$n=1$	2	3	4	5	6	7	8	9
0.05	2058.42	1238.36	846.14	565.41	309.64	170.36	125.97	104.40	92.23
0.1	10.22	31.32	21.51	29.24	27.66	23.53	29.64	40.07	63.41
0.2	0.99	1.12	1.18	1.50	2.11	2.73	3.72	5.09	6.32
0.4	1.07	1.26	1.44	1.64	1.91	2.08	2.22	2.34	2.47
0.8	1.23	1.39	1.54	1.70	1.80	1.96	2.10	2.25	2.39

them with less than ten directions. Kernel PCA moves each point to the correct source even when using only a small number of components.

To get some intuitive understanding in a low-dimensional case, figure 7 depicts the results of denoising a half circle and a square in the plane, using kernel PCA, a *nonlinear autoencoder*, *principal curves*, and linear PCA. The principal curves algorithm [34] iteratively estimates a curve capturing the structure of the data. The data are projected to the closest point on a curve which the algorithm tries to construct such that each point is the average of all data points projecting onto it. It can be shown that the only straight lines satisfying the latter are principal components, so principal curves are a generalization of the latter. The algorithm uses a smoothing parameter which is annealed during the iteration. In the nonlinear autoencoder algorithm, a ‘bottleneck’ 5-layer network is trained to reproduce the input values as outputs (i.e. it is used in autoassociative mode). The hidden unit activations in the third layer form a lower-dimensional representation of the data, closely related to PCA (see for instance [35]). Training is done by conjugate gradient descent. In all algorithms, parameter values were selected such that the best possible denoising result was obtained. The figure shows that on the closed square problem, kernel PCA does (subjectively) best, followed by principal curves and the nonlinear autoencoder; linear PCA fails completely. However, note that all algorithms except for kernel PCA actually provide an explicit one-dimensional parameterization of the data, whereas kernel PCA only provides us with a means of mapping points to their denoised versions (in this case, we used four kernel PCA features, and hence obtain a four-dimensional parameterization).

## A.2 Handwritten Digit Denoising

To test our approach on real-world data, we also applied the algorithm to the USPS database of handwritten digits (e.g. [36], [9]) of 7291 training patterns and 2007 test patterns (size  $16 \times 16$ ). For each of the ten digits, we randomly chose 300 examples from the training set and 50 examples from the test set. We used the method of section IV-B, with  $k(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|^2 / (0.5 \cdot 16^2))$ . The width 0.5 equals twice the average of the data’s variance in each dimension. In figure 8, we give two possible depictions of the eigenvectors found by kernel PCA, compared to those

found by linear PCA for the USPS set. The second row shows the approximate pre-images of the eigenvectors  $\mathbf{V}^k$ ,  $k = 2^0, \dots, 2^8$ , found by our iterative algorithm. In the third row each image is computed as follows: Pixel  $i$  is the projection of the  $\Phi$ -image of the  $i$ -th canonical basis vector in input space onto the corresponding eigenvector in features space (upper left  $\Phi(\mathbf{e}_1) \cdot \mathbf{V}^k$ , lower right  $\Phi(\mathbf{e}_{256}) \cdot \mathbf{V}^k$ ). In the linear case, both methods would simply yield the eigenvectors of linear PCA depicted in the first row; in this sense, they may be considered as generalized eigenvectors in input space. We see that the first eigenvectors are almost identical (except for arbitrary signs). However, we also see that eigenvectors in linear PCA start to focus on high-frequency structures already at smaller eigenvalue size. To understand this, note that in linear PCA we only have a maximum number of 256 eigenvectors, contrary to kernel PCA which gives us the number of training examples (here 3000) possible eigenvectors.

This also explains some of the results we found when working with the USPS set (figures 9 and 10). In these experiments, linear and kernel PCA were trained with the original data. To the test set, we added

- (i) additive Gaussian noise with zero mean and standard deviation  $\sigma = 0.5$ , or
- (ii) ‘speckle’ noise, where each pixel is flipped to black or white with probability  $p = 0.2$ .

For the noisy test sets, we computed the projections onto the first  $n$  linear and nonlinear components, and carried out reconstruction for each case. The results were compared by taking the mean squared distance of each reconstructed digit of the noisy test set to its original counterpart.

For the optimal number of components in linear and kernel PCA, our approach did better by a factor of 1.6 for the Gaussian noise, and 1.2 for the ‘speckle’ noise (the opti-

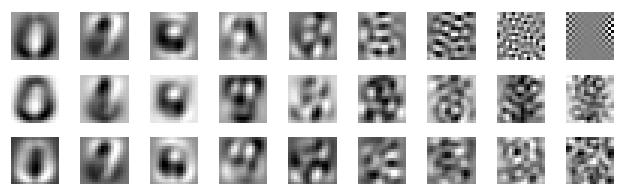


Fig. 8. Visualization of eigenvectors (see text). Depicted are the  $2^0, \dots, 2^8$ -th eigenvector (from left to right). First row: linear PCA, second and third row: different visualizations for kernel PCA.

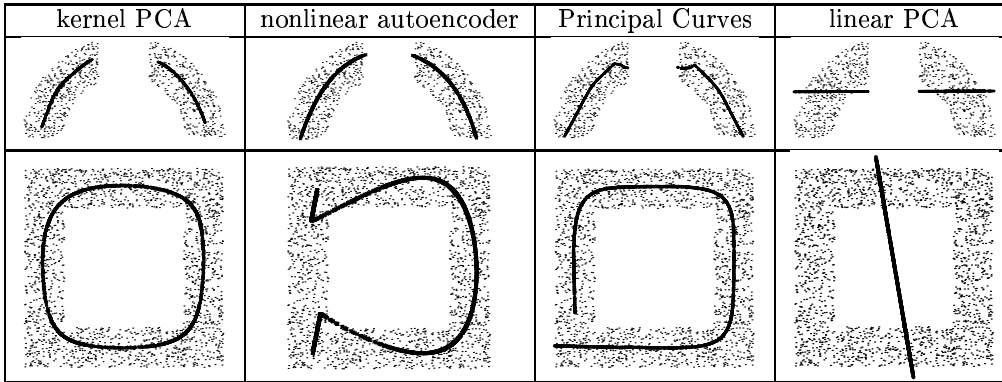


Fig. 7. Denoising in 2-d (see text). Depicted are the data set (small points) and its denoised version (big points, joining up to solid lines). For linear PCA, we used one component for reconstruction, as using two components, reconstruction is perfect and thus does not denoise. Note that all algorithms except for our approach have problems in capturing the circular structure in the bottom example (taken from [31]).

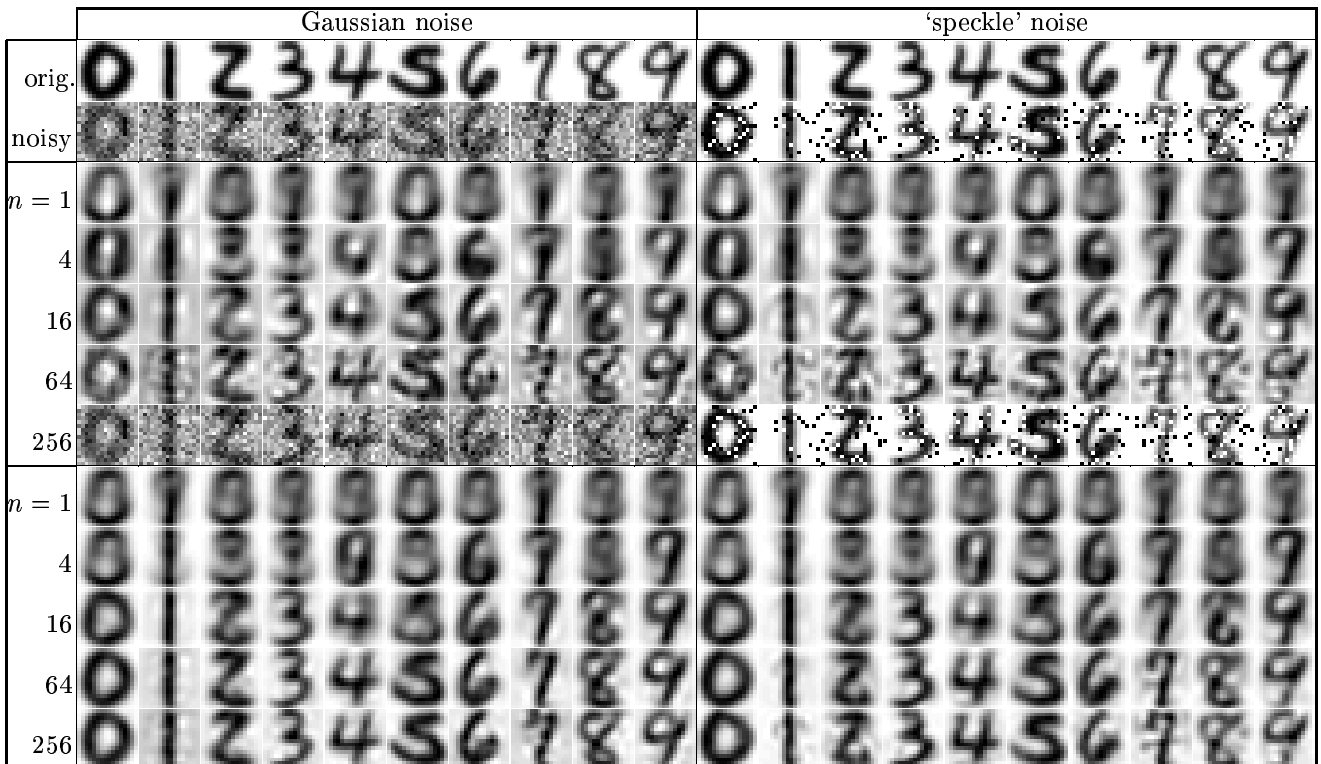


Fig. 9. De-noising of USPS data (see text). The left half shows: *top*: the first occurrence of each digit in the test set, *second row*: the upper digit with additive Gaussian noise ( $\sigma = 0.5$ ), *following five rows*: the reconstruction for linear PCA using  $n = 1, 4, 16, 64, 256$  components, and, *last five rows*: the results of our approach using the same number of components. In the right half we show the same but for 'speckle' noise with probability  $p = 0.2$ .

mal number of components were 32 in linear PCA, and 512 and 256 in kernel PCA, respectively). Taking identical numbers of components in both algorithms, kernel PCA becomes up to 8 times better than linear PCA. However, note that kernel PCA comes with a higher computational complexity.

### B. Speeding up Support Vector Decision Rules

As in section VI-A, we used the USPS handwritten digit database. We approximated the SV expansions (30) of ten binary classifiers, each trained to separate one digit from the rest. We used the Gaussian kernel  $k(\mathbf{x}, \mathbf{y}) =$

$\exp(-\|\mathbf{x} - \mathbf{y}\|^2 / (0.5 \cdot 16^2))$ , and the approximation techniques described in sections V-C and V-D.

The original SV system had on average 254 SVs per classifier. Tables II and III show the classification error results for approximation using the reduced set selection techniques described in section V-C, while table IV gives results for using the reduced set construction method described in V-D, for varying numbers of RS vectors. Shown in each line is the number of misclassified digits for each single classifier and the error of the combined 10-class machine. In all RS systems, the optimal SV threshold was re-computed on the training set.

TABLE II

NUMBERS OF TEST ERRORS FOR EACH BINARY RECOGNIZER, AND TEST ERROR RATES FOR 10-CLASS CLASSIFICATION USING THE RS METHOD OF SECTION V-C.1. TOP: NUMBERS OF SVs FOR THE ORIGINAL SV RBF-SYSTEM. BOTTOM, FIRST ROW: ORIGINAL SV SYSTEM, WITH 254 SVs ON AVERAGE; FOLLOWING ROWS: SYSTEMS WITH VARYING AVERAGE NUMBERS OF RS VECTORS. IN THE SYSTEM RSS- $n$ , EQUAL FRACTIONS OF SVs WERE REMOVED FROM EACH RECOGNIZER SUCH THAT ON AVERAGE,  $n$  RS VECTORS WERE LEFT.

digit	0	1	2	3	4	5	6	7	8	9	ave.
#SVs	219	91	316	309	288	340	213	206	304	250	254
	0	1	2	3	4	5	6	7	8	9	10-class
SV-254	16	13	30	17	32	22	11	12	26	17	4.4%
RSS-10	154	262	205	162	163	161	169	96	146	177	59.5%
RSS-15	171	107	200	130	153	158	173	93	143	177	47.7%
RSS-20	149	82	171	111	153	160	147	116	161	177	44.0%
RSS-25	120	82	143	141	188	161	143	96	152	177	45.0%
RSS-50	47	18	70	52	192	95	54	38	97	157	17.6%
RSS-75	23	15	36	30	65	47	21	29	56	41	7.0%
RSS-100	19	15	42	22	40	29	14	18	37	27	5.5%
RSS-125	15	13	29	26	32	20	15	14	32	19	4.5%
RSS-150	18	12	28	21	35	32	9	15	23	14	4.5%
RSS-175	17	13	28	22	31	25	12	15	23	18	4.5%
RSS-200	14	13	27	25	27	26	11	13	26	21	4.5%
RSS-225	15	13	26	23	30	27	11	14	25	20	4.2%
RSS-250	15	13	27	24	32	28	11	14	25	18	4.3%

TABLE III

NUMBERS OF TEST ERRORS FOR EACH BINARY RECOGNIZER, AND TEST ERROR RATES FOR 10-CLASS CLASSIFICATION USING THE RS METHOD OF SECTION V-C.2 ( $c_i = \alpha/|\alpha_i|$ ). FIRST ROW: ORIGINAL SV SYSTEM, WITH 254 SVs ON AVERAGE; FOLLOWING ROWS: SYSTEMS WITH VARYING AVERAGE NUMBERS OF RS VECTORS. IN THE SYSTEM RSS<sub>2</sub>- $n$ ,  $\lambda$  WAS ADJUSTED SUCH THAT THE AVERAGE NUMBER OF RS VECTORS LEFT WAS  $n$  (THE CONSTANT  $\lambda$ , GIVEN IN PARENTHESES, WAS CHOSEN SUCH THAT THE NUMBERS  $n$  WERE COMPARABLE TO TABLE II). THE RESULTS CAN BE FURTHER IMPROVED USING THE METHOD OF SECTION V-C.3 ( $c_i = \alpha/\max_j |\alpha_i^j|$ ). FOR INSTANCE, USING ABOUT 570 EXPANSION VECTORS (WHICH IS THE SAME NUMBER THAT WE GET WHEN TAKING THE UNION OF ALL SVs IN THE RSS<sub>2</sub> - 74 SYSTEM), THIS LED TO AN IMPROVED ERROR RATE OF 5.5%.

digit		0	1	2	3	4	5	6	7	8	9	10-class
SV-254		16	13	30	17	32	22	11	12	26	17	4.4%
RSS <sub>2</sub> -30	(4.00)	319	25	198	169	161	152	116	147	95	118	51.5%
RSS <sub>2</sub> -50	(3.34)	225	24	171	146	149	124	94	147	100	101	28.5%
RSS <sub>2</sub> -74	(2.55)	113	25	100	100	120	95	40	147	83	50	10.8%
RSS <sub>2</sub> -101	(1.73)	38	21	46	64	81	54	23	143	49	37	5.9%
RSS <sub>2</sub> -126	(1.06)	21	19	35	34	47	27	15	72	41	23	5.3%
RSS <sub>2</sub> -151	(0.62)	19	20	30	24	31	30	10	27	33	18	4.5%
RSS <sub>2</sub> -174	(0.33)	16	15	24	26	31	26	10	19	30	19	4.3%
RSS <sub>2</sub> -200	(0.13)	17	15	25	27	34	27	11	14	26	22	4.3%
RSS <sub>2</sub> -224	(0.04)	16	15	26	24	32	27	11	14	28	19	4.3%
RSS <sub>2</sub> -234	(0.02)	16	14	26	24	32	28	11	14	26	19	4.3%

For the method of section V-C, RSS- $n$  means, that for each binary classifier, we removed support vectors until on average  $n$  were left; for the method of section V-D, that we approximated each decision function using  $n$  vectors (i.e.  $n$  steps using the described iteration procedure). For large numbers  $n$ , i.e. small reductions of the decision functions complexity, the accuracy of the original system can be approached closely with both techniques. In Tables II and III, we see that removal of about 40% of the support vectors leaves the error practically unchanged. Reducing

the numbers of support vectors further can lead to large performance losses.

The reduced set construction method, which is computationally and conceptually more complex, performs better in this situation as it is able to utilize vectors different from the original support patterns in the expansion. To get a speedup by a factor of 10, we have to use a system with 25 RS vectors (RSC-25). For the method in table IV, the classification accuracy only drops moderately from 4.4% to 5.1%, which is still competitive with convolutional neural

TABLE IV

NUMBER OF TEST ERRORS FOR EACH BINARY RECOGNIZER, AND TEST ERROR RATES FOR 10-CLASS CLASSIFICATION USING THE RS CONSTRUCTION METHOD OF SECTION V-D. FIRST ROW: ORIGINAL SV SYSTEM, WITH 254 SVs ON AVERAGE (SEE ALSO TABLES II AND III); FOLLOWING ROWS: SYSTEMS WITH VARYING NUMBERS OF RS VECTORS (RSC- $n$  STANDS FOR  $n$  VECTORS CONSTRUCTED) PER BINARY RECOGNIZER, COMPUTED BY ITERATING ONE-TERM APPROXIMATIONS, SEPARATELY FOR EACH RECOGNIZER. LAST TWO ROWS: WITH A SUBSEQUENT GLOBAL GRADIENT DESCENT, THE RESULTS CAN BE FURTHER IMPROVED (SEE TEXT).

digit	0	1	2	3	4	5	6	7	8	9	10-class
SV-254	16	13	30	17	32	22	11	12	26	17	4.4%
RSC-10	26	13	45	49	35	54	22	24	39	24	7.1%
RSC-15	17	16	43	50	49	37	14	18	45	35	6.4%
RSC-20	27	11	38	30	35	43	12	16	30	25	5.6%
RSC-25	21	12	38	32	31	22	12	18	33	28	5.1%
RSC-50	18	10	33	28	32	23	12	15	35	27	5.0%
RSC-100	14	13	26	22	30	26	11	14	28	23	4.8%
RSC-150	13	14	28	32	27	24	12	14	29	26	4.7%
RSC-200	14	13	28	28	29	24	10	15	26	26	4.9%
RSC-250	12	13	26	26	32	25	11	14	26	24	4.6%
RSC <sub>2</sub> -25	14	14	31	22	30	23	11	14	26	17	4.7%
RSC <sub>2</sub> <sup>gd</sup> -25	16	13	32	19	31	26	11	15	25	18	5.0%



Fig. 11. Complete display of reduced set vectors constructed by the iterative approach of section V-D for  $n = 20$ , with coefficients (Top: recognizer of digit 0, ..., bottom: digit 9). Note that positive coefficients (roughly) correspond to positive examples in the classification problem.

networks on that data base [36]. Moreover, we can further improve this result by adding the second phase of the traditional RS algorithm, where a global gradient descent is performed in the space of all  $(\mathbf{z}_i, \beta_i)$  [24], [25] (computationally more expensive than the first phase by about two orders of magnitude): this led to an error rate of 4.7%. For the considered kernel, this is almost identical to the traditional RS method, which yielded 5.0% (for polynomial kernels, the latter method led to 4.3% at the same speedup

[24]). Note that the traditional RS method restarts the second phase many times to make sure that the global minimum of the cost function is actually found, which makes it plausible that the final results are similar.

Finally, figure 11 shows the RS-20 vectors of the 10 binary classifiers. As an aside, note that unlike the approach of [24], our algorithm produces images which do look meaningful (i.e. digit-like).



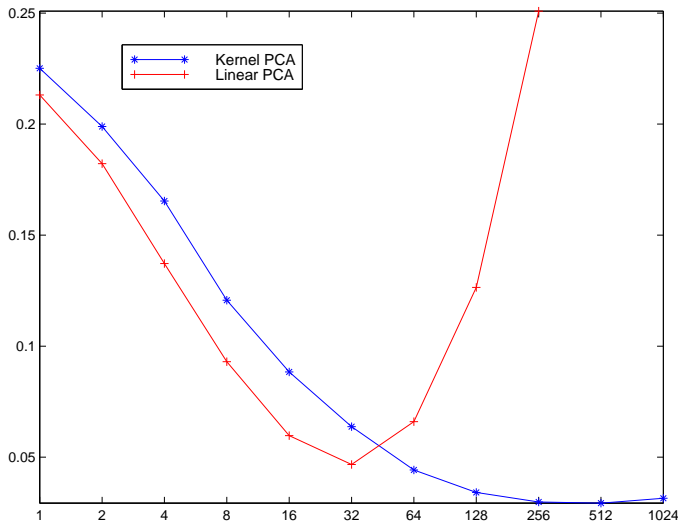


Fig. 10. Mean squared error of denoised images vs. # of features used, for kernel PCA and linear PCA. Kernel PCA exploits nonlinearities and has the potential to utilize more features to code structure rather than noise. Therefore, it outperforms linear PCA denoising if a sufficiently large number of features is used.

## VII. DISCUSSION

Algorithms utilizing Mercer kernels construct their solutions as expansions  $\Psi = \sum_{i=1}^{\ell} \alpha_i \Phi(\mathbf{x}_i)$  in terms of mapped input patterns. However, the map  $\Phi$  is often unknown, or too complex to provide any intuition about the solution  $\Psi$ . This has motivated our efforts to reduce the complexity of the expansion, summarized in this paper.

As an extreme case, we have first studied how to approximate  $\Psi$  by a single  $\Phi(\mathbf{z})$  (i.e. how to find an approximate pre-image of  $\Psi$ ), and proposed a fixed-point iteration algorithm to perform the task.

In situations where no good approximate pre-image exists, one can still reduce the complexity of  $\sum_{i=1}^{\ell} \alpha_i \Phi(\mathbf{x}_i)$  by expressing it as a sparser expansion  $\sum_{i=1}^m \beta_i \Phi(\mathbf{z}_i)$  ( $m < \ell$ ). We have proposed methods for computing the optimal coefficients  $\beta_i$  and for coming up with suitable patterns  $\mathbf{z}_i$  either by selecting among the  $\mathbf{x}_i$  or by iterating the above pre-image algorithm.

Both types of approximations are of theoretical interest for feature space methods; however, they also lead to practical applications. In this paper, we have considered two such applications, namely, the problem of statistical denoising via kernel PCA reconstruction, and the problem of speeding up SV decision rules. We address them in turn.

*Kernel PCA de-noising.* In de-noising experiments on real-world data, we obtained results significantly better than using linear PCA. Our interpretation of this finding is as follows. Linear PCA can extract at most  $N$  components, where  $N$  is the dimensionality of the data. Being a basis transform, all  $N$  components together fully describe the data. If the data are noisy, this implies that a certain fraction of the components will be devoted to the extraction of noise. Kernel PCA, on the other hand, allows the extraction of up to  $\ell$  features, where  $\ell$  is the number of training examples. Accordingly, kernel PCA can provide a

larger number of features carrying information about the structure in the data (in our experiments, we had  $\ell > N$ ). In addition, if the structure to be extracted is nonlinear, then linear PCA must necessarily fail, as we have illustrated with toy examples.

Open questions and problems include the choice of a suitable kernel for a given noise reduction problem, possibly in conjunction with the regularization properties of the kernel (e.g. [6]), the application of the approach to compression, and the comparison (and connection) to alternative nonlinear denoising methods (cf. [37]).

*Speeding up SV machines.* We have shown experimentally that our approximation algorithms can be used to speed up SV machines significantly. Note that in the Gaussian RBF case, the approximation can never be as good as the original, since the kernel matrix  $K_{ij} = (k(\mathbf{x}_i, \mathbf{x}_j))$  has full rank [29].

As in [25], good *RS construction* results were obtained even though the objective function did not decrease to zero (in our RS construction experiments, it was reduced by a factor of 2 to 20 in the first phase, depending on how many RS vectors were computed; the global gradient descent yielded another factor 2 – 3). We conjecture that this is due to the following: in classification, we are not interested in  $\|\Psi - \Psi'\|$ , but in  $\int \text{sgn}(\sum_{i=1}^{N_z} \alpha_i k(\mathbf{x}, \mathbf{x}_i) + b) - \text{sgn}(\sum_{j=1}^{N_z} \beta_j k(\mathbf{x}, \mathbf{z}_j) + \tilde{b}) dP(\mathbf{x})$ , where  $P$  is the underlying probability distribution of the patterns (cf. [38]). This is consistent with the fact that the performance of a RS SV classifier can be improved by re-computing an optimal threshold  $b$ .

The previous RS construction method [24], [25] can be used for any SV kernel; the new one is limited to  $k(\mathbf{x}, \mathbf{y}) = k(\|\mathbf{x} - \mathbf{y}\|^2)$ . However, it is fast, and it led to interpretable RS images and an interesting connection between clustering and approximation in feature spaces. It appears intriguing to pursue the question whether this connection could be exploited to form more general types of approximations of SV and kernel PCA expansions by making use of Gaussians of variable widths.

The *RS selection* methods, on the other hand, are applicable for any SV kernel. In our experiments, they led to worse reduction rates than RS construction; however, they are simpler and computationally faster. Among the first two RS selection methods, the one described in section V-C.1 was slightly superior at higher reductions; however, the one given in section V-C.2 is computationally cheaper since unlike the former, it does not remove the SVs one at a time and therefore it need not be iterated. Moreover, we found that it can be improved by simultaneously approximating several vectors, corresponding to the 10 binary recognizers in a digit recognition task.

The proposed methods are applicable to any feature space algorithm based on Mercer kernels. For instance, we could also speed up SV regression machines or kernel PCA feature extractors. Moreover, we expect further possibilities to open up in the future, as Mercer kernel methods are being applied in an increasing number of learning and signal processing problems.

## REFERENCES

- [1] S. Saitoh, *Theory of Reproducing Kernels and its Applications*, Longman Scientific & Technical, Harlow, England, 1988.
- [2] B. E. Boser, I. M. Guyon, and V. N. Vapnik, "A training algorithm for optimal margin classifiers," in *Proceedings of the 5th Annual ACM Workshop on Computational Learning Theory*, D. Haussler, Ed., Pittsburgh, PA, July 1992, pp. 144–152, ACM Press.
- [3] M. Aizerman, E. Braverman, and L. Rozonoer, "Theoretical foundations of the potential function method in pattern recognition learning," *Automation and Remote Control*, vol. 25, pp. 821 – 837, 1964.
- [4] V. Vapnik, *The Nature of Statistical Learning Theory*, Springer Verlag, New York, 1995.
- [5] B. Schölkopf, A. Smola, and K.-R. Müller, "Nonlinear component analysis as a kernel eigenvalue problem," *Neural Computation*, vol. 10, pp. 1299 – 1319, 1998.
- [6] R. C. Williamson, A. J. Smola, and B. Schölkopf, "Generalization performance of regularization networks and support vector machines via entropy numbers of compact operators," Tech. Rep. 19, NeuroCOLT, <http://www.neurocolt.com>, 1998.
- [7] N. Aronszajn, "Theory of reproducing kernels," *Transactions of the American Mathematical Society*, vol. 68, pp. 337 – 404, 1950.
- [8] G. Wahba, *Spline Models for Observational Data*, vol. 59 of *CBMS-NSF Regional Conference Series in Applied Mathematics*, SIAM, Philadelphia, 1990.
- [9] B. Schölkopf, *Support Vector Learning*, R. Oldenbourg Verlag, München, 1997, Doktorarbeit, TU Berlin.
- [10] F. Girosi, "An equivalence between sparse approximation and support vector machines," *Neural Computation*, vol. 10, no. 6, pp. 1455–1480, 1998.
- [11] T. J. Hastie and R. J. Tibshirani, *Generalized Additive Models*, vol. 43 of *Monographs on Statistics and Applied Probability*, Chapman & Hall, London, 1990.
- [12] D. Mattera, "personal communication," 1998.
- [13] S. Mika, "Nichtlineare Signalverarbeitung in Feature-Räumen," Diplomarbeit, Technische Universität Berlin, 1998.
- [14] K. Tsuda, "Support vector classifier with asymmetric kernel function," in *Proceedings ESANN*, M. Verleysen, Ed., Brussels, 1999, pp. 183 – 188, D Facto.
- [15] B. Schölkopf, P. Simard, A. Smola, and V. Vapnik, "Prior knowledge in support vector kernels," in *Advances in Neural Information Processing Systems 10*, M. Jordan, M. Kearns, and S. Solla, Eds., Cambridge, MA, 1998, pp. 640 – 646, MIT Press.
- [16] A. Smola and B. Schölkopf, "From regularization operators to support vector kernels," in *Advances in Neural Information Processing Systems 10*, M. Jordan, M. Kearns, and S. Solla, Eds., Cambridge, MA, 1998, pp. 343 – 349, MIT Press.
- [17] V. Vapnik, *Statistical Learning Theory*, Wiley, New York, 1998.
- [18] B. Carl and I. Stephani, *Entropy, compactness, and the approximation of operators*, Cambridge University Press, Cambridge, UK, 1990.
- [19] R. Williamson, A. Smola, and B. Schölkopf, "Entropy numbers, operators and support vector kernels," in *Advances in Kernel Methods — Support Vector Learning*, B. Schölkopf, C. Burges, and A. Smola, Eds., pp. 127 – 144, MIT Press, Cambridge, MA, 1999.
- [20] B. Schölkopf, J. Shawe-Taylor, A. J. Smola, and R. C. Williamson, "Kernel-dependent support vector error bounds," in *Proceedings ICANN*, 1999, to appear.
- [21] C. J. C. Burges, "Geometry and invariance in kernel based methods," in *Advances in Kernel Methods — Support Vector Learning*, B. Schölkopf, C. Burges, and A. Smola, Eds., Cambridge, MA, 1999, pp. 89 – 116, MIT Press.
- [22] C. J. C. Burges, "A tutorial on support vector machines for pattern recognition," *Data Mining and Knowledge Discovery*, vol. 2, no. 2, pp. 1–47, 1998.
- [23] A. Smola and B. Schölkopf, "A tutorial on support vector regression," <http://svm.first.gmd.de>, 1998.
- [24] C. J. C. Burges, "Simplified support vector decision rules," in *Proceedings, 13th Intl. Conf. on Machine Learning*, L. Saitta, Ed., San Mateo, CA, 1996, pp. 71–77, Morgan Kaufmann.
- [25] C. J. C. Burges and B. Schölkopf, "Improving the accuracy and speed of support vector learning machines," in *Advances in Neural Information Processing Systems 9*, M. Mozer, M. Jordan, and T. Petsche, Eds., Cambridge, MA, 1997, pp. 375–381, MIT Press.
- [26] B. Schölkopf, P. Knirsch, A. Smola, and C. Burges, "Fast approximation of support vector kernel expansions, and an interpretation of clustering as approximation in feature spaces," in *Mustererkennung 1998 — 20. DAGM-Symposium*, P. Levi, M. Schanz, R.-J. Ahlers, and F. May, Eds., Berlin, 1998, Informatik aktuell, pp. 124 – 132, Springer.
- [27] T. Frieß and R. F. Harrison, "Linear programming support vector machines for pattern classification and regression estimation; and the SR algorithm: improving speed and tightness of VC bounds in SV algorithms," Research Report 706, University of Sheffield, Dept. Automatic Control and Systems Engineering, 1998.
- [28] E. Osuna and F. Girosi, "Reducing run-time complexity in support vector machines," in *Advances in Kernel Methods — Support Vector Learning*, B. Schölkopf, C. Burges, and A. Smola, Eds., pp. 271 – 283, MIT Press, Cambridge, MA, 1999.
- [29] C. A. Micchelli, "Interpolation of scattered data: distance matrices and conditionally positive definite functions," *Constructive Approximation*, vol. 2, pp. 11–22, 1986.
- [30] B. Schölkopf, S. Mika, A. Smola, G. Rätsch, and K.-R. Müller, "Kernel PCA pattern reconstruction via approximate pre-images," in *Proceedings of the 8th International Conference on Artificial Neural Networks*, L. Niklasson, M. Bodén, and T. Ziemke, Eds., Berlin, 1998, Perspectives in Neural Computing, pp. 147 – 152, Springer Verlag.
- [31] S. Mika, B. Schölkopf, A. Smola, K.-R. Müller, M. Scholz, and G. Rätsch, "Kernel PCA and de-noising in feature spaces," in *Advances in Neural Information Processing Systems 11*, 1999.
- [32] B. Schölkopf, K. Sung, C. Burges, F. Girosi, P. Niyogi, T. Poggio, and V. Vapnik, "Comparing support vector machines with Gaussian kernels to radial basis function classifiers," *IEEE Trans. Sign. Processing*, vol. 45, pp. 2758 – 2765, 1997.
- [33] J. M. Buhmann, "Data clustering and learning," in *The Handbook of Brain Theory and Neural Networks*, M. A. Arbib, Ed., pp. 278–281, MIT Press, 1995.
- [34] T. Hastie and W. Stuetzle, "Principal curves," *JASA*, vol. 84, pp. 502 – 516, 1989.
- [35] K. I. Diamantaras and S. Y. Kung, *Principal Component Neural Networks*, Wiley, New York, 1996.
- [36] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. J. Jackel, "Backpropagation applied to handwritten zip code recognition," *Neural Computation*, vol. 1, pp. 541 – 551, 1989.
- [37] S. Mallat and Z. Zhang, "Matching Pursuit in a time-frequency dictionary," *IEEE Transactions on Signal Processing*, vol. 41, pp. 3397–3415, 1993.
- [38] K. L. Blackmore, R. C. Williamson, and I. M. Y. Mareels, "Decision region approximation by polynomials or neural networks," *IEEE Transactions in Information Theory*, vol. 43, pp. 903 – 907, 1997.



**Bernhard Schölkopf** received an M.Sc. in mathematics and the Lionel Cooper Memorial Prize from the University of London in 1992, followed in 1994 by the Diplom in physics from the Eberhard-Karls-Universität, Tübingen (Germany), with a thesis written at the Max-Planck-Institute for biological cybernetics. Three years later, he obtained a Ph.D. in computer science from the Technical University Berlin. His thesis on Support Vector Learning won the prize of the German Association for Computer Science (GI) 1997. He has worked at AT&T Bell Labs (with V. Vapnik) and at the Australian National University. At present, he is a tenured researcher at GMD FIRST, Berlin. His scientific interests include machine learning and perception.



**Sebastian Mika** is a doctoral student at GMD FIRST, Berlin. He received the Diplom in computer science from the Technical University of Berlin in 1998. His scientific interests are in the fields of Machine Learning and Kernel methods.



**Gunnar Rätsch** is a doctoral student at GMD FIRST, Berlin. He received the Diplom in computer science from the University of Potsdam in 1998, along with the prize for the best student of the faculty of Natural Sciences. His scientific interests are in the fields of Boosting and Kernel methods.



**Chris Burges** is a Distinguished Member of Technical Staff at Bell Laboratories, Lucent Technologies. Educated as a physicist, he joined AT&T in 1986 and worked on network performance and routing algorithms. He moved to applied neural network research in 1990 and has worked on handwriting and machine print recognition and speaker identification. For the last several years he has concentrated on the theory and application of SVMs.



**Alexander Smola** is a Postdoc at GMD FIRST, Berlin. He received the Diplom in physics from the Technical University of Munich in 1996. Two years later, he earned a Ph.D. in computer science from the Technical University of Berlin, with a thesis on Learning with Kernels. During his studies, which were supported by the Stiftung Maximilaneum, he spent a year at AT&T Bell Labs and at the Collegio Ghislieri in Pavia (Italy). In addition, he has spent two spells at the Australian National University. His scientific goal is to make generalization error bounds of statistical learning theory applicable in practice.



**Philipp Knirsch** is a computer science student at the University of Tübingen, Germany. He has worked for various companies in the IT area during the last 10 years, including but not limited to AT&T Bell Labs, Hewlett-Packard, 21TORR medienDesign GmbH, Festo AG and the Max-Planck-Institute. He has worked in several areas of computer science, such as computer graphics, support vector learning machines, network programming, system administration, code optimization, databases, compression technology, wavelet analysis and others.

His current work and interest are in the fields of high performance clustered network servers, real time presentation of visual stimuli on low end Linux workstations, and optimizations and improvements of reduced set vector methods for support vector machines.



**Klaus-Robert Müller** received the Diplom degree in mathematical physics 1989 and the Ph.D. in theoretical computer science in 1992, both from University of Karlsruhe, Germany. From 1992 to 1994 he was a Postdoc at the Research Institute for Computer Architecture and Software Technology of the German National Research Center for Information Technology (GMD FIRST) in Berlin. From 1994 to 1995 he was a European Community STP Research Fellow at University of Tokyo in Prof.

Amari's Lab. Since 1995 he has a tenure position and is department head of the intelligent data analysis (IDA) group at GMD FIRST in Berlin. The IDA group has twice won the price for the best scientific and technical project within GMD (1996 and 1998). He is currently lecturing at Humboldt University and Technical University Berlin. He has worked on statistical physics and statistical learning theory of neural networks and time-series analysis. His present interests are expanded to support vector learning machines, boosting, non-stationary blind separation techniques and recently also medical data analysis (MEG, EEG).