

Generalization Performance of Regularization Networks and Support Vector Machines *via* Entropy Numbers of Compact Operators*

Robert C. Williamson[†] Alex J. Smola^{†¶}
Bob.Williamson@anu.edu.au smola@first.gmd.de

Bernhard Schölkopf^{†¶}
bs@first.gmd.de

[†]Department of Engineering, Australian National University,
Canberra, ACT 0200, Australia

[¶]GMD FIRST, Rudower Chaussee 5,
12489 Berlin, Germany

March 11, 1999

Abstract

We derive new bounds for the generalization error of kernel machines, such as support vector machines and related regularization networks by obtaining new bounds on their covering numbers. The proofs make use of a viewpoint that is apparently novel in the field of statistical learning theory. The hypothesis class is described in terms of a linear operator mapping from a possibly infinite dimensional unit ball in feature space into a finite dimensional space. The covering numbers of the class are then determined via the entropy numbers of the operator. These numbers, which characterize the degree of compactness of the operator, can be bounded in terms of the eigenvalues of an integral operator induced by the kernel function used by the machine. As a consequence we are able to theoretically explain the effect of the choice of kernel function on the generalization performance of support vector machines.

Index terms: ϵ -entropy, covering numbers, statistical learning theory, support vector machines, linear operators.

*Parts of this work have been submitted to EUROCOLT'99 (European Conference on Computational Learning Theory) and have also appeared (in summary form) in *Advances in Kernel Methods*, MIT Press, 1999.

1 Introduction

In this paper we give new bounds on the covering numbers for kernel machines. This leads to improved bounds on their generalization performance. Kernel machines perform a mapping from input space into a feature space (see e.g. [1, 33]), construct regression functions or decision boundaries based on this mapping, and use constraints in feature space for capacity control. Support Vector (SV) machines, which have recently been proposed as a new class of learning algorithms solving problems of pattern recognition, regression estimation, and operator inversion [54] are a well known example of this class. We will use SV machines as our model of choice to show how bounds on the covering numbers can be obtained. We outline the relatively standard methods one can then use to hence bound their generalization performance. SV machines, like most kernel based methods, possess the nice property of defining the feature map in a manner that allows its computation implicitly at little additional computational cost. Our reasoning also applies to similar algorithms such as regularization networks [16] or certain unsupervised learning algorithms [42]. Let us now take a closer look at SV machines. Central to them are two ideas: capacity control by maximizing margins, and the use of nonlinear kernel functions.

Capacity control. In order to perform pattern recognition using linear hyperplanes, often a maximum margin of separation between the classes is sought for, as this leads to good generalization ability independent of the dimensionality [55, 54, 44]. It can be shown that for separable training data

$$(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m) \in \mathbb{R}^d \times \{\pm 1\}, \quad (1)$$

this is achieved by minimizing $\|\mathbf{w}\|_2$ subject to the constraints $y_j(\langle \mathbf{w}, \mathbf{x}_j \rangle + b) \geq 1$ for $j = 1, \dots, m$, and some $b \in \mathbb{R}$. The decision function then takes the form

$$f(\mathbf{x}) = \text{sgn}(\langle \mathbf{w}, \mathbf{x} \rangle + b). \quad (2)$$

Similarly, a linear regression

$$f(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b \quad (3)$$

can be estimated from data

$$(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m) \in \mathbb{R}^d \times \mathbb{R} \quad (4)$$

by finding the flattest function which approximates the data within some margin of error: in this case, one minimizes $\|\mathbf{w}\|_2$ subject to $|f(\mathbf{x}_j) - y_j| \leq \varepsilon$, where the parameter $\varepsilon > 0$ plays the role of the margin, albeit not in the space of the inputs \mathbf{x} , but in that of the outputs y . The analogy to pattern recognition is somewhat loose, and there exist alternative ways of introducing a margin (e.g. in the space \mathbb{R}^m of all outputs y_1, \dots, y_m , [60]).

In both cases, generalizations for the nonseparable or nonrealizable case exist, using various types of cost functions [14, 54, 47].

Nonlinear kernels. In order to apply the above reasoning to a rather general class of *nonlinear* functions, one can use kernels computing dot products in high-dimensional spaces nonlinearly related to input space [1, 10]. Under certain conditions on a kernel k , to be stated below (Theorem 4), there exists a nonlinear map Φ into a reproducing kernel Hilbert space F (see e.g. [39]) such that k computes the dot product in F , i.e.

$$k(\mathbf{x}, \mathbf{y}) = \langle \Phi(\mathbf{x}), \Phi(\mathbf{y}) \rangle_F. \quad (5)$$

Given any algorithm which can be expressed in terms of dot products exclusively, one can thus construct a nonlinear version of it by substituting a kernel for the dot product. Examples of such machines include SV pattern recognition [10], SV regression estimation [54], and kernel principal component analysis [42].

By using the kernel trick for SV machines, the maximum margin idea is thus extended to a large variety of nonlinear function classes (e.g. radial basis function networks, polynomial networks, neural networks), which in the case of regression estimation comprise functions written as kernel expansions

$$f(\mathbf{x}) = \sum_{i=1}^m \alpha_i k(\mathbf{x}_i, \mathbf{x}) + b, \quad (6)$$

with $\alpha_j \in \mathbb{R}$, $j = 1, \dots, m$. It has been noticed that different kernels can be characterized by their regularization properties [49]: SV machines are regularization networks minimizing the regularized risk $R_{reg}[f] = R_{emp}[f] + \frac{\lambda}{2} \|Pf\|^2$, (with a regularization parameter $\lambda \geq 0$, and a regularization operator P) over the set of functions of the form (6), provided that k and P are interrelated by $k(\mathbf{x}_s, \mathbf{x}_t) = \langle (Pk)(\mathbf{x}_s, \cdot), (Pk)(\mathbf{x}_t, \cdot) \rangle$. To this end, k is chosen as Green's function of P^*P where P^* is the adjoint of P .

This provides insight into the regularization properties of SV kernels. However, it does not completely settle the issue of how to select a kernel for a given learning problem, and how using a specific kernel might influence the performance of a SV machine.

1.1 Outline of the paper

In the present work, we show that properties of the spectrum of the kernel can be used to make statements about the generalization error of the associated class of learning machines. Unlike in previous SV learning studies, the kernel is no longer merely a means of broadening the class of functions used, e.g. by making a nonseparable dataset separable in a feature space nonlinearly related to input space. Rather, we now view it as a constructive handle by which we can control the generalization error.

A key feature of the present paper is the manner in which we *directly* bound the covering numbers of interest rather than making use of a combinatorial dimension (such as the VC-dimension or the fat-shattering dimension) and subsequent application of a general result relating such dimensions to covering

numbers. We bound covering numbers directly by viewing the relevant class of functions as the image of a unit ball under a particular compact operator. A general overview of the method is given in Section 3.

The remainder of the paper is organized as follows. We start by introducing notation and definitions (Section 2). Section 4 formulates generalization error bounds in terms of covering numbers. Section 5 contains the main result bounding entropy numbers in terms of the spectrum of a given kernel. The results in this paper rest on a connection between covering numbers of function classes and entropy numbers of suitably defined operators. In particular we derive an upper bound on the entropy numbers in terms of the size of the weight vector in feature space and the eigenvalues of the kernel used. Section 6 shows how to make use of kernels such as $k(x) = e^{-x^2}$ which do not have a discrete spectrum. Section 7 presents some results on the entropy numbers obtained for given rates of decay of eigenvalues and 8 shows how to extend the results to several dimensions. The concluding section (Section 9) indicates how the various results in the paper can be glued together in order to obtain overall bounds on the generalization error. Most of the examples we provide for the calculation of eigenvalues are for translation invariant kernels (i.e. convolutional kernels); this is merely for convenience — the general theory is not restricted to such kernels.

We do not present a single master generalization error theorem for three key reasons: 1) the only novelty in the paper lies in the computation of covering numbers themselves; 2) the particular statistical result one needs to use depends on the specific problem situation; 3) many of the results obtained are in a form which, whilst quite amenable to ready computation on a computer, do not provide much direct insight by merely looking at them, except perhaps in the asymptotic sense, and finally 4) some applications (such as classification) where further quantities like margins are estimated in a data dependent fashion, need an additional luckiness argument [45] to apply the bounds.

Thus although our goal has been theorems, we are ultimately forced to resort to a computer to make use of our results. This is not necessarily a disadvantage — it is both a strength and a weakness of Structural Risk Minimization (SRM) [56] that a good generalization error bound is both necessary and sufficient to make the method work well. It is our expectation that the refined (and significantly more tight) covering number bounds obtainable by our methods will be exploitable in SRM algorithms — they could be used for example for model selection. If one is running a computer program anyway, there is little point in expending a large effort to make the generalization error bounds directly consumable in a pencil and paper sense.

2 Definitions and Notation

For $d \in \mathbb{N}$, \mathbb{R}^d denotes the d -dimensional space of vectors $\mathbf{x} = (x_1, \dots, x_d)^T$. We define spaces ℓ_p^d as follows: as vector spaces, they are identical to \mathbb{R}^d , in addition, they are endowed with p -norms: for $0 < p < \infty$, $\|\mathbf{x}\|_{\ell_p^d} := \|\mathbf{x}\|_p =$

$(\sum_{j=1}^d |x_j|^p)^{1/p}$; for $p = \infty$, $\|\mathbf{x}\|_{\ell_\infty^d} := \|\mathbf{x}\|_\infty = \max_{j=1,\dots,d} |x_j|$. Note that a different normalization of the the ℓ_p^d norm is used in some papers in learning theory (e.g. [52]). For $0 < p < \infty$, $\ell_p = \ell_p^\infty$. We use the shorthand sequence notation $(x_j)_j = (x_1, x_2, \dots)$.

Given m points $\mathbf{x}_1, \dots, \mathbf{x}_m \in \ell_p^d$, we use the shorthand $\mathbf{X}^m = (\mathbf{x}_1, \dots, \mathbf{x}_m)$.

Suppose \mathcal{F} is a class of functions $f: \mathbb{R}^d \rightarrow \mathbb{R}$. The ℓ_∞^d norm *with respect to* \mathbf{X}^m of $f \in \mathcal{F}$ is defined as $\|f\|_{\ell_\infty^d} := \max_{i=1,\dots,m} |f(\mathbf{x}_i)|$. Likewise $\|f\|_{\ell_p^{\mathbf{X}^m}} = \|(f(\mathbf{x}_1), \dots, f(\mathbf{x}_m))\|_{\ell_p^m}$.

Given some set \mathcal{X} with a σ -algebra, a measure μ on \mathcal{X} , some $1 \leq p < \infty$ and a function $f: \mathcal{X} \rightarrow \mathbb{R}$ we define $\|f\|_{L_p(\mathcal{X}, \mathbb{R})} := (\int |f(x)|^p d\mu(x))^{1/p}$ if the integral exists and $\|f\|_{L_\infty(\mathcal{X}, \mathbb{R})} := \text{ess sup}_{x \in \mathcal{X}} |f(x)|$. For $1 \leq p \leq \infty$, we let $L_p(\mathcal{X}, \mathbb{R}) := \{f: \mathcal{X} \rightarrow \mathbb{R}: \|f\|_{L_p(\mathcal{X}, \mathbb{R})} < \infty\}$. We let $L_p(\mathcal{X}) := L_p(\mathcal{X}, \mathbb{R})$.

Let $\mathfrak{L}(E, F)$ be the set of all bounded linear operators T between the normed spaces $(E, \|\cdot\|_E)$ and $(F, \|\cdot\|_F)$, i.e. operators such that the image of the (closed) unit ball

$$U_E := \{x \in E: \|x\|_E \leq 1\} \tag{7}$$

is bounded. The smallest such bound is called the *operator norm*,

$$\|T\| := \sup_{x \in U_E} \|Tx\|_F. \tag{8}$$

The n th *entropy number* of a set $M \subset E$, for $n \in \mathbb{N}$, is

$$\epsilon_n(M) := \inf\{\epsilon > 0 \quad : \quad \text{there exists an } \epsilon\text{-cover for } M \text{ in } E \\ \text{containing } n \text{ or fewer points}\} \tag{9}$$

The *entropy numbers* of an operator $T \in \mathfrak{L}(E, F)$ are defined as

$$\epsilon_n(T) := \epsilon_n(T(U_E)). \tag{10}$$

Note that $\epsilon_1(T) = \|T\|$, and that $\epsilon_n(T)$ certainly is well defined for all $n \in \mathbb{N}$ if T is a *compact operator*, i.e. if $T(U_E)$ is precompact, i.e. if for any $\epsilon > 0$ there exists a finite cover of $T(U_E)$ with open ϵ balls on \mathcal{X} . The *dyadic entropy numbers* of an operator are defined by

$$e_n(T) := \epsilon_{2^{n-1}}(T), \quad n \in \mathbb{N}; \tag{11}$$

similarly, the dyadic entropy numbers of a set are defined from its entropy numbers. A very nice introduction to entropy numbers of operators is [13]. The ϵ -*covering number* of \mathcal{F} with respect to the metric d denoted $\mathcal{N}(\epsilon, \mathcal{F}, d)$ is the smallest number of elements of an ϵ -cover for \mathcal{F} using the metric d .

In this paper, E and F will always be *Banach spaces*, i.e. complete normed spaces (for instance ℓ_p^d spaces). In some cases, they will be *Hilbert spaces* H , i.e. Banach spaces endowed with a dot product $\langle \cdot, \cdot \rangle_H$ giving rise to its norm via $\|x\|_H = \sqrt{\langle x, x \rangle_H}$.

By \log and \ln , we denote the logarithms to base 2 and e , respectively. By i , we denote the imaginary unit $i = \sqrt{-1}$, k will always be a kernel, and d and

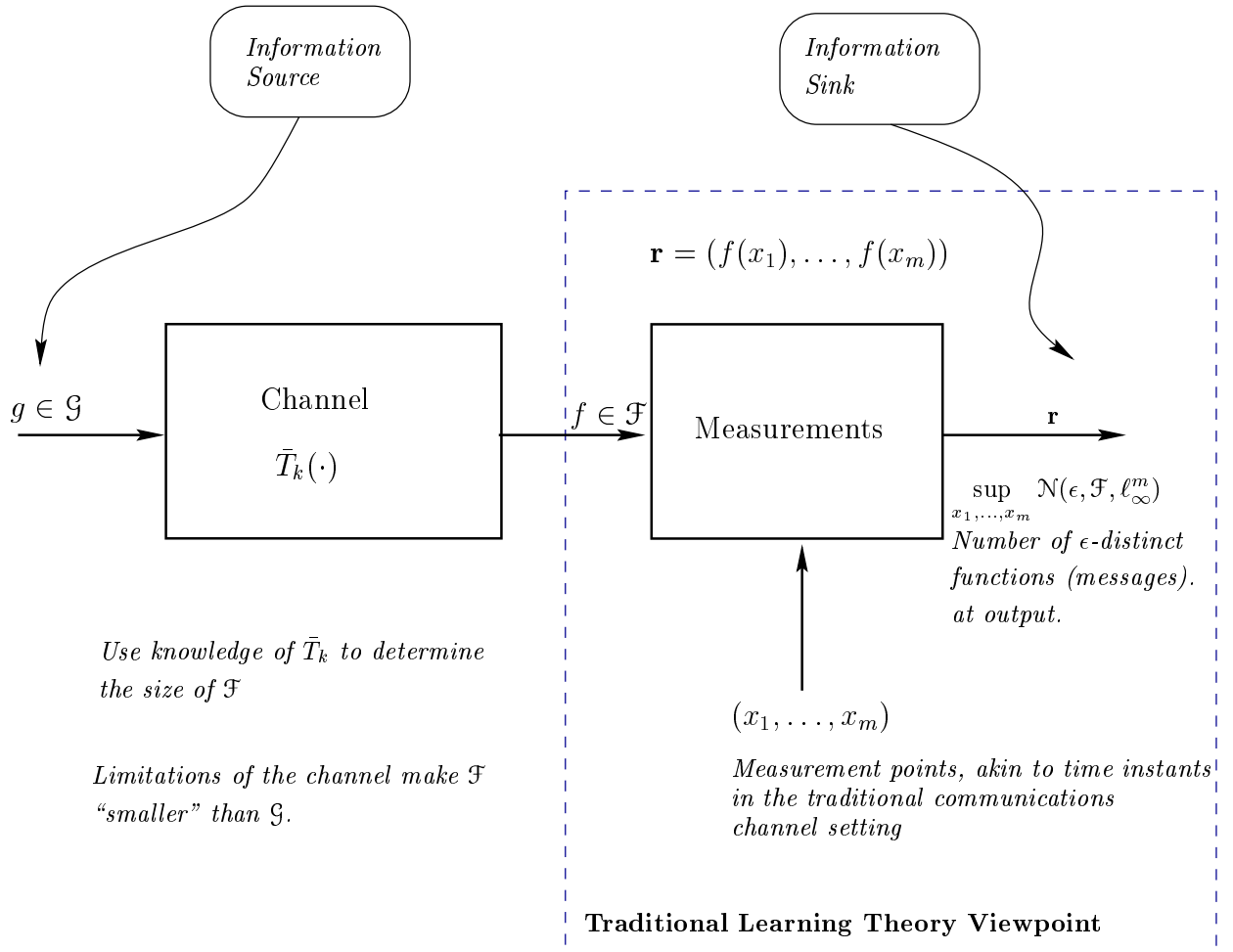


Figure 1: Schematic picture of the new viewpoint.

m will be the input dimensionality and the number of examples

$$(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m) \in \mathbb{R}^d \times \mathbb{R}, \quad (12)$$

respectively. We will map the input data into a feature space via a mapping Φ . We let $\tilde{\mathbf{x}} := \Phi(\mathbf{x})$.

3 Operator Theory Methods for Entropy Numbers

In this section we briefly explain the new viewpoint implicit in the present paper. With reference to Figure 1, consider the traditional viewpoint in statistical learning theory. One is given a class of functions \mathcal{F} , and the generalization performance attainable using \mathcal{F} is determined via the covering numbers of \mathcal{F} . More precisely, for some set \mathcal{X} , and $\mathbf{x}_i \in \mathcal{X}$ for $i = 1, \dots, m$, define the ϵ -Growth

function of the function class \mathcal{F} on \mathcal{X} as

$$\mathcal{N}^m(\epsilon, \mathcal{F}) := \sup_{\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathcal{X}} \mathcal{N}(\epsilon, \mathcal{F}, \ell_\infty^{\mathbf{X}^m}), \quad (13)$$

where $\mathcal{N}(\epsilon, \mathcal{F}, \ell_\infty^{\mathbf{X}^m})$ is the ϵ -covering number of \mathcal{F} with respect to $\ell_\infty^{\mathbf{X}^m}$. Many generalization error bounds can be expressed in terms of $\mathcal{N}^m(\epsilon, \mathcal{F})$. An example is given in the following section.

The key novelty in the present work solely concerns the manner in which the covering numbers are computed. Traditionally, appeal has been made to a result such as the so-called Sauer’s lemma (originally due to Vapnik and Chervonenkis). In the case of function learning, a generalization called the VC-dimension of real valued functions, or a variation due to Pollard (called the pseudo-dimension), or or a scale-sensitive generalization of that (called the fat-shattering dimension) is used to bound the covering numbers. These results reduce the computation of $\mathcal{N}^m(\epsilon, \mathcal{F})$ to the computation of a single “dimension-like” quantity. An overview of these various dimensions, some details of their history, and some examples of their computation can be found in [5].

In the present work, we view the class \mathcal{F} as being induced by an operator \bar{T}_k depending on some kernel function k . Thus \mathcal{F} is the image of a “base class” \mathcal{G} under \bar{T}_k . The analogy implicit in the picture is that the quantity that matters is the number of ϵ -distinguishable messages obtainable at the information sink. (Recall the equivalence up to a constant factor of 2 in ϵ of packing and covering numbers.) In a typical communications problem, one tries to maximize the number of distinguishable messages (per unit time), in order to maximize the information transmission rate. But from the point of view of the receiver, the job is made easier the *smaller* the number of distinct messages that one needs to be concerned with decoding. The significance of the picture is that the kernel in question is exactly the kernel that is used, for example, in support vector machines. As a consequence, the determination of $\mathcal{N}^m(\epsilon, \mathcal{F})$ can be done in terms of properties of the operator \bar{T}_k . The latter thus plays a constructive role in controlling the complexity of \mathcal{F} and hence the difficulty of the learning task. We believe that the new viewpoint in itself is potentially very valuable, perhaps more so than the specific results in the paper. A further exploitation of the new viewpoint can be found in [60, 41, 50, 48]. There are in fact a variety of ways to define exactly what is meant by \bar{T}_k , and we have deliberately not been explicit in the picture. We make use of one particular \bar{T}_k in this paper. A slightly different approach is taken in [60].

We conclude this section with some historical remarks.

The concept of the metric entropy of a set has been around for some time. It seems to have been introduced by Pontriagin and Schnirelmann [36] and was studied in detail by Kolmogorov and others [27]. The use of metric entropy to say something about linear operators was developed independently by several people. Prosser [37] appears to have been the first to make the idea explicit. He determined the effect of an operator’s spectrum on its entropy numbers. In particular, he proved a number of results concerning the asymptotic rate of decrease of the entropy numbers in terms of the asymptotic behavior of the

eigenvalues. A similar result is actually implicit in section 22 of Shannon’s famous paper [43], where he considered the effect of different convolution operators on the entropy of an ensemble. Prosser’s paper [37] led to a handful of papers (see e.g. [38, 22, 3, 29]) which studied various convolutional operators. A connection between Prosser’s ϵ -entropy of an operator and Kolmogorov’s ϵ -entropy of a stochastic process was shown in [2]. Independently, another group of mathematicians including Carl and Stephani [13] studied covering numbers [53] and later entropy numbers [35] in the context of operator ideals. (They seem to be unaware of Prosser’s work — see e.g. [11, p. 136].)

Connections between the local theory of Banach spaces and uniform convergence of empirical means has been noted before (e.g. [34]). More recently Gurvits [21] has obtained a result relating the Rademacher type of a Banach space to the fat-shattering dimension of linear functionals on that space and hence via the key result in [4] to the covering numbers of the induced class. We will make further remarks concerning the relationship between Gurvits’ approach and ours in [60]; for now let us just note that the equivalence of the type of an operator (or of the space it maps to), and the rate of decay of its entropy numbers has been (independently) shown by Kolchinskii [25, 26] and Defant and Junge [15, 23]. Note that the exact formulation of their results differs. Kolchinskii was motivated by probabilistic problems not unlike ours.

4 Generalization Bounds via Uniform Convergence

The generalization performance of learning machines can be bounded via uniform convergence results as in [57, 56]. A recent review can be found in [5]; see also [30]. The key thing about these results is the role of the covering numbers of the hypothesis class — the focus of the present paper. Results for both classification and regression are now known. For the sake of concreteness, we quote below a result suitable for regression which was proved in [4]. For results on classifier performance in terms of covering numbers see [8]. Let $P_m(f) := \frac{1}{m} \sum_{i=1}^m f(\mathbf{x}_i)$ denote the *empirical mean* of f on the sample $\mathbf{x}_1, \dots, \mathbf{x}_m$.

Lemma 1 (Alon, Ben–David, Cesa–Bianchi, and Haussler, 1997) *Let \mathcal{F} be a class of functions from \mathcal{X} into $[0, 1]$ and let P be a distribution over \mathcal{X} . Then, for all $\epsilon > 0$ and all $m \geq \frac{2}{\epsilon^2}$,*

$$\Pr \left\{ \sup_{f \in \mathcal{F}} |P_m(f) - P(f)| > \epsilon \right\} \leq 12m \cdot \mathbf{E} \left[\mathcal{N} \left(\frac{\epsilon}{6}, \mathcal{F}, \ell_\infty^{\bar{\mathbf{X}}^{2m}} \right) \right] e^{-\epsilon^2 m / 36} \quad (14)$$

where \Pr denotes the probability w.r.t. the sample $\mathbf{x}_1, \dots, \mathbf{x}_m$ drawn i.i.d. from P , and \mathbf{E} the expectation w.r.t. a second sample $\bar{\mathbf{X}}^m = (\bar{\mathbf{x}}_1, \dots, \bar{\mathbf{x}}_{2m})$ also drawn i.i.d. from P .

In order to use this lemma one usually makes use of the fact that for any P ,

$$\mathbf{E} \left[\mathcal{N}(\epsilon, \mathcal{F}, \ell_\infty^{\bar{\mathbf{X}}^m}) \right] \leq \mathcal{N}^m(\epsilon, \mathcal{F}). \quad (15)$$

The above result can be used to give a generalization error result by applying it to the loss-function induced class. The following Lemma, which is an improved version of [9, Lemma 17], is useful in this regard (a similar result was obtained by [6]):

Lemma 2 *Denote \mathcal{F} a set of functions from \mathcal{X} to $[a, b]$ with $a < b$, $a, b \in \mathbb{R}$ and $l : \mathbb{R} \rightarrow [0, \infty)$ a loss function. Let $\mathbf{z} := (\mathbf{x}_i, y_i)_{i=1}^m$, $l_f|_{\mathbf{z}_j} := l(f(\mathbf{x}_j) - y_j)$, $l_f|_{\mathbf{z}} := (l_f|_{\mathbf{z}_j})_{j=1}^m$, $l_{\mathcal{F}}|_{\mathbf{z}} := \{l_f|_{\mathbf{z}} : f \in \mathcal{F}\}$ and $\mathcal{N}(\epsilon, l|_{\mathbf{z}}) := \mathcal{N}(\epsilon, l_{\mathcal{F}}|_{\mathbf{z}}, \ell_{\infty}^{\mathbf{z}})$. Then the following two statements hold:*

1. *Suppose l satisfies the Lipschitz-condition*

$$l(\xi) - l(\xi') \leq C|\xi - \xi'| \text{ for all } \xi, \xi' \in [a - b, b - a]. \quad (16)$$

Then for all $\epsilon > 0$

$$\max_{\mathbf{z} \in (\mathcal{X} \times [a, b])^m} \mathcal{N}(\epsilon, l|_{\mathbf{z}}) \leq \max_{\mathbf{x} \in \mathcal{X}^m} \mathcal{N}\left(\frac{\epsilon}{C}, \mathcal{F}|_{\mathbf{x}}, \ell_{\infty}^{\mathbf{x}}\right) \quad (17)$$

and

$$\max_{\mathbf{z} \in (\mathcal{X} \times [a, b])^m} \mathcal{N}(\epsilon, l|_{\mathbf{z}}) \leq \max_{\mathbf{x} \in \mathcal{X}^m} \mathcal{N}\left(\frac{\epsilon m}{C}, \mathcal{F}|_{\mathbf{x}}, \ell_1^{\mathbf{x}}\right). \quad (18)$$

2. *Suppose that for some $C, \tilde{C} > 0$, l satisfies the “approximate Lipschitz-condition”*

$$l(\xi) - l(\xi') \leq \max(C|\xi - \xi'|, \tilde{C}) \text{ for all } \xi, \xi' \in [a - b, b - a] \quad (19)$$

then for all $\epsilon > \tilde{C}/C$

$$\max_{\mathbf{z} \in (\mathcal{X} \times [a, b])^m} \mathcal{N}(\epsilon, l|_{\mathbf{z}}) \leq \max_{\mathbf{x} \in \mathcal{X}^m} \mathcal{N}\left(\frac{\epsilon}{C}, \mathcal{F}|_{\mathbf{x}}, \ell_{\infty}^{\mathbf{x}}\right). \quad (20)$$

Proof We show that, for any sequence \mathbf{z} of (\mathbf{x}, y) pairs in $\mathcal{X} \times [a, b]$ and any functions f and g , if the restrictions of f and g to \mathbf{x} are close, then the restrictions of l_f and l_g to \mathbf{z} are close. Thus, given a cover of $\mathcal{F}|_{\mathbf{x}}$ we can construct a cover of $l_{\mathcal{F}}|_{\mathbf{z}}$ that is no bigger. For part 1 we get:

$$\begin{aligned} \frac{1}{m} \left| \sum_{j=1}^m l(g(\mathbf{x}_j) - y_j) - l(f(\mathbf{x}_j) - y_j) \right| &\leq \frac{1}{m} \sum_{j=1}^m |l(g(\mathbf{x}_j) - y_j) - l(f(\mathbf{x}_j) - y_j)| \\ &\leq \frac{1}{m} \sum_{j=1}^m C|g(\mathbf{x}_j) - f(\mathbf{x}_j)| \\ &= \frac{C}{m} \|g(\mathbf{X}^m) - f(\mathbf{X}^m)\|_{\ell_1^m} \\ &\leq C \|g(\mathbf{X}^m) - f(\mathbf{X}^m)\|_{\ell_{\infty}^m}. \end{aligned}$$

In the second case we proceed similarly

$$\begin{aligned} \frac{1}{m} \left| \sum_{j=1}^m l(g(\mathbf{x}_j) - y_j) - l(f(\mathbf{x}_j) - y_j) \right| &\leq \frac{C}{m} \sum_{j=1}^m \max(|g(\mathbf{x}_j) - f(\mathbf{x}_j)|, \tilde{C}/C) \\ &\leq C\epsilon \quad \text{for } \epsilon \geq \tilde{C}/C. \end{aligned}$$

■

The second case can be useful, when the exact form of the cost function is not known, happens to be discontinuous or is badly behaved in some other way.¹ It shows how down to a scale \tilde{C}/C statements about the covering numbers of the loss-function induced class can be made. Applying the result above to polynomial loss leads to the following corollary:

Corollary 3 *Let the assumptions be as above in lemma 2. Then for loss functions of type*

$$l(\eta) = \frac{1}{p} \eta^p \quad \text{with } p > 1 \quad (21)$$

we have $C = (b - a)^{(p-1)}$, in particular $C = (b - a)$ for $p = 2$ and therefore

$$\max_{\mathbf{z} \in (X \times [a, b])^m} \mathcal{N}(\epsilon, l|\mathbf{z}) \leq \max_{\mathbf{x} \in X^m} \mathcal{N} \left(\frac{\epsilon}{(b - a)^{p-1} X C}, \mathcal{F}|\mathbf{x} \right) \quad (22)$$

One can readily combine the uniform convergence results with the above results to get overall bounds on generalization performance. We do not explicitly state such a result here since the particular uniform convergence result needed depends on the exact set-up of the learning problem. A typical uniform convergence result takes the form

$$P^m \{ \sup_f |R_{emp}(f) - R(f)| > \epsilon \} \leq c_1(m) \mathcal{N}^m(\epsilon, \mathcal{F}) e^{-\epsilon^\beta m / c_2}. \quad (23)$$

Even the exponent in (23) depends on the setting: In regression β can be set to 1, however in agnostic learning [24] in general $\beta = 2$, except if the class is convex in which case it can be set to 1 [31]. Since our primary interest is in determining $\mathcal{N}^m(\epsilon, \mathcal{F})$ we will not try to summarize the large body of results now available on uniform convergence and generalization error.

These generalization bounds are typically used by setting the right hand side equal to δ and solving for $m = m(\epsilon, \delta)$ (which is called the sample complexity). Another way to use these results is as a learning curve bound $\bar{\epsilon}(\delta, m)$ where

$$P^m \{ \sup_{f \in \mathcal{F}} |R_{emp}(f) - R(f)| > \bar{\epsilon}(\delta, m) \} \leq \delta.$$

We note here that the determination of $\bar{\epsilon}(\delta, m)$ is quite convenient in terms of e_n , the dyadic entropy number associated with the covering number $\mathcal{N}^m(\epsilon, \mathcal{F})$

¹The two cases could be combined into one by writing the conditions in terms of the modulus of continuity. For the sake of clarity, however, we refrained from doing so.

in (23). Setting the right hand side of (23) equal to δ , we have

$$\begin{aligned}
\delta &= c_1(m)\mathcal{N}^m(\epsilon, \mathcal{F})e^{-\epsilon^\beta m/c_2} \\
\Rightarrow \log\left(\frac{\delta}{c_1(m)}\right) + \frac{\epsilon^\beta m}{c_2 \ln 2} &= \log \mathcal{N}^m(\epsilon, \mathcal{F}) \\
\Rightarrow e_{\lceil \log\left(\frac{\delta}{c_1(m)}\right) + \frac{\epsilon^\beta m}{c_2 \ln 2} + 1 \rceil} &\leq \epsilon.
\end{aligned} \tag{24}$$

Thus $\bar{\epsilon}(\delta, m) = \min\{\epsilon: (24) \text{ holds}\}$. Thus the use of ϵ_n or e_n (which will arise naturally from our techniques) is in fact a convenient thing to do for finding learning curves.

5 Entropy Numbers for Kernel Machines

In the following we will mainly consider machines where the mapping into feature space is defined by Mercer kernels $k(\mathbf{x}, \mathbf{y})$ as they are easier to deal with using functional analytic methods. (More general kernels are considered in [48].) Such machines have become very popular due to the success of SV machines. Nonetheless in Appendix A we will show how a more direct approach could be taken towards upper-bounding entropy numbers.

5.1 Mercer's Theorem, Feature Spaces and Scaling

Our goal is to make statements about the shape of the image of the input space \mathcal{X} under the feature map $\Phi(\cdot)$. We will make use of Mercer's theorem. The version stated below is a special case of the theorem proven in [28, p. 145]. In the following we will assume (\mathcal{X}, μ) to be a finite measure space, i.e. $\mu(\mathcal{X}) < \infty$.

Theorem 4 (Mercer) *Suppose $k \in L_\infty(\mathcal{X}^2)$ is a symmetric kernel (hence $k(x, x') = k(x', x)$) such that the integral operator $T_k : L_2(\mathcal{X}) \rightarrow L_2(\mathcal{X})$,*

$$T_k f(\cdot) := \int_{\mathcal{X}} k(\cdot, \mathbf{y})f(\mathbf{y})d\mu(\mathbf{y}) \tag{25}$$

is positive. Let $\psi_j \in L_2(\mathcal{X})$ be the eigenfunction of T_k associated with the eigenvalue $\lambda_j \neq 0$ and normalized such that $\|\psi_j\|_{L_2} = 1$. Suppose ψ_j is continuous for all $j \in \mathbb{N}$.

1. $(\lambda_j(T))_j \in \ell_1$.
2. $\psi_j \in L_\infty(\mathcal{X})$ and $\sup_j \|\psi_j\|_{L_\infty} < \infty$.
3. $k(\mathbf{x}, \mathbf{y}) = \sum_{j \in \mathbb{N}} \lambda_j \psi_j(\mathbf{x})\psi_j(\mathbf{y})$ holds for all (\mathbf{x}, \mathbf{y}) , where the series converges absolutely and uniformly for all (\mathbf{x}, \mathbf{y}) .

We will call a kernel satisfying the conditions of this theorem a *Mercer kernel*. From statement 2 of Mercer's theorem there exists some constant $C_k \in \mathbb{R}^+$ depending on $k(\cdot, \cdot)$ such that

$$|\psi_j(\mathbf{x})| \leq C_k \text{ for all } j \in \mathbb{N} \text{ and } \mathbf{x} \in \mathcal{X}. \tag{26}$$

Note that if \mathcal{X} is compact and k is continuous, then ψ_j is continuous (cf. e.g. [7, p.270]). Alternatively, if k is translation invariant, then ψ_j are scaled cosine functions and thus continuous. Moreover from statement 3 it follows that $k(\mathbf{x}, \mathbf{y})$ corresponds to a dot product in ℓ_2 i.e. $k(\mathbf{x}, \mathbf{y}) = \langle \Phi(\mathbf{x}), \Phi(\mathbf{y}) \rangle_{\ell_2}$ with

$$\begin{aligned} \Phi : \mathcal{X} &\rightarrow \ell_2 \\ \Phi : \mathbf{x} &\mapsto (\phi_j(\mathbf{x}))_j := (\sqrt{\lambda_j} \psi_j(\mathbf{x}))_j \end{aligned} \quad (27)$$

for all $\mathbf{x} \in \mathcal{X}$. In the following we will (without loss of generality) assume the sequence of $(\lambda_j)_j$ be sorted in nonincreasing order. From the argument above, one can see that $\Phi(\mathcal{X})$ lives not only in ℓ_2 but in an axis parallel parallelepiped with lengths $2C_k \sqrt{\lambda_j}$.

We remark that the measure μ need have nothing to do with the distribution of examples. We make use of result 3 of the above theorem to bound the covering numbers of the class of functions implemented using SV machines. The specific bounds *will* depend on μ since that will affect the ψ_j , C_k and the $(\lambda_j)_j$. The question of the optimal μ to use and how it may be chosen if one knows P (the distribution from which the \mathbf{x}_j are drawn) is not considered here. In all cases considered in this paper we will in fact take μ to be Lebesgue measure.

It will be useful to consider maps that map $\Phi(\mathcal{X})$ into balls of some radius R centered at the origin. The following proposition shows that the class of all these maps is determined by elements of ℓ_2 and the sequence of eigenvalues $(\lambda_j)_j$.

Proposition 5 (Mapping $\Phi(\mathcal{X})$ into ℓ_2) *Let S be the diagonal map*

$$\begin{aligned} S : \mathbb{R}^{\mathbb{N}} &\rightarrow \mathbb{R}^{\mathbb{N}} \\ S : (x_j)_j &\mapsto S(x_j)_j = (s_j x_j)_j \text{ with } s_j \in \mathbb{R}. \end{aligned} \quad (28)$$

Then S maps $\Phi(\mathcal{X})$ into a ball of finite radius R_S centered at the origin if and only if $(\sqrt{\lambda_j} s_j)_j \in \ell_2$.

Proof

(\Leftarrow) Suppose $(s_j \sqrt{\lambda_j})_j \in \ell_2$ and let $R_S^2 := C_k^2 \|(s_j \sqrt{\lambda_j})_j\|_{\ell_2}^2 < \infty$. For any $\mathbf{x} \in \mathcal{X}$,

$$\|S\Phi(\mathbf{x})\|_{\ell_2}^2 = \sum_{j \in \mathbb{N}} s_j^2 \lambda_j |\psi_j(\mathbf{x})|^2 \leq \sum_{j \in \mathbb{N}} s_j^2 \lambda_j C_k^2 = R_S^2. \quad (29)$$

Hence $S\Phi(\mathcal{X}) \subseteq \ell_2$.

(\Rightarrow) Suppose $(s_j \sqrt{\lambda_j})_j$ is not in ℓ_2 . Hence the sequence $(A_n)_n$ with $A_n := \sum_{j=1}^n s_j^2 \lambda_j$ is unbounded. Now define

$$a_n(\mathbf{x}) := \sum_{j=1}^n s_j^2 \lambda_j |\psi_j(\mathbf{x})|^2. \quad (30)$$

Then $\|a_n(\cdot)\|_{L_1(\mathcal{X})} = A_n$ due to the normalization condition on ψ_j . However, as $\mu(\mathcal{X}) < \infty$ there exists a set $\tilde{\mathcal{X}}$ of nonzero measure such that

$$a_n(\mathbf{x}) \geq \frac{A_n}{\mu(\mathcal{X})} \quad \text{for all } \mathbf{x} \in \tilde{\mathcal{X}}. \quad (31)$$

Combining the left side of (29) with (30) we obtain $\|S\Phi(\mathbf{x})\|_{\ell_2}^2 \geq a_n(\mathbf{x})$ for all $n \in \mathbb{N}$ and all \mathbf{x} . Since $a_n(\mathbf{x})$ is unbounded for a set $\tilde{\mathcal{X}}$ with nonzero measure in \mathcal{X} , we can see that $S\Phi(\mathcal{X}) \not\subset \ell_2$. ■

Once we know that $\Phi(\mathcal{X})$ is contained in the parallelepiped described above we can use this result to construct a mapping \hat{A} from the unit ball in ℓ_2 to an ellipsoid \mathcal{E} such that $\Phi(\mathcal{X}) \subset \mathcal{E}$ as in the following diagram.

$$\begin{array}{ccc} \mathcal{X} & \xrightarrow{\Phi} & \Phi(\mathcal{X}) \subset \ell_2 & \xrightarrow{\hat{A}^{-1}} & U_{\ell_2} \subset \ell_2 & (32) \\ & & \cap & \nearrow \hat{A} & \\ & & \ell_2 \supset \mathcal{E} & & \end{array}$$

The operator \hat{A} will be useful for computing the entropy numbers of concatenations of operators. (Knowing the inverse will allow us to compute the forward operator, and that can be used to bound the covering numbers of the class of functions, as shown in the next subsection.) We thus seek an operator $\hat{A} : \ell_2 \rightarrow \ell_2$ such that

$$\hat{A}^{-1}\Phi(\mathcal{X}) \subset U_{\ell_2}. \quad (33)$$

This means that $\mathcal{E} := AU_{\ell_2}$ will be such that $\Phi(\mathcal{X}) \subset \mathcal{E}$. The latter can be ensured by constructing \hat{A} such that

$$\hat{A}: (x_j)_j \mapsto (R_{\hat{A}} \cdot a_j \cdot x_j)_j \text{ with } R_{\hat{A}}, a_j \in \mathbb{R}^+ \quad (34)$$

where C_k and a_j are chosen with respect to a specific kernel and where $R_{\hat{A}} := C_k \|(\sqrt{\lambda_j}/a_j)_j\|_{\ell_2}$. From Proposition 5 it follows that all those operators \hat{A} for which $R_{\hat{A}} < \infty$ will satisfy (33). We call such scaling (inverse) operators *admissible*.

5.2 Entropy Numbers

The next step is to compute the entropy numbers of the operator \hat{A} and use this to obtain bounds on the entropy numbers for kernel machines like SV machines. We will make use of the following theorem due to Gordon, König and Schütt [17, p. 226] (stated in the present form in [13, p. 17]).

Theorem 6 *Let $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_j \geq \dots \geq 0$ with $\sigma_i \in \mathbb{R}_0^+$ be a non-increasing sequence of non-negative numbers and let*

$$D\mathbf{x} = (\sigma_1 x_1, \sigma_2 x_2, \dots, \sigma_j x_j, \dots) \quad (35)$$

for $\mathbf{x} = (x_1, x_2, \dots, x_j, \dots) \in \ell_p$ be the diagonal operator from ℓ_p into itself,

generated by the sequence $(\sigma_j)_j$, where $1 \leq p \leq \infty$. Then for all $n \in \mathbb{N}$,

$$\sup_{j \in \mathbb{N}} n^{-\frac{1}{j}} (\sigma_1 \sigma_2 \cdots \sigma_j)^{\frac{1}{j}} \leq \epsilon_n(D) \leq 6 \sup_{j \in \mathbb{N}} n^{-\frac{1}{j}} (\sigma_1 \sigma_2 \cdots \sigma_j)^{\frac{1}{j}}. \quad (36)$$

We can exploit the freedom in choosing \hat{A} to minimize an entropy number as the following corollary shows. This will be a key ingredient of the calculation of the covering numbers for SV classes, as shown below.²

Corollary 7 (Entropy numbers for $\Phi(\mathcal{X})$) *Let $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be a Mercer kernel and let \hat{A} be defined by (34). Then there exists an operator A such that for all $n \in \mathbb{N}$*

$$\epsilon_n(A: \ell_2 \rightarrow \ell_2) \leq \inf_{(a_s)_s: (\sqrt{\lambda_s}/a_s)_s \in \ell_2} \sup_{j \in \mathbb{N}} 6C_k \left\| \left(\sqrt{\lambda_s}/a_s \right)_s \right\|_{\ell_2} n^{-\frac{1}{j}} (a_1 a_2 \cdots a_j)^{\frac{1}{j}}. \quad (37)$$

This result follows immediately by identifying D and A and exploiting the freedom that we still have in choosing a particular operator A among the class of admissible ones.

As already described in Section 1 the hypotheses that a SV machine generates can be expressed as $\langle \mathbf{w}, \tilde{\mathbf{x}} \rangle + b$ where both \mathbf{w} and $\tilde{\mathbf{x}}$ are defined in the feature space $\mathcal{S} = \text{span}(\Phi(\mathcal{X}))$ and $b \in \mathbb{R}$. The kernel trick as introduced by [1] was then successfully employed in [10] and [14] to extend the Optimal Margin Hyperplane classifier to what is now known as the SV machine. We deal with the “+ b ” term in Section 9; for now we consider the class

$$\mathcal{F}_\Lambda := \{f_{\mathbf{w}}: \mathbf{x} \mapsto \langle \mathbf{w}, \mathbf{x} \rangle: \mathbf{x} \in \mathcal{S}, \|\mathbf{w}\| \leq \Lambda\} \subseteq \mathbb{R}^{\mathcal{S}}.$$

Note that \mathcal{F}_Λ depends implicitly on k since \mathcal{S} does.

We seek the ℓ_∞^m covering numbers for the class \mathcal{F}_Λ induced by the kernel in terms of the parameter Λ which is the inverse of the size of the margin in feature space, or equivalently, the size of the weight vector in feature space as defined by the dot product in \mathcal{S} (see [55, 54] for details). In the following we will call such hypothesis classes with length constraint on the weight vectors in feature space *SV classes*. Let T be the operator $T = S_{\tilde{\mathbf{X}}_m} \Lambda$ where $\Lambda \in \mathbb{R}$ and the operator $S_{\tilde{\mathbf{X}}_m}$ is defined by

$$\begin{aligned} S_{\tilde{\mathbf{X}}_m}: \ell_2 &\rightarrow \ell_\infty^m \\ S_{\tilde{\mathbf{X}}_m}: \mathbf{w} &\mapsto (\langle \tilde{\mathbf{x}}_1, \mathbf{w} \rangle, \dots, \langle \tilde{\mathbf{x}}_m, \mathbf{w} \rangle). \end{aligned} \quad (38)$$

with $\tilde{\mathbf{x}}_j \in \Phi(\mathcal{X})$ for all j . The following theorem is useful when computing entropy numbers in terms of T and A . It is originally due to Maurey, and was extended by Carl [12].

²[20] show that for Mercer kernels the minimization as required in (37) is well defined, i.e. there exists some operator A such that the infimum of $\epsilon_n(\hat{A})$ over all \hat{A} is obtained.

Theorem 8 (Carl and Stephani [13, p. 246]) *Let $S \in \mathcal{L}(H, \ell_\infty^m)$ where H is a Hilbert space. Then there exists a constant $c > 0$ such that for all $m \in \mathbb{N}$, and $1 \leq j \leq m$*

$$e_n(S) \leq c \|S\| \left(n^{-1} \log \left(1 + \frac{m}{n} \right) \right)^{1/2}. \quad (39)$$

An alternative proof of this result (given in [60]) provides a small explicit value for the constant: $c \leq 103$. However there is reason to believe that c should be 1.86, the constant obtainable for identity maps from ℓ_2^m into ℓ_∞^m .³

The restatement of Theorem 8 in terms of $\epsilon_{2n-1} = e_n$ will be useful in the following. Under the assumptions above we have

$$\epsilon_n(S) \leq c \|S\| \left((\log n + 1)^{-1} \log \left(1 + \frac{m}{\log n + 1} \right) \right)^{1/2}. \quad (40)$$

Now we can combine the bounds on entropy numbers of A and $S_{\mathbf{X}^m}$ to obtain bounds for SV classes. First we need the following lemma.

Lemma 9 (Carl and Stephani [13, p. 11]) *Let E, F, G be Banach spaces, $R \in \mathcal{L}(F, G)$, and $S \in \mathcal{L}(E, F)$. Then, for $n, t \in \mathbb{N}$,*

$$\epsilon_{nt}(RS) \leq \epsilon_n(R) \epsilon_t(S) \quad (41)$$

$$\epsilon_n(RS) \leq \epsilon_n(R) \|S\| \quad (42)$$

$$\epsilon_n(RS) \leq \epsilon_n(S) \|R\|. \quad (43)$$

Note that the latter two inequalities follow directly from the fact that $\epsilon_1(R) = \|R\|$ for all $R \in \mathcal{L}(F, G)$.

Theorem 10 (Bounds for SV classes) *Let k be a Mercer kernel, let Φ be induced via (27) and let $T := S_{\mathbf{X}^m} \Lambda$ where $S_{\mathbf{X}^m}$ is given by (38) and $\Lambda \in \mathbb{R}^+$. Let A be defined as in corollary 7 and suppose $\tilde{\mathbf{x}}_j = \Phi(\mathbf{x}_j)$ for $j = 1, \dots, m$. Then the entropy numbers of T satisfy the following inequalities:*

$$\epsilon_n(T) \leq c \|A\| \Lambda \log^{-1/2} n \log^{1/2} \left(1 + \frac{m}{\log n} \right) \quad (44)$$

$$\epsilon_n(T) \leq 6\Lambda \epsilon_n(A) \quad (45)$$

$$\epsilon_{nt}(T) \leq 6c\Lambda \log^{-1/2} n \log^{1/2} \left(1 + \frac{m}{\log n} \right) \epsilon_t(A)$$

where c is defined as in Lemma 8.

This result gives several options for bounding $\epsilon_n(T)$. We shall see in examples later that the best inequality to use depends on the rate of decay of the

³In fact (39) is just one of three terms to be considered when bounding $e_n(S)$. However, the other two terms are merely improvements of (39) for the case of small n (since $e_n(S) \leq \|S\|$ by definition) and for n in the order of m or larger. This leads to rates decaying exponentially, i.e. $e_n = O(2^{-n/m})$, however this case is not interesting for learning theoretical studies, since it corresponds to overly complex models. See [60] for details.

eigenvalues of k . The result gives effective bounds on $\mathcal{N}^m(\epsilon, \mathcal{F}_\Lambda)$ since

$$\epsilon_n(T: \ell_2 \rightarrow \ell_\infty^m) \leq \epsilon_0 \Rightarrow \mathcal{N}^m(\epsilon_0, \mathcal{F}_\Lambda) \leq n.$$

Proof We will use the following factorization of T to upper bound $\epsilon_n(T)$.

$$\begin{array}{ccc} U_{\ell_2} \subset \ell_2 & \xrightarrow{T} & \ell_\infty^m \\ \downarrow \Lambda & \nearrow S_{\Phi(\mathbf{x}^m)} & \uparrow S_{(A^{-1}\Phi(\mathbf{x}^m))} \\ \Lambda U_{\ell_2} \subset \ell_2 & \xrightarrow{A} & \Lambda \mathcal{E} \subset \ell_2 \end{array} \quad (46)$$

The top arrow in the diagram follows from the definition of T . The fact that remainder commutes stems from the fact that since A is diagonal, it is self-adjoint and so for any $\tilde{\mathbf{x}} \in \mathcal{S}$

$$\langle \mathbf{w}, \tilde{\mathbf{x}} \rangle = \langle \mathbf{w}, AA^{-1}\tilde{\mathbf{x}} \rangle = \langle A\mathbf{w}, A^{-1}\tilde{\mathbf{x}} \rangle. \quad (47)$$

Instead of computing the entropy number of $T = S_{\tilde{\mathbf{x}}^m} \Lambda$ directly, which is difficult or wasteful, as the the bound on $S_{\tilde{\mathbf{x}}^m}$ does not take into account that $\tilde{\mathbf{x}} \in \mathcal{E}$ but just makes the assumption of $\tilde{\mathbf{x}} \in \rho U_{\ell_2}$ for some $\rho > 0$, we will represent T as $S_{(A^{-1}\tilde{\mathbf{x}}^m)} AA$. This is more efficient as we constructed A such that $A^{-1}\Phi(\mathcal{X}) \subseteq U_{\ell_2}$ filling a larger proportion of it than $\frac{1}{\rho}\Phi(\mathcal{X})$ does.

By construction of A and due to the Cauchy-Schwarz inequality we have $\|S_{A^{-1}\tilde{\mathbf{x}}^m}\| \leq 1$. Thus applying lemma 9 to the factorization of T and using Theorem 8 proves the theorem. \blacksquare

As we shall see in Section 7, one can give asymptotic rates of decay for $\epsilon_n(A)$. (In fact we give non-asymptotic results with explicitly evaluable constants.) It is thus of some interest to give overall asymptotic rates of decay of $\epsilon_n(T)$ in terms of the order of $\epsilon_n(A)$.

Lemma 11 (Rate bounds on ϵ_n) *Let k be a Mercer kernel and suppose A is the scaling operator associated with it as defined by (34).*

1. If $\epsilon_n(A) = O(\log^{-\alpha} n)$ for some $\alpha > 0$ then for fixed m ,

$$\epsilon_n(T) = O(\log^{-(\alpha+1/2)} n). \quad (48)$$

2. If $\log \epsilon_n(A) = O(\log^{-\beta} n)$ for some $\beta > 0$ then for fixed m ,

$$\log \epsilon_n(T) = O(\log^{-\beta} n). \quad (49)$$

This Lemma shows that in the first case, Maurey's result (theorem 8) allows an improvement in the exponent of the entropy number of T , whereas in the second, it affords none (since the entropy numbers decay so fast anyway). The Maurey result may still help in that case though for nonasymptotic n .

Proof From theorem 8 we know that $\epsilon_n(S) = O(\log^{-1/2} n)$. Now use (41), splitting the index n in the following way:

$$n = n^\tau n^{(1-\tau)} \text{ with } \tau \in (0, 1). \quad (50)$$

For the first case this yields

$$\begin{aligned} \epsilon_n(T) &= O(\log^{-1/2} n^\tau) O(\log^{-\alpha} j^{\tau-1}) \\ &= \tau^{-1/2} (1-\tau)^{-\alpha} O(\log^{-(\alpha+1/2)} n) = O(\log^{-(\alpha+1/2)} n). \end{aligned} \quad (51)$$

In the second case we have

$$\log \epsilon_n(T) = \log \left((\tau^{-1/2}) O(\log^{-1/2} n) \right) + (1-\tau)^{-\beta} O(\log^{-\beta} n) = O(\log^{-\beta} n). \quad (52)$$

In a nutshell we can always obtain rates of convergence better than those due to Maurey's theorem because we are not dealing with *arbitrary* mappings into infinite dimensional spaces. In fact, for logarithmic dependency of $\epsilon_n(T)$ on n , the effect of the kernel is so strong that it completely dominates the $1/\sqrt{n}$ behavior for arbitrary Hilbert spaces. An example of such a kernel is $k(x, y) = \exp(-(x-y)^2)$; see Proposition 15 and also Section 6 for the discretization question.

We conclude this section by remarking that instead of the a priori bounds developed here, one may often obtain better bounds in terms of the empirical Gram matrix. A sketch of this reasoning is provided in Appendix A.

6 Discrete Spectra of Convolution Operators

The results presented above show that if one knows the eigenvalue sequence $(\lambda_i)_i$ of a compact operator, one can bound its entropy numbers. Whilst it is always possible to assume that the *data* fed into a SV machine have bounded support, the same can not be said of the kernel $k(\cdot, \cdot)$; a commonly used kernel is $k(x, y) = \exp(-(x-y)^2)$ which has noncompact support. The induced integral operator

$$(T_k f)(x) = \int_{-\infty}^{\infty} k(x, y) f(y) dy \quad (53)$$

then has a continuous spectrum (a nondenumerable infinity of eigenvalues) and thus T_k is not compact [7, p.267]. The question arises: can we make use of such kernels in SV machines and still obtain generalization error bounds of the form developed above? Note that by a theorem of Widom [59], the eigenvalue decay of any convolution operator defined on a compact set via a kernel having compact support can decay no faster than $\lambda_j = \Omega(e^{-j^2})$ and thus if one seeks very rapid decay of eigenvalues (with concomitantly small entropy numbers), one must use convolution kernels with noncompact support.

We will resolve these issues in the present section. Before doing so, let us first consider the case that $\text{supp} k \subseteq [-a, a]$ for some $a < \infty$. Suppose further

that the data points \mathbf{x}_j satisfy $\mathbf{x}_j \in [-b, b]$ for all j . If $k(\cdot, \cdot)$ is a convolution kernel (i.e. $k(x, y) = k(x - y, 0)$ which allows us to write with some abuse of notation $k(x - y) := k(x - y, 0)$), then the SV hypothesis $h_k(\cdot)$ can be written

$$h_k(x) := \sum_{j=1}^m \alpha_j k(x, \mathbf{x}_j) = \sum_{j=1}^m \alpha_j k_v(x, \mathbf{x}_j) =: h_{k_v}(x) \quad (54)$$

for $v \geq 2(a + b)$ where $k_v(\cdot)$ is the v -periodic extension of $k(\cdot)$ (analogously $k_v(x - y) := k_v(x - y, 0)$):

$$k_v(x) := \sum_{j=-\infty}^{\infty} k(x - jv). \quad (55)$$

We now relate the eigenvalues of T_{k_v} to the Fourier transform of $k(\cdot)$. We do so for the case of $d = 1$ and then state the general case later.

Lemma 12 *Let $k: \mathbb{R} \rightarrow \mathbb{R}$ be a symmetric convolution kernel, let $K(\omega) = F[k(x)](\omega)$ denote the Fourier transform of $k(\cdot)$ (see (63)) and k_v denote the v -periodical kernel derived from k (also assume that k_v exists). Then k_v has a representation as a Fourier series with $\omega_0 := \frac{2\pi}{v}$ and*

$$\begin{aligned} k_v(x - y) &= \sum_{j=-\infty}^{\infty} \frac{\sqrt{2\pi}}{v} K(j\omega_0) e^{ij\omega_0 x} \\ &= \frac{\sqrt{2\pi}}{v} K(0) + \sum_{j=1}^{\infty} \frac{2}{v} \sqrt{2\pi} K(j\omega_0) \cos(j\omega_0(x - y)). \end{aligned} \quad (56)$$

Moreover the eigenvalues λ_j of T_{k_v} satisfy $\lambda_j = \sqrt{2\pi} K(j\omega_0)$ for $j \in \mathbb{Z}$ and $C_k = \sqrt{\frac{2}{v}}$.

Proof Clearly the Fourier series coefficients K_j of k_v exist (as k_v exists) with

$$K_j := \frac{1}{\sqrt{v}} \int_{v/2}^{v/2} e^{-ij\omega_0 x} k_v(x) dx$$

and therefore by the definition of k_v and the existence of $K(\omega)$ we conclude

$$\begin{aligned} K_j &= \frac{1}{\sqrt{v}} \int_{v/2}^{v/2} \sum_{j=-\infty}^{\infty} e^{-ij\omega_0 x} k(x - jv) \\ &= \frac{1}{\sqrt{v}} \sum_{j=-\infty}^{\infty} \int_{v/2}^{v/2} e^{-ij\omega_0 x} k(x - jv) = \sqrt{\frac{2\pi}{v}} K(j\omega_0). \end{aligned}$$

This and the fact that $\{x \mapsto v^{-1/2} e^{ij\omega_0 x}; j \in \mathbb{Z}\}$ forms an orthogonal basis in $L_2([-\frac{v}{2}, \frac{v}{2}], \mathbb{C})$ proves (56). (Note that from $k(x) = k(-x)$ we conclude $\overline{K(\omega)} = K(-\omega)$). Furthermore, we are interested in real valued basis functions for $k(x - y)$. The functions

$$\begin{aligned} \psi_0(x) &:= \frac{1}{\sqrt{v}} \\ \psi_j(x) &:= \sqrt{\frac{2}{v}} \cos(j\omega_0 x) \text{ and } \psi_{-j}(x) := \sqrt{\frac{2}{v}} \sin(j\omega_0 x) \text{ for all } j \in \mathbb{N} \end{aligned} \quad (57)$$

form an eigensystem of the integral operator defined by k_v with the corresponding eigenvalues $\sqrt{2\pi}K(j\omega_0)$. Finally one can see that $C_k = \sqrt{\frac{2}{v}}$ by computing the max over $j \in \mathbb{N}$ and $x \in [-v/2, v/2]$. ■

Thus even though T_k may not be compact, T_{k_v} can be (if $(K(j\omega_0))_{j \in \mathbb{N}} \subset \ell_2$ for example). The above lemma can be applied whenever we can form $k_v(\cdot)$ from $k(\cdot)$. Clearly $k(x) = O(x^{-(1+\epsilon)})$ for some $\epsilon > 0$ suffices to ensure the sum in (55) converges.

Let us now consider how to choose v . Note that the Riemann-Lebesgue lemma tells us that for integrable $k(\cdot)$ of bounded variation (surely any kernel one would use would satisfy that assumption), one has $K(\omega) = O(1/\omega)$. There is an tradeoff in choosing v in that for large enough ω , $K(\omega)$ is a decreasing function of ω (at least as fast as $1/\omega$) and thus by Lemma 12, $\lambda_j = \sqrt{2\pi}K(2\pi j/v)$ is an increasing function of v . This suggests one should choose a small value of v . But a small v will lead to high empirical error (as the kernel “wraps around” and its localization properties are lost) and large C_k . There are several approaches to picking a value of v . One obvious one is to *a priori* pick some $\tilde{\epsilon} > 0$ and choose the smallest v such that $|k(x) - k_v(x)| \leq \tilde{\epsilon}$ for all $x \in [-v/2, v/2]$. Thus one would obtain a hypothesis $h_{k_v}(x)$ uniformly within $C\tilde{\epsilon}$ of $h_k(x)$ where $\sum_{j=1}^m |\alpha_j| \leq C$.

Remark 13 *The above Lemma can be readily extended to d dimensions. Assume $k(\mathbf{x})$ is v -periodic in each direction ($\mathbf{x} = (x_1, \dots, x_d)$), we get*

$$\lambda_{\mathbf{j}} = (2\pi)^{\frac{d}{2}} K(\omega_0 \mathbf{j}) = (2\pi)^{\frac{d}{2}} K(\omega_0 \|\mathbf{j}\|) \quad (58)$$

for radially symmetric k and finally for the eigenfunctions $C_k = (2/v)^{\frac{d}{2}}$.

Finally it is worth explicitly noting how the choice of a different bandwidth of the kernel, i.e. letting $k^{(\sigma)}(\mathbf{x}) := \sigma^d k(\sigma \mathbf{x})$, affects the eigenspectrum of the corresponding operator. We have $K^{(\sigma)}(\boldsymbol{\omega}) = K(\boldsymbol{\omega}/\sigma)$, hence scaling a kernel by σ means more densely spaced eigenvalues in the spectrum of the integral operator $T_{k^{(\sigma)}}$.

In conclusion: in order to obtain a discrete spectrum one needs to use a periodic kernel. For a given problem one can always periodize a nonperiodic kernel in a way that changes the final hypothesis in an arbitrarily small way. One can then make use of the results of the present paper.

7 Covering Numbers for Given Decay Rates

In this section we will show how the asymptotic behavior of $\epsilon_n(A: \ell_2 \rightarrow \ell_2)$, where A is the scaling operator introduced before, depends on the eigenvalues of T_k .

A similar analysis has been carried out by Prosser [37], in order to compute the entropy numbers of integral operators. However all of his operators mapped into $L_2(\mathcal{X}, \mathbb{C})$. Furthermore, whilst our propositions are stated as asymptotic

results as his were, the proofs actually give non-asymptotic information with explicit constants.

Note that we need to sort the eigenvalues in a nonincreasing manner because of the requirements in corollary 7. If the eigenvalues were unsorted one could obtain far too small numbers in the geometrical mean of $\lambda_1, \dots, \lambda_j$. Many one-dimensional kernels have nondegenerate systems of eigenvalues in which case it is straightforward to explicitly compute the geometrical means of the eigenvalues as will be shown below. Note that whilst all of the examples below are for convolution kernels, i.e. $k(x, y) = k(x - y)$, there is nothing in the formulations of the propositions themselves that requires this. When we consider the d -dimensional case we shall see that with rotationally invariant kernels, degenerate systems of eigenvalues are generic. In section 8.2 we will show how to systematically deal with that case.

Let us consider the special case where $(\lambda_j)_j$ decays asymptotically with some polynomial or exponential degree. In this case we can choose a sequence $(a_j)_j$ for which we can evaluate (37) explicitly. In what follows, by the eigenvalues of a kernel k we mean the eigenvalues of the induced integral operator T_k .

Proposition 14 (Polynomial Decay) *Let k be a Mercer kernel with eigenvalues $\lambda_j = O(j^{-(\alpha+1)})$ for some $\alpha > 0$. Then for any $\delta \in (0, \alpha/2)$ we have*

$$\epsilon_n(A: \ell_2 \rightarrow \ell_2) = O(\ln^{-\frac{\alpha}{2} + \delta} n). \quad (59)$$

The rate obtained is tight within logarithmic factors and can be bounded from below by $\Omega(\ln^{-\frac{\alpha}{2}} n)$.

An example of such a kernel is $k(x) = e^{-x}$. The proof can be found in the appendix.

The next theorem covers a wide range of practically used kernels, namely those with exponential polynomial decay in their eigenvalues. For instance the Gaussian kernel $k(x) = e^{-x^2}$ has exponential quadratic decay in λ_i . The ‘‘damped harmonic oscillator’’ kernel $k(x) = \frac{1}{1+x^2}$ is another example, this time with just exponential decay in its eigenvalues.

Proposition 15 (Exponential–Polynomial Decay) *Suppose k is a Mercer kernel with $\lambda_j = O(e^{-\alpha j^p})$ for some $\alpha, p > 0$. Then*

$$\ln \epsilon_n^{-1}(A: \ell_2 \rightarrow \ell_2) = O(\ln^{\frac{p}{p+1}} n). \quad (60)$$

The rate is tight.

See the appendix for a proof. (A more precise, but rather more complex, calculation is given in [20].) Whilst this theorem gives the guarantees on the learning rates of estimators using such types of kernels (which is theoretically pleasing and leads to desirable sample complexity rates), it may not always be wise to use the theoretically obtained bounds. Instead, one should take advantage of the estimates based on an analysis of the distribution of the training data since the rates obtained by the latter may turn out to be far superior wrt. the theoretical predictions (cf. sec. 6 and [41]).

8 Higher Dimensions

Things get a somewhat more complicated in higher dimensions. For simplicity, we will restrict ourselves to translation invariant kernels in what follows.

There are basically two ways that can be pursued for constructing kernels in $\mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ with $d > 1$ if no particular assumptions on the data we are dealing with are made. Firstly one could construct kernels by

$$k(\mathbf{x} - \mathbf{y}) = k(x_1 - y_1) \times \cdots \times k(x_d - y_d). \quad (61)$$

This choice will usually lead to preferred directions in input space as the kernels are not rotationally invariant in general. The second approach consists in setting

$$k(\mathbf{x} - \mathbf{y}) = k(\|\mathbf{x} - \mathbf{y}\|_{\ell_2}). \quad (62)$$

This approach also leads to translationally invariant kernels which are also rotationally invariant. In the following we will exploit this approach to compute regularization operators and corresponding Green's functions. It is quite straightforward, however, to generalize our exposition to the rotational asymmetric case. Now let us define the basic ingredients needed for the further calculations.

8.1 Basic Tools

The d -dimensional Fourier transform is defined by

$$F : L_2(\mathbb{R}^d) \rightarrow L_2(\mathbb{R}^d) \text{ with } F[f](\boldsymbol{\omega}) := \frac{1}{(2\pi)^{d/2}} \int_{\mathbb{R}^d} e^{-i\langle \boldsymbol{\omega}, \mathbf{x} \rangle} f(\mathbf{x}) d\mathbf{x}. \quad (63)$$

Then its inverse transform is given by

$$F^{-1} : L_2(\mathbb{R}^d) \rightarrow L_2(\mathbb{R}^d) \text{ with } F^{-1}[f](\mathbf{x}) = \frac{1}{(2\pi)^{d/2}} \int_{\mathbb{R}^d} e^{i\langle \boldsymbol{\omega}, \mathbf{x} \rangle} f(\boldsymbol{\omega}) d\boldsymbol{\omega}. \quad (64)$$

F can be shown to be an isometry on $L_2(\mathbb{R}^d)$.

Now introduce regularization operators P defined by

$$\langle Pf, Pg \rangle := \int_{\text{supp } P(\boldsymbol{\omega})} \frac{\overline{F[f](\boldsymbol{\omega})} F[g](\boldsymbol{\omega})}{P(\boldsymbol{\omega})} d\boldsymbol{\omega} \quad (65)$$

for some nonnegative function $P(\boldsymbol{\omega})$ converging to 0 for $\|\boldsymbol{\omega}\| \rightarrow \infty$. It can be shown [49] that for a kernel to be a Green's function of P^*P , i.e.

$$\langle Pk(\mathbf{x}), Pk(\mathbf{x} - \mathbf{x}_0) \rangle = k(\mathbf{x}_0), \quad (66)$$

we need $F[k](\boldsymbol{\omega}) = P(\boldsymbol{\omega})$. For radially symmetric functions, i.e. $f(\mathbf{x}) = f(\|\mathbf{x}\|_2)$, we can explicitly carry out the integration on the sphere to obtain Fourier transform which is also radially symmetric (see e.g. [51, 32]), namely

$$F[f](\|\boldsymbol{\omega}\|) = \omega^{-\nu} H_\nu[r^\nu f(r)](\|\boldsymbol{\omega}\|), \quad (67)$$

where $\nu := \frac{1}{2}d - 1$ and H_ν is the Hankel transform over the positive real line. The latter is defined by

$$H_\nu[f](\omega) := \int_0^\infty r f(r) J_\nu(\omega r) dr. \quad (68)$$

Here J_ν is the Bessel function of the first kind defined by

$$J_\nu(r) := r^\nu 2^{-\nu} \sum_{j=0}^{\infty} \frac{(-1)^j r^{2j}}{2^{2j} j! (j + \nu + 1)!}. \quad (69)$$

Note that $H_\nu = H_\nu^{-1}$, i.e. $f = H_\nu[H_\nu[f]]$ (in L_2) due to the Hankel inversion theorem [51].

8.2 Degenerate Systems

Computing the Fourier transform for a given kernel k gives us the continuous spectrum. As pointed out in Section 6, we are interested in the discrete spectrum of integral kernels defined on \mathcal{X} . This means that the eigenvalues are defined on the grid $\omega_0 \mathbb{Z}^d$ with $\omega_0 = 2\pi/v$. Assuming $k(\mathbf{x})$ is rotationally invariant, so is $K(\boldsymbol{\omega})$ and therefore also the eigenvalues $\lambda_{\mathbf{j}} = (2\pi)^{\frac{d}{2}} K(\mathbf{j}\omega_0)$ as shown in Lemma 12. Consequently we have degeneracies in the point spectrum of the integral operator given by k (or k_v respectively) as all $\mathbf{j}\omega_0$ with equal length will have the same eigenvalue. In order to deal with this case efficiently we slightly modify Theorem 6 for our purposes. The following theorem allows proper account to be taken of the multiplicity of eigenvalues, and thus allows the straight-forward calculation of the sought for entropy numbers.

Theorem 16 *Let $(s_t)_t \in \mathbb{N}^{\mathbb{N}_0}$ be an increasing sequence with $s_0 = 1$ and $(\sigma_j)_j \in \mathbb{R}^{\mathbb{N}}$ be a non-increasing sequence of non-negative numbers with*

$$\sigma_{s_j} < \sigma_{s_{\bar{j}}} \text{ for } j < \bar{j} \text{ and } \sigma_j = \sigma_{s_t} \text{ for } s_{t-1} < j \leq s_t$$

and let

$$D\mathbf{x} = (\sigma_1 x_1, \sigma_2 x_2, \dots, \sigma_j x_j, \dots) \quad (70)$$

for $\mathbf{x} = (x_1, x_2, \dots, x_j, \dots) \in \ell_p$ be the diagonal operator from ℓ_p into itself, generated by the sequence $(\sigma_j)_j$, where $1 \leq p \leq \infty$. Then for all $n \in \mathbb{N}$,

$$\sup_{t \in \mathbb{N}} n^{-\frac{1}{s_t}} (\sigma_1 \sigma_2 \cdots \sigma_{s_t})^{\frac{1}{s_t}} \leq \epsilon_n(D) \leq 6 \sup_{t \in \mathbb{N}} n^{-\frac{1}{s_t}} (\sigma_1 \sigma_2 \cdots \sigma_{s_t})^{\frac{1}{s_t}}. \quad (71)$$

See the appendix for a proof.

This theorem allows us to obtain a similar result to corollary 7.

Corollary 17 (Entropy numbers for degenerate systems)

Let $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be a Mercer kernel and let A be defined by (34) with the additional restriction that the coefficients a_j have to match the degeneracy of

λ_j , i.e. $a_{s_j} \geq a_{s_{\bar{j}}}$ for $j < \bar{j}$ and $a_j = a_{s_t}$ for $s_{t-1} < j \leq s_t$. Then⁴ one can choose A such that

$$\epsilon_n(A: \ell_2 \rightarrow \ell_2) \leq \inf_{(a_j)_{j: (\sqrt{\lambda_j}/a_j)_j \in \ell_2}} \sup_{t \in \mathbb{N}} 6C_k \left\| \left(\sqrt{\lambda_j}/a_j \right)_j \right\|_{\ell_2} n^{-\frac{1}{s_t}} (a_1 a_2 \dots a_{s_t})^{\frac{1}{s_t}} \quad (72)$$

This result by itself may not appear too useful. However it is in fact exactly what we need for the degenerate case (it is slightly tighter than the original statement, as the supremum effectively has to be carried out only over a subset of \mathbb{N}). Finally we have to compute the degree of multiplicity that occurs for different indices \mathbf{j} . For this purpose consider shells of radius r in \mathbb{R}^d centered at the origin, i.e. rS^{d-1} , which contain a nonzero number of elements of \mathbb{Z}^d . Denote the corresponding radii by r_j and let $n(r_j, d)$ be the number of elements on these shells. Observe that $n(r, d) \neq 0$ only when $r^2 \in \mathbb{N}$. Thus

$$\begin{aligned} n(r, d) &:= |\mathbb{Z}^d \cap rS^{d-1}| \\ N(r, d) &:= \sum_{\{0 \leq \rho \leq r: \rho^2 \in \mathbb{N}\}} n(\rho, d). \end{aligned} \quad (73)$$

The determination of $n(r, d)$ is a classical problem which is completely solved by the use of the θ -series. (see e.g. [19]):

Theorem 18 (Occupation numbers of shells) *Let the formal power series $\theta(x)$ be defined by*

$$\theta(x) := \sum_{j=-\infty}^{\infty} x^{j^2} = 1 + 2 \sum_{j=1}^{\infty} x^{j^2}. \quad (74)$$

Then

$$(\theta(x))^d = \sum_{j=1}^{\infty} n(\sqrt{j}, d) x^j. \quad (75)$$

This theorem allows one to readily compute $n(r, d)$ exactly; see the appendix for some Maple code to do so. (Note that whilst there do exist closed form asymptotic approximate formulae for $n(r, d)$ [19, p. 155], they are inordinately complicated and of little use for our purposes.)

We can now construct an index of the eigenvalues which satisfies the required ordering (at least for nonincreasing functions $K(\omega)$) and we get the following result:

Corollary 19 (Radially Symmetric Systems on a Lattice)

Let $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be a Mercer kernel with eigenvalues given by a radially symmetric nonincreasing function on a lattice, i.e. $\lambda_{\mathbf{j}} = \lambda(\|\mathbf{j}\|)$ with $\mathbf{j} \in \mathbb{Z}^d$ and let A be defined by (34) with the additional restriction that the coefficients $a_{\mathbf{i}}$

⁴ See [20] for a proof that (72) is well defined for kernels with summable eigenvalues.

have to match the degeneracy of $\lambda_{\mathbf{j}}$, i.e. $a_{\mathbf{j}} = a(\|\mathbf{j}\|)$. Then

$$\epsilon_n(A: \ell_2 \rightarrow \ell_2) \leq \inf_{(a_{\mathbf{j}})_{\mathbf{j}}: \left(\frac{\sqrt{\lambda_{\mathbf{j}}}}{a_{\mathbf{j}}}\right)_{\mathbf{j}} \in (\ell_2)^d} \sup_{t \in \mathbb{N}} 6C_k \left\| \left(\frac{\sqrt{\lambda_{\mathbf{j}}}}{a_{\mathbf{j}}}\right)_{\mathbf{j}} \right\|_{(\ell_2)^d} n^{-\frac{1}{N(r_t, d)}} \left(\prod_{q=1}^t a(r_q)^{n(r_q, d)} \right)^{\frac{1}{N(r_t, d)}}. \quad (76)$$

Note that this result, although it may seem straightforward, cannot be obtained from corollary 7 directly as there the sup would have to be carried out over \mathbb{N} instead of $(N(r_t, d))_t$. The different formulation allows us to compute bounds on the entropy numbers more easily.

8.3 Bounds for Kernels in \mathbb{R}^d

Let us conclude this section with some examples of the eigenvalue sequences for kernels typically used in SV machines. These can then be used to evaluate the right hand side in corollary 19. Recall that $\nu = \frac{d}{2} - 1$. First we have to compute the Fourier/Hankel transform for the kernels.

Example 20 (Gaussian RBFs) For Gaussian rbfs in d dimensions we have $k(r) = \sigma^{-d} e^{-\frac{r^2}{2\sigma^2}}$ and correspondingly

$$\begin{aligned} F[k](\omega) &= \omega^{-\nu} \sigma^{-d} H_{\nu} \left[r^{\nu} e^{-\frac{r^2}{2\sigma^2}} \right] (\omega) \\ &= \omega^{-\nu} \sigma^{2(\nu+1)-d} \omega^{\nu} e^{-\frac{\omega^2 \sigma^2}{2}} \\ &= e^{-\frac{\omega^2 \sigma^2}{2}}. \end{aligned}$$

Example 21 (Exponential RBFs) In the case of $k(r) = e^{-ar}$ we get

$$\begin{aligned} F[k](\omega) &= \omega^{-\nu} H_{\nu} \left[r^{\nu} e^{-ar} \right] (\omega) \\ &= \omega^{-\nu} 2^{\nu+1} \omega^{\nu} a \pi^{-\frac{1}{2}} \left(\nu + \frac{3}{2} \right) (a^2 + \omega^2)^{-\nu - \frac{3}{2}} \\ &= 2^{\frac{d}{2}} a \pi^{-\frac{1}{2}} \left(\frac{d}{2} + 1 \right) \frac{1}{(a^2 + \omega^2)^{\frac{d+1}{2}}} \end{aligned}$$

i.e. in the case of $d = 1$ we recover the damped harmonic oscillator (in frequency domain). In general we get a decay in terms of the eigenvalues like $\omega^{-(d+1)}$. Moreover we can conclude from this that the Fourier transform of k , viewed itself as a kernel, i.e. $k(r) = (1 + r^2)^{-\frac{d+1}{2}}$, yields the initial kernel as its corresponding power spectrum in Fourier domain.

Example 22 (Damped Harmonic Oscillator) Another way to generalize the harmonic oscillator, this time in a way, that k does not depend on the dimensionality d is to set $k(r) = \frac{1}{a^2 + r^2}$. Following [58, section 13.6] we get

$$F[k](\omega) = \omega^{-\nu} H_{\nu} \left[\frac{r^{\nu}}{a^2 + r^2} \right] (\omega)$$

$$= \omega^{-\nu} a^\nu K_\nu(\omega a)$$

where K_ν is the Bessel function of the second kind, defined by (see [51])

$$K_\nu(x) = \int_0^\infty e^{-x \cosh t} \cosh(\nu t) dt. \quad (77)$$

It is possible to upper bound $F[k]$ via

$$K_\nu(x) = \sqrt{\frac{\pi}{2x}} e^{-x} \left[\sum_{j=0}^{p-1} (2x)^{-j} \frac{? (\nu + j + \frac{1}{2})}{j! (\nu - j + \frac{1}{2})} + \theta \cdot (2x)^{-p} \frac{? (\nu + p + \frac{1}{2})}{j! (\nu - p + \frac{1}{2})} \right] \quad (78)$$

with $p > \nu - \frac{1}{2}$ and $\theta \in [0, 1]$ [18, eq. 8.451.6]). As one can see the term in the brackets $[\cdot]$ converges to 1 for $x \rightarrow \infty$ and we get exponential decay of the eigenvalues.

Using Theorem 18, Corollary 19 and Remark 13 one may compute the entropy numbers numerically for a particular kernel and a particular set of parameters. This may seem unsatisfactory from a theoretician's point of view. However, as the ultimate goal is to use the obtained bounds for model selection, it is desirable to obtain as tight bounds (especially in the constants) as possible. Hence if much more precise bounds can be obtained by some not too expensive numerical calculation it is definitely worth while to use those instead of a theoretically nice but not sufficiently tight upper bound. The computational effort to calculate these quantities is typically negligible in comparison to training the actual learning machine.

Notwithstanding the above, in order to give a feeling for the effect of the decay of the Fourier transform of the kernel on the entropy numbers of the A operator, we conclude with the following general result, the proof of which is relegated to the appendix.

Proposition 23 (Polynomial exponential decay in \mathbb{R}^d) For kernels $k(\cdot, \cdot)$ in $\mathbb{R}^d \times \mathbb{R}^d$ with $\lambda(\omega) = O(e^{-\alpha \|\omega\|^p})$ with $\alpha, p > 0$ the following bound on the entropy number of the corresponding scaling operator holds.

$$\ln \epsilon_n^{-1}(A: \ell_2 \rightarrow \ell_2) = O(\ln^{\frac{p}{p+d}} n)$$

9 Conclusions

We have shown how to connect properties known about mappings into feature spaces with bounds on the covering numbers. Our reasoning relied on the fact that this mapping exhibits certain decay properties to ensure rapid convergence and a constraint on the size of the weight vector in feature space. This means that the corresponding algorithms have to restrict exactly this quantity to ensure good generalization performance. This is exactly what is done in Support Vector machines.

The actual application of our results, perhaps for model selection using structural risk minimization, is somewhat involved, but is certainly doable — see the references below. Here we outline one possible path. As said before, the viewpoint in this paper is new, and perhaps there will be refinements soon forthcoming which would make the codification of our existing results into a single generalization bound premature.

9.1 A Possible Procedure to use the Results of this Paper

Choose k and σ The kernel k may be chosen for a variety of reasons, which we have nothing additional to say about here. The choice of σ should take account of the discussion in Section 6.

Choose the period v of the kernel One suggested procedure is outlined in Section 6.

Bound $\epsilon_n(A)$ This can be done using Corollary 7 (for the case $d = 1$) or Corollary 17 or 19 for the case $d > 1$. Some examples of this sort of calculation are given in Section 7.

Bound $\epsilon_n(T)$ Using Theorem 10.

Take account of the “+ b ” The key observation is that given a class \mathcal{F} with known $\mathcal{N}^m(\epsilon, \mathcal{F})$, one can bound $\mathcal{N}^m(\epsilon, \mathcal{F}^+)$ as follows. (Here $\mathcal{F}^+ := \{f + b : f \in \mathcal{F}, b \in \mathbb{R}\}$.) Suppose V_ϵ is an ϵ -cover for \mathcal{F} and elements of \mathcal{F}^+ are uniformly bounded by B (this implies a limit on $|b|$ as well as a uniform bound on elements of \mathcal{F}). Then

$$V_\epsilon^+ := \bigcup_{j=-B/\epsilon}^{B/\epsilon} V_\epsilon + j\epsilon$$

is an ϵ -cover for \mathcal{F}^+ and thus $\mathcal{N}^m(\epsilon, \mathcal{F}^+) \leq \frac{2B}{\epsilon} \mathcal{N}^m(\epsilon, \mathcal{F})$. Observe that this will only be “noticeable” for classes \mathcal{F} with very slowly growing covering numbers (polynomial in $1/\epsilon$).

Take account of the loss function using Lemma 2 for example.

Plug into a uniform convergence result See the pointers to the literature and the example in Section 4.

9.2 Further Work

The operator-theoretic viewpoint introduced in this paper seems fruitful. The overall bounds for SV classes can, via a somewhat involved argument, be considerably simplified [20]. The general approach can be applied to various other learning machines such as convex combinations of basis functions and multilayer networks [48]. When combined with an appropriate statistical argument [46] they yield bounds on generalization performance that appear to work well for model

selection [41]. The methods can also be applied to some problems of unsupervised learning [50].

The results of the present paper hinge on the measurement of the size of the weight vector \mathbf{w} by a ℓ_2 norm. In [60] we show the effect of different norms for measuring the size of \mathbf{w} , as well as presenting a number of related results.

We expect that further refinements and extensions to these techniques will continue to yield valuable results.

Acknowledgements

The authors thank Peter Bartlett, Bernd Carl, André Elisseeff, Ying Guo, John Shawe–Taylor, and Anja Westerhoff for helpful discussions and comments. This work was supported by the Australian Research Council, a grant of the DFG (# Ja 379/71) and Neurocolt II. This work would not have been possible had the European cup final not been held in London in 1996 immediately prior to COLT96.

A Empirical Bounds

Instead of theoretically determining the shape of $\Phi(\mathcal{X})$ *a priori* one could use the training and/or test data to empirically estimate its shape and use this quantity to compute an operator B_{emp} in place of A to (32). In this subsection we will sketch a possible approach — the full development of these ideas requires considerable further work; see [41].

For instance assume that we are given an m -sample of data points $\mathbf{X} := \{\mathbf{x}_1, \dots, \mathbf{x}_m\} \subset \mathcal{X}$, not necessarily only from the training set but perhaps also comprising unlabelled test samples, drawn from the same distribution P . Now suppose we could estimate the first j radii (e.g. in a manner similar to the radius estimate in [40]) $\{r_1, \dots, r_j\}$ of an ellipsoid enclosing $\Phi(\mathbf{X})$ with probability say $1 - \eta$.⁵ Denote by $\{\mathbf{e}_1, \dots, \mathbf{e}_{j-1}\} \subset \ell_2$ a set of orthogonal vectors pointing in the directions given by the radii and P_j the projector onto $(\text{span}\{\mathbf{e}_1, \dots, \mathbf{e}_{j-1}\})$. Note that we only have to be sure that with probability $1 - \eta$ the data lies inside the ellipsoid and that we need no statement on the precision of the estimate of the radii — this makes a big difference in terms of the volume. Then with probability $1 - \eta$ we could upper bound the covering numbers of a scaling operator B_{emp} by making use of corollary 7. Due to the ellipsoid condition the following inequality holds for all \mathbf{x}_s :

$$\sum_{t=1}^{j-1} \frac{\langle \mathbf{e}_t, \Phi(\mathbf{x}_t) \rangle^2}{r_t^2} \leq 1 \quad (79)$$

and moreover $\|P_j \Phi(\mathbf{x}_t)\|_{\ell_2} \leq r_j$ for all $1 \leq t \leq M$. Hence for an operator B_{emp}^{-1} scaling the first $j-1$ directions $\mathbf{e}_1, \dots, \mathbf{e}_{j-1}$ by r_t^{-1} and the rest by r_j^{-1} , $B_{\text{emp}}^{-1} \Phi(\mathbf{X})$ would still be enclosed in $\sqrt{2}U_{\ell_2}$. Hence we have a similar situation as in the case where we explicitly computed all eigenvalues analytically. Setting $r_t := r_j$ for $t > j$ and applying corollary 7 leads to the following upper bound on the entropy of B_{emp}

$$\epsilon_n(B_{\text{emp}}) \leq 6\sqrt{2} \sup_{j \in \mathbb{N}} n^{-\frac{1}{j}} (r_1 r_2 \cdots r_j)^{\frac{1}{j}} \quad (80)$$

The aim is now to find the maximum n for which the estimate will not break down yet (for we have the liberty of distributing the covering numbers arbitrarily between the shrinkage operator B_{emp} and the actual evaluation operator $S_{\mathbf{X}^m}$ as shown in section 5). In other words we are looking for that particular value of j where \sup_j is taken on for the smallest radius estimated. Ignoring the fact that $n \in \mathbb{N}$ for a moment we arrive at the following equation:

$$n^{-\frac{1}{j}} (r_1 r_2 \cdots r_j)^{\frac{1}{j}} = n^{-\frac{1}{j+1}} (r_1 r_2 \cdots r_j \cdot r_j)^{\frac{1}{j+1}} \quad (81)$$

⁵This for instance could be done in a way similar to Kernel-PCA [49] by computing the eigensystem (λ_i, α_{ij}) of the Gram matrix $k_{ij} = k(x_i, x_j)$. Then, possibly after some ordering, $\sqrt{\lambda_i} \geq r_i = \sqrt{\lambda_i} \max_j |\alpha_{ij}|$.

Solving for n yields and taking $n \in \mathbb{N}$ into account yields

$$n_j = \left\lceil \frac{r_1 r_2 \cdots r_j}{r_j^j} \right\rceil \text{ and therefore } \epsilon_{n_j+1}(B_{\text{emp}}) \leq 6\sqrt{2}r_j \quad (82)$$

This calculation is valid as n_j is a nondecreasing function of j : $\frac{n_{j+1}}{n_j} = \left(\frac{r_j}{r_{j+1}}\right)^j$. If this assumption failed to be true one would have to redefine B_{emp} to scale only the first \tilde{j} directions for which this happened to hold — it gains one little to scale in directions where the decay rate is too slow.

Instead of taking real data (which may be expensive to get) we also could upper bound the first j radii by a Monte–Carlo method, once we can bound the set \mathcal{X} . This is also useful when no analytic expansion in terms of eigenvalues of the operator can be obtained or where it would be too tedious to obtain explicitly. In cases with a sufficient amount of computational power available this may even be a more practical and faster way than computing the spectrum given by k analytically. The latter, at least in order to obtain optimal bounds, would have to be done for each learning problem anew. The method proposed here would obviate the need for such detailed theoretical calculations which may be impractical to carry out in some instances.

B Proofs of Results in Section 7

Proof (Proposition 14) The proof uses Corollary 7. Since $\lambda_j = O(j^{-\alpha-1})$ there exists some $\beta \in \mathbb{R}_+$ with $\lambda_j \leq \beta^2 j^{-\alpha-1}$. In this case all sequences $(a_j)_j = (j^{-\frac{\tau}{2}})_j$ with $0 < \tau < \alpha$ lead to an admissible scaling property. One has

$$\left\| \left(\frac{\sqrt{\lambda_j}}{a_j} \right)_j \right\|_{\ell_2} = \beta \left\| \left(j^{\frac{\tau-\alpha-2}{2}} \right)_j \right\|_{\ell_2} = \beta \sqrt{\zeta(\alpha - \tau + 1)} \quad (83)$$

where $\zeta(\cdot)$ is Riemann’s zeta function. Moreover one can bound $\zeta(\cdot)$ by

$$x + \gamma \leq \zeta\left(1 + \frac{1}{x}\right) \leq x + 1 \quad (84)$$

where γ is Euler’s constant. The next step is to evaluate the expression

$$(a_1 a_2 \cdots a_j)^{\frac{1}{j}} = \left(\prod_{s=1}^j s^{-\frac{\tau}{2}} \right)^{\frac{1}{j}} = (j!)^{-\frac{\tau}{2j}} = ?(j+1)^{-\frac{\tau}{2}} \quad (85)$$

The Gamma function $?(x)$ can be bounded as follows; for $j > 1$,

$$\ln j - 1 \leq \frac{1}{j} \ln ?(j+1) \leq \ln j. \quad (86)$$

Hence one may bound $\epsilon_n(A)$

$$\epsilon_n(A) \geq C_k \beta \inf_{\tau \in (0, \alpha)} \sup_{j \in \mathbb{N}} n^{-\frac{1}{j}} \left(\frac{1}{\alpha - \tau} + \gamma \right)^{\frac{1}{2}} j^{-\frac{\tau}{2}} \quad (87)$$

$$\epsilon_n(A) \leq 6C_k \beta \inf_{\tau \in (0, \alpha)} \sup_{j \in \mathbb{N}} n^{-\frac{1}{j}} \left(\frac{1}{\alpha - \tau} + 1 \right)^{\frac{1}{2}} e^{\frac{\tau}{2}} j^{-\frac{\tau}{2}} \quad (88)$$

In order to avoid unneeded technicalities we will replace $\sup_{j \in \mathbb{N}}$ by $\sup_{j \in [1, \infty)}$. This is no problem when computing the upper bound, but it is an issue for the lower bound. However, $j^{-\frac{\tau}{2}}$ on $[1, \infty)$ is within a constant factor of $2^{-\frac{\tau}{2}}$ of its corresponding values on the integer domain \mathbb{N} , the biggest discrepancy being at $[1, 2]$.⁶ Thus we may safely ignore the concern. Next we compute

$$\sup_{j \in [1, \infty)} n^{-\frac{1}{j}} j^{-\frac{\tau}{2}} = \sup_{j \in [1, \infty)} e^{-\frac{1}{j} \ln n - \frac{\tau}{2} \ln j} = \left(\frac{2e \ln n}{\tau} \right)^{-\frac{\tau}{2}}. \quad (89)$$

The maximum of the argument is obtained for $j = \frac{2 \ln n}{\tau}$, hence (89) holds for all $\ln n \geq \frac{\tau}{2}$, which is fine since we want to compute bounds on $\epsilon_n(A)$ as $n \rightarrow \infty$. For the lower bounds on $\epsilon_n(A)$ we obtain

$$\begin{aligned} \epsilon_n(A) &\geq C_k \beta (2e)^{-\frac{\tau}{2}} \inf_{\tau \in (0, \alpha)} \left(\frac{1}{\alpha - \tau} + \gamma \right)^{\frac{1}{2}} \left(\frac{2 \ln n}{\tau} \right)^{-\frac{\tau}{2}} \\ &\geq C_k \beta (2e) 2^{-\frac{\tau}{2}} \inf_{\tau \in (0, \alpha)} \left(\frac{1}{\alpha - \tau} + \gamma \right)^{\frac{1}{2}} \inf_{\tau \in (0, \alpha)} \left(\frac{2 \ln n}{\tau} \right)^{-\frac{\tau}{2}} \\ &= C_k \beta (2e) 2^{-\frac{\tau}{2}} \left(\frac{1}{\alpha} + \gamma \right)^{\frac{1}{2}} \left(\frac{2 \ln n}{\alpha} \right)^{-\frac{\alpha}{2}}. \end{aligned} \quad (90)$$

This shows that $\epsilon_n(A)$ is always bounded from below by $\Omega \left(\ln^{-\frac{\alpha}{2}} n \right)$. Computation of the upper bound is slightly more effort, since one has to evaluate

$$\epsilon_n(A) \leq 6C_k \beta \inf_{\tau \in (0, \alpha)} \left(\frac{1}{\alpha - \tau} + 1 \right)^{\frac{1}{2}} \left(\frac{2 \ln n}{\tau} \right)^{-\frac{\tau}{2}}. \quad (91)$$

Clearly for any fixed $\tau \in (0, \alpha)$ we are able to obtain a rate of $\epsilon_n(A) = O \left(\ln^{-\frac{\tau}{2}} n \right)$, thus the theorem follows. For practical purposes a good approximation of the inf can be found as $\frac{1}{\alpha - \tau} = \ln(2 \ln n)$ by computing the derivative of the argument in (91) wrt. τ and dropping all terms independent of τ and n . However, numerical minimization of (91) is more advisable when small values of $\epsilon_n(A)$ are crucial. ■

For the proof of Proposition 15 we need the following standard Lemma:

⁶One may show [41] that $a_{j^*+1} \leq \sup_{j \in \mathbb{N}} n^{-\frac{1}{j}} (a_1, \dots, a_j)^{\frac{1}{j}} \leq a_{j^*}$ for that particular j^* where $\sup_{j \in \mathbb{N}}$ is actually obtained. Hence the maximum quotient a_{j+1}/a_j , which in the present case is $2^{-\frac{\tau}{2}}$, determines the value by which the bound has to be lowered in order to obtain a true lower bound.

Lemma 24 (Summation and Integration in \mathbb{R}^1) Suppose $f : \mathbb{R} \rightarrow \mathbb{R}$ is an integrable nonincreasing function. Then the following inequality holds for any $a \in \mathbb{Z}$

$$\int_a^\infty f(x)dx \leq \sum_{n=a}^\infty f(n) \leq \int_{a-1}^\infty f(x)dx. \quad (92)$$

Proof The proof relies on the fact that

$$f(n) \geq \int_n^{n+1} f(n)dn \geq f(n+1)$$

due to the monotonicity of f and a decomposition of the integral $\int_0^\infty = \sum_{n=0}^\infty \int_n^{n+1}$. The lemma is a direct consequence thereof. \blacksquare

Proof (Proposition 15) Since $\lambda_j = O(e^{-\alpha j^p})$ there exists some $\beta \in \mathbb{R}^+$ with $\lambda_j \leq \beta^2 e^{-\alpha j^p}$. Similarly to before we now use a series $(a_j)_j = e^{-\tau/2j^p}$. Then by applying lemma 24 we have that for any $\tau \in [0, \alpha)$,

$$\left\| \left(\frac{\sqrt{\lambda_j}}{a_j} \right)_j \right\|_{\ell_2} = \beta \left(\sum_{j=0}^\infty e^{(\tau-\alpha)j^p} \right)^{\frac{1}{2}} \begin{cases} \leq \beta \sqrt{1 + \frac{\Gamma(1/p)}{p(\alpha-\tau)^{1/p}}} \\ \geq \beta \sqrt{\frac{\Gamma(1/p)}{p(\alpha-\tau)^{1/p}}} \end{cases} \quad (93)$$

Next we have to apply a similar bound to the product of the first j diagonal entries of the scaling operator A .

$$(a_1 a_2 \dots a_j)^{\frac{1}{j}} = e^{-\frac{1}{2j}\tau \sum_{s=1}^j s^p} \begin{cases} \geq e^{-\frac{\tau}{2(p+1)}j^p} \\ \leq e^{-\frac{\tau}{2(p+1)}j^p + \frac{\tau}{2j(p+1)}} \leq e^{-\frac{\tau}{2(p+1)}j^p + \frac{\tau}{2(p+1)}} \end{cases} \quad (94)$$

The last inequality holds since $j \in \mathbb{N}$. Next we have to compute $\sup_{j \in \mathbb{N}} n^{-\frac{1}{j}} e^{-\frac{\tau}{2(p+1)}j^p} = \sup_{j \in \mathbb{N}} e^{-\frac{1}{j} \ln n - \frac{\tau}{2(p+1)}j^p}$. Differentiation of the exponent wrt. j leads to

$$j^{-2} \ln n - \frac{\tau p}{2(p+1)} j^{p-1} = 0 \Rightarrow j^{-1} \ln n = \frac{\tau p}{2(p+1)} j^p \quad (95)$$

and thus

$$\sup_{j \in [1, \infty)} n^{-\frac{1}{j}} e^{-\frac{\tau}{2(p+1)}j^p} = e^{-\left(\frac{\tau}{2}\right)^{1/(p+1)} \left(\frac{p+1}{p} \ln n\right)^{\frac{p}{p+1}}}. \quad (96)$$

Replacing the domain from $\sup_{j \in \mathbb{N}}$ to $\sup_{j \in [1, \infty)}$ is not a problem when it comes to computing upper bounds on $\epsilon_n(A)$. As for the lower bounds, again, a similar reasoning like in the previous proof would have to be applied.⁷ However, it is left out for the sake of clarity. Thus $\epsilon_n(A)$ can be bounded from below as

⁷As in the previous theorem, the problem reduces to bounding the quotient a_{j^*+1}/a_{j^*} where j^* is the variable for which $\sup_{j \in \mathbb{N}}$ is obtained. However, here the quotient can only be bounded by $e^{-\frac{\tau p}{2}j^{p-1}}$. Fortunately this is of lower order than the remaining terms, hence it will not change the *rate* of the lower bounds.

follows

$$\begin{aligned}
\epsilon_n(A) &\geq C_k \beta \inf_{\tau \in (0, \alpha)} \sqrt{\frac{\Gamma(1/p)}{p(\alpha-\tau)^{1/p}}} e^{-\left(\frac{\tau}{2}\right)^{1/(p+1)} \left(\frac{p+1}{p} \ln n\right)^{\frac{p}{p+1}}} \\
&\geq C_k \beta \inf_{\tau \in (0, \alpha)} \sqrt{\frac{\Gamma(1/p)}{p(\alpha-\tau)^{1/p}}} \inf_{\tau \in (0, \alpha)} e^{-\left(\frac{\tau}{2}\right)^{1/(p+1)} \left(\frac{p+1}{p} \ln n\right)^{\frac{p}{p+1}}} \\
&= C_k \beta \sqrt{\frac{\Gamma(1/p)}{p\alpha^{1/p}}} e^{-\left(\frac{\alpha}{2}\right)^{1/(p+1)} \left(\frac{p+1}{p} \ln n\right)^{\frac{p}{p+1}}} \tag{97}
\end{aligned}$$

Thus a lower bound on the rate of $\log \epsilon_n$ is $O(\log^{\frac{p}{p+1}} n)$. Moreover, for the upper bound we obtain

$$\epsilon_n(A) \leq 6C_k \beta \inf_{\tau \in (0, \alpha)} \sqrt{1 + \frac{\Gamma(1/p)}{p(\alpha-\tau)^{1/p}}} e^{-\left(\frac{\tau}{2}\right)^{1/(p+1)} \left(\frac{p+1}{p} \ln n\right)^{\frac{p}{p+1}} + \frac{\tau}{2j(p+1)}} \tag{98}$$

One could evaluate (98) numerically. However, it can be seen that for any fixed $\tau \in (0, \alpha)$ the rate of $\log \epsilon_n(A)$ can be bounded by $O(\log^{\frac{p}{p+1}} n)$, which shows that the obtained rates are tight. ■

C Proof of Theorem 16

Proof The first part of the inequality follows directly from theorem 6 as it is a weaker statement than the original one. We prove the second part by closely mimicking the proof in [13, p. 17]. We define

$$\delta(n) := 8 \sup_{t \in \mathbb{N}} n^{-\frac{1}{st}} (\sigma_1 \sigma_2 \cdots \sigma_{s_t})^{\frac{1}{st}} \tag{99}$$

and show that for all n there is an index s_j with $\sigma_{s_j+1} \leq \frac{\delta(n)}{4}$. For this purpose choose an index r such that $n \leq 2^{s_j+1}$ and thus $1 \leq 2n^{-1/(s_j+1)}$. Moreover we have

$$\sigma_{s_j+1} \leq (\sigma_1 \sigma_2 \cdots \sigma_{s_j+1})^{\frac{1}{s_j+1}} \tag{100}$$

because of the monotonicity of $(\sigma_j)_j$ and finally

$$\sigma_{s_j+1} \leq 2n^{-1/(s_j+1)} (\sigma_1 \sigma_2 \cdots \sigma_{s_j+1})^{\frac{1}{s_j+1}}. \tag{101}$$

Using the definition of $\delta(n)$ we thus conclude $\sigma_{s_j+1} \leq \delta(n)/4$. If this happens to be the case for σ_1 we have $\epsilon_n(D) \leq \sigma_1$ which proves the theorem.

If this is not the case there exists an index s_j such that $\sigma_{s_j+1} \leq \delta(n)/4 < \sigma_{s_j}$. Hence the corresponding sectional operator

$$\begin{aligned}
D_{s_j} : \ell_p &\rightarrow \ell_p \text{ with} \\
D_{s_j}(x_1, x_2, \dots, x_{s_j}, x_{s_j+1}, \dots) &= (\sigma_1 x_1, \sigma_2 x_2, \dots, \sigma_{s_j} x_{s_j}, 0, 0, \dots) \tag{102}
\end{aligned}$$

is of rank s_j and the image $D_{s_j}(U_p)$ of the closed unit ball U_p of ℓ_p is isometric to the subset $D^{(s_j)}(U_p^{(s_j)})$ of $\ell_p^{s_j}$. In any case $D_{s_j}(U_p)$ is a precompact subset of ℓ_p . So let y_1, y_2, \dots, y_N be a maximal system of elements in $D_{s_j}(U_p)$ with

$$\|y_j - y_{\bar{j}}\| > \delta(n)/2 \text{ for } j \neq \bar{j}. \quad (103)$$

The maximality of this system guarantees that

$$D_{s_j}(U_p) \subseteq \bigcup_{j=1}^N \left\{ y_j + \frac{\delta(n)}{2} U_p \right\} \quad (104)$$

and thus $\epsilon_N(D_{s_j}) \leq \delta(n)/2$. In order to get an estimate for $\epsilon_N(D)$ we split the operator D into two parts $D = (D - D_{s_j}) + D_{s_j}$ which allows us to bound

$$\epsilon_N(D) \leq \|D - D_{s_j}\| + \epsilon_N(D_{s_j}). \quad (105)$$

Using $\|D - D_{s_j}\| = \sigma_{s_j+1} \leq \delta(n)/4$ and the bound on $\epsilon_N(D_{s_j})$ we arrive at

$$\epsilon_N(D) \leq \frac{3}{4} \delta(n). \quad (106)$$

The final step is to show that $N \leq n$ as then by substituting in the definition of $\delta(n)$ into (106) yields the result. This is again achieved by a comparison of volumes. Consider the sets $\{y_j + (\delta(n)/4)U_p^{s_j}\}$ as subsets of the space $\ell_p^{s_j}$ which is possible since $y_j \in D_{s_j}(U_p)$ and $D_{s_j}(U_p) = D^{(s_j)}(U_p^{s_j})$. These sets are obviously pairwise disjoint. On the other hand we have

$$\bigcup_{j=1}^N \left\{ y_j + \frac{\delta(n)}{4} U_p^{s_j} \right\} \subseteq D^{(s_j)}(U_p^{s_j}) + \frac{\delta(n)}{4} U_p^{s_j} \subseteq 2D^{(s_j)}(U_p^{s_j}) \quad (107)$$

as $\delta(n)/4 < \sigma_1$. Now a comparison of the d -dimensional Euclidean volumes vol_d provides

$$N \left(\frac{\delta(n)}{4} \right)^{s_j} \text{vol}_{s_j}(U_p^{s_j}) \leq 2^{s_j} \sigma_1 \sigma_2 \cdots \sigma_{s_j} \text{vol}_{s_j}(U_p^{s_j}) \quad (108)$$

and therefore $N \leq (8/\delta(n))^{s_j} \sigma_1 \sigma_2 \cdots \sigma_{s_j}$. Using the definition of $\delta(n)$ this yields $N \leq n$. ■

D Proof of Proposition 23

Proof We will completely ignore the fact that we are actually dealing with a countable set of eigenvalues on a lattice and replace all summations by integrals without further worry. Of course this is not accurate but still will give us the correct rates for the entropy numbers.

Denote $1/\Lambda := (2\pi/v)^{\frac{d}{2}}$ the size of a unit cell, i.e. $\Lambda = (v/(2\pi))^{\frac{d}{2}}$ the

density of lattice points in frequency space as given in section 6. Then we get for infinitesimal volumes dV and numbers of points dN in frequency space

$$dV = S_{d-1}r^{d-1}dr \text{ and therefore } dN = \Lambda S_{d-1}r^{d-1}dr \quad (109)$$

(here S_{d-1} denotes the volume of the $d-1$ dimensional unit sphere) leading to

$$N(r, d) = \frac{1}{d}\Lambda S_{d-1}r^d. \quad (110)$$

We introduce a scaling operator whose eigenvalues decay like $a(\omega) = e^{-\frac{\tau}{2}\|\omega\|^p}$ for $\tau \in [0, \alpha]$. It is straightforward to check that all these values lead to both useful and admissible scaling operators. Now we will estimate the separate terms in (76).

$$\begin{aligned} \left\| \left(\frac{\sqrt{\lambda_{\mathbf{i}}}}{a_{\mathbf{i}}} \right)_{\mathbf{i}} \right\|_{\ell_2}^2 &\approx \int dN(\omega) \frac{\lambda(\omega)}{a^2(\omega)} = S_{d-1}\Lambda \int_0^\infty r^{d-1} \beta^2 e^{-(\alpha-\tau)\|\omega\|^p} \\ &= S_{d-1}\Lambda \beta^2 (\alpha - \tau)^{-\frac{d}{p}} \tau \left(\frac{d}{p} \right) p^{-1} \end{aligned} \quad (111)$$

Next we have

$$\ln \left(n^{-\frac{1}{N(r,d)}} \right) = -\frac{d}{\Lambda S_{d-1}r_0^d} \ln n \quad (112)$$

and

$$\begin{aligned} \ln \left(a_1 \cdot a_2 \cdots a_{N(r,d)} \right)^{\frac{1}{N(r,d)}} &= -\frac{d}{\Lambda S_{d-1}r_0^d} \sum_{j=1}^{N(r,d)} \ln a_j \\ &\approx dr^{-d} \int_0^r \omega^{d-1} \ln a(\omega) d\omega \\ &= -dr^{-d} \int_0^r \omega^{d-1} \frac{\tau}{2} \omega^p d\omega = -\frac{\tau}{2} \frac{d}{d+p} r^p. \end{aligned} \quad (113)$$

This leads to

$$\epsilon_n \leq 6C_k \beta \sqrt{\frac{S_{d-1}\Lambda \Gamma\left(\frac{d}{p}\right)}{p}} \inf_{\tau \in [0, \alpha]} (\alpha - \tau)^{-\frac{d}{2p}} \sup_{r \in \mathbb{R}^+} \exp \left(-\frac{d}{\Lambda S_{d-1}r^d} \ln n - \frac{\tau}{2} \frac{d}{d+p} r^p \right). \quad (115)$$

Computing the $\sup_{r \in \mathbb{R}^+}$ yields

$$r = \left(\frac{2}{\tau \Lambda S_{d-1}} \frac{(d+p)d}{p} \ln n \right)^{\frac{1}{d+p}} \quad (116)$$

and therefore

$$\epsilon_n \leq 6C_k \beta \sqrt{\frac{S_{d-1}\Lambda \Gamma\left(\frac{d}{p}\right)}{p}} \inf_{\tau \in [0, \alpha]} (\alpha - \tau)^{-\frac{d}{2p}} \exp \left(-\left(\frac{\tau}{2}\right)^{\frac{d}{d+p}} \left(\frac{(d+p)d}{p} \frac{\ln n}{\Lambda S_{d-1}} \right)^{\frac{p}{d+p}} \right). \quad (117)$$

Already from this expression one can observe the rate bounds on ϵ_n . What remains to be done is to compute the \inf_τ . This can be done by differentiating (117) w.r.t. τ . Define

$$T_n := \left(\frac{(d+p)d}{p} \frac{\ln n}{\Lambda S_{d-1}} \right)^{\frac{p}{d+p}} \quad (118)$$

which leads to the optimality condition on τ

$$(\alpha - \tau)\tau^{-\frac{p}{d+p}} = \frac{d+p}{2T_n p} 2^{\frac{d}{d+p}} \text{ with } \tau \in (0, \alpha] \quad (119)$$

which can be solved numerically. ■

E Maple code to compute $n(r,d)$

```
# This code defines a function t where
# t(m,d) is number of points on a sphere of radius^2=m from Z^d
h:=n->eval('if'(isolve(m^2=n,m)=NULL,0,'if'(n=0,1,2)),1):
powseries[powcreate](theta(n)=h(n)):
t:=(m,d)->
  coeff(convert(powseries[tpsform](powseries[evalpow](theta^d),
    x,m+1),polynom),x,m):
```

References

- [1] M. A. Aizerman, E. M. Braverman, and L. I. Rozonoér. Theoretical foundations of the potential function method in pattern recognition learning. *Automation and Remote Control*, 25:821–837, 1964.
- [2] S. Akashi. Characterization of ϵ -entropy in Gaussian processes. *Kodai Mathematical Journal*, 9:58–67, 1986.
- [3] S. Akashi. The asymptotic behaviour of ϵ -entropy of a compact positive operator. *Journal of Mathematical Analysis and Applications*, 153:250–257, 1990.
- [4] N. Alon, S. Ben-David, N. Cesa-Bianchi, and D. Haussler. Scale-sensitive Dimensions, Uniform Convergence, and Learnability. *Journal of the Association for Computing Machinery*, 44(4):615–631, 1997.
- [5] M. Anthony. Probabilistic analysis of learning in artificial neural networks: The PAC model and its variants. *Neural Computing Surveys*, 1:1–47, 1997. <http://www.icsi.berkeley.edu/~jagota/NCS>.
- [6] M. Anthony and P. Bartlett. *A Theory of Learning in Artificial Neural Networks*. Cambridge University Press, 1999.

- [7] R. Ash. *Information Theory*. Interscience Publishers, New York, 1965.
- [8] P. Bartlett and J. Shawe-Taylor. Generalization performance of support vector machines and other pattern classifiers. In B. Schölkopf, C.J.C. Burges, and A.J. Smola, editors, *Advances in Kernel Methods — Support Vector Learning*, pages 43–54, Cambridge, MA, 1999. MIT Press.
- [9] P.L. Bartlett, P. Long, and R.C. Williamson. Fat-Shattering and the Learnability of Real-Valued Functions. *Journal of Computer and System Sciences*, 52(3):434–452, 1996.
- [10] B. E. Boser, I. M. Guyon, and V. N. Vapnik. A training algorithm for optimal margin classifiers. In D. Haussler, editor, *5th Annual ACM Workshop on COLT*, pages 144–152, Pittsburgh, PA, 1992. ACM Press.
- [11] B. Carl. Entropy numbers of diagonal operators with an application to eigenvalue problems. *Journal of Approximation Theory*, 32:135–150, 1981.
- [12] B. Carl. Inequalities of Bernstein-Jackson-type and the degree of compactness of operators in Banach spaces. *Annales de l’Institut Fourier*, 35(3):79–118, 1985.
- [13] B. Carl and I. Stephani. *Entropy, compactness, and the approximation of operators*. Cambridge University Press, Cambridge, UK, 1990.
- [14] C. Cortes and V. Vapnik. Support vector networks. *Machine Learning*, 20:273 – 297, 1995.
- [15] M. Defant and M. Junge. Characterization of weak type by the entropy distribution of r -nuclear operators. *Studia Mathematica*, 107(1):1–14, 1993.
- [16] F. Girosi, M. Jones, and T. Poggio. Priors, stabilizers and basis functions: From regularization to radial, tensor and additive splines. A.I. Memo No. 1430, MIT, 1993.
- [17] Y. Gordon, H. König, and C. Schütt. Geometric and probabilistic estimates for entropy and approximation numbers of operators. *Journal of Approximation Theory*, 49:219–239, 1987.
- [18] I.S. Gradshteyn and I.M. Ryzhik. *Table of integrals, series, and products*. Academic Press, New York, 1981.
- [19] E. Grosswald. *Representation of Integers as Sums of Squares*. Springer, N. Y., 1985.
- [20] Y. Guo, P.L. Bartlett, J. Shawe-Taylor, and R.C. Williamson. Covering numbers for support vector machines. Submitted to COLT99, February 1999.
- [21] L. Gurvits. A note on a scale-sensitive dimension of linear bounded functionals in Banach spaces. In M. Li and A. Maruoka, editors, *Algorithmic Learning Theory ALT-97*, LNAI-1316, pages 352–363, Berlin, 1997. Springer.

- [22] D. Jagerman. ε -entropy and approximation of bandlimited functions. *SIAM Journal on Applied Mathematics*, 17(2):362–377, 1969.
- [23] M. Junge and M. Defant. Some estimates of entropy numbers. *Israel Journal of Mathematics*, 84:417–433, 1993.
- [24] M. J. Kearns, R. E. Schapire, and L. M. Sellie. Toward efficient agnostic learning. *Machine Learning*, 17(2):115–141, 1994.
- [25] V.I. Kolchinskiĭ. Operators of type p and metric entropy. *Teoriya Veroyatnostei i Matematicheskaya Statistika*, 38:69–76, 135, 1988. (In Russian. MR 89j:60007).
- [26] V.I. Kolchinskiĭ. Entropic order of operators in Banach spaces and the central limit theorem. *Theory of Probability and its Applications*, 36(2):303–315, 1991.
- [27] A.N. Kolmogorov and V.M. Tihomirov. ε -entropy and ε -capacity of sets in functional spaces. *American Mathematical Society Translations, Series 2*, 17:277–364, 1961.
- [28] H. König. *Eigenvalue Distribution of Compact Operators*. Birkhäuser, Basel, 1986.
- [29] T. Koski, L.-E. Persson, and J. Peetre. ε -entropy ε -rate, and interpolation spaces revisited with an application to linear communication channels. *Journal of Mathematical Analysis and Applications*, 186:265–276, 1994.
- [30] S.R. Kulkarni, G. Lugosi, and S.S. Venkatesh. Learning pattern classification — a survey. *IEEE Transactions on Information Theory*, 44(6):2178–2206, 1998.
- [31] W.S. Lee, P.L. Bartlett, and R.C. Williamson. The importance of convexity in learning with squared loss. *IEEE Transactions on Information Theory*, 44(5):1974–1980, 1998.
- [32] C. Müller. *Analysis of Spherical Symmetries in Euclidean Spaces*, volume 129 of *Applied Mathematical Sciences*. Springer, New York, 1997.
- [33] N.J. Nilsson. *Learning machines: Foundations of Trainable Pattern Classifying Systems*. McGraw–Hill, 1965.
- [34] A. Pajor. *Sous-espaces ℓ_n^1 des espaces de Banach*. Hermann, Paris, 1985.
- [35] A. Pietsch. *Operator ideals*. North-Holland, Amsterdam, 1980.
- [36] L.S. Pontriagin and L.G. Schnirelmann. Sur une propriété métrique de la dimension. *Annals of Mathematics*, 33:156–162, 1932.
- [37] R.T. Prosser. The ε -Entropy and ε -Capacity of Certain Time-Varying Channels. *Journal of Mathematical Analysis and Applications*, 16:553–573, 1966.

- [38] R.T. Prosser and W.L. Root. The ε -entropy and ε -capacity of certain time-invariant channels. *Journal of Mathematical Analysis and its Applications*, 21:233–241, 1968.
- [39] S. Saitoh. *Theory of Reproducing Kernels and its Applications*. Longman Scientific & Technical, Harlow, England, 1988.
- [40] B. Schölkopf, C. Burges, and V. Vapnik. Extracting support data for a given task. In U. M. Fayyad and R. Uthurusamy, editors, *Proceedings, First International Conference on Knowledge Discovery & Data Mining*, Menlo Park, 1995. AAAI Press.
- [41] B. Schölkopf, J. Shawe-Taylor, A.J. Smola, and R.C. Williamson. Generalization bounds via eigenvalues of the gram matrix. Submitted to COLT99, February 1999.
- [42] B. Schölkopf, A. Smola, and K.-R. Müller. Nonlinear component analysis as a kernel eigenvalue problem. *Neural Computation*, 10:1299 – 1319, 1998.
- [43] C.E. Shannon. A mathematical theory of communication. *Bell System Technical Journal*, 27:379–423, 623–656, 1948.
- [44] J. Shawe-Taylor, P.L. Bartlett, R.C. Williamson, and M. Anthony. A framework for structural risk minimisation. In *Proceedings of the 9th Annual Conference on Computational Learning Theory*, pages 68–76. ACM Press, 1996.
- [45] J. Shawe-Taylor, P.L. Bartlett, R.C. Williamson, and M. Anthony. Structural risk minimization over data-dependent hierarchies. *IEEE Transactions on Information Theory*, 44(5):1926–1940, 1998.
- [46] J. Shawe-Taylor and R.C. Williamson. Generalization performance of classifiers in terms of observed covering numbers. In *Proceedings of EUROCOLT'99*, 1999.
- [47] A. J. Smola and B. Schölkopf. On a kernel-based method for pattern recognition, regression, approximation and operator inversion. *Algorithmica*, 22:211–231, 1998.
- [48] A.J. Smola, A. Elisseeff, B. Schölkopf, and R.C. Williamson. Entropy numbers for convex combinations and multilayer networks. Submitted to COLT99, February 1999.
- [49] A.J. Smola, B. Schölkopf, and K.-R. Müller. The connection between regularization operators and support vector kernels. *Neural Networks*, 11:637–649, 1998.
- [50] A.J. Smola, R.C. Williamson, S. Mika, and B. Schölkopf. Regularized principal manifolds. In *Proceedings of EUROCOLT'99*, 1999.
- [51] I. H. Sneddon. *The Use of Integral Transforms*. McGraw-Hill, New York, 1972.

- [52] M. Talagrand. The Glivenko–Cantelli problem, ten years later. *Journal of Theoretical Probability*, 9(2):371–384, 1996.
- [53] H. Triebel. Interpolationseigenschaften von Entropie- und Durchmesseridealen kompakter Operatoren. *Studia Mathematica*, 34:89–107, 1970.
- [54] V. Vapnik. *The Nature of Statistical Learning Theory*. Springer, N.Y., 1995.
- [55] V. Vapnik and A. Chervonenkis. *Theory of Pattern Recognition [in Russian]*. Nauka, Moscow, 1974. (German Translation: W. Vapnik & A. Tscherwonkis, *Theorie der Zeichenerkennung*, Akademie-Verlag, Berlin, 1979).
- [56] V.N. Vapnik. *Estimation of Dependences from Empirical Data*. Springer, New York, 1982.
- [57] V.N. Vapnik and A.Ya. Chervonenkis. Necessary and sufficient conditions for the uniform convergence of means to their expectations. *Theory of Probability and its Applications*, 26(3):532–553, 1981.
- [58] G.N. Watson. *A Treatise on the Theory of Bessel Functions*. Cambridge University Press, Cambridge, UK, 2 edition, 1958.
- [59] H. Widom. Asymptotic behaviour of eigenvalues of certain integral operators. *Archive for Rational Mechanics and Analysis*, 17:215–229, 1964.
- [60] R.C. Williamson, B. Schölkopf, and A.J. Smola. A Maximum Margin Miscellany. Typescript, March 1999.