

to appear in Journal of Complexity

# Learning with generalization capability by kernel methods of bounded complexity

Věra Kůrková<sup>1,2</sup>

*Institute of Computer Science, Academy of Sciences of the Czech Republic  
Pod Vodárenskou věží 2, 182 07, Prague 8, Czech Republic*

Marcello Sanguineti<sup>1,3,\*</sup>

*Department of Communications, Computer, and System Sciences (DIST)  
University of Genoa, Via Opera Pia 13, 16145 Genova, Italy*

---

## Abstract

Learning from data with generalization capability is studied in the framework of minimization of regularized empirical error functionals over nested families of hypothesis sets with increasing model complexity. For Tikhonov's regularization with kernel stabilizers, minimization over restricted hypothesis sets containing for a fixed integer  $n$  only linear combinations of all  $n$ -tuples of kernel functions is investigated. Upper bounds are derived on the rate of convergence of suboptimal solutions from such sets to the optimal solution achievable without restrictions on model complexity. The bounds are of the form  $1/\sqrt{n}$  multiplied by a term that depends on the size of the sample of empirical data, the vector of output data, the Gram matrix of the kernel with respect to the input data, and the regularization parameter.

*Key words:* supervised learning, generalization, model complexity, kernel methods, minimization of regularized empirical errors, upper bounds on rates of approximate optimization

---

## 1 Introduction

A key property of systems performing intelligent computing, such as feature extraction, pattern recognition, semantic web realization, and classification, is learning ability. The goal of supervised learning is to adjust the parameters of a computational model so that it approximates to a desired accuracy a functional relationship between inputs and outputs by learning from a set of examples, i.e., a sample  $\mathbf{z} = \{(x_i, y_i) \in \Omega \times \mathfrak{R}, i = 1, \dots, m\}$  of  $m$  input/output pairs of *empirical data*. It is desirable that a model trained on a sample of empirical data also has a *generalization capability*, i.e., it is able to satisfactorily process new data, which were not used for learning. To endow a model with this capability, one needs some global knowledge of the desired input/output functional relationship, such as smoothness or lack of high-frequency oscillations.

In statistical learning theory [9,45], learning from empirical data is modelled as minimization of a functional, called *empirical error*. For a *sample*  $\mathbf{z}$  of data and a *loss function*  $V : \mathfrak{R}^2 \rightarrow [0, +\infty)$ , the empirical error  $\mathcal{E}_{\mathbf{z},V}$  is defined as  $\mathcal{E}_{\mathbf{z},V}(f) = \frac{1}{m} \sum_{i=1}^m V(f(x_i), y_i)$ , where  $f$  belongs to a function space, called *hypothesis space*, over which such a minimization is performed.

Mathematical modeling of generalization requires some *prior information* on the behavior of potential solutions. Such information is already expressed by the choice of a hypothesis space, over which the empirical error is minimized. It can be further specified by restricting minimization of the empirical error to a subset of the hypothesis space (containing only functions with some desired behavior). Alternatively, one can add to the empirical error a term penalizing undesired properties, or combine these two approaches. The first method is an application to learning of Ivanov's regularization, the second one of Tikhonov's, and the third one of Miller's [6, pp. 68-78].

*Tikhonov's regularization* [43,44], which was introduced into learning theory by

---

\* Corresponding author.

*Email addresses:* vera@cs.cas.cz (Věra Kůrková), marcello@dist.unige.it (Marcello Sanguineti).

<sup>1</sup> Collaboration between V. K. and M. S. was supported by the 2004-2006 Scientific Agreement among University of Genoa, National Research Council of Italy, and Czech Academy of Sciences, project "Learning from Data by Neural Networks and Kernel Methods: an Approach Based on Approximate Optimization".

<sup>2</sup> Partially supported by project 1ET100300419 "Intelligent Models, Algorithms, Methods, and Tools for Semantic Web Realization" of the program "Information Society" of the National Research Program of the Czech Republic.

<sup>3</sup> Partially supported by a PRIN grant from the Italian Ministry of University and Research.

Poggio and Girosi [20,35,36], leads to minimization over the whole hypothesis space of the *regularized empirical error functional*, defined as the sum of two functionals  $\mathcal{E}_{\mathbf{z},V} + \gamma\Psi$ . The first one, the empirical error  $\mathcal{E}_{\mathbf{z},V}$ , enforces closeness to the sample  $\mathbf{z}$  of empirical data, whereas  $\Psi$ , called *stabilizer*, expresses requirements on the global behavior of the desired input/output functional relationship. The *regularization parameter*  $\gamma$  controls the trade-off between fitting to empirical data and penalizing undesired behavior.

A large class of hypothesis spaces can be studied in the framework of the theory of Hilbert spaces of a special type, called *reproducing kernel Hilbert spaces* (RKHSs). Norms on such spaces often play the role of measures of various types of oscillations of input/output mappings. RKHSs were formally defined by Aronszajn [2], but their theory employs work by Schönberg [41], as well as many classical results on kernels and positive definite functions. RKHSs were introduced into applications closely related to learning by Parzen [33] and Wahba [47], and into learning theory by Cortes and Vapnik [8] and Girosi [19].

The *Representer Theorem* [10, p. 42], [18,20,26,35,37,39] states that for Tikhonov’s regularization with a stabilizer defined as a strictly increasing function of the norm on an RKHS, the problem of minimization of the regularized empirical error over the whole space has a unique solution of the form of a linear combination of the  $m$ -tuple of the kernel functions, which are parameterized by the input data vector  $\mathbf{x} = (x_1, \dots, x_m)$ . In particular, for a stabilizer equal to the square of the norm on an RKHS, the vector  $\mathbf{c}$  of the coefficients of the linear combination is obtained as the solution of the well-posed linear system of equations  $(\gamma m\mathcal{I} + \mathcal{K}[\mathbf{x}])\mathbf{c} = \mathbf{y}$ , where  $\mathcal{I}$  is the  $m \times m$  identity matrix,  $\mathcal{K}[\mathbf{x}]$  is the Gram matrix of the kernel  $K$  with respect to  $\mathbf{x}$ , and  $\mathbf{y} = (y_1, \dots, y_m)$  is the output data vector.

A paradigmatic example of a kernel is the Gaussian kernel, for which the solution given by the Representer Theorem has the form of an input/output function of a Gaussian radial-basis-function network with  $m$  units centered at the input data  $x_1, \dots, x_m$  [18]. The coefficients of the linear combination play the role of output weights of such a network. On the basis of this interpretation of the Representer Theorem, in [20, p. 219] it was argued that “the regularization principles lead to approximation schemes that are equivalent to networks with one layer of hidden units.”

The Representer Theorem was used to design a learning algorithm (see, e.g., [10, p. 42] and [37, pp. 538-539]) that requires one to solve the linear system of equations  $(\gamma m\mathcal{I} + \mathcal{K}[\mathbf{x}])\mathbf{c} = \mathbf{y}$ . An advantage of this algorithm is that it gives the best possible solution of the task of fitting a function to a given sample of empirical data and satisfying a global property describable in terms of a condition on smoothness that can be modelled in terms of a kernel.

However, practical applications of this algorithm are limited by the rate of convergence of iterative methods solving the system of equations and by the size of the condition number of the matrix  $\gamma m\mathcal{I} + \mathcal{K}[\mathbf{x}]$ . For some methods, the computational requirements for solving such a system grow polynomially with the size  $m$  of the sample (e.g., for the Gaussian elimination and  $m$  large enough, they grow as  $m^3/3$  [32, p. 175]). For some data and kernels, keeping the condition number of  $\gamma m\mathcal{I} + \mathcal{K}[\mathbf{x}]$  small requires a large value of the regularization parameter  $\gamma$ , which may cause poor fit to the empirical data.

The learning algorithm based on the Representer Theorem uses a computational model of complexity determined by the size  $m$  of the sample of data, and does not allow any flexibility in choosing the inner parameters of the computational units (as they are set equal to the input data).

In this paper, we investigate suboptimal solutions of the problems of minimization of regularized empirical error functionals over hypothesis sets corresponding to kernel models with limited complexity and flexible choice of parameters. We derive upper bounds on the rates of convergence of sequences of suboptimal solutions achievable by minimization over hypothesis sets formed by linear combinations of at most  $n$  kernel functions (either with arbitrary parameters or with parameters drawn from the data set) to the optimal solution given by the Representer Theorem. The upper bounds are of the form  $1/\sqrt{n}$  multiplied by a term that depends on the size  $m$  of the sample, the  $l_2$ -norm of the vector  $\mathbf{y} = (y_1, \dots, y_m)$  of output data, the minimal and the maximal eigenvalues of the Gram matrix  $\mathcal{K}[\mathbf{x}]$  of the kernel with respect to the input data, and the regularization parameter  $\gamma$ .

We state conditions on the sample, the kernel and the regularization parameter, under which the term multiplying  $1/\sqrt{n}$  is “small” and so suboptimal solutions converge “quickly” to the optimal one. Under such conditions, kernel methods with bounded model complexity provide good approximations to the best possible solution of the learning task. As our estimates are not merely asymptotic, they can be applied to any bound on model complexity. For the Gaussian kernel we derive an upper bound of the form  $\frac{3(1+\gamma)y_{\max}^2}{\gamma^2\sqrt{n}}$ , where  $y_{\max}$  is the maximum of the absolute values of output data.

The paper is organized as follows. Section 2 introduces concepts concerning minimization of functionals and Tikhonov’s regularization applied to learning from data with RKHSs as hypothesis spaces. Section 3 states the Representer Theorem and explores the condition numbers of the matrices used in algorithms based on this theorem. Section 4 develops tools for investigating approximate optimization over hypothesis sets with bounded model complexity. Section 5 describes continuity and convexity properties of regularized empirical error functionals with various types of loss functions and estimates rates of convergence of sequences of suboptimal solutions to the problem of learn-

ing by kernel methods with increasing model complexity. Section 6 illustrates the estimates on RKHSs defined by convolution kernels. Section 7 is a brief discussion. An Appendix describes properties of RKSHs and illustrates them by examples of kernels and types of oscillations measured by squares of norms defined by such kernels.

## 2 Tikhonov's regularization of the learning problem in reproducing kernel Hilbert spaces

By a normed linear space  $(X, \|\cdot\|)$  we mean a real normed linear space.  $\mathfrak{R}$  denotes the set of real numbers.

Let  $M$  be a subset of  $X$  and  $\Phi : X \rightarrow \mathfrak{R}$  be a functional. Using standard notation [15], we denote by

$$(M, \Phi)$$

the problem of minimizing  $\Phi$  over  $M$ ;  $M$  is called *hypothesis set*.

By  $\operatorname{argmin}(M, \Phi) = \{g \in M : \Phi(g) = \inf_{g \in M} \Phi(g)\}$  is denoted the set of *minimum points* of the problem  $(M, \Phi)$  and for any  $\varepsilon > 0$ ,  $\operatorname{argmin}_\varepsilon(M, \Phi) = \{g \in M : \Phi(g) < \inf_{g \in M} \Phi(g) + \varepsilon\}$  is the set of  $\varepsilon$ -near *minimum points* of  $(M, \Phi)$ . A minimum point of  $(M, \Phi)$  is called a *solution* of the problem  $(M, \Phi)$ . A sequence  $\{g_n\}$  of elements of  $M$  is called  $\Phi$ -*minimizing over  $M$*  if  $\lim_{n \rightarrow \infty} \Phi(g_n) = \inf_{g \in M} \Phi(g)$ .

Let  $\Omega$  be a set and  $\mathbf{z} = \{(x_i, y_i) \in \Omega \times \mathfrak{R}, i = 1, \dots, m\}$  an  $m$ -tuple of input/output pairs of data, called a *sample*. A standard approach to learning from empirical data [9,45] is based on minimization of the *empirical error* functional (also called the *empirical risk* functional), defined for any  $f$  in the hypothesis set as

$$\mathcal{E}_V(f) = \mathcal{E}_{\mathbf{z}, V}(f) = \frac{1}{m} \sum_{i=1}^m V(f(x_i), y_i),$$

where  $V : \mathfrak{R}^2 \rightarrow [0, \infty)$  satisfying for all  $y \in \mathfrak{R}$ ,  $V(y, y) = 0$  is called a *loss function*. When the sample  $\mathbf{z}$  is clear from the context, we write merely  $\mathcal{E}_V$  instead of  $\mathcal{E}_{\mathbf{z}, V}$ .

The most common loss function is the *square loss*, defined as

$$V(f(x), y) = (f(x) - y)^2.$$

In this paper, we mostly focus on the empirical error defined using the square loss, for which we merely write  $\mathcal{E}$ . So we let

$$\mathcal{E}(f) = \frac{1}{m} \sum_{i=1}^m (f(x_i) - y_i)^2.$$

Other common loss functions are the *absolute value loss*  $V(f(x), y) = |f(x) - y|$  and *Vapnik's  $\varepsilon$ -insensitive loss*  $V(f(x), y) = \max(|f(x) - y| - \varepsilon, 0)$ .

*Tikhonov's regularization* replaces the problem

$$(M, \mathcal{E}_V)$$

with the problem

$$(M, \mathcal{E}_V + \gamma\Psi),$$

where  $\Psi$  is a functional called *stabilizer* and  $\gamma > 0$  is a *regularization parameter* [43,44].

An important class of stabilizers are squares of norms on reproducing kernel Hilbert spaces (RKHSs). Such stabilizers often enable one to penalize high oscillations of various types. For a set  $\Omega$  and a symmetric positive semidefinite function  $K : \Omega \times \Omega \rightarrow \mathfrak{R}$ , called *kernel*, we denote by  $(\mathcal{H}_K(\Omega), \|\cdot\|_K)$  the RKHS defined by  $K$  (see Appendix). The squared norm  $\|\cdot\|_K^2$  is used as a stabilizer instead of  $\|\cdot\|_K$  for technical reasons, as the square of the norm on any Hilbert space is a uniformly convex functional (see Proposition 4.1 (iii)); this implies uniqueness of the solution of the regularized problem (see, e.g., [14, p. 10], [10, pp. 27, 42]) and convergence of minimizing sequences to this solution [31].

Using  $\|\cdot\|_K^2$  as a stabilizer, the regularized empirical error functional with a loss function  $V$  and a regularization parameter  $\gamma$  has the form

$$\mathcal{E}_{V,\gamma,K}(f) = \frac{1}{m} \sum_{i=1}^m V(f(x_i), y_i) + \gamma \|f\|_K^2.$$

As in the case of the empirical error, when the square loss is employed in the regularized empirical error we use the simplified notation

$$\mathcal{E}_{\gamma,K}(f) = \frac{1}{m} \sum_{i=1}^m (f(x_i) - y_i)^2 + \gamma \|f\|_K^2.$$

Thus we denote by

$$(\mathcal{H}_K(\Omega), \mathcal{E}_{\gamma, K}).$$

the problem of minimizing over the RKHS  $\mathcal{H}_K(\Omega)$  the regularized empirical error with the square loss and the stabilizer  $\|\cdot\|_K^2$ .

### 3 The Representer Theorem

Existence, uniqueness and an explicit formula describing the solution of the problem  $(\mathcal{H}_K(\Omega), \mathcal{E}_{\gamma, K})$  of minimizing over the whole RKHS the regularized empirical error with the square loss and the stabilizer  $\|\cdot\|_K^2$  are given by the *Representer Theorem*. For a kernel  $K$ , a positive integer  $m$ , and a vector  $\mathbf{x} = (x_1, \dots, x_m) \in \Omega^m$  of input data, we denote by  $\mathcal{K}[\mathbf{x}]$  the  $m \times m$  matrix defined as

$$\mathcal{K}[\mathbf{x}]_{ij} = K(x_i, x_j),$$

which is called the *Gram matrix of the kernel  $K$  with respect to the vector  $\mathbf{x}$* . We denote by  $\mathcal{I}$  the  $m \times m$  identity matrix.

**Theorem 3.1 (Representer Theorem)** *Let  $\Omega$  be a nonempty set,  $K : \Omega \times \Omega \rightarrow \mathfrak{R}$  a kernel,  $m$  a positive integer,  $\mathbf{x} = (x_1, \dots, x_m) \in \Omega^m$ ,  $\mathbf{y} = (y_1, \dots, y_m) \in \mathfrak{R}^m$ , and  $\gamma > 0$ . Then the problem  $(\mathcal{H}_K(\Omega), \mathcal{E}_{\gamma, K})$  has the unique solution*

$$g^o = \sum_{i=1}^m c_i K_{x_i}, \tag{1}$$

where  $\mathbf{c} = (c_1, \dots, c_m)$  is the unique solution of the well-posed linear system

$$(\gamma m \mathcal{I} + \mathcal{K}[\mathbf{x}])\mathbf{c} = \mathbf{y}. \tag{2}$$

The Representer Theorem was originally proven in [23]. An elegant proof using functional derivatives was given in [37, pp. 538-539], while for Mercer kernels a more sophisticated argument based on the Mercer Theorem was provided in [10, p. 42]. In [26] it was derived from the theory of inverse problems. Inspection of proofs shows that for any differentiable loss function  $V$ , the solution is of the form  $g^o = \sum_{i=1}^m c_i K_{x_i}$ . However, when  $V$  is not a polynomial of degree 2, the equation to be solved to compute the coefficients  $c_1, \dots, c_m$  is nonlinear [19, p. 1473]. A weaker form of the Representer Theorem, without a formula for

computing the coefficients  $c_1, \dots, c_m$ , even holds for an arbitrary loss function  $V$  and a stabilizer of the form  $\psi(\|\cdot\|_K)$ , where  $\psi : [0, +\infty) \rightarrow \mathfrak{R}$  is a strictly increasing function [39].

The Representer Theorem was exploited to design algorithms for learning from data (see, e.g., [10, p. 42] and [37, pp. 538-539]). However, its applications are limited by the rates of convergence of iterative methods solving the linear system of equations (2) and by the size of the condition number of the matrix  $\gamma m \mathcal{I} + \mathcal{K}[\mathbf{x}]$ .

We recall that the *condition number* of a nonsingular  $m \times m$  matrix  $\mathcal{A}$  with respect to a norm  $\|\cdot\|$  on  $\mathfrak{R}^m$  is defined as

$$\text{cond}(\mathcal{A}) = \|\mathcal{A}\| \|\mathcal{A}^{-1}\|,$$

where  $\|\mathcal{A}\|$  denotes the norm of  $\mathcal{A}$  as a linear operator on  $(\mathfrak{R}^m, \|\cdot\|)$ . We denote by  $\lambda_{\max}(\mathcal{A})$  and  $\lambda_{\min}(\mathcal{A})$  the maximal and minimal eigenvalues of the matrix  $\mathcal{A}$ , respectively.

It is easy to check that for any norm  $\|\cdot\|$  on  $\mathfrak{R}^m$  and any  $m \times m$  nonsingular matrix  $\mathcal{A}$ ,  $\text{cond}(\mathcal{A}) \geq \frac{|\lambda_{\max}(\mathcal{A})|}{|\lambda_{\min}(\mathcal{A})|}$  and for any symmetric nonsingular  $m \times m$  matrix  $\mathcal{A}$ ,  $\text{cond}_2(\mathcal{A}) = \frac{|\lambda_{\max}(\mathcal{A})|}{|\lambda_{\min}(\mathcal{A})|}$ , where  $\text{cond}_2(\mathcal{A})$  denotes the condition number of  $\mathcal{A}$  with respect to the  $l_2$ -norm on  $\mathfrak{R}^m$ .

To simplify the notation, we write  $\lambda_{\max}$  instead of  $\lambda_{\max}(\mathcal{K}[\mathbf{x}])$  and similarly for  $\lambda_{\min}$ . As  $\mathcal{K}[\mathbf{x}]$  is positive semidefinite, all its eigenvalues are nonnegative [32, p. 7]. As  $\lambda$  is an eigenvalue of  $\mathcal{K}[\mathbf{x}]$  if and only if  $\gamma m + \lambda$  is an eigenvalue of  $\gamma m \mathcal{I} + \mathcal{K}[\mathbf{x}]$ , we have

$$\text{cond}_2(\gamma m \mathcal{I} + \mathcal{K}[\mathbf{x}]) = \frac{\gamma m + \lambda_{\max}}{\gamma m + \lambda_{\min}} \leq \frac{\lambda_{\max}}{\lambda_{\min}} = \text{cond}_2(\mathcal{K}[\mathbf{x}]) \quad (3)$$

and

$$\text{cond}_2(\gamma m \mathcal{I} + \mathcal{K}[\mathbf{x}]) \leq 1 + \frac{\lambda_{\max}}{\gamma m}. \quad (4)$$

Equation (3) shows that when  $\text{cond}_2(\mathcal{K}[\mathbf{x}])$  is sufficiently small, good conditioning of  $\gamma m \mathcal{I} + \mathcal{K}[\mathbf{x}]$  is guaranteed for any value of  $\gamma$ . However, for large values of  $m$  the matrix  $\mathcal{K}[\mathbf{x}]$  might be ill-conditioned. For example, when the data are uniformly distributed over an interval, then the probability that  $\mathcal{K}[\mathbf{x}]$  is ill-conditioned increases with  $m$  (see [12, Theorem 2.2] and [13, Theorem 5.1]). On the other hand, equation (4) shows that  $\lim_{\gamma \rightarrow \infty} \text{cond}_2(\gamma m \mathcal{I} + \mathcal{K}[\mathbf{x}]) = 1$  and thus the regularization parameter  $\gamma$  can always be chosen such that  $\text{cond}_2(\gamma m \mathcal{I} + \mathcal{K}[\mathbf{x}])$  is close to 1. But good conditioning of  $\gamma m \mathcal{I} + \mathcal{K}[\mathbf{x}]$  is



not the only requirement for  $\gamma$ , as its value must also allow a good fit to the empirical data and thus it cannot be too large. Existence of a value of  $\gamma$  guaranteeing a good fit to data as well as good conditioning depends on the rate of convergence of the condition number of  $\gamma m\mathcal{I} + \mathcal{K}[\mathbf{x}]$  to 1. The smaller  $\frac{\lambda_{\max}}{m}$ , the faster such convergence. The problem of choosing  $\gamma$  in order to minimize the expected error was investigated in [11].

When a value of  $\gamma$  guaranteeing both a small condition number and a good fit to the empirical data cannot be found, algorithms for learning from data that differ from the one based on the Representer Theorem have to be applied. A variety of learning algorithms have been developed in the field of neurocomputing. Typically, such algorithms operate on networks of lower model complexity than the algorithm based on the Representer Theorem. The number of computational units in such networks is either set in advance or adjusted during learning, but, typically, it is much smaller than the size  $m$  of the sample used as a training set. Moreover, the values of the computational units' parameters (which are called *centroids* in the case of RBF networks) are not set equal to the input vectors from the data sample, but are searched for during learning.

#### 4 Minimization of functionals over hypothesis sets with bounded model complexity

In this section, we derive tools for estimating rates of convergence of suboptimal solutions over computational models with  $n$  units (the case of interest is  $n < m$ ) to the optimal solution given by the Representer Theorem. Such suboptimal solutions can be studied in terms of optimization over nested families of subsets of RKHSs formed by linear combinations of all  $n$ -tuples of kernel functions chosen from the sets  $\{K_x : x \in \Omega\}$  or  $\{K_{x_1}, \dots, K_{x_m}\}$ .

For a subset  $G$  of a linear space, let  $span_n G = \{\sum_{i=1}^n w_i g_i : w_i \in \mathfrak{R}, g_i \in G\}$  denote the set of linear combinations of all  $n$ -tuples of elements of  $G$ . The optimal solution to the problem  $(\mathcal{H}_K(\Omega), \mathcal{E}_{\gamma, K})$  described by the Representer Theorem is an element of  $span_m G_{K, \mathbf{x}} \subseteq span_m G_K$ , where  $G_{K, \mathbf{x}} = \{K_{x_1}, \dots, K_{x_m}\}$  and  $G_K = \{K_x : x \in \Omega\}$ . The set  $span_m G_K$  can be interpreted as the set of all input/output functions of a computational model with one hidden layer of  $m$  computational units computing functions from  $G_K$ . In particular, for the Gaussian kernel the solution has the form of an input/output function of a Gaussian radial-basis-function (RBF) network with  $m$  computational units [20].

To compare the optimal solution given by the Representer Theorem with suboptimal ones that can be obtained by minimization of  $\mathcal{E}_{\gamma, K}$  over restricted hypothesis sets (containing only linear combinations of all  $n$ -tuples of elements

of the set  $G_K$  or  $G_{K,\mathbf{x}}$ ), we shall employ a version of the Maurey-Jones-Barron Theorem [3,22,34], reformulated in [24] in terms of a norm called  $G$ -variation.

We recall that the *Minkowski functional* of a subset  $M$  of a linear space  $X$ , denoted by  $p_M$ , is defined for every  $f \in X$  as  $p_M(f) = \inf\{\lambda \in \mathfrak{R}_+ : f/\lambda \in M\}$ . If  $M$  is a subset of a normed linear space  $(X, \|\cdot\|)$ , we denote by  $\text{cl } M$  its *closure* with respect to the topology generated by  $\|\cdot\|$ , i.e.,  $\text{cl } M = \{f \in X : (\forall \varepsilon > 0) (\exists g \in M) \|f - g\| < \varepsilon)\}$ .

$G$ -variation norm, denoted by  $\|\cdot\|_G$ , is defined for a subset  $G$  of a normed linear space  $(X, \|\cdot\|)$  as the Minkowski functional of the closure of the convex hull of the set  $G \cup -G$ . So for every  $f \in X$  we have

$$\|f\|_G = \inf\{c > 0 : f/c \in \text{cl conv}(G \cup -G)\}.$$

For properties of  $G$ -variation, see [24,25,27,28,30].

Maurey-Jones-Barron's theorem stated in terms of  $G$ -variation [24,25] gives for a Hilbert space  $(X, \|\cdot\|)$ , its bounded subset  $G$  with  $s_G = \sup_{g \in G} \|g\|$ , and every  $f \in X$ , the following upper bound on the rate of approximation of  $f$  by  $\text{span}_n G$ :  $\|f - \text{span}_n G\| \leq \sqrt{\frac{(s_G \|f\|_G)^2 - \|f\|^2}{n}}$ .

Taking advantage of this upper bound, we shall estimate rates of convergence of suboptimal solutions over  $\text{span}_n G$  to the optimal solution of the problem  $(X, \Phi)$  of minimization of a continuous functional  $\Phi$  over a normed linear space  $X$ .

A functional  $\Phi : X \rightarrow \mathfrak{R}$  is *continuous* at  $f \in X$  if for any  $\varepsilon > 0$ , there exists  $\eta > 0$  such that  $\|f - g\| < \eta$  implies  $|\Phi(f) - \Phi(g)| < \varepsilon$ . A *modulus of continuity* of  $\Phi$  at  $f$  is a function  $\omega : [0, +\infty) \rightarrow [0, +\infty)$  defined as  $\omega(a) = \sup\{|\Phi(f) - \Phi(g)| : \|f - g\| \leq a\}$ .

$\Phi$  is *convex* on a convex set  $M \subseteq X$  if for all  $h, g \in M$  and all  $\lambda \in [0, 1]$ , we have  $\Phi(\lambda h + (1 - \lambda)g) \leq \lambda\Phi(h) + (1 - \lambda)\Phi(g)$  and it is *uniformly convex* if there exists a non-negative function  $\delta : \mathfrak{R}_+ \rightarrow \mathfrak{R}_+$  such that  $\delta(0) = 0$ ,  $\delta(t) > 0$  for all  $t > 0$ , and for all  $h, g \in M$  and all  $\lambda \in [0, 1]$ ,  $\Phi(\lambda h + (1 - \lambda)g) \leq \lambda\Phi(h) + (1 - \lambda)\Phi(g) - \lambda(1 - \lambda)\delta(\|h - g\|)$ . Any such function  $\delta$  is called a *modulus of convexity* of  $\Phi$  [31]<sup>4</sup>.

Next proposition states some elementary properties of uniformly convex func-

<sup>4</sup> The terminology is not unified: some authors use the term “strictly uniformly convex” instead of “uniformly convex”, while they reserve the term “uniformly convex” for the case where  $\delta : \mathfrak{R}_+ \rightarrow \mathfrak{R}_+$  merely satisfies  $\delta(0) = 0$  and  $\delta(t_0) > 0$  for some  $t_0 > 0$  (see, e.g., [46] and [14, p. 10]).

tionals and moduli of convexity.

**Proposition 4.1** *Let  $(X, \|\cdot\|)$  be a normed linear space,  $M \subseteq X$  convex, and  $\Phi$  a uniformly convex functional on  $M$  with a modulus of convexity  $\delta$ . Then the following hold:*

- (i) *if  $\Psi$  is convex on  $M$  and  $\gamma > 0$ , then  $\Psi + \gamma\Phi$  is uniformly convex on  $M$  with a modulus of convexity  $\gamma\delta$ ;*
- (ii) *if  $g^\circ \in \operatorname{argmin}(M, \Phi)$ , then for every  $g \in M$ ,  $\delta(\|g - g^\circ\|) \leq \Phi(g) - \Phi(g^\circ)$ ;*
- (iii) *if  $(X, \|\cdot\|)$  is a Hilbert space, then the functional  $\|\cdot\|^2 : X \rightarrow \mathfrak{R}$  is uniformly convex with a modulus of convexity  $\delta(t) = t^2$ .*

**Proof.** (i) follows directly from the definitions.

(ii) By the definition of uniformly convex functional, for every  $\lambda \in [0, 1]$  we have  $\lambda(1 - \lambda)\delta(\|g - g^\circ\|) \leq \lambda\Phi(g) + (1 - \lambda)\Phi(g^\circ) - \Phi(\lambda g + (1 - \lambda)g^\circ)$ . As  $\Phi(g^\circ) \leq \Phi(\lambda g + (1 - \lambda)g^\circ)$ , we get  $\lambda(1 - \lambda)\delta(\|g - g^\circ\|) \leq \lambda\Phi(g) + (1 - \lambda)\Phi(g^\circ) - \Phi(g^\circ) = \lambda(\Phi(g) - \Phi(g^\circ))$ . Hence  $(1 - \lambda)\delta(\|g - g^\circ\|) \leq \Phi(g) - \Phi(g^\circ)$  for every  $\lambda \in [0, 1]$ . So we obtain  $\delta(\|g - g^\circ\|) \leq \Phi(g) - \Phi(g^\circ)$ .

(iii) For every  $h, g \in X$  and every  $\lambda \in [0, 1]$ , we have  $\|\lambda h + (1 - \lambda)g\|^2 \leq \lambda\|h\|^2 + (1 - \lambda)\|g\|^2 - \lambda(1 - \lambda)\|h - g\|^2$  and thus  $\delta(t) = t^2$  is a modulus of convexity of  $\|\cdot\|^2$ .  $\square$

Next theorem gives upper bounds on rates of convergence of suboptimal solutions over  $\operatorname{span}_n G$  to the optimal solution of the problem  $(X, \Phi)$  of minimization of a continuous functional  $\Phi$  over a Hilbert space  $X$ . The estimates are formulated in terms of moduli of continuity and convexity of the functional to be minimized.

**Theorem 4.2** *Let  $(X, \|\cdot\|)$  be a Hilbert space,  $G$  its bounded subset,  $s_G = \sup_{g \in G} \|g\|$ ,  $\Phi : X \rightarrow (-\infty, +\infty]$  a functional,  $g^\circ \in \operatorname{argmin}(X, \Phi)$ ,  $\Phi$  continuous at  $g^\circ$  with a modulus of continuity  $\alpha$ ,  $\{\varepsilon_n\}$  a sequence of positive real numbers,  $g_n \in \operatorname{argmin}_{\varepsilon_n}(\operatorname{span}_n G, \Phi)$ , and  $a = (s_G \|g^\circ\|_G)^2 - \|g^\circ\|^2$ . Then, for every positive integer  $n$  the following estimates hold:*

- (i)  $\inf_{g \in \operatorname{span}_n G} \Phi(g) - \Phi(g^\circ) \leq \alpha\left(\sqrt{\frac{a}{n}}\right)$ ;
- (ii) *if  $\|g^\circ\|_G < \infty$  and  $\lim_{n \rightarrow \infty} \varepsilon_n = 0$ , then  $\{g_n\}$  is a  $\Phi$ -minimizing sequence and  $\Phi(g_n) - \Phi(g^\circ) \leq \alpha\left(\sqrt{\frac{a}{n}}\right) + \varepsilon_n$ ;*
- (iii) *if  $\Phi$  is uniformly convex with a modulus of convexity  $\delta$ , then  $\delta(\|g_n - g^\circ\|) \leq \alpha\left(\sqrt{\frac{a}{n}}\right) + \varepsilon_n$ .*

**Proof.** (i) For every positive integer  $n$  and every  $\varepsilon > 0$ , choose an  $\varepsilon$ -near best approximation  $f_n^\varepsilon$  of  $g^\circ$  in  $\operatorname{span}_n G$ . So  $\|g^\circ - f_n^\varepsilon\| < \|g^\circ - \operatorname{span}_n G\| + \varepsilon$ . As  $f_n^\varepsilon \in \operatorname{span}_n G$ , we have  $\inf_{g \in \operatorname{span}_n G} \Phi(g) - \Phi(g^\circ) \leq \Phi(f_n^\varepsilon) - \Phi(g^\circ)$ . Estimating

the right-hand side of this inequality in terms of the modulus of continuity  $\alpha$  of  $\Phi$  at  $g^o$ , we obtain  $\inf_{g \in \text{span}_n G} \Phi(g) - \Phi(g^o) \leq \alpha(\|f_n^\varepsilon - g^o\|) \leq \alpha(\|g^o - \text{span}_n G\| + \varepsilon)$ . By the upper bound from Maurey-Jones-Barron's theorem reformulated in terms of  $G$ -variation we get

$$\inf_{g \in \text{span}_n G} \Phi(g) - \Phi(g^o) \leq \alpha\left(\sqrt{\frac{a}{n}} + \varepsilon\right). \quad (5)$$

Infimizing (5) over  $\varepsilon$  we obtain (i).

(ii) By the definition of  $\varepsilon_n$ -near minimum point, we have  $\Phi(g_n) - \Phi(g^o) \leq \inf_{g \in \text{span}_n G} \Phi(g) - \Phi(g^o) + \varepsilon_n$ . So, by item (i) we get

$$\Phi(g_n) - \Phi(g^o) \leq \alpha\left(\sqrt{\frac{a}{n}}\right) + \varepsilon_n. \quad (6)$$

If  $\|g^o\|_G$  is finite and  $\lim_{n \rightarrow \infty} \varepsilon_n = 0$ , then the right-hand side of (6) converges to zero and so  $\{g_n\}$  is  $\Phi$ -minimizing.

(iii) By item (i), the definition of  $\varepsilon_n$ -near minimum point, and Proposition 4.1 (iii), we have  $\delta(\|g_n - g^o\|) \leq \Phi(g_n) - \Phi(g^o) < \inf_{g \in \text{span}_n G} \Phi(g) - \Phi(g^o) + \varepsilon_n \leq \alpha\left(\sqrt{\frac{a}{n}}\right) + \varepsilon_n$ .  $\square$

Theorem 4.2 can be also obtained as a corollary of [29, Theorem 4.2], which applies to other types of regularization, too, such as Ivanov's one. However, the direct argument used here is much simpler than the proof of [29, Theorem 4.2].

## 5 Suboptimal solutions over kernel models with bounded complexity

In this section, we derive estimates of rates of convergence of suboptimal solutions of the problems  $(\text{span}_n G_K, \mathcal{E}_{\gamma, K})$  to the optimal solution  $g^o$  given by the Representer Theorem for the problem  $(\mathcal{H}_K(\Omega), \mathcal{E}_{\gamma, K})$ . In contrast to the optimal solution  $g^o$ , which is a linear combination of the representer  $K_{x_1}, \dots, K_{x_m}$  determined by the sample  $\mathbf{x} = (x_1, \dots, x_m)$  of input data, suboptimal solutions are formed by linear combinations of *arbitrary  $n$ -tuples* of elements of  $G_K = \{K_x : x \in \Omega\}$ . In applications, a proper  $n$ -tuple together with coefficients of the linear combination can be adjusted by a suitable nonlinear programming algorithm (see, e.g., [1,7,21]).

To employ Theorem 4.2 to estimate rates of approximate minimization of

regularized empirical error functionals with kernel stabilizers, we need upper bounds on the moduli of continuity and convexity of these functionals. The next proposition describes convexity and continuity properties of regularized empirical error functionals with various loss functions.

**Proposition 5.1** *Let  $\Omega$  be a nonempty set,  $K : \Omega \times \Omega$  a kernel,  $s_K = \sup_{x \in \Omega} \sqrt{K(x, x)}$ ,  $\gamma > 0$ ,  $m$  a positive integer,  $\mathbf{x} = (x_1, \dots, x_m) \in \Omega^m$ ,  $\mathbf{y} = (y_1, \dots, y_m) \in \mathbb{R}^m$ ,  $y_{\min} = \min\{|y_i| : i = 1, \dots, m\}$ , and  $V : \mathbb{R} \rightarrow \mathbb{R}$  a loss function. Then the following hold:*

- (i) *if for every  $i = 1, \dots, m$  the functions  $V(\cdot, y_i) : \mathbb{R} \rightarrow \mathbb{R}$  are convex, then  $\mathcal{E}_{V, \gamma, K}$  is uniformly convex on  $\mathcal{H}_K(\Omega)$  with a modulus of convexity  $\delta(t) = \gamma t^2$ ;*
- (ii) *if  $V$  is either the square or the absolute value loss function, then at every  $f \in \mathcal{H}_K(\Omega)$  the functional  $\mathcal{E}_{V, \gamma, K}$  is continuous with a modulus of continuity bounded from above by the quadratic function  $\beta(t) = b_2 t^2 + b_1 t$ , where for the square loss  $b_2 = s_K^2 + \gamma$  and  $b_1 = 2(\|f\|_K (s_K^2 + \gamma) + y_{\min} s_K)$ , while for the absolute value loss,  $b_2 = \gamma$  and  $b_1 = s_K + 2\gamma\|f\|_K$ ;*
- (iii) *if  $V$  is the square loss function, then there exists a unique minimum point  $g^\circ$  of the problem  $(\mathcal{H}_K(\Omega), \mathcal{E}_{V, \gamma, K})$  and for every  $f \in \mathcal{H}_K(\Omega)$*

$$\|f - g^\circ\|_K^2 \leq \frac{\mathcal{E}_{V, \gamma, K}(f) - \mathcal{E}_{V, \gamma, K}(g^\circ)}{\gamma}.$$

**Proof.** (i) It is easy to show that for such loss functions the empirical error functional  $\mathcal{E}_V = 1/m \sum_{i=1}^m V(f(x_i), y_i)$  is convex, and so the statement follows from Proposition 4.1 (i) and (iii).

(ii) For the square loss, by inequality (A.1) we obtain  $|\mathcal{E}_{V, \gamma, K}(f) - \mathcal{E}_{V, \gamma, K}(g)| = \left| \frac{1}{m} \sum_{i=1}^m ((f(x_i) - y_i)^2 - (g(x_i) - y_i)^2) + \gamma (\|f\|_K^2 - \|g\|_K^2) \right| \leq \left| \frac{1}{m} \sum_{i=1}^m (f(x_i) - g(x_i)) (f(x_i) + g(x_i) - 2y_i) \right| + \gamma \| \|f\|_K - \|g\|_K \| (\|f\|_K + \|g\|_K) \leq \sup_{x \in \Omega} |f(x) - g(x)| (\sup_{x \in \Omega} |f(x) + g(x)| + 2y_{\min}) + \gamma \|f - g\|_K (\|f\|_K + \|g\|_K).$

Let  $t > 0$  and  $f, g$  be such that  $\|f - g\|_K \leq t$ . Then  $|\mathcal{E}_{V, \gamma, K}(f) - \mathcal{E}_{V, \gamma, K}(g)| \leq t s_K (\|s_K \|f + g\|_K + 2y_{\min}) + t \gamma (\|f\|_K + \|g\|_K) \leq t s_K (2\|f\|_K s_K + t s_K + 2y_{\min}) + \gamma t (2\|f\|_K + t) \leq t^2 (s_K^2 + \gamma) + 2t(\|f\|_K s_K^2 + y_{\min} s_K + \gamma \|f\|_K)$ . Thus,  $\|f - g\|_K < t$  implies  $|\mathcal{E}_{V, \gamma, K}(f) - \mathcal{E}_{V, \gamma, K}(g)| \leq \beta(t) = b_2 t^2 + b_1 t$ , where  $b_2 = s_K^2 + \gamma$  and  $b_1 = 2(\|f\|_K (s_K^2 + \gamma) + y_{\min} s_K)$ .

Similarly, for the absolute value loss we have  $|\mathcal{E}_{V, \gamma, K}(f) - \mathcal{E}_{V, \gamma, K}(g)| = \left| \frac{1}{m} \sum_{i=1}^m (|f(x_i) - y_i| - |g(x_i) - y_i|) + \gamma (\|f\|_K - \|g\|_K) \right| \leq \sup_{x \in \Omega} |f(x) - g(x)| + \gamma \| \|f\|_K - \|g\|_K \| (\|f\|_K + \|g\|_K) \leq s_K \|f - g\|_K + \gamma \|f - g\|_K (\|f\|_K + \|g\|_K)$ . If  $\|f - g\|_K \leq t$ , then  $|\mathcal{E}_{V, \gamma, K}(f) - \mathcal{E}_{V, \gamma, K}(g)| \leq s_K t + t \gamma (\|f\|_K + \|g\|_K) \leq s_K t + t \gamma (t + 2\|f\|_K)$ . Hence  $|\mathcal{E}_{V, \gamma, K}(f) - \mathcal{E}_{V, \gamma, K}(g)| \leq \beta(t) = b_2 t^2 + b_1 t$ , where  $b_2 = \gamma$  and  $b_1 = s_K + 2\gamma\|f\|_K$ .

(iii) The existence of a unique minimum point  $g^\circ$  follows from the Representer Theorem. By Proposition 4.1 (i), (ii), and (iii), for every  $f \in \mathcal{H}_K(\Omega)$  we have  $\gamma \|f - g^\circ\|_K^2 \leq |\mathcal{E}_{V,\gamma,K}(f) - \mathcal{E}_{V,\gamma,K}(g^\circ)|$ .  $\square$

The assumptions of Proposition 5.1 (i) are satisfied by both the square loss and the absolute value loss. So these two loss functions determine uniformly convex functionals  $\mathcal{E}_{V,\gamma,K}$  with quadratic moduli of convexity. Their moduli of continuity at any  $f \in \mathcal{H}_K(\Omega)$  are bounded from above by the quadratic function  $\beta(t) = b_2 t^2 + b_1 t$ , where for both losses  $b_2$  depends on  $\gamma$  and for the square loss, also on  $s_K$ , while  $b_1$  depends on  $\gamma$ ,  $s_K$ ,  $\|f\|_K$  and for the square loss, also on  $y_{\min}$ . The larger the regularization parameter  $\gamma$ , the larger the coefficients of the quadratic function bounding the moduli of continuity. Generally, the modulus of continuity of  $\mathcal{E}_{V,\gamma,K}$  depends on the moduli of continuity of the functions  $V(\cdot, y_i)$ ,  $i = 1, \dots, m$ .

To simplify the formulas, in the following we assume that  $y_{\min} = \min\{|y_i| : i = 1, \dots, m\} = 0$ . Note that although the next theorem holds for any positive integer  $n$ , it is useful only for  $n < m$  since by the Representer Theorem, the minimum point of  $\mathcal{E}_{\gamma,K}$  over  $\text{span}_m G_K$  is equal to the minimum point over the whole space  $\mathcal{H}_K(\Omega)$ .

**Theorem 5.2** *Let  $\Omega$  be a nonempty set,  $K : \Omega \times \Omega \rightarrow \mathfrak{R}$  a kernel,  $s_K = \sup_{x \in \Omega} \sqrt{K(x, x)}$ ,  $m$  a positive integer,  $\mathbf{x} = (x_1, \dots, x_m) \in \Omega^m$ ,  $\mathbf{y} = (y_1, \dots, y_m) \in \mathfrak{R}^m$ ,  $\min\{|y_i| : i = 1, \dots, m\} = 0$ ,  $g^\circ = \sum_{i=1}^m c_i K_{x_i}$  the unique solution of  $(\mathcal{H}_K(\Omega), \mathcal{E}_{\gamma,K})$ ,  $\{\varepsilon_n\}$  a sequence of positive real numbers such that  $\lim_{n \rightarrow \infty} \varepsilon_n = 0$ , and  $\{g_n\}$  a sequence of  $\varepsilon_n$ -near minimum points of  $(\text{span}_n G_K, \mathcal{E}_K)$ . Let  $a = (s_K \|g^\circ\|_{G_K})^2 - \|g^\circ\|_K^2$ ,  $u = (s_K^2 + \gamma)a$ , and  $v = 2(s_K^2 + \gamma)\|g^\circ\|_K \sqrt{a}$ . Then, for every positive integer  $n$  the following estimates hold:*

- (i)  $\inf_{g \in \text{span}_n G_K} \mathcal{E}_{\gamma,K}(g) - \mathcal{E}_{\gamma,K}(g^\circ) \leq \frac{u}{n} + \frac{v}{\sqrt{n}}$ ;
- (ii)  $\mathcal{E}_{\gamma,K}(g_n) - \mathcal{E}_K(g^\circ) \leq \frac{u}{n} + \frac{v}{\sqrt{n}} + \varepsilon_n$ ;
- (iii)  $\|g_n - g^\circ\|_K^2 \leq \frac{1}{\gamma} \left( \frac{u}{n} + \frac{v}{\sqrt{n}} + \varepsilon_n \right)$ ;
- (iv)  $\sup_{x \in \Omega} |g_n(x) - g^\circ(x)|^2 \leq \frac{s_K^2}{\gamma} \left( \frac{u}{n} + \frac{v}{\sqrt{n}} + \varepsilon_n \right)$ .

**Proof.** (i) Combining Theorem 4.2 (i) with Proposition 5.1 (ii), we get  $\inf_{g \in \text{span}_n G_K} \mathcal{E}_{\gamma,K}(g) - \mathcal{E}_{\gamma,K}(g^\circ) \leq \beta\left(\sqrt{\frac{a}{n}}\right)$ , where  $\beta(t) = (s_K^2 + \gamma)(t^2 + 2\|g^\circ\|_K t)$ , which gives for  $\inf_{g \in \text{span}_n G_K} \mathcal{E}_{\gamma,K}(g) - \mathcal{E}_{\gamma,K}(g^\circ)$  the upper bound  $(s_K^2 + \gamma) \left( \frac{a}{n} + 2\|g^\circ\|_K \sqrt{\frac{a}{n}} \right) = \frac{u}{n} + \frac{v}{\sqrt{n}}$ .

Similarly, item (ii) follows from Theorem 4.2 (ii) and Proposition 5.1 (ii), item (iii) follows from (ii) and Proposition 5.1 (iii), and item (iv) from (iii) and inequality (A.1).  $\square$

Thus when  $u$  and  $v$  are not too large, it is possible to choose  $n$  small enough so that a computational model with  $n$  units is implementable and a suboptimal solution over such a model approximates well the optimal solution given by the Representer Theorem.

Only two terms in the above formulas defining  $u$  and  $v$  cannot be derived directly from the data sample  $\mathbf{z}$ , the kernel  $K$  and the regularization parameter  $\gamma$ : the values of the two norms of the optimal solution  $g^o$ , i.e., its  $G_K$ -variation and its norm  $\|\cdot\|_K$ . The next proposition estimates these two values in terms of the size  $m$  of the sample, the regularization parameter  $\gamma$ , the  $l_2$ -norm of the output vector  $\mathbf{y}$ , and the maximal and minimal eigenvalues,  $\lambda_{\max}$  and  $\lambda_{\min}$ , of the Gram matrix  $\mathcal{K}[\mathbf{x}]$  of the kernel  $K$  with respect to the input data vector  $\mathbf{x}$ . The  $l_1$ - and  $l_2$ -norm on  $\Re^m$  are denoted by  $\|\cdot\|_1$  and  $\|\cdot\|_2$ , respectively.

The estimates in the rest of the paper (Proposition 5.3, Theorem 5.4, and Corollaries 6.1 and 6.2) involve an upper bound on  $\|g^o\|_{G_K}$ , which is also an upper bound on  $\|g^o\|_{G_{K,\mathbf{x}}}$ . Thus, all these estimates can be applied also to approximate solutions over hypothesis sets formed by functions from  $\text{span}_n G_{K,\mathbf{x}}$ . Such solutions are obtained when  $n$  representers are chosen from the set  $G_{K,\mathbf{x}}$ , as, e.g., in [42], where approximation techniques were proposed that reduce the Gram matrix  $\mathcal{K}[\mathbf{x}]$  to a sparse matrix of lower rank.

**Proposition 5.3** *Let  $\Omega$  be a nonempty set,  $K : \Omega \times \Omega \rightarrow \Re$  a kernel,  $s_K = \sup_{x \in \Omega} \sqrt{K(x,x)}$ ,  $\gamma > 0$ ,  $m$  a positive integer,  $\mathbf{x} = (x_1, \dots, x_m) \in \Omega^m$ ,  $\mathbf{y} = (y_1, \dots, y_m) \in \Re^m$ ,  $g^o = \sum_{i=1}^m c_i K_{x_i}$  the unique solution of  $(\mathcal{H}_K(\Omega), \mathcal{E}_{\gamma,K})$ . Then the following estimates hold:*

- (i)  $\|g^o\|_{G_K} \leq \frac{\sqrt{m}\|\mathbf{y}\|_2}{\gamma m + \lambda_{\min}}$ ;
- (ii)  $\|g^o\|_K \leq \frac{\sqrt{\lambda_{\max}}\|\mathbf{y}\|_2}{\gamma m + \lambda_{\min}}$ ;
- (iii)  $s_K^2 \|g^o\|_{G_K}^2 - \|g^o\|_K^2 \leq \frac{(s_K^2 m - \lambda_{\min})\|\mathbf{y}\|_2^2}{(\gamma m + \lambda_{\min})^2}$ .

**Proof.** (i) From the Representer Theorem, the definition of  $G_K$ -variation, and the Cauchy-Schwartz inequality it follows that

$$\|g^o\|_{G_K} \leq \sum_{i=1}^m |c_i| = \|\mathbf{c}\|_1 \leq \sqrt{m} \|\mathbf{c}\|_2, \quad (7)$$

where  $\mathbf{c} = (\gamma m \mathcal{I} + \mathcal{K}[\mathbf{x}])^{-1} \mathbf{y}$ . By the definition of the norm of an operator,  $\|\mathbf{c}\|_2 \leq \|(\gamma m \mathcal{I} + \mathcal{K}[\mathbf{x}])^{-1}\|_2 \|\mathbf{y}\|_2$ . As  $(\gamma m \mathcal{I} + \mathcal{K}[\mathbf{x}])^{-1}$  is symmetric and positive definite, its  $l_2$ -norm is equal to its maximal eigenvalue, i.e.,  $\frac{1}{\gamma m + \lambda_{\min}}$ . So we have

$$\|\mathbf{c}\|_2 \leq \frac{\|\mathbf{y}\|_2}{\gamma m + \lambda_{\min}} \quad (8)$$

and thus  $\|g^o\|_{G_K} \leq \frac{\sqrt{m}\|\mathbf{y}\|_2}{\gamma m + \lambda_{\min}}$ .

(ii) By the Representer Theorem,  $\|g^o\|_K^2 = \left\langle \sum_{i=1}^m c_i K_{x_i}, \sum_{j=1}^m c_j K_{x_j} \right\rangle_K = \sum_{i,j=1}^m c_i c_j K(x_i, x_j) = \mathbf{c}^T \mathcal{K}[\mathbf{x}] \mathbf{c}$ , where  $\mathbf{c}^T$  denotes the transpose of the vector  $\mathbf{c}$ . As  $\lambda_{\min} \|\mathbf{c}\|_2^2 \leq \mathbf{c}^T \mathcal{K}[\mathbf{x}] \mathbf{c} \leq \lambda_{\max} \|\mathbf{c}\|_2^2$  [32, p. 21], we have

$$\lambda_{\min} \|\mathbf{c}\|_2^2 \leq \|g^o\|_K^2 \leq \lambda_{\max} \|\mathbf{c}\|_2^2. \quad (9)$$

Thus by (8),  $\|g^o\|_K \leq \frac{\sqrt{\lambda_{\max}}\|\mathbf{y}\|_2}{\gamma m + \lambda_{\min}}$ .

(iii) By (7), (8), and (9), we obtain

$$s_K^2 \|g^o\|_{G_K}^2 - \|g^o\|_K^2 \leq s_K^2 m \|\mathbf{c}\|_2^2 - \lambda_{\min} \|\mathbf{c}\|_2^2 \leq (s_K^2 m - \lambda_{\min}) \|\mathbf{c}\|_2^2 \leq \frac{(s_K^2 m - \lambda_{\min}) \|\mathbf{y}\|_2^2}{(\gamma m + \lambda_{\min})^2}.$$

□

As both  $\lambda_{\min}$  and  $\lambda_{\max}$  are nonnegative, we can further simplify as follows the upper bounds from Proposition 5.3:

$$(i) \quad \|g^o\|_{G_K} \leq \frac{\|\mathbf{y}\|_2}{\gamma \sqrt{m}}, \quad (10)$$

$$(ii) \quad \|g^o\|_K \leq \frac{\sqrt{\lambda_{\max}}\|\mathbf{y}\|_2}{\gamma m}, \quad (11)$$

$$(iii) \quad s_K^2 \|g^o\|_{G_K}^2 - \|g^o\|_K^2 \leq \frac{s_K^2 \|\mathbf{y}\|_2^2}{\gamma^2 m}. \quad (12)$$

Combining Proposition 5.3 with Theorem 5.2 and inequalities (10)-(12), we shall derive upper bounds on rates of convergence of approximate solutions of the problems  $(\text{span}_n G_K, \mathcal{E}_{\gamma, K})$  to the solution of the problem  $(\mathcal{H}_K(\Omega), \mathcal{E}_{\gamma, K})$  in terms of  $s_K$ ,  $m$ ,  $\gamma$ ,  $\|\mathbf{y}\|_2$ ,  $\lambda_{\min}$ , and  $\lambda_{\max}$ .

**Theorem 5.4** *Let  $\Omega$  be a nonempty set,  $K : \Omega \times \Omega \rightarrow \mathfrak{R}$  a kernel,  $s_K = \sup_{x \in \Omega} \sqrt{K(x, x)}$ ,  $\gamma > 0$ ,  $m$  a positive integer,  $\mathbf{x} = (x_1, \dots, x_m) \in \Omega^m$ ,  $\mathbf{y} = (y_1, \dots, y_m) \in \mathfrak{R}^m$ ,  $\min\{|y_i| : i = 1, \dots, m\} = 0$ ,  $g^o = \sum_{i=1}^m c_i K_{x_i}$  the unique solution of  $(\mathcal{H}_K(\Omega), \mathcal{E}_{\gamma, K})$ ,  $\{\varepsilon_n\}$  a sequence of positive real numbers, and  $\{g_n\}$  a sequence of  $\varepsilon_n$ -near minimum points of  $(\text{span}_n G_K, \mathcal{E}_{\gamma, K})$ . Let*

$$\bar{u} = \left( s_K^2 + \gamma \right) \frac{(s_K^2 m - \lambda_{\min}) \|\mathbf{y}\|_2^2}{(\gamma m + \lambda_{\min})^2} \leq \left( s_K^2 + \gamma \right) \frac{s_K^2 \|\mathbf{y}\|_2^2}{\gamma^2 m} \quad \text{and}$$



$$\bar{v} = 2 \left( s_K^2 + \gamma \right) \frac{\sqrt{\lambda_{\max}} \|\mathbf{y}\|_2}{(\gamma m + \lambda_{\min})^2} \sqrt{(s_K^2 m - \lambda_{\min})} \|\mathbf{y}\|_2 \leq 2 \left( s_K^2 + \gamma \right) \frac{\sqrt{\lambda_{\max} s_K} \|\mathbf{y}\|_2^2}{\gamma^2 m^{3/2}}.$$

Then, for every positive integer  $n$  the following estimates hold:

- (i)  $\inf_{g \in \text{span}_n G_K} \mathcal{E}_{\gamma, K}(g) - \mathcal{E}_{\gamma, K}(g^o) \leq \frac{\bar{u}}{n} + \frac{\bar{v}}{\sqrt{n}};$
- (ii)  $\mathcal{E}_{\gamma, K}(g_n) - \mathcal{E}_K(g^o) \leq \frac{\bar{u}}{n} + \frac{\bar{v}}{\sqrt{n}} + \varepsilon_n;$
- (iii)  $\|g_n - g^o\|_K^2 \leq \frac{1}{\gamma} \left( \frac{\bar{u}}{n} + \frac{\bar{v}}{\sqrt{n}} + \varepsilon_n \right);$
- (iv)  $\sup_{x \in \Omega} |g_n(x) - g^o(x)|^2 \leq \frac{s_K^2}{\gamma} \left( \frac{\bar{u}}{n} + \frac{\bar{v}}{\sqrt{n}} + \varepsilon_n \right).$

Thus, to obtain a good approximation of the solution of  $(\mathcal{H}_K(\Omega), \mathcal{E}_{\gamma, K})$  given by the Representer Theorem by a suboptimal solution computable by a model with at most  $n < m$  computational units, both  $\frac{\bar{u}}{n}$  and  $\frac{\bar{v}}{\sqrt{n}}$  have to be sufficiently small for some  $n$ , for which models with  $n$  computational units computing functions from  $G_K$  are implementable.

## 6 Estimates for convolution kernels

In this section, we illustrate the estimates given in Theorem 5.4 by examples of RSH with  $\Omega = \mathfrak{R}^d$  and convolution kernels. Let  $K(u, v) = \psi(\|u - v\|)$  be a convolution kernel, where  $\psi : \mathfrak{R} \rightarrow [0, 1]$  is monotonically decreasing and satisfies  $\psi(0) = 1$  (this includes the Gaussian kernel). The following corollary estimates rates of convergence of suboptimal solutions for input/output pairs of data  $(x_1, y_1), \dots, (x_m, y_m)$  for which the input data are sufficiently separated so that there exists  $t \in [0, 1]$  such that for all distinct  $i, j \in \{1, \dots, m\}$ ,  $\psi(\|x_i - x_j\|) \leq t$ .

**Corollary 6.1** *Let  $K : \mathfrak{R}^d \times \mathfrak{R}^d \rightarrow \mathfrak{R}$  be a kernel such that  $K(s, t) = \psi(\|s - t\|)$  with  $\psi : \mathfrak{R} \rightarrow [0, 1]$  monotonically decreasing, satisfying  $\psi(0) = 1$ , and such that for all distinct  $i, j \in \{1, \dots, m\}$ ,  $\psi(\|x_i - x_j\|) \leq t$  for some  $t > 0$ . Let  $\gamma > 0$ ,  $m$  be a positive integer,  $\mathbf{x} = (x_1, \dots, x_m) \in \mathfrak{R}^{dm}$ ,  $\mathbf{y} = (y_1, \dots, y_m) \in \mathfrak{R}^m$ ,  $\min\{|y_i| : i = 1, \dots, m\} = 0$ ,  $g^o = \sum_{i=1}^m c_i K_{x_i}$  the unique solution of  $(\mathcal{H}_K(\mathfrak{R}^d), \mathcal{E}_K)$ ,  $\{\varepsilon_n\}$  a sequence of positive real numbers, and  $\{g_n\}$  a sequence of  $\varepsilon_n$ -near minimum points of  $(\text{span}_n G_K, \mathcal{E}_{\gamma, K})$ . Let*

$$\hat{u} = (1 + \gamma) \frac{\|\mathbf{y}\|_2^2}{\gamma^2 m} \quad \text{and}$$

$$\hat{v} = 2(1 + \gamma) \frac{\sqrt{1 + (m-1)t} \|\mathbf{y}\|_2^2}{\gamma^2 m^{3/2}}.$$

Then, for every positive integer  $n$  the following estimates hold:

- (i)  $\inf_{g \in \text{span}_n G_K} \mathcal{E}_{\gamma, K}(g) - \mathcal{E}_{\gamma, K}(g^o) \leq \frac{\hat{u}}{n} + \frac{\hat{v}}{\sqrt{n}}$ ;
- (ii)  $\mathcal{E}_{\gamma, K}(g_n) - \mathcal{E}_K(g^o) \leq \frac{\hat{u}}{n} + \frac{\hat{v}}{\sqrt{n}} + \varepsilon_n$ ;
- (iii)  $\|g_n - g^o\|_K^2 \leq \frac{1}{\gamma} \left( \frac{\hat{u}}{n} + \frac{\hat{v}}{\sqrt{n}} + \varepsilon_n \right)$ ;
- (iv)  $\sup_{x \in \Omega} |g_n(x) - g^o(x)|^2 \leq \frac{1}{\gamma} \left( \frac{\hat{u}}{n} + \frac{\hat{v}}{\sqrt{n}} + \varepsilon_n \right)$ .

**Proof.** As  $s_K = 1$  and  $\lambda_{\max} \leq \|\mathcal{K}[\mathbf{x}]\|_1 = \max_{j=1, \dots, m} \sum_{i=1}^m |K[\mathbf{x}]_{i,j}|$  [32, pp. 6, 21-23], we have  $\lambda_{\max} \leq 1 + (m-1)t$ . Hence the estimates (i)-(iv) follow from Theorem 5.4 with  $\bar{u} = \hat{u}$  and  $\bar{v} \leq \hat{v}$ .

□

Bounding from above the right-hand-side of the estimates from Corollary 6.1 in terms of the maximum of the absolute values of output data, we obtain the following corollary.

**Corollary 6.2** *Let  $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  be a kernel such that  $K(s, t) = \psi(\|s-t\|)$  with  $\psi : \mathbb{R} \rightarrow [0, 1]$  monotonically decreasing, satisfying  $\psi(0) = 1$ , and such that for all distinct  $i, j \in \{1, \dots, m\}$ ,  $\psi(\|x_i - x_j\|) \leq t$  for some  $t > 0$ . Let  $\gamma > 0$ ,  $m$  be a positive integer,  $\mathbf{x} = (x_1, \dots, x_m) \in \mathbb{R}^{dm}$ ,  $\mathbf{y} = (y_1, \dots, y_m) \in \mathbb{R}^m$ ,  $\min\{|y_i| : i = 1, \dots, m\} = 0$ ,  $y_{\max} = \max\{|y_i| : i = 1, \dots, m\}$ ,  $g^o = \sum_{i=1}^m c_i K_{x_i}$  the unique solution of  $(\mathcal{H}_K(\mathbb{R}^d), \mathcal{E}_K)$ ,  $\{\varepsilon_n; n = 1, \dots, m\}$  positive real numbers,  $\{g_n : n = 1, \dots, m\}$   $\varepsilon_n$ -near minimum points of  $(\text{span}_n G_K, \mathcal{E}_{\gamma, K})$ , and  $b = \frac{3(1+\gamma)y_{\max}^2}{\gamma^2}$ .*

Then, for every positive integer  $n \leq m$  the following estimates hold:

- (i)  $\inf_{g \in \text{span}_n G_K} \mathcal{E}_{\gamma, K}(g) - \mathcal{E}_{\gamma, K}(g^o) \leq \frac{b}{\sqrt{n}}$ ;
- (ii)  $\mathcal{E}_{\gamma, K}(g_n) - \mathcal{E}_K(g^o) \leq \frac{b}{\sqrt{n}} + \varepsilon_n$ ;
- (iii)  $\|g_n - g^o\|_K^2 \leq \frac{1}{\gamma} \left( \frac{b}{\sqrt{n}} + \varepsilon_n \right)$ ;
- (iv)  $\sup_{x \in \Omega} |g_n(x) - g^o(x)|^2 \leq \frac{1}{\gamma} \left( \frac{b}{\sqrt{n}} + \varepsilon_n \right)$ .

**Proof.** As  $\|\mathbf{y}\|_2^2 \leq m y_{\max}^2$ , we have  $\frac{\hat{u}}{n} + \frac{\hat{v}}{\sqrt{n}} \leq \frac{(1+\gamma)y_{\max}^2}{\gamma^2} \left( \frac{1}{n} + \frac{2\sqrt{1+(m-1)t}}{\sqrt{mn}} \right)$ , which for  $t \in [0, 1]$  and  $n \leq m$  is bounded from above by  $\frac{(1+\gamma)y_{\max}^2}{\gamma^2} \left( \frac{1}{n} + \frac{2}{\sqrt{n}} \right) \leq \frac{3(1+\gamma)y_{\max}^2}{\gamma^2 \sqrt{n}}$ . Hence the estimates (i)-(iv) follow from Corollary 6.1. □

So, when  $\gamma$  is not too small and  $y_{\max}$  is not too large, Corollary 6.2 guarantees a good approximation of the optimal solution by suboptimal ones.

In particular for the Gaussian kernel, the minimum of the regularized empirical error functional over the set of functions computable by Gaussian radial-basis function networks with  $n$  computational units approximates the global mini-

mum over the whole RKHS within  $\frac{b}{\sqrt{n}}$ , where  $b = \frac{3(1+\gamma)y_{\max}^2}{\gamma^2}$ . For example, for  $\gamma = 0.5$  we have  $\frac{1+\gamma}{\gamma^2} = 6$  and so  $b = 18y_{\max}^2$ .

## 7 Discussion

We have compared two approaches to learning from data with generalization capability, both modeling learning as a minimization of the empirical error functional with the square loss function regularized by the square of a norm on an RKHS, but differing in the hypothesis sets over which minimization is performed. The first approach, which is based on the Representer Theorem, considers minimization of the regularized empirical error over the whole RKHS, whereas the second one over its subset formed by functions computable by linear combinations of  $n$  computational units defined by the kernel.

We have derived upper bounds on the errors of approximation of the optimal solution by the suboptimal ones obtainable with  $n$  increasing. We have shown that when the absolute values of output data are not too large and the regularization parameter is not too small, suboptimal solutions approximate the optimal one within an accuracy  $\frac{c}{\sqrt{n}}$  with  $c$  moderate. In such cases, algorithms operating on models with  $n$  computational units can approximate the optimal solution quite well. Hence, when the solution of the system of linear equations described in the Representer Theorem is not computationally feasible or when the system is ill-conditioned, models with bounded complexity provide a useful and quite accurate alternative to the learning algorithms based on the Representer Theorem. For convolution kernels on  $\mathbb{R}^d \times \mathbb{R}^d$  the upper bounds from Corollaries 6.1 and 6.2 do not depend on the number  $d$  of variables, so the approximation of the optimal solution by such models does not exhibit the curse of dimensionality [4].

Minimization over a set of parameters of a chosen model is a nonlinear programming problem [35, p. 1489], which can be solved by iterative methods such as gradient descent [7, pp. 103-106, 173-174] (possibly with additive stochastic terms to avoid local minima, due to the nonconvexity of  $\mathcal{E}_{\gamma,K}$  as a function of the parameters), genetic algorithms [21], and simulated annealing [1].

## Acknowledgments

The authors thank Martin Holeňa, Paul Kainen, Per Kullstam, and Andrew Vogt for fruitful comments and discussions.

## A Appendix

A Reproducing Kernel Hilbert Space (RKHS) is a Hilbert space  $X$  formed by functions defined on a nonempty set  $\Omega$  such that for every  $u \in \Omega$  the evaluation functional  $\mathcal{F}_u$ , defined for any  $f \in X$  as  $\mathcal{F}_u(f) = f(u)$ , is bounded [2,5,10].

RKHSs can be characterized in terms of *kernels*, which are *symmetric positive semidefinite* functions  $K : \Omega \times \Omega \rightarrow \mathfrak{R}$ , i.e., functions satisfying for all positive integers  $m$ , all  $(w_1, \dots, w_m) \in \mathfrak{R}^m$ , and all  $(u_1, \dots, u_m) \in \Omega^m$ ,

$$\sum_{i,j=1}^m w_i w_j K(u_i, u_j) \geq 0.$$

By the Riesz Representation Theorem [17, p. 200], for every  $u \in \Omega$  there exists a unique element  $K_u \in X$ , called the *representer* of  $u$ , such that  $\mathcal{F}_u(f) = \langle f, K_u \rangle$  for all  $f \in X$  (this property is called the *reproducing property*). It is easy to check that the function  $K : \Omega \times \Omega$  defined for all  $u, v \in \Omega$  as  $K(u, v) = \langle u, v \rangle$  is a kernel.

On the other hand, every kernel  $K : \Omega \times \Omega \rightarrow \mathfrak{R}$  generates an RKHS  $\mathcal{H}_K(\Omega)$  that is the completion of the linear span of the set  $\{K_u : u \in \Omega\}$ , with the inner product defined as  $\langle K_u, K_v \rangle_K = K(u, v)$  and the induced norm  $\|\cdot\|_K$  (see, e.g., [2] and [5, p. 81]).

By the reproducing property and the Cauchy-Schwartz inequality, for every  $f \in \mathcal{H}_K(\Omega)$  and every  $u \in \Omega$  we have  $|f(u)| = |\langle f, K_u \rangle_K| \leq \|f\|_K \sqrt{K(u, u)} \leq s_K \|f\|_K$ , where  $s_K = \sup_{u \in \Omega} \sqrt{K(u, u)}$ . Thus for every kernel  $K$ , we have

$$\sup_{u \in \Omega} |f(u)| \leq s_K \|f\|_K. \quad (\text{A.1})$$

A paradigmatic example of a kernel on  $\mathfrak{R}^d$  is the *Gaussian kernel*  $K : \mathfrak{R}^d \times \mathfrak{R}^d \rightarrow \mathfrak{R}$ , defined as  $K(u, v) = \exp(-\|u - v\|^2)$ . Other examples of kernels are  $K(u, v) = \exp(-\|u - v\|)$ ,  $K(u, v) = \langle u, v \rangle^p$  (*homogeneous polynomial* of degree  $p$ ), where  $\langle \cdot, \cdot \rangle$  is any inner product on  $\mathfrak{R}^d$ ,  $K(u, v) = (1 + \langle u, v \rangle)^p$  (*inhomogeneous polynomial* of degree  $p$ ), and  $K(u, v) = (a^2 + \|u - v\|^2)^{-\alpha}$ , with  $\alpha > 0$  [10, p. 38].

The role of  $\|\cdot\|_K^2$  as a stabilizer can be illustrated by two examples of classes of kernels. The first one is formed by *Mercer kernels*, i.e., continuous, symmetric, and positive definite functions  $K : \Omega \times \Omega \rightarrow \mathfrak{R}$ , where  $\Omega \subset \mathfrak{R}^d$  is compact. For a Mercer kernel  $K$ ,  $\|f\|_K^2$  can be expressed using eigenvectors and eigenvalues of the compact linear operator  $L_K : \mathcal{L}_2(\Omega) \rightarrow \mathcal{C}(\Omega)$  defined for every  $f \in \mathcal{L}_2(\Omega)$  as  $L_K(f)(x) = \int_{\Omega} K(x, u) f(u) du$ , where  $\mathcal{L}_2(\Omega)$  and  $\mathcal{C}(\Omega)$  denote the spaces

of square integrable and of continuous functions on  $\Omega$ , respectively. By the Mercer Theorem [10, p. 34]

$$\|f\|_K^2 = \sum_{i=1}^{\infty} \frac{c_i^2}{\lambda_i},$$

where the  $\lambda_i$ 's are the eigenvalues of  $L_K$  and the  $c_i$ 's are the coefficients of the representation  $f = \sum_{i=1}^{\infty} c_i \phi_i$ , where  $\{\phi_i\}$  is the orthonormal basis of  $\mathcal{H}_K(\Omega)$  formed by the eigenvectors of  $L_K$ .

Note that the sequence  $\{\lambda_i\}$  is either finite or convergent to zero (for  $K$  smooth enough, the convergence to zero is rather fast [16, p. 1119]). Thus, the stabilizer  $\|\cdot\|_K^2$  penalizes functions for which the sequences of coefficients  $\{c_i\}$  do not converge to zero sufficiently quickly. So the functional  $\|\cdot\|_K^2$  plays the role of a high-frequency filter.

The second class of kernels illustrating the role of  $\|\cdot\|_K^2$  as a stabilizer consists of *convolution kernels*, i.e., kernels defined on  $\mathfrak{R}^d \times \mathfrak{R}^d$  such that  $K(x, y) = k(x - y)$ , for which the Fourier transform  $\tilde{k}$  of  $k$  is positive. For such kernels, the value of the stabilizer at any  $f \in \mathcal{H}_K(\Omega)$  can be expressed as

$$\|f\|_K^2 = \frac{1}{(2\pi)^{d/2}} \int_{\mathfrak{R}^d} \frac{\tilde{f}(\omega)^2}{\tilde{k}(\omega)} d\omega \quad (\text{A.2})$$

(see, e.g., [19] and [40, p. 97]). So the function  $1/\tilde{k}$  plays a role analogous to that of the sequence  $\{1/\lambda_i\}$  in the case of a Mercer kernel. For example, the Gaussian kernel is a convolution kernel with a positive Fourier transform.

Another example of a convolution kernel with a positive Fourier transform is  $K(u, v) = k(u - v)$ , where  $k(t) = \exp(-a \|t\|)$ ,  $\tilde{k}(\omega) = 2^{d/2} a \pi^{-1/2} \Gamma(d/2 + 1) (a^2 + \|\omega\|^2)^{-(d+1)/2}$  [40, p. 107] and  $\Gamma$  denotes the gamma function, defined for  $s > 0$  as  $\Gamma(s) = \int_0^{\infty} \exp(-r) r^{s-1} dr$ . In this case, the rate of decay of  $\tilde{k}(\omega)$  is of the order of  $\|\omega\|^{-(d+1)}$ . In particular, for  $d = 1$  and  $a = 1$  one gets  $K(u, v) = k(u - v) = \exp(-|u - v|)$ . Since  $\Gamma(1) = 1$ ,  $\Gamma(1/2) = \sqrt{\pi}$ , and  $\Gamma(s + 1) = s \Gamma(s)$ , we have  $\tilde{k}(\omega) = (\sqrt{2\pi}(1 + \omega^2))^{-1}$ . Thus  $\|f\|_K^2 = 1/2\pi \int_{\mathfrak{R}} \tilde{f}(\omega)^2 (\sqrt{2\pi}(1 + \omega^2)) d\omega = 1/\sqrt{2\pi} \int_{\mathfrak{R}} f(\omega)^2 d\omega + 1/\sqrt{2\pi} \int_{\mathfrak{R}} \omega^2 f(\omega)^2 d\omega$ . As  $\tilde{f}' = \omega \tilde{f}(\omega)$  and  $\int_{\mathfrak{R}} f(t)^2 dt = 1/2\pi \int_{\mathfrak{R}} f(\omega)^2 d\omega$ , by Parseval's formula [38, p. 172] we have  $\|f\|_K^2 = \sqrt{2\pi} (\|f\|_{\mathcal{L}_2}^2 + \|f'\|_{\mathcal{L}_2}^2)$ , where  $\|\cdot\|_{\mathcal{L}_2}$  denotes the  $\mathcal{L}_2$ -norm. So, as pointed out in [19], in this case the norm on the RKHS is equal to the Sobolev norm.

For more details on kernels and their role in learning theory, see, e.g., [40].

## References

- [1] E. Aarts and J. Korst. *Simulated Annealing and Boltzmann Machines: A Stochastic Approach to Combinatorial Optimization and Neural Computing*. John Wiley & Sons, New York, NY, 1989.
- [2] N. Aronszajn. Theory of reproducing kernels. *Transactions of AMS*, 68:337-404, 1950.
- [3] A. R. Barron. Universal approximation bounds for superpositions of a sigmoidal function. *IEEE Trans. on Information Theory*, 39:930-945, 1993.
- [4] R. Bellman. *Dynamic Programming*. Princeton University Press, Princeton, New Jersey, 1957.
- [5] C. Berg, J. P. R. Christensen, and P. Ressel. *Harmonic Analysis on Semigroups*. Springer-Verlag, New York, 1984.
- [6] M. Bertero. Linear inverse and ill-posed problems. *Advances in Electronics and Electron Physics*, 75:1-120, 1989.
- [7] D. P. Bertsekas. *Nonlinear Programming*. Athena Scientific, Belmont, MA, 1999.
- [8] C. Cortes and V. Vapnik, Support vector networks. *Machine Learning*, 20:1-25, 1995.
- [9] N. Cristianini and J. Shawe-Taylor. *An Introduction to Support Vector Machines and Other Kernel-Based Learning Methods*. Cambridge University Press, 2000.
- [10] F. Cucker and S. Smale. On the mathematical foundations of learning. *Bulletin of AMS*, 39:1-49, 2001.
- [11] F. Cucker and S. Smale. Best choices for regularization parameters in learning theory: on the bias-variance problem. *Foundations of Computational Mathematics*, 2:413-428, 2002.
- [12] J. A. Cuesta-Albertos and M. Wschebor. Some remarks on the condition number of a real random square matrix. *J. of Complexity*, 19:548-554, 2003.
- [13] J. Demmel. The geometry of ill-conditioning. *J. of Complexity*, 3:201-229, 1987.
- [14] A. L. Dontchev. *Perturbations, Approximations and Sensitivity Analysis of Optimal Control Systems*. Lecture Notes in Control and Information Sciences, vol. 52. Springer-Verlag, Berlin Heidelberg, 1983.
- [15] A. L. Dontchev and T. Zolezzi. *Well-Posed Optimization Problems*. Lecture Notes in Mathematics, vol. 1543. Springer-Verlag, Berlin Heidelberg, 1993.
- [16] N. Dunford, J. T. Schwartz: *Linear Operators. Part II: Spectral Theory*. Interscience Publishers, New York, NY, 1963.
- [17] A. Friedman. *Modern Analysis*. Dover, New York, 1982.

- [18] F. Girosi. Regularization theory, Radial Basis Functions and networks. In *From Statistics to Neural Networks. Theory and Pattern Recognition Applications - Subseries F, Computer and Systems Sciences* (V. Cherkassky, J. H. Friedman, and H. Wechsler, Eds.), pp. 166-187. Springer-Verlag, 1994.
- [19] F. Girosi. An equivalence between sparse approximation and support vector machines. *Neural Computation*, 10:1455-1480, 1998.
- [20] F. Girosi, M. Jones, and T. Poggio. Regularization theory and neural networks architectures. *Neural Computation*, 7:219-269, 1995.
- [21] D. E. Goldberg. *Genetic Algorithms in Search, Optimization, and Machine Learning*. Addison-Wesley, Reading, MA, 1989.
- [22] L. K. Jones. A simple lemma on greedy approximation in Hilbert space and convergence rates for projection pursuit regression and neural network training. *Ann. of Statistics*, 20:608-613, 1992.
- [23] G. S. Kimeldorf and G. Wahba. A correspondence between Bayesian estimation on stochastic processes and smoothing by splines. *Ann. Math. Statist.*, 41:495-502, 1970.
- [24] V. Kůrková. Dimension-independent rates of approximation by neural networks. In *Computer-Intensive Methods in Control and Signal Processing: The Curse of Dimensionality* (K. Warwick and M. Kárný, Eds.), pp. 261-270. Birkhäuser, Boston, 1997.
- [25] V. Kůrková. High-dimensional approximation by neural networks. In *Advances in Learning Theory: Methods, Models and Applications* (J. Stuykens et al., Eds.), pp. 69-88. IOS Press, Amsterdam, 2003.
- [26] V. Kůrková. Learning from data as an inverse problem. *COMPSTAT 2004 - Proceedings in Computational Statistics* (J. Antoch, Ed.), pp. 1377-1384. Physica-Verlag/Springer, Heidelberg, 2004.
- [27] V. Kůrková and M. Sanguineti. Bounds on rates of variable-basis and neural-network approximation. *IEEE Trans. on Information Theory*, 47:2659-2665, 2001.
- [28] V. Kůrková and M. Sanguineti. Comparison of worst case errors in linear and neural network approximation. *IEEE Trans. on Information Theory*, 48:264-275, 2002.
- [29] V. Kůrková and M. Sanguineti. Error estimates for approximate optimization by the extended Ritz method. *SIAM J. on Optimization*, to appear.
- [30] V. Kůrková, P. Savický, and K. Hlaváčková. Representations and rates of approximation of real-valued Boolean functions by neural networks. *Neural Networks*, 11:651-659, 1998.
- [31] E. S. Levitin and B. T. Polyak. Convergence of minimizing sequences in conditional extremum problems. *Sov. Math., Dokl.*, 7:764-767, 1966. English translation: *Dokl. Akad. Nauk SSSR*, 168:997-1000, 1966.

- [32] J. M. Ortega. *Numerical Analysis: A Second Course*. SIAM, Philadelphia, 1990 (reprint of the 1972 edition by Academic Press, New York, NY).
- [33] E. Parzen. An approach to time series analysis. *Ann. Math. Statist.*, 32:951-989, 1961.
- [34] Pisier, G.. Remarques sur un résultat non publié de B. Maurey. *Séminaire d'Analyse Fonctionnelle* 1980-81, Exposé no. V, pp. V.1-V.12. École Polytechnique, Centre de Mathématiques, Palaiseau, France.
- [35] T. Poggio and F. Girosi. Networks for approximation and learning. *Proc. IEEE* 78:1481-1497, 1990.
- [36] T. Poggio and F. Girosi. Regularization algorithms for learning that are equivalent to multilayer networks. *Science* 247:978-982, 1990.
- [37] T. Poggio and S. Smale. The mathematics of learning: dealing with data. *Notices of the AMS* 50:536-544, 2003.
- [38] W. Rudin. *Functional Analysis*. McGraw-Hill, New York, N.Y., 1973.
- [39] B. Schölkopf, R. Herbrich, A. J. Smola, and R. C. Williamson. A generalized Representer Theorem. *Lecture Notes in Artificial Intelligence* (Proc. of the Annual Conference on Computational Learning Theory - COLT), pp. 416-426. Springer-Verlag, London, 2001.
- [40] B. Schölkopf and A. J. Smola. *Learning With Kernels – Support Vector Machines, Regularization, Optimization and Beyond*. MIT Press, Cambridge, MA, 2002.
- [41] I. J. Schönberg. Metric spaces and completely monotone functions. *Ann. of Math.*, 39:811-841, 1938.
- [42] A. J. Smola and B. Schölkopf. Sparse greedy matrix approximation for machine learning, *Proc. 17th International Conf. on Machine Learning*, pp. 911–918. Morgan Kaufmann, San Francisco, CA, 2000.
- [43] A. N. Tikhonov. Solutions of incorrectly formulated problems and the regularization method. *Soviet Math. Dokl.*, 4:1035-1038, 1963.
- [44] A. N. Tikhonov and V. Y. Arsenin. *Solutions of Ill-posed Problems*. W.H. Winston, Washington, D.C., 1977.
- [45] V. Vapnik. *Statistical Learning Theory*. Wiley, Chichester, GB, 1998.
- [46] A. A. Vladimirov, Yu. E. Nesterov, and Yu. N. Chekanov. On uniformly convex functionals. *Vestnik Moskovskogo Universiteta. Seriya 15 - Vychislitel'naya Matematika i Kibernetika*, 3:12-23, 1978. (English translation: *Moscow University Computational Mathematics and Cybernetics*, pp. 10-21, 1979).
- [47] G. Wahba. *Splines Models for Observational Data*. CBMS-NSF Regional Conference Series in Applied Mathematics, vol. 59. SIAM, Philadelphia, PA, 1990.