

# Adaptation for regularization operators in learning theory

Andrea Caponnetto <sup>a b c</sup> Yuan Yao <sup>d</sup>

<sup>a</sup> *C.B.C.L., McGovern Institute, Massachusetts Institute of Technology, Bldg. 46-5155, , 77 Massachusetts Avenue, Cambridge, MA 02139*

<sup>b</sup> *D.I.S.I., Università di Genova, Via Dodecaneso 35, 16146 Genova, Italy*

<sup>c</sup> *Department of Computer Science, University of Chicago, 1100 East 58th Street, Chicago, IL 60637*

<sup>d</sup> *Department of Mathematics, University of California, Berkeley, CA 94720*

**Abstract**

We consider learning algorithms induced by regularization methods in the regression setting. We show that previously obtained error bounds for these algorithms using a-priori choices of the regularization parameter, can be attained using a suitable a-posteriori choice based on validation. In particular, these results prove adaptation of the rate of convergence of the estimators to the minimax rate induced by the "effective dimension" of the problem. We also show universal consistency for this class methods.

This report describes research done at the Center for Biological & Computational Learning, which is in the McGovern Institute for Brain Research at MIT, as well as in the Dept. of Brain & Cognitive Sciences, and which is affiliated with the Computer Sciences & Artificial Intelligence Laboratory (CSAIL).

This research was sponsored by grants from: Office of Naval Research (DARPA) Contract No. MDA972-04-1-0037, Office of Naval Research (DARPA) Contract No. N00014-02-1-0915, National Science Foundation (ITR/SYS) Contract No. IIS-0112991, National Science Foundation (ITR) Contract No. IIS-0209289, National Science Foundation-NIH (CRCNS) Contract No. EIA-0218693, National Science Foundation-NIH (CRCNS) Contract No. EIA-0218506, and National Institutes of Health (Conte) Contract No. 1 P20 MH66239-01A1.

Additional support was provided by: Central Research Institute of Electric Power Industry (CRIEPI), Daimler-Chrysler AG, Compaq/Digital Equipment Corporation, Eastman Kodak Company, Honda R&D Co., Ltd., Industrial Technology Research Institute (ITRI), Komatsu Ltd., Eugene McDermott Foundation, Merrill-Lynch, NEC Fund, Oxygen, Siemens Corporate Research, Inc., Sony, Sumitomo Metal Industries, and Toyota Motor Corporation.

This work was supported by the NSF grant 0325113. The first author is also supported by the FIRB Project ASTAA and the IST Programme of the European Community, under the PASCAL Network of Excellence, IST-2002-506778.

## 1. INTRODUCTION

We show that previous results in [2] about rates of convergence for regularization methods using a-priori choices of the regularization parameter, can be attained using a suitable a-posteriori choice based on validation. We also show universal consistency for this class methods. The framework for semi-supervised statistical learning theory is the same one considered in [2]. The algorithms we consider are based on the formalism of *regularization methods* for linear ill-posed inverse problems in their classical setting (see for example [11] for general reference). Some popular algorithms from this class are: regularized least-squares, truncated SVD, Landweber method and  $\nu$ -method.

The paper is organized as follows. In Section 2 we focus on a-priori choices of the regularization parameter for regularization methods. Theorem 1 shows universal consistency for a large class of choice rules, and Theorem 2 shows specific rates of convergence under suitable prior assumptions (parameterized by the constants  $r$ ,  $s$ ,  $C_r$  and  $D_s$ ) on the unknown probability measure  $\rho$ . Unlabelled data are added to the training set in order to improve the rates for a certain range of the parameters  $r$  and  $s$ .

In Section 3 we consider a validation technique for the a-posteriori choice of the regularization parameter. Theorem 3 shows how error bounds for the estimators  $f_{\tilde{\mathbf{z}}, \lambda}$ , with a-priori choices of  $\lambda$ , can be transferred to the estimators  $f_{\mathbf{z}^{\text{tot}}}$  which use the validation examples  $\mathbf{z}^v$  in  $\mathbf{z}^{\text{tot}} = (\tilde{\mathbf{z}}, \mathbf{z}^v)$  to determine  $\lambda$ . The subsequent corollaries are applications of Theorem 3 to the choices of  $\lambda$  described in Section 2.

In Sections 4 and 5 we give the proofs of the results stated in the previous Sections, using some lemmas from [2].

## 2. A-PRIORI CHOICE OF THE REGULARIZATION PARAMETER.

We consider the setting of semi-supervised statistical learning. We assume that  $Y \subset [-M, M]$  and we let the supervised part of the training set be equal to

$$\mathbf{z} = (z_1, \dots, z_m),$$

with  $z_i = (x_i, y_i)$  drawn i.i.d. according to the probability measure  $\rho$  over  $Z = X \times Y$ . Moreover we assume that the unsupervised part of the training set is  $(x_{m+1}^u, \dots, x_{\tilde{m}}^u)$ , with  $x_i^u$  drawn i.i.d. according to the marginal probability measure over  $X$ ,  $\rho_X$ . For sake of brevity we also introduce the complete training set

$$\tilde{\mathbf{z}} = (\tilde{z}_1, \dots, \tilde{z}_{\tilde{m}}),$$

with  $\tilde{z}_i = (\tilde{x}_i, \tilde{y}_i)$ , where we introduced the compact notations  $\tilde{x}_i$  and  $\tilde{y}_i$ , defined by

$$\tilde{x}_i = \begin{cases} x_i & \text{if } 1 \leq i \leq m, \\ x_i^u & \text{if } m < i \leq \tilde{m}, \end{cases}$$

and

$$\tilde{y}_i = \begin{cases} \frac{\tilde{m}}{m} y_i & \text{if } 1 \leq i \leq m, \\ 0 & \text{if } m < i \leq \tilde{m}. \end{cases}$$

It is clear that, in the supervised setting, the semi-supervised part of the training set is missing, whence  $\tilde{m} = m$  and  $\tilde{\mathbf{z}} = \mathbf{z}$ .

In the following we will study the generalization properties of a class of estimators  $f_{\tilde{\mathbf{z}}, \lambda}$  belonging to the *hypothesis space*  $\mathcal{H}$ : the RKHS of functions on  $X$  induced by the bounded Mercer kernel  $K$  (in the following  $\kappa = \sup_{x \in X} K(x, x)$ ). The learning algorithms that we consider, have the general form

$$(1) \quad f_{\tilde{\mathbf{z}}, \lambda} = G_\lambda(T_{\tilde{\mathbf{x}}}) g_{\mathbf{z}},$$

where  $T_{\tilde{\mathbf{x}}} \in \mathcal{L}(\mathcal{H})$  is given by,

$$T_{\tilde{\mathbf{x}}} f = \frac{1}{\tilde{m}} \sum_{i=1}^{\tilde{m}} K_{\tilde{x}_i} \langle K_{\tilde{x}_i}, f \rangle_{\mathcal{H}},$$

$g_{\mathbf{z}} \in \mathcal{H}$  is given by,

$$g_{\mathbf{z}} = \frac{1}{\tilde{m}} \sum_{i=1}^{\tilde{m}} K_{\tilde{x}_i} \tilde{y}_i = \frac{1}{m} \sum_{i=1}^m K_{x_i} y_i,$$

and the *regularization parameter*  $\lambda$  lays in the range  $(0, \kappa]$ . We will often use the shortcut notation  $\hat{\lambda} = \frac{\lambda}{\kappa}$ .

The functions  $G_\lambda : [0, \kappa] \rightarrow \mathbb{R}$ , which select the *regularization method*, will be characterized in terms of the constants  $A$  and  $B_r$  in  $[0, +\infty]$ , defined as follows

$$(2) \quad A = \sup_{\lambda \in (0, \kappa]} \sup_{\sigma \in [0, \kappa]} |(\sigma + \lambda)G_\lambda(\sigma)|$$

$$(3) \quad B_r = \sup_{t \in [0, r]} \sup_{\lambda \in (0, \kappa]} \sup_{\sigma \in [0, \kappa]} |1 - G_\lambda(\sigma)\sigma| \sigma^t \lambda^{-t}, \quad r > 0.$$

Finiteness of  $A$  and  $B_r$  (with  $r$  over a suitable range) are standard in the literature of ill-posed inverse problems (see for reference [11]). Regularization methods have been recently studied in the context of learning theory in [12, 8, 7, 9, 1].

The main results of the paper, Theorems 1 and 2, describe the convergence rates of  $f_{\tilde{\mathbf{z}}, \lambda}$  to the *target function*  $f_{\mathcal{H}}$ . Here, the target function is the “best” function which can be arbitrarily well approximated by elements of our hypothesis space  $\mathcal{H}$ . More formally,  $f_{\mathcal{H}}$  is the projection of the regression function  $f_\rho(x) = \int_{\mathcal{Y}} y d\rho_x(y)$  onto the closure of  $\mathcal{H}$  in  $\mathcal{L}^2(X, \rho_X)$ .

The convergence rates in Theorem 2, will be described in terms of the constants  $C_r$  and  $D_s$  in  $[0, +\infty]$  characterizing the probability measure  $\rho$ . These constants can be described in terms of the integral operator  $L_K : \mathcal{L}^2(X, \rho_X) \rightarrow \mathcal{L}^2(X, \rho_X)$  of kernel  $K$ . Note that the same integral operator is denoted by  $T$ , when seen as a bounded operator from  $\mathcal{H}$  to  $\mathcal{H}$ .

The constants  $C_r$  characterize the conditional distributions  $\rho_{|x}$  through  $f_{\mathcal{H}}$ , they are defined as follows

$$(4) \quad C_r = \begin{cases} \kappa^r \|L_K^{-r} f_{\mathcal{H}}\|_\rho & \text{if } f_{\mathcal{H}} \in \text{Im } L_K^r \\ +\infty & \text{if } f_{\mathcal{H}} \notin \text{Im } L_K^r \end{cases}, \quad r > 0.$$

Finiteness of  $C_r$  is a common *source condition* in the inverse problems literature (see [11] for reference). This type of condition has been introduced in the statistical learning literature in [6, 16, 3, 15, 4].

The constants  $D_s$  characterize the marginal distribution  $\rho_X$  through the *effective dimension*  $\mathcal{N}(\lambda) = \text{Tr}[T(T + \lambda)^{-1}]$ , they are defined as follows

$$(5) \quad D_s = 1 \vee \sup_{\lambda \in (0, 1]} \sqrt{\mathcal{N}(\lambda)\lambda^s}, \quad s \in (0, 1].$$

Finiteness of  $D_s$  was implicitly assumed in [3, 4].

The next theorem shows (strong) universal consistency (in probability) for the estimators  $f_{\tilde{\mathbf{z}}, \lambda}$  under mild assumptions on the choice of  $\lambda$ . The function  $|x|_+$ , appearing in the text of Theorem 1, is the “positive part” of  $x$ , that is  $\frac{x+|x|}{2}$ .

**Theorem 1.** *Let  $\{\tilde{\mathbf{z}}_m\}_{m=1}^\infty$  be a sequence of training sets composed of  $m$  labelled examples drawn i.i.d. from a probability measure  $\rho$  over  $Z$ , and  $\tilde{m}_m - m \geq 0$  unlabelled examples drawn i.i.d. from the marginal measure of  $\rho$  over  $X$ . Let the regularization parameter choice,  $\lambda_m : \mathbb{N} \rightarrow (0, \kappa]$ , fulfill the conditions*

$$(6) \quad \lim_{m \rightarrow \infty} \lambda_m = 0,$$

$$(7) \quad \lim_{m \rightarrow \infty} \sqrt{m}\lambda_m = \infty.$$

Then, if  $B_{\bar{r}} < +\infty$  for some  $\bar{r} > 0$ , it holds <sup>1</sup>

$$\lim_{m \rightarrow \infty} \|f_{\tilde{\mathbf{z}}_m, \lambda_m} - f_{\mathcal{H}}\|_{\rho} \stackrel{P}{=} 0.$$

Theorem 2 below is a restatement in a slightly modified form of Theorem 2 in [2]. In particular the introduction of the parameter  $q > 1$  will be useful when we will merge this result with Theorem 3 in the proof of Corollary 2.

**Theorem 2.** *Let  $r > 0$ ,  $s \in (0, 1]$  and  $\alpha \in [0, |2 - 2r - s|_+]$ . Furthermore, let  $m$  and  $\lambda$  satisfy the constraints  $\lambda \leq \|T\|$  and*

$$(8) \quad \dot{\lambda} = q \left( \frac{4D_s \log \frac{6}{\delta}}{\sqrt{m}} \right)^{\frac{2}{2r+s+t_1}},$$

for some  $q \geq 1$ ,  $\delta \in (0, 1/3)$  and  $t_1$  defined in eq. (10). Finally, assume  $\tilde{m} \geq 4 \vee m \dot{\lambda}^{-\alpha}$ . Then, with probability greater than  $1 - 3\delta$ , it holds

$$\|f_{\tilde{\mathbf{z}}, \lambda} - f_{\mathcal{H}}\|_{\rho} \leq q^r E_r \left( \frac{4D_s \log \frac{6}{\delta}}{\sqrt{m}} \right)^{\frac{2r-t_2}{2r+s+t_1}},$$

where

$$(9) \quad E_r = C_r (30A + 2(3+r)B_r + 1) + 9MA,$$

$$(10) \quad t_1 = |2 - 2r - s|_+ - \alpha,$$

$$(11) \quad t_2 = |1 - 2r - 2s - t_1|_+.$$

The proofs of the above Theorems is postponed to Section 4.

### 3. ADAPTATION.

In this section we show the adaptation properties of the estimators obtained by a suitable data-dependent choice of the regularization parameter. The main results of this section are obtained assuming that

$$(12) \quad f_{\mathcal{H}} = f_{\rho},$$

this is true for every  $\rho$  when the underlying kernel  $K$  is *universal* (see [17]). In fact for this class of kernels the RKHS  $\mathcal{H}$  is always dense in  $\mathcal{L}^2(X, \rho_X)$ . The Gaussian kernel is a popular instance of a kernel in this family.

Let the *validation set*

$$\mathbf{z}^v = (z_1^v, \dots, z_{m^v}^v),$$

be composed of  $m^v$  labelled examples  $z_i^v = (x_i^v, y_i^v)$  drawn i.i.d. from the probability measure  $\rho$  over  $Z = X \times Y$ . The validation set  $\mathbf{z}^v$  is, by assumption, independent of the training set  $\tilde{\mathbf{z}}$ , and these two sets define the *learning set*

$$\mathbf{z}^{\text{tot}} = (\tilde{\mathbf{z}}, \mathbf{z}^v),$$

which represents the total input of the adaptive learning algorithm. Following the notations of the previous Section, we let  $\tilde{m}$  be the total number of examples in  $\tilde{\mathbf{z}}$ , and  $m$  the number of its labelled examples.

<sup>1</sup>We say that the sequence of random variables  $\{X_m\}_{m \in \mathbb{N}}$  converges in probability to the random variable  $X$  (and we write  $\lim_{m \rightarrow \infty} X_m \stackrel{P}{=} X$  or  $X_m \xrightarrow{P} X$ ), if for every  $\epsilon > 0$ ,  $\lim_{m \rightarrow \infty} \mathbb{P}[|X_m - X| \geq \epsilon] = 0$ . This is equivalent to say that, for every  $\delta \in (0, 1)$ ,  $\mathbb{P}[|X_m - X| \geq \epsilon(m, \delta)] \leq \delta$ , with  $\lim_{m \rightarrow \infty} \epsilon(m, \delta) = 0$ .

Now let us explain how  $\mathbf{z}^v$  is used for the choice of  $\lambda$ . We consider the finite set of positive reals  $\Lambda_m$  depending on  $m$ , the number of labelled examples in  $\tilde{\mathbf{z}}$ , and the data-dependent choice for the regularization parameter is

$$(13) \quad \hat{\lambda}_{\mathbf{z}^v} = \operatorname{argmin}_{\lambda \in \Lambda_m} \frac{1}{m^v} \sum_{i=1}^{m^v} (T_M f_{\tilde{\mathbf{z}}, \lambda}(x_i^v) - y_i^v)^2,$$

where the truncation operator  $T_M : \mathcal{L}^2(X, \rho_X) \rightarrow \mathcal{L}^2(X, \rho_X)$  is defined by

$$T_M f(x) = (|f(x)| \wedge M) \operatorname{sign} f(x).$$

The final learning estimator, whose adaptation properties are investigated in this Section, is defined as follows

$$(14) \quad f_{\mathbf{z}^{\text{tot}}} = T_M f_{\tilde{\mathbf{z}}, \hat{\lambda}_{\mathbf{z}^v}}.$$

Theorem 3 below is the main result of this Section and shows an important property of the estimator  $f_{\mathbf{z}^{\text{tot}}}$ . It will be used to extend to  $f_{\mathbf{z}^{\text{tot}}}$  convergence results similar to the ones obtained in the previous Section.

**Theorem 3.** *Let  $\rho$ ,  $K$ ,  $m$ ,  $\tilde{m}$ ,  $m^v$ ,  $\Lambda_m$ ,  $\delta \in (0, 1)$ ,  $\epsilon > 0$  and  $\lambda_m \in \Lambda_m$  be such that with probability greater than  $1 - \delta$ , it holds*

$$\|f_{\tilde{\mathbf{z}}, \lambda_m} - f_\rho\|_\rho \leq \epsilon.$$

*Then, with probability greater than  $1 - 2\delta$ , it holds*

$$\|f_{\mathbf{z}^{\text{tot}}} - f_\rho\|_\rho \leq \hat{\epsilon},$$

*with*

$$\hat{\epsilon}^2 = 2\epsilon^2 + \frac{80M^2}{m^v} \log \frac{2|\Lambda_m|}{\delta}.$$

The proof of Theorem 3 is postponed to Section 5.

The first corollary of Theorem 3 proves universal consistency for the estimators  $f_{\mathbf{z}^{\text{tot}}}$  under mild assumptions on the cardinalities of the grids  $\Lambda_m$  and validation sets  $\mathbf{z}^v$ .

**Corollary 1.** *Let  $K$  be a universal kernel,  $Q$  be a constant greater than 1, and define*

$$(15) \quad \Lambda_m = \{\kappa, \kappa Q^{-1}, \dots, \kappa Q^{-|\Lambda_m|+1}\},$$

*with*

$$(16) \quad |\Lambda_m| = \omega(1).$$

*Moreover let  $\{\mathbf{z}_m^{\text{tot}}\}_{m=1}^\infty$  be a sequence of learning sets drawn according to a probability measure  $\rho$  over  $Z$ . Assume  $\mathbf{z}_m^{\text{tot}} = (\tilde{\mathbf{z}}_m, \mathbf{z}_m^v)$ , with the training sets  $\tilde{\mathbf{z}}_m$  composed of  $m$  labelled examples and  $\tilde{m}_m - m \geq 0$  unlabelled examples, and  $\mathbf{z}_m^v$  the validation sets composed by  $m_m^v = \omega(\log |\Lambda_m|)$  examples. Then, if  $B_{\bar{r}} < +\infty$  for some  $\bar{r} > 0$ , it holds*

$$\lim_{m \rightarrow \infty} \|f_{\mathbf{z}_m^{\text{tot}}} - f_\rho\|_\rho \stackrel{P}{=} 0.$$

*Proof.* The result is a corollary of theorems 1 and 3. The universality of  $K$  enforces the equality (12) (see [17]). Condition (16) implies that the regularization parameter  $\lambda_m = \kappa Q^{-(\lceil \log \log m \rceil \wedge |\Lambda_m|)}$ , which belongs to  $\Lambda_m$ , fulfills the assumptions (6) and (7). Hence, using the assumption on  $m_m^v$ , we get that for every  $\delta \in (0, 1)$ , with probability greater than  $1 - 2\delta$

$$\|f_{\mathbf{z}_m^{\text{tot}}} - f_\rho\|_\rho^2 \leq o(1) + \frac{80M^2}{\omega(\log |\Lambda_m|)} \log \frac{2|\Lambda_m|}{\delta} \rightarrow 0.$$

□

The second corollary proves explicit rates for the convergence of  $f_{\mathbf{z}^{\text{tot}}}$  to  $f_\rho$  over specific prior classes defined in term of finiteness of the constants  $C_r$  and  $D_s$ . The main assumption is the requirement  $m^\vee \geq m/\log m$ . Since this constrain can be fulfilled still being  $m^\vee$  asymptotically negligible with respect to  $m$ , the rates (expressed in terms of  $m$ ) that are obtained in the second part of the corollary are minimax optimal over the corresponding priors (see [4]).

**Corollary 2.** *Let  $K$  be a universal kernel. Consider a learning set  $\mathbf{z}^{\text{tot}}$  with  $m^\vee \geq \frac{m}{\log m}$  and  $\tilde{m} \geq 4 \vee m^{1+\eta}$ , for some constants  $\eta \geq 0$ ,  $r > 0$ , and  $s \in (0, 1]$ . Define  $\Lambda_m$  as in eq. (15) with  $Q$  an arbitrary constant greater than 1 and*

$$(17) \quad \frac{\eta}{\alpha} \log_Q m + 1 \leq |\Lambda_m| \leq m,$$

with  $\alpha$  defined by eq. (20).

Moreover assume that for some  $\delta \in (0, 1/6)$ ,  $m$  is large enough that it holds

$$(18) \quad Q \left( \frac{4D_s \log \frac{6}{\delta}}{\sqrt{m}} \right)^{\frac{2\eta}{\alpha}} \leq \kappa^{-1} \|T\|.$$

Then, with probability greater than  $1 - 6\delta$

$$(19) \quad \|f_{\mathbf{z}^{\text{tot}}} - f_\rho\|_\rho \leq 4(Q^r E_r D_s + 3M) \log(6m/\delta) m^{-\frac{1}{2} \frac{2r-t_2}{2r+s+t_1}},$$

where  $E_r$ ,  $t_1$  and  $t_2$  are the constants defined in equations (9), (10) and (11) substituting

$$(20) \quad \alpha = |2 - 2r - s|_+ \wedge \frac{\eta}{1 + \eta} (2r + s + |2 - 2r - s|_+).$$

In particular, if  $r + s \geq \frac{1}{2}$  and  $\eta = \frac{|2-2r-s|_+}{2r+s}$ , and assuming

$$(21) \quad 2 \log_Q m + 1 \leq |\Lambda_m| \leq m,$$

and

$$Q \left( \frac{4D_s \log \frac{6}{\delta}}{\sqrt{m}} \right)^{\frac{2}{2r+s}} \leq \kappa^{-1} \|T\|,$$

with probability greater than  $1 - 6\delta$ , it holds

$$\|f_{\mathbf{z}^{\text{tot}}} - f_\rho\|_\rho \leq 4(Q^r E_r D_s + 3M) \log(6m/\delta) m^{-\frac{1}{2} \frac{2r}{2r+s}}.$$

*Proof.* The result is a corollary of theorems 2 and 3. The universality of  $K$  enforces the equality (12) (see [17]).

First, from equations (20) and (10), by simple algebra we get

$$\frac{\eta}{\alpha} = \frac{1}{2r + s + t_1}.$$

Therefore condition (18) is equivalent to

$$\dot{\lambda}_q = q \left( \frac{4D_s \log \frac{6}{\delta}}{\sqrt{m}} \right)^{\frac{2}{2r+s+t_1}} \leq \kappa^{-1} \|T\| \quad \forall q \in [1, Q],$$

and condition  $\dot{\lambda} \leq \kappa^{-1} \|T\|$  in the text of Theorem 2 is verified by  $\dot{\lambda}_q$  for every  $q \in [1, Q]$ . Moreover, since  $D_s \geq 1$  and  $\delta \leq 1/6$ , for every  $q \in [1, Q]$  we can write

$$\tilde{m} \geq 4 \vee m^{1+\eta} = 4 \vee m \left( m^{-\frac{2}{\alpha}} \right)^{-\alpha} \geq 4 \vee m \dot{\lambda}_q^{-\alpha},$$

which shows that also the other assumption of Theorem 2 is verified.

Hence, by Theorem 2 we get that for every  $q \in [1, Q]$ , with probability greater than  $1 - 3\delta$ , it holds

$$\|f_{\mathbf{z}, \lambda} - f_\rho\|_\rho \leq \epsilon = Q^r E_r \left( \frac{4D_s \log \frac{6}{\delta}}{\sqrt{m}} \right)^{\frac{2r-t_2}{2r+s+t_1}}.$$

The next step is verifying that for some  $\bar{q} \in [1, Q]$ ,  $\lambda_{\bar{q}} = \kappa \dot{\lambda}_{\bar{q}} \in \Lambda_m$ , and hence applying Theorem 3.

In fact, from definition (15), assumption (17) and Proposition 5, it is clear that

$$\min \Lambda_m \leq \kappa m^{-\frac{\eta}{\alpha}} \leq \lambda_1 \leq \lambda_{\bar{q}} \leq \lambda_Q \leq \|T\| \leq \kappa = \max \Lambda_m,$$

for some  $\bar{q}$ .

Applying Theorem 3, we get that with probability greater than  $1 - 6\delta$  it holds

$$\|f_{\mathbf{z}}^{\text{tot}} - f_{\rho}\|_{\rho} \leq \hat{\epsilon},$$

with, using again condition (17) and the assumption  $m^{\vee} \geq \frac{m}{\log m}$ , the chain of inequalities

$$\begin{aligned} \hat{\epsilon} &= \left( \epsilon^2 + \frac{80M^2}{m^{\vee}} \log \frac{6|\Lambda_m|}{\delta} \right)^{\frac{1}{2}} \\ &\leq \left( \epsilon^2 + \frac{80M^2}{m} \log^2 \frac{6m}{\delta} \right)^{\frac{1}{2}} \leq \epsilon + \frac{12M}{\sqrt{m}} \log \frac{6m}{\delta} \\ &\leq 4Q^r E_r D_s \log(6/\delta) m^{-\frac{1}{2} \frac{2r-t_2}{2r+s+t_1}} + 12M \log(6m/\delta) m^{-\frac{1}{2}} \\ &\leq 4(Q^r E_r D_s + 3M) \log(6m/\delta) m^{-\frac{1}{2} \frac{2r-t_2}{2r+s+t_1}}, \end{aligned}$$

which concludes the proof of the first part of the Corollary.

The second part of the Corollary is an instantiation of the previous result. In fact by equations (20) and (10), the assumption  $\eta = \frac{|2-2r-s|_+}{2r+s}$  implies  $\alpha = |2-2r-s|_+$  and  $t_1 = 0$ . Moreover from the assumption  $r+s \geq \frac{1}{2}$  and eq. (11) we get  $t_2 = 0$ , and noticing that  $\frac{\eta}{\alpha} = \frac{1}{2r+s} \leq 2$ , it is clear that condition (21) implies condition (17).  $\square$

#### 4. PROOFS OF THEOREMS 1 AND 2

In this section we give the proof of Theorems 1 and 2. We use various propositions taken [2], which we state without proof.

**4.1.** Before proving Theorem 1, we begin showing some preliminary propositions. The first one is a technical result about sequences of real numbers.

**Proposition 1.** *Let  $\{a_i\}_{i \in \mathbb{N}}$  and  $\{b_i\}_{i \in \mathbb{N}}$  be two non-increasing sequences of reals in the interval  $(0, 1)$  with*

$$\begin{aligned} \lim_{i \rightarrow \infty} a_i &= 0, \\ \lim_{i \rightarrow \infty} b_i &= 0. \end{aligned}$$

*Then there exists a sequence  $\{c_i\}_{i \in \mathbb{N}}$  of reals in the interval  $(0, 1)$  such that, defining  $d_i = \log c_i / \log b_i$ , the following properties hold,*

- i)  $\{d_i\}_{i \in \mathbb{N}}$  is a non-increasing sequence of positive reals.*
- ii)  $\{c_i\}_{i \in \mathbb{N}}$  is a non-increasing sequence of positive reals, with*

$$\begin{aligned} c_i &\geq a_i \quad \forall i \in \mathbb{N}, \\ \lim_{i \rightarrow \infty} c_i &= 0. \end{aligned}$$

*Proof.* We consider the sequence  $\{c_i\}_{i \in \mathbb{N}}$  of positive numbers constructed by the recursive rule

$$\begin{aligned} c_1 &= a_1, \\ c_{i+1} &= a_{i+1} \vee (b_{i+1})^{\frac{\log c_i}{\log b_i}}. \end{aligned}$$

Let us prove point *i)* by induction.

Since by assumption  $a_1$  and  $b_1$  belong to  $(0, 1)$ , by construction  $d_1 = \frac{\log c_1}{\log b_1} = \frac{\log a_1}{\log b_1} > 0$ . Now, for  $i \geq 1$  assume  $d_i > 0$ , then by construction, either  $c_{i+1} = (b_{i+1})^{d_i}$ , and hence  $d_{i+1} = \frac{\log c_{i+1}}{\log b_{i+1}} = d_i > 0$ , or  $c_{i+1} = (b_{i+1})^{d_{i+1}} = a_{i+1} \geq (b_{i+1})^{d_i}$ , and hence, since  $a_{i+1}$  and  $b_{i+1}$  belong to  $(0, 1)$ , it holds

$$\begin{aligned} d_{i+1} &= \frac{\log a_{i+1}}{\log b_{i+1}} > 0, \\ (b_{i+1})^{d_{i+1}} \geq (b_{i+1})^{d_i} &\Rightarrow d_{i+1} \leq d_i. \end{aligned}$$

Let us now prove point *ii*).

First, by construction  $c_i \geq a_i > 0$ . Moreover, again by construction, either  $c_{i+1} = a_{i+1}$ , and hence,

$$c_{i+1} = a_{i+1} \leq a_i \leq c_i,$$

or  $c_{i+1} = (b_{i+1})^{d_i}$  and hence, since  $d_i > 0$  by point *i*), it holds

$$c_{i+1} = (b_{i+1})^{d_i} \leq (b_i)^{d_i} = c_i.$$

Therefore the sequence  $\{c_i\}_{i \in \mathbb{N}}$  is non-increasing and  $c_i \leq c_1 = a_1 < 1$ .

Finally, we prove that  $\lim_i c_i = 0$ .

Let us assume there exists an infinite increasing sequence of naturals  $\{i(k)\}_{k \in \mathbb{N}}$ , such that

$$c_{i(k)} = a_{i(k)} \quad \forall k \in \mathbb{N}.$$

Since, by assumption,  $\lim_i a_i = 0$ , then  $\lim_k c_{i(k)} = 0$ . Therefore, since we already proved that  $\{a_i\}_{i \in \mathbb{N}}$  is non-increasing,  $\lim_i c_i = 0$ . Which proves the Proposition, if  $\{i(k)\}_{k \in \mathbb{N}}$  exists.

If  $\{i(k)\}_{k \in \mathbb{N}}$  does not exist, by construction, there exists  $I \in \mathbb{N}$  such that

$$c_{i+1} = (b_{i+1})^{d_i} \quad \forall i \geq I.$$

Therefore, recalling the definition of  $d_i$ , by induction, it follows

$$c_i = (b_i)^{d_I} \quad \forall i > I.$$

Recalling that  $d_I > 0$  and  $\lim_i b_i = 0$ , the relation above proves that, also in this case,  $\lim_i c_i = 0$ . □

The next proposition introduces the functions  $f_\lambda^{\text{tr}}$  and shows some simple results related to them.

**Proposition 2.** For any  $\lambda > 0$  let the truncated function  $f_\lambda^{\text{tr}}$  be defined by

$$(22) \quad f_\lambda^{\text{tr}} = P_\lambda f_\mathcal{H}$$

where  $P_\lambda$  is the orthogonal projector in  $\mathcal{L}^2(X, \rho_X)$  defined by

$$(23) \quad P_\lambda = \Theta_\lambda(L_K),$$

with

$$(24) \quad \Theta_\lambda(\sigma) = \begin{cases} 1 & \text{if } \sigma \geq \lambda, \\ 0 & \text{if } \sigma < \lambda. \end{cases}$$

Then the function  $a : (0, \kappa] \rightarrow \mathbb{R}$ , defined by

$$(25) \quad a(\lambda) = \|f_\lambda^{\text{tr}} - f_\mathcal{H}\|_\rho,$$

is non-decreasing and fulfills the following properties

$$(26) \quad 0 \leq a(\lambda) \leq M \quad \forall \lambda \in (0, \kappa],$$

$$(27) \quad \lim_{\lambda \rightarrow 0} a(\lambda) = 0.$$

*Proof.* Recall that the self-adjoint integral operator  $L_K$  has a countable eigensystem  $\{(\lambda_i, \phi_i)\}_{i=1}^{\infty}$  with positive eigenvalues decreasing to zero (see [5]). Moreover  $L_K^{\frac{1}{2}}$  is an isometry between  $\mathcal{L}^2(X, \rho_X)$  and  $\mathcal{H}$  (again, see [5]). Therefore, since  $f_{\mathcal{H}}$  is the projection of  $f_{\rho}$  over the closure of  $\mathcal{H}$  in  $\mathcal{L}^2(X, \rho_X)$ , it holds

$$f_{\mathcal{H}} = \sum_{i=1}^{\infty} |\langle f_{\rho}, \phi_i \rangle_{\rho}|^2 \phi_i.$$

Hence, by the definition of  $f_{\lambda}^{\text{tr}}$ , and recalling that  $Y \subset [-M, M]$ , we get

$$0 \leq a(\lambda)^2 = \sum_{\lambda_i < \lambda} |\langle f_{\rho}, \phi_i \rangle_{\rho}|^2 \leq \sum_{i=1}^{\infty} |\langle f_{\rho}, \phi_i \rangle_{\rho}|^2 \leq \|f_{\rho}\|_{\rho}^2 \leq M^2.$$

Monotonicity and convergence to zero for  $a(\lambda)$  follow from the relation above by standard arguments on convergent series of positive numbers.  $\square$

The next proposition is used in the proof of Theorem 1.

**Proposition 3.** *Let  $\bar{r}$  be a positive number. Then, there exists a function*

$$R : (0, 1] \rightarrow (0, \bar{r}]$$

such that

$$(28) \quad \kappa^{R(\dot{\lambda})} \left\| L_K^{-R(\dot{\lambda})} P_{\dot{\lambda}} f_{\mathcal{H}} \right\|_{\rho} \leq 4M, \quad \forall \dot{\lambda} \in (0, 1],$$

$$(29) \quad \lim_{\dot{\lambda} \rightarrow 0} \dot{\lambda}^{R(\dot{\lambda})} = 0.$$

*Proof.* Let  $\{\lambda_i, \phi_i\}$  be the eigensystem of the positive compact operator  $L_K$  (we also use the shortcut notation  $\dot{\lambda}_i = \kappa^{-1} \lambda_i$ ). First, if the range of  $L_K$  is finite dimensional, the choice  $R(\dot{\lambda}) = \bar{r}$  fulfills trivially the required conditions. Second, from definition (25), it is clear that if the sequence  $\{a(\lambda_i)\}_i$  has only a finite number of positive elements,  $f_{\mathcal{H}}$  belongs to the finite dimensional range of the projector  $P_{\dot{\lambda}}$ , for some positive  $\bar{\lambda}$ , and the choice  $R(\dot{\lambda}) = \bar{r}$  is again a trivial solution.

Therefore in the following we assume  $\lambda_i > 0$  and  $a(\lambda_i) > 0$  for every  $i \in \mathbb{N}$ . Moreover, from Proposition 5,  $\lambda_i \leq \kappa$ , and by eq. (26),  $a(\lambda) \leq M$ . Hence we can apply Proposition 1 to the non-increasing sequences  $\{a_i\}_i$  and  $\{b_i\}_i$  defined by

$$a_i = \frac{a(\lambda_i)}{2M},$$

$$b_i = \frac{\lambda_i}{2\kappa}.$$

The function  $R$  is defined in terms of the sequence  $\{d_i\}_i$  constructed in Proposition 1 as follows

$$R(\lambda) = \begin{cases} \bar{r} & \text{if } \lambda_1 < \lambda \leq 1, \\ \bar{r} d_i / (\bar{r} \vee d_1) & \text{if } \lambda_{i+1} < \lambda \leq \lambda_i, \quad i \geq 1. \end{cases}$$

Equality (29) can be proved, recalling that by Proposition 1  $c_i = b_i^{d_i} \leq \dot{\lambda}_i^{d_i}$  goes to zero as  $i \rightarrow \infty$ , and hence

$$\lim_{\dot{\lambda} \rightarrow 0} \dot{\lambda}^{R(\dot{\lambda})} = \lim_{i \rightarrow \infty} \dot{\lambda}_i^{R(\dot{\lambda}_i)} = \lim_{i \rightarrow \infty} \left( (2b_i)^{d_i} \right)^{\bar{r}/(\bar{r} \vee d_1)} \leq 2^{\bar{r}} \left( \lim_{i \rightarrow \infty} c_i \right)^{\bar{r}/(\bar{r} \vee d_1)} = 0.$$

Since by Proposition 1  $\{d_i\}_i$  is a sequence of non-increasing positives, then  $R$  is non-decreasing. Therefore, defining  $f_i = \langle f_{\mathcal{H}}, \phi_i \rangle_{\rho}$ , we can write

$$\begin{aligned}
\kappa^{2R(\lambda)} \left\| L_K^{-R(\lambda)} P_{\lambda} f_{\mathcal{H}} \right\|_{\rho}^2 &= \sum_{\lambda_i \geq \lambda} f_i^2 \lambda_i^{-2R(\lambda)} \leq \sum_i f_i^2 \lambda_i^{-2R(\lambda_i)} \\
&= \sum_i f_i^2 \left( (2b_i)^{d_i} \right)^{-\bar{r}/(\bar{r} \vee d_1)} \\
\left( b_i^{d_i} = c_i < 1 \right) &\leq \sum_i f_i^2 c_i^{-1} \\
\left( c_i > a_i \right) &\leq \sum_i f_i^2 a_i^{-1} = 2M \sum_i f_i^2 a(\lambda_i)^{-1} \\
&= 2M \sum_{k=0}^{\infty} \sum_{2^{-k-1} < a(\lambda_i)/M \leq 2^{-k}} f_i^2 a(\lambda_i)^{-1} \\
&\leq 2 \sum_{k=0}^{\infty} 2^{k+1} \sum_{2^{-k-1} < a(\lambda_i)/M \leq 2^{-k}} f_i^2 \\
\left( a(\lambda)^2 = \sum_{\lambda_i \leq \lambda} f_i^2 \right) &\leq 4M^2 \sum_{k=0}^{\infty} 2^{-k} = 8M^2,
\end{aligned}$$

which proves inequality (28) and concludes the proof.  $\square$

We now state four propositions from [2]. The first one introduces the empirical and ideal estimators least-squares  $f_{\bar{\mathbf{z}}, \lambda}^{\text{ls}}$  and  $f_{\lambda}^{\text{ls}}$ .

**Proposition 4.** *Assume  $\lambda \leq \|T\|$  and*

$$(30) \quad \lambda \tilde{m} \geq 16\kappa \mathcal{N}(\lambda) \log^2 \frac{6}{\delta},$$

for some  $\delta \in (0, 1)$ . Then, with probability greater than  $1 - \delta$ , it holds

$$\left\| (T + \lambda)^{\frac{1}{2}} (f_{\bar{\mathbf{z}}, \lambda}^{\text{ls}} - f_{\lambda}^{\text{ls}}) \right\|_{\mathcal{H}} \leq 8 \left( M + \sqrt{\kappa \frac{m}{\tilde{m}}} \left\| f_{\lambda}^{\text{ls}} \right\|_{\mathcal{H}} \right) \left( \frac{2}{m} \sqrt{\frac{\kappa}{\lambda}} + \sqrt{\frac{\mathcal{N}(\lambda)}{m}} \right) \log \frac{6}{\delta}$$

where

$$(31) \quad f_{\bar{\mathbf{z}}, \lambda}^{\text{ls}} = (T_{\bar{\mathbf{x}}} + \lambda)^{-1} g_{\bar{\mathbf{z}}},$$

$$(32) \quad f_{\lambda}^{\text{ls}} = (T + \lambda)^{-1} L_K f_{\mathcal{H}}.$$

*Proof.* See Proposition 1 in [2].  $\square$

The second one gives two simple properties of the operator  $T$  and the effective dimension  $\mathcal{N}(\lambda)$ .

**Proposition 5.** *For every probability measure  $\rho_X$  and  $\lambda > 0$ , it holds*

$$\|T\| \leq \kappa,$$

and

$$\lambda \mathcal{N}(\lambda) \leq \kappa.$$

*Proof.* See Proposition 2 in [2].  $\square$

The other two propositions from [2] estimate two different terms which appear in the proofs of Theorems 1 and 1. The symbol  $\lfloor x \rfloor$  in the text below represents the greater integer less or equal to  $x$ .

**Proposition 6.** Let  $f$  belong to  $\text{Im } L_K^r$  for some  $r > 0$ . Then, if  $\lambda \in (0, \kappa]$ , it holds

$$\left\| \sqrt{T} (G_\lambda(T_{\bar{\mathbf{x}}}) T_{\bar{\mathbf{x}}} - \text{Id}) P_\lambda f \right\|_{\mathcal{H}} \leq B_r \|L_K^r f\|_\rho (1 + \sqrt{\gamma})(2 + r\gamma\lambda^{\frac{3}{2}-r} + \gamma^\eta)\lambda^r,$$

where  $P_\lambda$  is defined in eq. (23), and

$$(33) \quad \begin{aligned} \gamma &= \lambda^{-1} \|T - T_{\bar{\mathbf{x}}}\|, \\ \eta &= |r - \frac{1}{2}| - \lfloor |r - \frac{1}{2}| \rfloor. \end{aligned}$$

*Proof.* See Proposition 6 in [2]. □

**Proposition 7.** Let the operator  $\Omega_\lambda$  be defined by

$$(34) \quad \Omega_\lambda = \sqrt{T} G_\lambda(T_{\bar{\mathbf{x}}}) (T_{\bar{\mathbf{x}}} + \lambda)(T + \lambda)^{-\frac{1}{2}}.$$

Then, if  $\lambda \in (0, \kappa]$ , it holds

$$\|\Omega_\lambda\| \leq (1 + 2\sqrt{\gamma}) A,$$

with  $\gamma$  defined in eq. (33).

*Proof.* See Proposition 7 in [2]. □

We finally need the following probabilistic inequality based on a result of [14], see also Th. 3.3.4 of [18]. We report it without proof.

**Proposition 8.** Let  $(\Omega, \mathcal{F}, P)$  be a probability space and  $\xi$  be a random variable on  $\Omega$  taking value in a real separable Hilbert space  $\mathcal{K}$ . Assume that there are two positive constants  $H$  and  $\sigma$  such that

$$\begin{aligned} \|\xi(\omega)\|_{\mathcal{K}} &\leq \frac{H}{2} \quad \text{a.s.}, \\ \mathbb{E}[\|\xi\|_{\mathcal{K}}^2] &\leq \sigma^2, \end{aligned}$$

then, for all  $m \in \mathbb{N}$  and  $0 < \delta < 1$ ,

$$(35) \quad \mathbb{P}_{(\omega_1, \dots, \omega_m) \sim P^m} \left[ \left\| \frac{1}{m} \sum_{i=1}^m \xi(\omega_i) - \mathbb{E}[\xi] \right\|_{\mathcal{K}} \leq 2 \left( \frac{H}{m} + \frac{\sigma}{\sqrt{m}} \right) \log \frac{2}{\delta} \right] \geq 1 - \delta.$$

We are now ready to prove Theorem 1.

**Proof of Theorem 1.** Let us consider the expansion

$$\begin{aligned} \sqrt{T}(f_{\bar{\mathbf{z}}, \lambda} - f_{\mathcal{H}}) &= \sqrt{T}(G_\lambda(T_{\bar{\mathbf{x}}}) g_{\mathbf{z}} - f_\lambda^{\text{tr}}) + \sqrt{T}(f_\lambda^{\text{tr}} - f_{\mathcal{H}}) \\ &= \Omega_\lambda (T + \lambda)^{\frac{1}{2}} (f_{\bar{\mathbf{z}}, \lambda}^{\text{ls}} - f_{\bar{\mathbf{z}}', \lambda}^{\text{ls}}) + \sqrt{T}(G_\lambda(T_{\bar{\mathbf{x}}}) T_{\bar{\mathbf{x}}} - \text{Id}) f_\lambda^{\text{tr}} + \sqrt{T}(f_\lambda^{\text{tr}} - f_{\mathcal{H}}) \\ &= \Omega_\lambda \left( (T + \lambda)^{\frac{1}{2}} (f_{\bar{\mathbf{z}}, \lambda}^{\text{ls}} - f_\lambda^{\text{ls}}) + (T + \lambda)^{\frac{1}{2}} (f_\lambda^{\text{ls}} - \bar{f}_\lambda^{\text{ls}}) + (T + \lambda)^{\frac{1}{2}} (\bar{f}_\lambda^{\text{ls}} - f_{\bar{\mathbf{z}}', \lambda}^{\text{ls}}) \right) \\ &\quad + \sqrt{T}(G_\lambda(T_{\bar{\mathbf{x}}}) T_{\bar{\mathbf{x}}} - \text{Id}) f_\lambda^{\text{tr}} + \sqrt{T}(f_\lambda^{\text{tr}} - f_{\mathcal{H}}) \end{aligned}$$

where the operator  $\Omega_\lambda$  is defined by equation (34), the ideal RLS estimators are  $f_\lambda^{\text{ls}} = (T + \lambda)^{-1} T_M f_{\mathcal{H}}$  and  $\bar{f}_\lambda^{\text{ls}} = (T + \lambda)^{-1} T f_\lambda^{\text{tr}}$ , and  $f_{\bar{\mathbf{z}}', \lambda}^{\text{ls}} = (T_{\bar{\mathbf{x}}} + \lambda)^{-1} T_{\bar{\mathbf{x}}} f_\lambda^{\text{tr}}$  is the RLS estimator constructed by the training set

$$\bar{\mathbf{z}}' = ((\tilde{x}_1, f_\lambda^{\text{tr}}(\tilde{x}_1)) \dots, (\tilde{x}_m, f_\lambda^{\text{tr}}(\tilde{x}_m))).$$

Hence we get the following decomposition,

$$(36) \quad \|f_{\bar{\mathbf{z}}, \lambda} - f_{\mathcal{H}}\|_\rho \leq D(\bar{\mathbf{z}}, \lambda) \left( S^{\text{ls}}(\bar{\mathbf{z}}, \lambda) + R(\lambda) + \bar{S}^{\text{ls}}(\bar{\mathbf{z}}, \lambda) \right) + P(\bar{\mathbf{z}}, \lambda) + P^{\text{tr}}(\lambda),$$

with

$$\begin{aligned}
(37) \quad S^{\text{ls}}(\tilde{\mathbf{z}}, \lambda) &= \left\| (T + \lambda)^{\frac{1}{2}} (f_{\tilde{\mathbf{z}}, \lambda}^{\text{ls}} - f_{\lambda}^{\text{ls}}) \right\|_{\mathcal{H}}, \\
\bar{S}^{\text{ls}}(\tilde{\mathbf{z}}, \lambda) &= \left\| (T + \lambda)^{\frac{1}{2}} (\bar{f}_{\tilde{\mathbf{z}}, \lambda}^{\text{ls}} - \bar{f}_{\lambda}^{\text{ls}}) \right\|_{\mathcal{H}}, \\
D(\tilde{\mathbf{z}}, \lambda) &= \|\Omega_{\lambda}\|, \\
P(\tilde{\mathbf{z}}, \lambda) &= \left\| \sqrt{T} (G_{\lambda}(T_{\tilde{\mathbf{x}}}) T_{\tilde{\mathbf{x}}} - \text{Id}) f_{\lambda}^{\text{tr}} \right\|_{\mathcal{H}}, \\
P^{\text{tr}}(\lambda) &= \|f_{\lambda}^{\text{tr}} - f_{\mathcal{H}}\|_{\rho}, \\
R(\lambda) &= \left\| (T + \lambda)^{\frac{1}{2}} (\bar{f}_{\lambda}^{\text{ls}} - f_{\lambda}^{\text{ls}}) \right\|_{\mathcal{H}}.
\end{aligned}$$

Terms  $S^{\text{ls}}$  and  $\bar{S}^{\text{ls}}$  will be estimated by Proposition 4, terms  $P^{\text{tr}}$  and  $R$  by Proposition 2, term  $D$  by Proposition 7, and finally term  $P$  by Propositions 6, 3 and 2.

**Step 1:** Estimate of  $S^{\text{ls}}$ . Since  $L_{\kappa}^{\frac{1}{2}}$  is an isometry between  $\mathcal{L}^2(X, \rho_X)$  and  $\mathcal{H}$  (see [5]), we obtain

$$(38) \quad \|f_{\lambda}^{\text{ls}}\| \leq \left\| \sqrt{T} (T + \lambda)^{-1} \right\| \left\| L_{\kappa}^{\frac{1}{2}} f_{\mathcal{H}} \right\|_{\mathcal{H}} \leq \frac{\|f_{\mathcal{H}}\|_{\rho}}{\sqrt{\lambda}} \leq \frac{M}{\sqrt{\lambda}}.$$

Now, let  $\delta$  be an arbitrary real in  $(0, 1)$ . From the assumptions on  $\lambda_m$ , for large enough  $m$ , we have

$$\begin{aligned}
\lambda_m \sqrt{m} &\geq 4\kappa \log \frac{6}{\delta}, \\
\lambda_m &\leq \|T\|.
\end{aligned}$$

Hence, by Proposition 5, for large enough  $m$ , the assumptions of Proposition 4 are verified, and we get that with probability greater than  $1 - \delta$

$$\begin{aligned}
S^{\text{ls}}(\tilde{\mathbf{z}}_m, \lambda_m) &\leq 8 \left( M + \sqrt{\kappa \frac{m}{\tilde{m}_m}} \|f_{\lambda_m}^{\text{ls}}\|_{\mathcal{H}} \right) \left( \frac{2}{m} \sqrt{\frac{\kappa}{\lambda_m}} + \sqrt{\frac{\mathcal{N}(\lambda_m)}{m}} \right) \log \frac{6}{\delta} \\
(\text{Prop.5, eq. (38)}) &\leq 8M \left( 1 + \sqrt{\frac{\kappa}{\lambda_m}} \right) \frac{1}{\sqrt{m}} \left( 2\sqrt{\frac{\kappa}{\lambda_m m}} + \sqrt{\frac{\kappa}{\lambda_m}} \right) \log \frac{6}{\delta} \\
(\lambda_m \leq \kappa, m \geq 4) &\leq \frac{32M\kappa}{\lambda_m \sqrt{m}} \log \frac{6}{\delta} \rightarrow 0.
\end{aligned}$$

Hence it holds

$$(39) \quad \lim_{m \rightarrow \infty} S^{\text{ls}}(\tilde{\mathbf{z}}_m, \lambda_m) \stackrel{P}{=} 0.$$

**Step 2:** Estimate of  $\bar{S}^{\text{ls}}$ . This term can be estimated observing that  $\tilde{\mathbf{z}}'$  is a training set of  $\tilde{m}$  supervised samples drawn i.i.d. from the probability measure  $\rho'$  with marginal  $\rho_X$  and conditional  $\rho'_{|x}(y) = \delta(y - f_{\lambda}^{\text{tr}}(x))$ . Therefore the regression function induced by  $\rho'$  is  $f_{\rho'} = f_{\lambda}^{\text{tr}}$ , and the support of  $\rho'$  is included in  $X \times [-M', M']$ , with  $M' = \sup_{x \in X} f_{\rho'}(x) \leq \sqrt{\kappa} \|f_{\lambda}^{\text{tr}}\|_{\mathcal{H}}$ . Reasoning as in the analysis of  $S^{\text{ls}}$ , we obtain that, for every  $\delta \in (0, 1)$  and large enough  $m$ , with probability greater than  $1 - \delta$  it holds

$$\begin{aligned}
\bar{S}^{\text{ls}}(\tilde{\mathbf{z}}_m, \lambda_m) &\leq 8 \left( M' + \sqrt{\kappa} \|f_{\lambda_m}^{\text{tr}}\|_{\mathcal{H}} \right) \left( \frac{2}{\tilde{m}_m} \sqrt{\frac{\kappa}{\lambda_m}} + \sqrt{\frac{\mathcal{N}(\lambda_m)}{\tilde{m}_m}} \right) \log \frac{6}{\delta} \\
(\text{Prop.5}) &\leq 16\sqrt{\kappa} \|f_{\lambda_m}^{\text{tr}}\|_{\mathcal{H}} \frac{1}{\sqrt{m}} \left( 2\sqrt{\frac{\kappa}{\lambda_m m}} + \sqrt{\frac{\kappa}{\lambda_m}} \right) \log \frac{6}{\delta} \\
(m \geq 4) &\leq \frac{32\kappa}{\sqrt{\lambda_m m}} \left\| P_{\lambda_m} L_{\kappa}^{-\frac{1}{2}} P_{\lambda_m} f_{\mathcal{H}} \right\|_{\rho} \log \frac{6}{\delta} \leq \frac{32\kappa M}{\lambda_m \sqrt{m}} \log \frac{6}{\delta} \rightarrow 0.
\end{aligned}$$

Hence it holds

$$(40) \quad \lim_{m \rightarrow \infty} \bar{S}^{\text{ls}}(\tilde{\mathbf{z}}_m, \lambda_m) \stackrel{P}{=} 0.$$

**Step 3:** Estimate of  $P^{\text{tr}}$ . By definition (25),  $P^{\text{tr}}(\lambda) = a(\lambda)$ . Hence from eq. (27)

$$(41) \quad \lim_{m \rightarrow \infty} P^{\text{tr}}(\lambda_m) = \lim_{m \rightarrow \infty} a(\lambda_m) = \lim_{\lambda \rightarrow 0} a(\lambda) = 0,$$

where we used the assumption (6).

**Step 4:** Estimate of  $R$ . Since from the definitions of  $f_\lambda^{\text{ls}}$  and  $\bar{f}_\lambda^{\text{ls}}$ ,

$$R(\lambda) = \left\| (T + \lambda)^{-\frac{1}{2}} T (\bar{f}_\lambda^{\text{ls}} - f_\lambda^{\text{ls}}) \right\|_{\mathcal{H}} \leq \left\| \sqrt{T} (f_\lambda^{\text{tr}} - f_{\mathcal{H}}) \right\|_{\mathcal{H}} \leq P^{\text{tr}}(\lambda),$$

from (41) we get

$$(42) \quad \lim_{m \rightarrow \infty} R(\lambda_m) = 0.$$

**Step 5:** Estimate of  $D$ . In order to estimate  $D(\tilde{\mathbf{z}}, \lambda)$ , we have first to estimate the quantity  $\gamma = \gamma(\tilde{\mathbf{z}}, \lambda)$  (see definition (33)) appearing in the Proposition 7. Our estimate for  $\gamma(\tilde{\mathbf{z}}, \lambda)$  follows from Proposition 8 applied to the random variable  $\xi : X \rightarrow \mathcal{L}_{\text{HS}}(\mathcal{H})$  defined by

$$\xi(x)[\cdot] = \lambda^{-1} K_x \langle K_x, \cdot \rangle_{\mathcal{H}}.$$

We can set  $H = \frac{2\kappa}{\lambda}$  and  $\sigma = \frac{H}{2}$ , and obtain that, for every  $\delta \in (0, 1)$  and  $m \geq 4$ , with probability greater than  $1 - \delta$

$$\gamma(\tilde{\mathbf{z}}_m, \lambda) \leq \lambda^{-1} \|T - T_{\tilde{\mathbf{x}}}\|_{\text{HS}} \leq \frac{2}{\lambda} \left( \frac{2\kappa}{\tilde{m}_m} + \frac{\kappa}{\sqrt{\tilde{m}_m}} \right) \log \frac{2}{\delta} \leq 4 \frac{\kappa}{\lambda \sqrt{m}} \log \frac{2}{\delta} =: \epsilon(m, \lambda, \delta).$$

From the expression of  $\epsilon(m, \lambda, \delta)$  we see that, by the assumption (7), for every  $\delta \in (0, 1)$ ,

$$\lim_{m \rightarrow \infty} \epsilon(m, \lambda_m, \delta) = 0,$$

and hence,

$$(43) \quad \lim_{m \rightarrow \infty} \gamma(\tilde{\mathbf{z}}_m, \lambda_m) \stackrel{P}{=} 0.$$

Finally, from eq. (43) and Proposition 7 we find

$$(44) \quad D(\tilde{\mathbf{z}}_m, \lambda_m) \leq \left(1 + 2\sqrt{\gamma(\tilde{\mathbf{z}}_m, \lambda_m)}\right) A \stackrel{P}{\leq} 3A.$$

**Step 6:** Estimate of  $P$ . First, notice that by the definition (3), WLOG we can assume  $\bar{r} < \frac{1}{2}$ . Moreover by condition (6), we can assume  $m$  large enough that  $\lambda_m \leq \kappa$ . We consider the function  $R$  introduced by Proposition 3, and apply Proposition 6, with  $f = P_{\lambda_m} f_{\mathcal{H}}$  and  $r_m = R(\kappa^{-1} \lambda_m) \leq \bar{r}$ , getting

$$P(\tilde{\mathbf{z}}_m, \lambda_m) \leq B_{\bar{r}} \left(1 + \gamma(\tilde{\mathbf{z}}_m, \lambda_m)^{\frac{1}{2}}\right) \left(2 + r_m \gamma(\tilde{\mathbf{z}}_m, \lambda_m) + \gamma(\tilde{\mathbf{z}}_m, \lambda_m)^{\frac{1}{2} - r_m}\right) \kappa^{r_m} \|L_K^{-r_m} P_{\lambda_m} f_{\mathcal{H}}\|_{\rho} (\kappa^{-1} \lambda_m)^{r_m}.$$

This result together with eq. (43), and recalling that by Proposition 3 and assumption (6), the sequence  $\{r_m\}_m$  verifies the two conditions

$$\begin{aligned} \kappa^{r_m} \|L_K^{-r_m} P_{\lambda_m} f_{\mathcal{H}}\|_{\rho} &\leq 4M \quad \forall m, \\ \lim_{m \rightarrow \infty} (\kappa^{-1} \lambda_m)^{r_m} &= 0, \end{aligned}$$

proves that

$$(45) \quad \lim_{m \rightarrow \infty} P(\tilde{\mathbf{z}}_m, \lambda_m) \stackrel{P}{=} 0.$$

The proof of the Theorem is completed considering the limit  $m \rightarrow \infty$  of estimate (36), and using equations (39), (40), (41), (42), (44) and (45).  $\square$

**4.2.** Before showing the proof of Theorem 2, we state two propositions from [2] which describe properties of the functions  $f_\lambda^{\text{tr}}$  and  $f_\lambda^{\text{ls}}$  (defined in eq. (22) and eq. (32) respectively) when  $f_{\mathcal{H}} \in \text{Im } L_K^r$ .

**Proposition 9.** *Let  $f_{\mathcal{H}} \in \text{Im } L_K^r$  for some  $r > 0$ . Then, the following estimates hold,*

$$\begin{aligned} \|f_\lambda^{\text{tr}} - f_{\mathcal{H}}\|_\rho &\leq \lambda^r \|L_K^{-r} f_{\mathcal{H}}\|_\rho, \\ \|f_\lambda^{\text{tr}}\|_{\mathcal{H}} &\leq \begin{cases} \lambda^{-\frac{1}{2}+r} \|L_K^{-r} f_{\mathcal{H}}\|_\rho & \text{if } r \leq \frac{1}{2}, \\ \kappa^{-\frac{1}{2}+r} \|L_K^{-r} f_{\mathcal{H}}\|_\rho & \text{if } r > \frac{1}{2}. \end{cases} \end{aligned}$$

*Proof.* See Proposition 3 in [2].  $\square$

**Proposition 10.** *Let  $f_{\mathcal{H}} \in \text{Im } L_K^r$  for some  $r > 0$ . Then, the following estimates hold,*

$$\begin{aligned} \|f_\lambda^{\text{ls}} - f_{\mathcal{H}}\|_\rho &\leq \lambda^r \|L_K^{-r} f_{\mathcal{H}}\|_\rho, \quad \text{if } r \leq 1 \\ \|f_\lambda^{\text{ls}}\|_{\mathcal{H}} &\leq \begin{cases} \lambda^{-\frac{1}{2}+r} \|L_K^{-r} f_{\mathcal{H}}\|_\rho & \text{if } r \leq \frac{1}{2}, \\ \kappa^{-\frac{1}{2}+r} \|L_K^{-r} f_{\mathcal{H}}\|_\rho & \text{if } r > \frac{1}{2}. \end{cases} \end{aligned}$$

*Proof.* See Proposition 1 in [2].  $\square$

We are now ready to prove Theorem 2.

**Proof of Theorem 2.** We consider the same decomposition (see equations (36) and (37)) for  $\|f_{\bar{z},\lambda} - f_{\mathcal{H}}\|_\rho$  that we used in the proof of Theorem 1.

Terms  $S^{\text{ls}}$  and  $\bar{S}^{\text{ls}}$  will be estimated by Proposition 4, term  $D$  by Proposition 7, term  $P$  by Proposition 6 and finally terms  $P^{\text{tr}}$  and  $R$  by Proposition 9.

Let us begin with the estimates of  $S^{\text{ls}}$  and  $\bar{S}^{\text{ls}}$ . First observe that, by Proposition 5, it holds

$$\dot{\lambda} \leq \kappa^{-1} \|T\| \leq 1,$$

therefore, since by assumption  $\tilde{m} \geq m \dot{\lambda}^{-|2-2r-s|_++t_1} \geq m \dot{\lambda}^{-|1-2r|_++t_1}$ , we get,

$$\dot{\lambda} \tilde{m} \geq \dot{\lambda}^{-|1-2r|_++1+t_1} m \geq \dot{\lambda}^{2r+t_1} m.$$

Moreover, by eq. (8) and definition (5), we find

$$\dot{\lambda}^{2r+t_1} m = 16q^{2r+s+t_1} D_s^2 \dot{\lambda}^{-s} \log^2 \frac{6}{\delta} \geq 16\mathcal{N}(\lambda) \log^2 \frac{6}{\delta},$$

hence the hypothesis (30) in the text of Proposition 4 is verified.

Regarding the estimate of  $S^{\text{ls}}$ . Applying Proposition 4 and recalling that by assumption  $\tilde{m} \geq m \dot{\lambda}^{-|2-2r-s|_++t_1} \geq m \dot{\lambda}^{-|1-2r|_++t_1}$  and from Proposition 10,  $\sqrt{\kappa} \|f_\lambda^{\text{ls}}\|_{\mathcal{H}} \leq C_r \dot{\lambda}^{-|\frac{1}{2}-r|_+}$ , we get that with probability greater than  $1 - \delta$

$$\begin{aligned} (46) \quad S^{\text{ls}}(\bar{z}, \lambda) &\leq 8 \left( M + \sqrt{\frac{m}{\tilde{m}}} C_r \dot{\lambda}^{-|\frac{1}{2}-r|_+} \right) \left( \frac{2}{m} \sqrt{\frac{\kappa}{\lambda}} + \sqrt{\frac{\mathcal{N}(\lambda)}{m}} \right) \log \frac{6}{\delta} \\ &\leq 8(M + \dot{\lambda}^{-\frac{t_1}{2}} C_r) \frac{1}{\sqrt{m}} \left( \frac{2}{\sqrt{m\lambda}} + \frac{D_s}{\sqrt{\lambda^s}} \right) \log \frac{6}{\delta} \\ (eq. (8)) &= 2q^{-r-\frac{s}{2}-\frac{t_1}{2}} (M + C_r) \dot{\lambda}^r \left( 1 + \frac{\dot{\lambda}^{-\frac{1}{2}(1-2r-2s-t_1)}}{2q^{r+\frac{s}{2}+\frac{t_1}{2}} D_s^2 \log \frac{6}{\delta}} \right) \\ (t_1 \geq 0, q \geq 1) &\leq 3q^{-r-\frac{s}{2}-\frac{t_1}{2}} (M + C_r) \dot{\lambda}^{r-\frac{t_2}{2}} \\ (q \geq 1) &\leq 3(M + C_r) \dot{\lambda}^{r-\frac{t_2}{2}}. \end{aligned}$$

The term  $\bar{S}^{\text{ls}}$  can be estimated observing that  $\bar{\mathbf{z}}'$  is a training set of  $\tilde{m}$  supervised samples drawn i.i.d. from the probability measure  $\rho'$  with marginal  $\rho_X$  and conditional  $\rho'_x(y) = \delta(y - f_\lambda^{\text{tr}}(x))$ . Therefore the regression function induced by  $\rho'$  is  $f_{\rho'} = f_\lambda^{\text{tr}}$ , and the support of  $\rho'$  is included in  $X \times [-M', M']$ , with  $M' = \sup_{x \in X} f_{\rho'}(x) \leq \sqrt{\kappa} \|f_\lambda^{\text{tr}}\|_{\mathcal{H}}$ . Again applying Proposition 4, we obtain that with probability greater than  $1 - \delta$  it holds

$$\begin{aligned}
(47) \quad \bar{S}^{\text{ls}}(\bar{\mathbf{z}}, \lambda) &\leq 8 \left( M' + \sqrt{\kappa} \|f_\lambda^{\text{tr}}\|_{\mathcal{H}} \right) \left( \frac{2}{\tilde{m}} \sqrt{\frac{\kappa}{\lambda}} + \sqrt{\frac{\mathcal{N}(\lambda)}{\tilde{m}}} \right) \log \frac{6}{\delta} \\
&\leq 16\sqrt{\kappa} \|f_\lambda^{\text{tr}}\|_{\mathcal{H}} \left( \frac{2}{\tilde{m}} \sqrt{\frac{\kappa}{\lambda}} + \sqrt{\frac{\mathcal{N}(\lambda)}{\tilde{m}}} \right) \log \frac{6}{\delta} \\
(\text{Prop.9}) \quad &\leq 16\sqrt{\frac{m}{\tilde{m}}} C_r \lambda^{-|\frac{1}{2}-r|_+} \left( \frac{2}{m} \sqrt{\frac{\kappa}{\lambda}} + \sqrt{\frac{\mathcal{N}(\lambda)}{m}} \right) \log \frac{6}{\delta} \\
&\leq 16C_r \frac{\lambda^{-\frac{t_1}{2}}}{\sqrt{m}} \left( \frac{2}{\sqrt{m\lambda}} + \frac{D_s}{\sqrt{\lambda^s}} \right) \log \frac{6}{\delta} \\
(\text{eq. (8)}) \quad &= 4q^{-r-\frac{s}{2}-\frac{t_1}{2}} C_r \lambda^r \left( 1 + \frac{\lambda^{-\frac{1}{2}(1-2r-2s-t_1)}}{2q^{r+\frac{s}{2}+\frac{t_1}{2}} D_s^2 \log \frac{6}{\delta}} \right) \\
(t_1 \geq 0, q \geq 1) \quad &\leq 6q^{-r-\frac{s}{2}-\frac{t_1}{2}} C_r \lambda^{r-\frac{t_2}{2}} \\
(q \geq 1) \quad &\leq 6C_r \lambda^{r-\frac{t_2}{2}}.
\end{aligned}$$

In order to get an upper bound for  $D$  and  $P$ , we have first to estimate the quantity  $\gamma = \gamma(\bar{\mathbf{z}}, \lambda)$  (see definition (33)) appearing in the Propositions 6 and 7. Our estimate for  $\gamma(\bar{\mathbf{z}}, \lambda)$  follows from Proposition 8 applied to the random variable  $\xi : X \rightarrow \mathcal{L}_{\text{HS}}(\mathcal{H})$  defined by

$$\xi(x)[\cdot] = \lambda^{-1} K_x \langle K_x, \cdot \rangle_{\mathcal{H}}.$$

We can set  $H = \frac{2\kappa}{\lambda}$  and  $\sigma = \frac{H}{2}$ , and obtain that with probability greater than  $1 - \delta$

$$\begin{aligned}
\gamma(\bar{\mathbf{z}}, \lambda) &\leq \lambda^{-1} \|T - T_{\bar{\mathbf{x}}}\|_{\text{HS}} \leq \frac{2}{\lambda} \left( \frac{2\kappa}{\tilde{m}} + \frac{\kappa}{\sqrt{\tilde{m}}} \right) \log \frac{2}{\delta} \leq 4 \frac{1}{\lambda \sqrt{\tilde{m}}} \log \frac{2}{\delta} \\
&\leq 4 \frac{\lambda^{|1-r-\frac{s}{2}|_+ - 1 - \frac{t_1}{2}}}{\sqrt{m}} \log \frac{2}{\delta} \leq \lambda^{|1-r-\frac{s}{2}|_+ - (1-r-\frac{s}{2})} \leq \lambda^{|r+\frac{s}{2}-1|_+} \leq 1,
\end{aligned}$$

where we used the assumption  $\tilde{m} \geq 4 \vee m \lambda^{-|2-2r-s|_+ + t_1}$  and the expression for  $\lambda$  in the text of the Theorem.

Hence, since  $\gamma(\bar{\mathbf{z}}, \lambda) \leq \lambda^{|r+\frac{s}{2}-1|_+}$ , from Proposition 7 we get

$$(48) \quad D(\bar{\mathbf{z}}, \lambda) \leq (1 + 2\sqrt{\gamma})A \leq 3A,$$

and from Proposition 6

$$\begin{aligned}
(49) \quad P(\bar{\mathbf{z}}, \lambda) &\leq B_r C_r (1 + \sqrt{\gamma}) (2 + r\gamma \lambda^{\frac{3}{2}-r} + \gamma^\eta) \lambda^r \\
&\leq 2B_r C_r (3 + r\gamma \lambda^{\frac{3}{2}-r}) \lambda^r \\
&\leq 2B_r C_r (3 + r\lambda^{|r+\frac{s}{2}-1|_+ + \frac{3}{2}-r}) \lambda^r \\
&\leq 2B_r C_r (3 + r\lambda^{\frac{s+1}{2}}) \lambda^r \leq 2B_r C_r (3 + r) \lambda^r.
\end{aligned}$$

Regarding terms  $P^{\text{tr}}$  and  $R$ . From Proposition 9 we get

$$(50) \quad P^{\text{tr}}(\lambda) \leq C_r \lambda^r,$$

and hence,

$$(51) \quad \begin{aligned} R(\lambda) &= \left\| (T + \lambda)^{-\frac{1}{2}} T (f_{\lambda}^{\text{ls}} - f_{\lambda}^{\text{ls}}) \right\|_{\mathcal{H}} \\ &\leq \left\| \sqrt{T} (f_{\lambda}^{\text{ls}} - f_{\lambda}^{\text{ls}}) \right\|_{\mathcal{H}} \leq P^{\text{tr}} \leq C_r \dot{\lambda}^r. \end{aligned}$$

The proof is completed by plugging inequalities (46), (47), (48), (49), (50) and (51) in (36), recalling the expression for  $\dot{\lambda}$ .  $\square$

## 5. PROOF OF THEOREM 3

The following result is due to [13], adapted to a suitable form used in this paper.

**Proposition 11.** *Let  $\{X_i\}_{i=1}^n$  be a set of real valued i.i.d. random variables with mean  $\mu$ ,  $|X_i| \leq B$  and  $\mathbb{E}[(X_i - \mu)^2] \leq \sigma^2$ , for all  $i \in \{1, \dots, n\}$ . Then for arbitrary  $\alpha > 0$ ,  $\epsilon > 0$ ,*

$$(52) \quad \mathbb{P} \left[ \frac{1}{n} \sum_{i=1}^n X_i - \mu \geq \alpha \sigma^2 + \epsilon \right] \leq e^{-\frac{6n\alpha\epsilon}{3+4\alpha B}},$$

and

$$(53) \quad \mathbb{P} \left[ \mu - \frac{1}{n} \sum_{i=1}^n X_i \geq \alpha \sigma^2 + \epsilon \right] \leq e^{-\frac{6n\alpha\epsilon}{3+4\alpha B}}.$$

*Proof.* It suffices to prove the one side inequality (52). For any  $s > 0$ ,

$$\begin{aligned} &\mathbb{P} \left[ \frac{1}{n} \sum_{i=1}^n X_i - \mu \geq \alpha \sigma^2 + \epsilon \right] = \mathbb{P} \left[ e^{\frac{s}{n} \sum_{i=1}^n (X_i - \mu)} \geq e^{s(\alpha \sigma^2 + \epsilon)} \right] \\ &\leq e^{-s\epsilon - s\alpha \sigma^2} \mathbb{E} e^{\frac{s}{n} \sum_{i=1}^n (X_i - \mu)}, \quad \text{by Markov inequality} \\ &= e^{-s\epsilon - s\alpha \sigma^2} \prod_{i=1}^n \mathbb{E} e^{\frac{s}{n} (X_i - \mu)}, \quad \text{by independence of } X_i \end{aligned}$$

Denote  $Z_i = X_i - \mu$ ,  $t = s/n$  and  $B_1 = 2B$ . Thus for those  $s$  such that  $sB < 3n/2$  (or equivalently  $B_1 t/3 < 1$ ),

$$\begin{aligned} \mathbb{E} e^{tZ_i} &= 1 + \sum_{k=1}^{\infty} \frac{t^k}{k!} \mathbb{E}[Z_i^k] \leq 1 + 0 + \sum_{k=2}^{\infty} \frac{t^k}{k!} B_1^{k-2} \sigma^2 \leq 1 + \frac{t^2 \sigma^2}{2} \sum_{k=0}^{\infty} \left( \frac{B_1 t}{3} \right)^k \\ &= 1 + \frac{3t^2 \sigma^2}{6 - 4Bt} \leq \exp \left( \frac{3t^2 \sigma^2}{6 - 4Bt} \right) = \exp \left( \frac{3s^2 \sigma^2}{n^2(6 - 4sB/n)} \right) \end{aligned}$$

whence

$$e^{-s\epsilon - s\alpha \sigma^2} \prod_{i=1}^n \mathbb{E} e^{\frac{s}{n} (X_i - \mu)} \leq e^{-s\epsilon} \exp \left\{ \frac{s\sigma^2}{n} \left( \frac{3s}{6 - 4sB/n} - n\alpha \right) \right\}.$$

Setting  $s = s_0 = \frac{6\alpha n}{3 + 4\alpha B}$  (one can check that  $s_0 B = \frac{6n\alpha B}{3 + 4\alpha B} < 3n/2$ ), we have  $\frac{3s_0}{6 - 4s_0 B/n} - n\alpha = 0$  and thus *r.h.s.*  $\leq e^{-s_0 \epsilon} = \exp \left( -\frac{6n\alpha\epsilon}{3 + 4\alpha B} \right)$ , which gives estimate (52).  $\square$

We are now ready to prove Theorem 3.

**Proof of Theorem 3.** The strategy of the proof is the following. Define

$$(54) \quad \lambda_m^* = \operatorname{argmin}_{\lambda \in \Lambda_m} \int_Z (T_M f_{\mathbf{z}, \lambda}(x) - y)^2 d\rho.$$

Notice that, since for every  $f \in \mathcal{L}^2(X, \rho_X)$ ,

$$\int_Z (f(x) - y)^2 d\rho = \|f - f_{\rho}\|_{\rho}^2 + \int_Z (f_{\rho}(x) - y)^2 d\rho,$$

definition (54) of  $\lambda_m^*$ , is equivalent to

$$\lambda_m^* = \operatorname{argmin}_{\lambda \in \Lambda_m} \|T_M f_{\tilde{\mathbf{z}}, \lambda} - f_\rho\|_\rho.$$

Now, from the equality above, the assumption of the Theorem, and recalling that  $f_\rho(x) \in Y \subset [-M, M]$ , we get that with probability greater than  $1 - \delta$  it holds

$$(55) \quad \|T_M f_{\tilde{\mathbf{z}}, \lambda_m^*} - f_\rho\|_\rho \leq \|T_M f_{\tilde{\mathbf{z}}, \lambda_m} - f_\rho\|_\rho \leq \|f_{\tilde{\mathbf{z}}, \lambda_m} - f_\rho\|_\rho \leq \epsilon.$$

We claim that for every  $\tilde{\mathbf{z}}, \lambda > 0$  and  $\delta \in (0, 1)$ , with probability greater than  $1 - \delta$  over the probability measure  $\rho^{m^\vee}$ , it holds

$$(56) \quad \|f_{\mathbf{z}^{\text{tot}}} - f_\rho\|_\rho^2 \leq 2 \|T_M f_{\tilde{\mathbf{z}}, \lambda_m^*} - f_\rho\|_\rho^2 + \frac{80M^2}{m^\vee} \log \frac{2|\Lambda_m|}{\delta}.$$

Estimates (55) and (56) together will complete the proof of the Theorem.

We now proceed to proving eq. (56). For  $i = 1, \dots, m^\vee$ , let us define the random variables

$$\xi_i^\lambda = (T_M f_{\tilde{\mathbf{z}}, \lambda}(x_i^\vee) - y_i^\vee)^2 - (f_\rho(x_i^\vee) - y_i^\vee)^2.$$

Clearly

$$\begin{aligned} |\xi_i^\lambda| &\leq 4M^2, \\ \mathbb{E}[\xi_i^\lambda] &= \int_Z (T_M f_{\tilde{\mathbf{z}}, \lambda}(x) - y)^2 d\rho - \int_Z (f_\rho(x) - y)^2 d\rho = \|T_M f_{\tilde{\mathbf{z}}, \lambda} - f_\rho\|_\rho^2, \\ \mathbb{E}[(\xi_i^\lambda)^2] &= \int_Z (T_M f_{\tilde{\mathbf{z}}, \lambda}(x) - f_\rho(x))^2 (T_M f_{\tilde{\mathbf{z}}, \lambda}(x) + f_\rho(x) - 2y)^2 d\rho \\ &\leq 16M^2 \|T_M f_{\tilde{\mathbf{z}}, \lambda} - f_\rho\|_\rho^2 = 16M^2 \mathbb{E}[\xi_i^\lambda]. \end{aligned}$$

Hence, using Proposition 11 with  $X_i = \xi_i^\lambda$ ,  $\mu = \mathbb{E}[\xi_i^\lambda]$ ,  $B = 4M^2$  and  $\sigma^2 = \mathbb{E}[(\xi_i^\lambda)^2] \leq 16M^2 \mu$ , we obtain that for all  $\lambda \in \Lambda_m$  with probability greater than  $1 - \delta$ ,

$$\frac{1}{m^\vee} \sum_{i=1}^{m^\vee} \xi_i^\lambda \leq (1 + \alpha') \mathbb{E}[\xi_i^\lambda] + \epsilon,$$

and

$$\mathbb{E}[\xi_i^\lambda] \leq \frac{1}{1 - \alpha'} \left( \frac{1}{m^\vee} \sum_{i=1}^{m^\vee} \xi_i^\lambda \right) + \frac{\epsilon}{1 - \alpha'},$$

where  $\alpha' = 16\alpha M^2$  and  $\epsilon = \frac{3 + \alpha'}{6\alpha m^\vee} \log \frac{2|\Lambda_m|}{\delta}$ . Therefore

$$\begin{aligned} \|f_{\mathbf{z}^{\text{tot}}} - f_\rho\|_\rho^2 &= \mathbb{E}[\xi_i^{\lambda_{\mathbf{z}^\vee}}] \leq \frac{1}{1 - \alpha'} \left( \frac{1}{m^\vee} \sum_{i=1}^{m^\vee} \xi_i^{\lambda_{\mathbf{z}^\vee}} \right) + \frac{\epsilon}{1 - \alpha'} \\ &\leq \frac{1}{1 - \alpha'} \left( \frac{1}{m^\vee} \sum_{i=1}^{m^\vee} \xi_i^{\lambda_m^*} \right) + \frac{\epsilon}{1 - \alpha'} \\ &\leq \frac{1 + \alpha'}{1 - \alpha'} \mathbb{E}[\xi_i^{\lambda_m^*}] + \frac{2\epsilon}{1 - \alpha'} \\ &= \frac{1 + \alpha'}{1 - \alpha'} \|T_M f_{\tilde{\mathbf{z}}, \lambda_m^*} - f_\rho\|_\rho^2 + \frac{2\epsilon}{1 - \alpha'}. \end{aligned}$$

Setting  $\alpha = 1/(48M^2)$ , this gives  $\alpha' = 1/3$  and

$$\|f_{\mathbf{z}^{\text{tot}}} - f_\rho\|_\rho^2 \leq 2 \|T_M f_{\tilde{\mathbf{z}}, \lambda_m^*} - f_\rho\|_\rho^2 + \frac{80M^2}{m^\vee} \log \frac{2|\Lambda_m|}{\delta},$$

which proves eq. (56), as desired.  $\square$

## ACKNOWLEDGEMENTS

The authors wish to thank Peter Bickel for an inspiring discussion leading to this work, Bo Li for pointing out reference [10] and E. De Vito, T. Poggio, L. Rosasco, S. Smale and A. Verri for useful discussions and suggestions.

## REFERENCES

- [1] F. Bauer, S. Pereverzev, and L. Rosasco. On regularization algorithms in learning theory. Preprint, 2005.
- [2] A. Caponnetto. Optimal rates for regularization operators in learning theory. Preprint, 2006.
- [3] A. Caponnetto and E. De Vito. Fast rates for regularized least-squares algorithm. Technical report, Massachusetts Institute of Technology, Cambridge, MA, April 2005. CBCL Paper#248/AI Memo#2005-013.
- [4] A. Caponnetto and E. De Vito. Optimal rates for regularized least-squares algorithm. 2005. *to appear in* Foundations of Computational Mathematics.
- [5] F. Cucker and S. Smale. On the mathematical foundations of learning. *Bull. Amer. Math. Soc. (N.S.)*, 39(1):1–49 (electronic), 2002.
- [6] E. De Vito, A. Caponnetto, and L. Rosasco. Model selection for regularized least-squares algorithm in learning theory. *Foundation of Computational Mathematics*, 5(1):59–85, February 2005.
- [7] E. De Vito, L. Rosasco, and A. Caponnetto. Discretization error analysis for tikhonov regularization. *to appear in* *Analisis and Applications*, 2005.
- [8] E. De Vito, L. Rosasco, A. Caponnetto, U. De Giovannini, and F. Odone. Learning from examples as an inverse problem. *Journal of Machine Learning Research*, 6:883–904, 2005.
- [9] E. De Vito, L. Rosasco, and A. Verri. Spectral methods for regularization in learning theory. Preprint, 2005.
- [10] S. Dudoit and M.J. van der Laan. Asymptotics of cross-validated risk estimation in estimator selection and performance assessment. *Statistical Methodology*, 2(2):131–154, 2001.
- [11] H. W. Engl, M. Hanke, and A. Neubauer. *Regularization of inverse problems*, volume 375 of *Mathematics and its Applications*. Kluwer Academic Publishers Group, Dordrecht, 1996.
- [12] T. Evgeniou, M. Pontil, and T. Poggio. Regularization networks and support vector machines. *Adv. Comp. Math.*, 13:1–50, 2000.
- [13] M. Hamers and M. Kohler. A bound on the expected maximal deviations of sample averages from their means. Preprint 2001-9, Mathematical Institute A, University of Stuttgart.
- [14] I. F. Pinelis and A. I. Sakhanenko. Remarks on inequalities for probabilities of large deviations. *Theory Probab. Appl.*, 30(1):143–148, 1985.
- [15] S. Smale and D. Zhou. Learning theory estimates via integral operators and their approximations. Preprint, Toyota Technological Institute, Chicago, 2005.
- [16] S. Smale and D.X. Zhou. Shannon sampling II : Connections to learning theory. *Appl. Comput. Harmonic Anal.* to appear.
- [17] I. Steinwart. On the influence of the kernel on the consistency of support vector machines. *Journal of Machine Learning Research*, 2:67–93, 2002.
- [18] V. Yurinsky. *Sums and Gaussian vectors*, volume 1617 of *Lecture Notes in Mathematics*. Springer-Verlag, Berlin, 1995.