

Nr. 22
12. January 2016

Leopold-Franzens-Universität Innsbruck



Preprint-Series: Department of Mathematics - Applied Mathematics

Nyström type subsampling analyzed as a regularized projection

Galyna Kriukova, Sergiy Pereverzyev Jr., Pavlo Tkackenko



APPLIEDMATHEMATICS

Technikerstraße 13 - 6020 Innsbruck - Austria
Tel.: +43 512 507 53803 Fax: +43 512 507 53898
<https://applied-math.uibk.ac.at>

Nyström type subsampling analyzed as a regularized projection

Galyna Kriukova¹ Sergiy Pereverzyev Jr.²
Pavlo Tkachenko¹

January 12, 2016

Abstract

In the statistical learning theory the Nyström type subsampling methods are considered as tools for dealing with big data. In this paper we consider Nyström subsampling as a special form of the projected Lavrentiev regularization, and study it using the approaches developed in the regularization theory. As a result, we prove that the same capacity independent learning rates that are guaranteed for standard algorithms running with quadratic computational complexity can be obtained with subquadratic complexity by the Nyström subsampling approach, provided that the subsampling size is chosen properly. We propose a priori rule for choosing the subsampling size and a posteriori strategy for dealing with uncertainty in the choice of it. The theoretical results are illustrated by numerical experiments.

1 Introduction

Regularization based kernel methods, such as kernel ridge regression (KRR), provide an effective framework for the supervised learning [12, 13]. However,

¹Johann Radon Institute for Computational and Applied Mathematics, Austrian Academy of Sciences, Altenbergerstraße 69, A-4040 Linz, Austria

²Institute of Mathematics, University of Innsbruck, Technikestraße 13, A-6020 Innsbruck, Austria

a standard implementation of these methods is infeasible when dealing with the so-called ‘‘Big Data’’.

The Big Data concept can be considered from different points of view. Here, by ‘‘Big Data’’, we mean data sets exceeding the computational capacity of conventional learning systems. For example, in KRR, one receives a training data set \mathbf{z} of N samples of the form $\mathbf{z} = \{(x_i, y_i)\}_{i=1}^N$, where each input $x_i \in X \subset \mathbb{R}^d$ is related to the output $y_i \in \mathbb{R}$ by an unknown target function $f^*: X \rightarrow \mathbb{R}$, and the goal is to approximate this function by the minimizer $f_{\mathbf{z}}^\alpha$ of the regularized empirical risk functional:

$$T_{\mathbf{z}}^\alpha(f) := \frac{1}{|\mathbf{z}|} \sum_{(x_i, y_i) \in \mathbf{z}} (f(x_i) - y_i)^2 + \alpha \|f\|_{\mathcal{H}_{\mathbf{K}}}^2. \quad (1)$$

Here, $\mathcal{H}_{\mathbf{K}}$ denotes the reproducing kernel Hilbert space (RKHS) generated by a kernel $\mathbf{K}: X \times X \rightarrow \mathbb{R}$, $|\mathbf{z}| = N$, and α is a regularization parameter.

By the representer theorem for RKHS [6], the minimizer of (1) is equal to

$$f_{\mathbf{z}}^\alpha = \sum_{x_i: (x_i, y_i) \in \mathbf{z}} c_i \mathbf{K}(\cdot, x_i),$$

where $\mathbf{c} = (c_i)_{i=1}^{|\mathbf{z}|} = (\mathbf{K} + \alpha |\mathbf{z}| \mathbf{I})^{-1} \mathbf{Y}$, $\mathbf{Y} = (y_i)_{i=1}^{|\mathbf{z}|}$, \mathbf{I} is the $|\mathbf{z}| \times |\mathbf{z}|$ diagonal identity matrix, and \mathbf{K} denotes the $|\mathbf{z}| \times |\mathbf{z}|$ kernel matrix with entries $\mathbf{K}_{ij} = \mathbf{K}(x_i, x_j)$.

Now, it is clear that KRR will suffer from at least quadratic computational complexity $\mathcal{O}(N^2)$ in the number of observations N , as this is the complexity of computing the kernel matrix \mathbf{K} . In the Big Data setting, where N is large, this is not acceptable. Therefore, learning schemes have been designed to avoid the computation of the exact minimizers $f_{\mathbf{z}}^\alpha$.

One family of such schemes, which we broadly refer to as the Nyström type subsampling, consists of methods replacing the kernel matrix \mathbf{K} with a smaller matrix obtained by column subsampling [15, 16]. This can also be interpreted as a restriction of the minimization of $T_{\mathbf{z}}^\alpha(f)$ to the space

$$\mathcal{H}_{\mathbf{K}}^{\mathbf{z}^\nu} := \{f | f = \sum_{x_i: (x_i, y_i) \in \mathbf{z}^\nu} c_i \mathbf{K}(\cdot, x_i), c_i \in \mathbb{R}\},$$

where $\mathbf{z}^\nu \subset \mathbf{z}$, and $|\mathbf{z}^\nu| = N_\nu \ll N$.

It can be shown [11] that the minimizer $f_{\mathbf{z}, \mathbf{z}^\nu}^\alpha$ of $T_{\mathbf{z}}^\alpha(f)$ over the space $\mathcal{H}_{\mathbf{K}}^{\mathbf{z}^\nu}$ has the form

$$f_{\mathbf{z}, \mathbf{z}^\nu}^\alpha = (\mathbf{P}_{\mathbf{z}^\nu} \mathbf{S}_{\mathbf{z}}^* \mathbf{S}_{\mathbf{z}} \mathbf{P}_{\mathbf{z}^\nu} + \alpha \mathbf{I})^{-1} \mathbf{P}_{\mathbf{z}^\nu} \mathbf{S}_{\mathbf{z}}^* \mathbf{Y}, \quad (2)$$

where $P_{\mathbf{z}^\nu}$ is the orthogonal projection operator with range $\mathcal{H}_K^{\mathbf{z}^\nu}$, $S_{\mathbf{z}}: \mathcal{H}_K \rightarrow \mathbb{R}^{|\mathbf{z}|}$ is the sampling operator, $S_{\mathbf{z}}f = (f(x_1), f(x_2), \dots, f(x_N))$, $x_i: (x_i, y_i) \in \mathbf{z}$, and $S_{\mathbf{z}}^*: \mathbb{R}^{|\mathbf{z}|} \rightarrow \mathcal{H}_K$ is the adjoint of $S_{\mathbf{z}}$. If the norm $\|\cdot\|_{\mathbb{R}^{|\mathbf{z}|}}$ is defined as $|\mathbf{z}|^{-1}$ times the Euclidean norm, then

$$S_{\mathbf{z}}^* \mathbf{u}(\cdot) = |\mathbf{z}|^{-1} \sum_{i=1}^{|\mathbf{z}|} u_i K(\cdot, x_i), \quad \mathbf{u} = (u_1, u_2, \dots, u_N).$$

Observe that the computational complexity of the minimization of $T_{\mathbf{z}}^\alpha(f)$ over the space $\mathcal{H}_K^{\mathbf{z}^\nu}$ is of order $\mathcal{O}(|\mathbf{z}| \cdot |\mathbf{z}^\nu|^2) = \mathcal{O}(N \cdot N_\nu^2)$. Of course, N_ν is expected to increase with N , such that a linear complexity in N seems impossible. Therefore, the main question about the Nyström type subsampling is the following: how big should N_ν be to incur no loss of the performance compared to the full kernel matrix \mathbf{K} ; or, that is the same, is it possible to realize the Nyström approach with a complexity that is subquadratic in the number of observations N without losing the performance?

A positive answer to this question has been recently given in [1, 11]. However, in [1], the error analysis is derived in a fixed design regression setting, such that x_i , $i = 1, 2, \dots, |\mathbf{z}|$, are assumed to be uniformly sampled, for example. The study [11] extends the results of [1] to a general statistical learning setting. At the same time, the analysis of [11] is fairly technical and lengthy. In particular, it is based on the assumptions describing the capacity of the hypothesis space \mathcal{H}_K with respect to the unknown distribution ρ_X from which $\{x_i\}_{i=1}^{|\mathbf{z}|}$ is assumed to be sampled.

In the present study, we are going to analyze the so-called plain Nyström approach as a particular case of the regularized projection scheme. Therefore, we will use some arguments developed in the regularization theory for analyzing regularized projection approximations [9, 10]. Instead of the assumption on the capacity of the solution space, these arguments rely on the assumption on the approximation power of the projection method induced by the projector such as $P_{\mathbf{z}^\nu}$ in (2). For the purpose of our study, the arguments developed in [9, 10] should be accompanied by the ones that take into account that in the context of learning, the regularized projection schemes, such as (2), operate only with noisy versions of the operators describing the learning tasks.

An analysis incorporating the above mentioned arguments is presented in the next section. Unlike [11], it gives capacity independent learning rates for the Nyström type subsampling. Moreover, it indicates a rather general

a priori choice of the subsampling size $|\mathbf{z}^\nu|$ that allows a subquadratic complexity without loss of the performance. Such a priori choice of $|\mathbf{z}^\nu|$ requires a knowledge of the regularity of the unknown target function with respect to \mathbf{K} and ρ_X . In Section 3, we consider a situation when such a priori knowledge is not accurate, and may lead to uncertain parameter $|\mathbf{z}^\nu|$. In Section 4, we discuss some simulations illustrating our theoretical results.

2 Approximation power, regularity and learning rate

A training data set $\mathbf{z} = \{(x_i, y_i)\}_{i=1}^{|\mathbf{z}|}$ is assumed to be sampled from the so-called sample space $Z = X \times Y$ endowed with a fixed but unknown probability distribution ρ , which can be factorized as $\rho(x, y) = \rho(y|x)\rho_X(x)$, where $\rho(\cdot|x)$ is the conditional distribution on $Y \subset \mathbb{R}$ given $x \in X$, and ρ_X is the so-called marginal distribution, from which the set of inputs $\{x_i\}_{i=1}^{|\mathbf{z}|}$ is supposed to be sampled.

A common assumption to simplify analysis is that $Y = [-D, D]$ for some $D > 0$. A weaker condition can be found in [3].

Given a training set $\mathbf{z} \subset Z$, the goal is to find an estimate $f = f_{\mathbf{z}}$ with a small expected risk

$$\mathcal{E}(f) = \int_{X \times Y} (f(x) - y)^2 d\rho(x, y).$$

Once we choose $\mathcal{H}_{\mathbf{K}}$ as the so-called hypothesis space, the best possible risk value is clearly

$$\inf_{f \in \mathcal{H}_{\mathbf{K}}} \mathcal{E}(f).$$

As in [11], we assume that there exists $f^\dagger \in \mathcal{H}_{\mathbf{K}}$ such that

$$\mathcal{E}(f^\dagger) = \min_{f \in \mathcal{H}_{\mathbf{K}}} \mathcal{E}(f).$$

To formulate our further assumptions we need some operators, which are traditionally used in the context of regression learning. At first we consider the space $L_2(X, \rho_X)$ of square integrable functions with respect to ρ_X equipped with the usual norm $\|\cdot\|_\rho = \|\cdot\|_{L_2(X, \rho_X)}$. It is well-known [5] that for $f, f^\dagger \in \mathcal{H}_{\mathbf{K}}$ we have

$$\mathcal{E}(f) - \mathcal{E}(f^\dagger) = \|f - f^\dagger\|_\rho^2. \quad (3)$$

It is also known that if the kernel \mathbf{K} is bounded then $\mathcal{H}_{\mathbf{K}}$ is continuously embedded in $L_2(X, \rho_X)$, such that the canonical embedding operator $\mathbf{J}_{\mathbf{K}} : \mathcal{H}_{\mathbf{K}} \rightarrow L_2(X, \rho_X)$ is continuous. Then we consider the adjoint operator $\mathbf{J}_{\mathbf{K}}^* : L_2(X, \rho_X) \rightarrow \mathcal{H}_{\mathbf{K}}$,

$$\mathbf{J}_{\mathbf{K}}^* f(\cdot) = \int_X \mathbf{K}(\cdot, x) f(x) d\rho_X(x),$$

and the covariance operator $\mathbf{C} = \mathbf{J}_{\mathbf{K}}^* \mathbf{J}_{\mathbf{K}} : \mathcal{H}_{\mathbf{K}} \rightarrow \mathcal{H}_{\mathbf{K}}$,

$$\mathbf{C}f(\cdot) = \int_X \mathbf{K}(\cdot, x) \langle f, \mathbf{K}_x \rangle_{\mathbf{K}} d\rho_X(x),$$

where $\mathbf{K}_x(\cdot) = \mathbf{K}(\cdot, x)$, and $\langle \cdot, \cdot \rangle_{\mathbf{K}}$ is the inner product in $\mathcal{H}_{\mathbf{K}}$.

The operator \mathbf{C} can be proved to be a positive trace class operator. Therefore, the operator $\mathbf{C}^{1/2} = \sqrt{\mathbf{C}}$ is well-defined and relates the norms of $f \in \mathcal{H}_{\mathbf{K}}$ in $\mathcal{H}_{\mathbf{K}}$ and $L_2(X, \rho_X)$ as follows

$$\|f\|_{\rho} = \|\mathbf{C}^{1/2} f\|_{\mathbf{K}} \quad (4)$$

where $\|\cdot\|_{\mathbf{K}} = \|\cdot\|_{\mathcal{H}_{\mathbf{K}}}$.

We will measure the approximation power of the projection method induced by the projector $\mathbf{P}_{\mathbf{z}^\nu}$ in terms of the quantity $\|\mathbf{C}^{1/2}(\mathbf{I} - \mathbf{P}_{\mathbf{z}^\nu})\|_{\mathcal{H}_{\mathbf{K}} \rightarrow \mathcal{H}_{\mathbf{K}}}$ that has been also studied in [11] (see Lemma 6 [11]). At the same time, such kind of measure is usual in studying regularized projection methods [9, 10], and in spirit of that studies we assume that there is $\beta > 0$ such that the following holds with probability $1 - \delta$

$$\Delta_m := \|\mathbf{C}^{1/2}(\mathbf{I} - \mathbf{P}_{\mathbf{z}^\nu})\|_{\mathcal{H}_{\mathbf{K}} \rightarrow \mathcal{H}_{\mathbf{K}}} \leq d_{\delta, \beta} m^{-\beta}, m = |\mathbf{z}^\nu|, \quad (5)$$

where $d_{\delta, \beta} = \mathcal{O}(\log^{\beta_1} \frac{1}{\delta})$ and β_1 is a positive number depending only on β .

Note, that a probabilistic character of the assumption (5) is natural, because in the plain Nyström approach the points forming \mathbf{z}^ν are sampled uniformly at random without replacement from the training set \mathbf{z} .

As we have already mentioned, in [11], the Nyström subsampling approach was studied under assumptions on the capacity of $\mathcal{H}_{\mathbf{K}}$. These assumptions are formulated in [11] with the use of the quantity $\mathcal{N}_{\infty}(\lambda) = \sup\{\mathcal{N}_x(\lambda), x \in X\}$, where $\mathcal{N}_x(\lambda) = \langle \mathbf{K}_x, (\mathbf{C} + \lambda \mathbf{I})^{-1} \mathbf{K}_x \rangle_{\mathbf{K}}$. If in spirit of Assumption 3 [11] we assume that $\mathcal{N}_{\infty}(\lambda) = \mathcal{O}(\lambda^{-\gamma})$, $0 < \gamma \leq 1$, then from Lemma 6 [11] it follows that our assumption (5) is satisfied with any $\beta \in (0, 1/2\gamma)$.

Our last assumption describes the regularity of f^\dagger in terms of source condition concept that is fairly standard in the regularization theory [8]. In the context of the learning theory this concept has been introduced in [2]. Within this concept, we assume that f^\dagger admits the representation

$$f^\dagger = \varphi(\mathbf{C})v^\dagger, v^\dagger \in \mathcal{H}_K, \|v^\dagger\|_K \leq R, \quad (6)$$

where the function φ is operator monotone on $[0, d]$, $d > \|\mathbf{C}\|_{\mathcal{H}_K \rightarrow \mathcal{H}_K}$, and such that $\varphi(0) = 0$ and φ^2 is a concave function.

As it has been shown in [9] an important implication of operator monotonicity is that there is a number d_φ depending only on φ such that for any self-adjoint operators C, C_1 with spectra in $[0, d]$ it holds

$$\|\varphi(C) - \varphi(C_1)\|_{\mathcal{H}_K \rightarrow \mathcal{H}_K} \leq d_\varphi \varphi(\|C - C_1\|_{\mathcal{H}_K \rightarrow \mathcal{H}_K}). \quad (7)$$

Moreover, as a consequence of the concavity of φ^2 we have (see Proposition 2 [9])

$$\|(\mathbf{I} - \mathbf{P}_{\mathbf{z}^\nu})\varphi(\mathbf{C})\|_{\mathcal{H}_K \rightarrow \mathcal{H}_K} \leq \varphi(\|\mathbf{C}^{1/2}(\mathbf{I} - \mathbf{P}_{\mathbf{z}^\nu})\|_{\mathcal{H}_K \rightarrow \mathcal{H}_K}^2). \quad (8)$$

Note that our assumption (6) generalizes Assumption 4 of [11], where only the case of operator monotone functions $\varphi(t) = t^s$, $0 < s \leq \frac{1}{2}$, has been studied.

In the sequel we extensively use the following bounds (see, e.g., [2]) that hold under the above assumptions with probability at least $1 - \delta$ and quantify the perturbation effect of random sampling:

$$\|\mathbf{C} - \mathbf{S}_{\mathbf{z}}^* \mathbf{S}_{\mathbf{z}}\|_{\mathcal{H}_K \rightarrow \mathcal{H}_K} \leq d_{1,\delta} |\mathbf{z}|^{-\frac{1}{2}}, \quad (9)$$

$$\|\mathbf{S}_{\mathbf{z}}^* \mathbf{S}_{\mathbf{z}} f - \mathbf{S}_{\mathbf{z}}^* \mathbf{Y}\|_K \leq d_{2,\delta} |\mathbf{z}|^{-\frac{1}{2}}, \quad (10)$$

where $d_{1,\delta}$ and $d_{2,\delta}$ are of order $\mathcal{O}(\log \frac{1}{\delta})$ and depend only on K and ρ .

The following capacity independent learning rates have been proven in [2] for KRR (1)

Theorem 1 ([2]). *Consider a sampling space $Z = X \times [-D, D]$, where the input space $X \subset \mathbb{R}^d$ is closed. Consider also a bounded and continuous kernel K defined on X . If minimizer f^\dagger of the expected risk $\mathcal{E}(f)$ over \mathcal{H}_K meets the assumption (6), then for $\alpha = \alpha_{\mathbf{z}} = \Theta^{-1}(|\mathbf{z}|^{-1/2})$, $\Theta(t) = \varphi(t)t$, we have with probability at least $1 - \delta$ that*

$$\|f^\dagger - f_{\mathbf{z}}^{\alpha_{\mathbf{z}}}\|_\rho = \mathcal{O}\left(\varphi(\Theta^{-1}(|\mathbf{z}|^{-1/2}))\sqrt{\Theta^{-1}(|\mathbf{z}|^{-1/2})\log \frac{1}{\delta}}\right). \quad (11)$$

Note that for $\varphi(t) = t^s$ the above theorem gives us the learning rate $\mathcal{O}\left(|\mathbf{z}|^{-\frac{s+\frac{1}{2}}{2(s+1)}}\right)$ that matches the result obtained in seminal paper by Smale and Zhou [14]. Moreover, for $\varphi(t) = t^s$ the rate (11) can be thought of as the limit case of the capacity dependent learning rate $\mathcal{O}\left(|\mathbf{z}|^{-\frac{(s+\frac{1}{2})\mu}{2s\mu+\mu+1}}\right)$ obtained in [3] under the assumptions that the eigenvalues λ_i of the covariance operator \mathbf{C} have polynomial decay $\lambda_i \asymp i^{-\mu}$ with $\mu > 1$.

Now we are going to prove that the same learning rate (11) can be achieved in Nyström type subsampling (2) if the approximation power of $\mathbf{P}_{\mathbf{z}^\nu}$ is high enough.

Theorem 2. *Assume the conditions of Theorem 1, and let (5) be satisfied. If the size $m = |\mathbf{z}^\nu|$ of a subsampling \mathbf{z}^ν is chosen such that*

$$\Delta_m \leq \sqrt{\Theta_{1/2}^{-1}(|\mathbf{z}|^{-1/2})}, \Theta_{1/2}(t) = \varphi(t)\sqrt{t},$$

then with probability at least $1 - \delta$ we have

$$\|f^\dagger - f_{\mathbf{z}, \mathbf{z}^\nu}^{\alpha_{\mathbf{z}, \mathbf{z}^\nu}}\|_\rho = \mathcal{O}\left(\varphi\left(\Theta^{-1}(|\mathbf{z}|^{-1/2})\right) \sqrt{\Theta^{-1}(|\mathbf{z}|^{-1/2}) \log^{\beta_2} \frac{1}{\delta}}\right), \quad (12)$$

where $\beta_2 = \max\{1, \beta_1\}$, and β_1 is the same as in (5).

Before proving this statement, we first comment on the computational complexity of Nyström approximation (2) with a subsampling size $|\mathbf{z}^\nu|$ chosen according to Theorem 2.

In view of the assumption (5) it is clear that the condition of the theorem can be satisfied with

$$|\mathbf{z}^\nu| \asymp [\Theta_{1/2}^{-1}(|\mathbf{z}|^{-1/2})]^{-\frac{1}{2\beta}}.$$

Let the assumption (6) be satisfied with

$$\varphi(t) = o(t^{\frac{1-\beta}{2\beta}}) \text{ as } t \rightarrow 0, \quad (13)$$

i.e. $\Theta_{1/2}(t) = o(t^{1/2\beta})$. Then

$$|\mathbf{z}|^{-\beta} = o(\Theta_{1/2}^{-1}(|\mathbf{z}|^{-1/2})) = o(|\mathbf{z}^\nu|^{-2\beta}),$$

which means that $|\mathbf{z}^\nu|^2 = o(|\mathbf{z}|)$ as $|\mathbf{z}| \rightarrow \infty$.

On the other hand, the computational complexity of (2) is of order $\mathcal{O}(|\mathbf{z}||\mathbf{z}^\nu|^2)$ (see, e.g. [11]), and under the condition (13) it is subquadratic, because $|\mathbf{z}||\mathbf{z}^\nu|^2 = o(|\mathbf{z}|^2)$.

Thus, under the conditions of Theorem 2 Nyström subsampling has the same learning rate as the one guaranteed by Theorem 1 for KRR based on the whole sample \mathbf{z} . Moreover, Theorem 2 allows an estimation of a regularity range, such as (13), for which the above mentioned learning rate can be achieved with subquadratic complexity. Note, that the condition (13) is automatically satisfied with $\beta \geq 1$, for example.

Proof of Theorem 2. It is known (see, e.g. [9]) that the following inequality holds true for functions φ mentioned in the assumption (6)

$$\sup_t |(1 - (\alpha + t)^{-1}t)\varphi(t)t^q| \leq h_{\varphi,q}\varphi(\alpha)\alpha^q, q \in [0, 1/2], \quad (14)$$

where $h_{\varphi,q}$ depends only on φ and q .

Note also that, by very definition, $\Theta_{1/2}(|\mathbf{z}|^{-1/2}) > \Theta(|\mathbf{z}|^{-1/2})$, and therefore

$$\Delta_m^2 = \Theta_{1/2}^{-1}(|\mathbf{z}|^{-1/2}) < \Theta^{-1}(|\mathbf{z}|^{-1/2}) = \alpha_{\mathbf{z}}. \quad (15)$$

Moreover, without loss of generality we can assume that $|\mathbf{z}|$ is so large that

$$\varphi(\max\{d_{1,\delta}, d_{2,\delta}\}|\mathbf{z}|^{-1/2}) < [\max\{d_{1,\delta}, d_{2,\delta}\}], \quad (16)$$

where $d_{1,\delta}, d_{2,\delta}$ are the numbers appearing in (9), (10). This is not a real restriction, because the left-hand side of (16) tends to zero as $|\mathbf{z}| \rightarrow \infty$. A direct implication of (16) is that with probability at least $1 - \delta$

$$\alpha_{\mathbf{z}} = \Theta^{-1}(|\mathbf{z}|^{-1/2}) > \max\{\|\mathbf{C} - \mathbf{C}_{\mathbf{z}}\|_{\mathcal{H}_K \rightarrow \mathcal{H}_K}, \|\mathbf{C}_{\mathbf{z}}f^\dagger - \mathbf{S}_{\mathbf{z}}^*\mathbf{Y}\|_K\}. \quad (17)$$

Consider the decomposition

$$f^\dagger - f_{\mathbf{z},\mathbf{z}^\nu}^{\alpha_{\mathbf{z}}} = \sigma_1 + \sigma_2 + \sigma_3, \quad (18)$$

where

$$\begin{aligned} \sigma_1 &= f^\dagger - \mathbf{P}_{\mathbf{z}^\nu}f^\dagger, \\ \sigma_2 &= \mathbf{P}_{\mathbf{z}^\nu}f^\dagger - (\alpha_{\mathbf{z}}\mathbf{I} + \mathbf{P}_{\mathbf{z}^\nu}\mathbf{C}_{\mathbf{z}}\mathbf{P}_{\mathbf{z}^\nu})^{-1}\mathbf{P}_{\mathbf{z}^\nu}\mathbf{C}_{\mathbf{z}}\mathbf{P}_{\mathbf{z}^\nu}f^\dagger, \\ \sigma_3 &= (\alpha_{\mathbf{z}}\mathbf{I} + \mathbf{P}_{\mathbf{z}^\nu}\mathbf{C}_{\mathbf{z}}\mathbf{P}_{\mathbf{z}^\nu})^{-1}(\mathbf{P}_{\mathbf{z}^\nu}\mathbf{C}_{\mathbf{z}}\mathbf{P}_{\mathbf{z}^\nu}f^\dagger - \mathbf{P}_{\mathbf{z}^\nu}\mathbf{S}_{\mathbf{z}}^*\mathbf{Y}), \end{aligned}$$

and we use notation $\mathbf{C}_z = \mathbf{S}_z^* \mathbf{S}_z$.

Now we are going to bound each term of (18). From (4)–(6) and (8) we have

$$\begin{aligned}
\|\sigma_1\|_\rho &= \|\mathbf{C}^{1/2}(\mathbf{I} - \mathbf{P}_{z^\nu})\varphi(\mathbf{C})v^\dagger\|_{\mathcal{K}} \\
&\leq R\|\mathbf{C}^{1/2}(\mathbf{I} - \mathbf{P}_{z^\nu})\|_{\mathcal{H}_{\mathcal{K}} \rightarrow \mathcal{H}_{\mathcal{K}}}\|(\mathbf{I} - \mathbf{P}_{z^\nu})\varphi(\mathbf{C})\|_{\mathcal{H}_{\mathcal{K}} \rightarrow \mathcal{H}_{\mathcal{K}}} \\
&\leq R\Delta_m\varphi(\Delta_m^2) = R\Theta_{1/2}(\Delta_m^2) \\
&\leq R\Theta_{1/2}(\Theta_{1/2}^{-1}(|z|^{-1/2})) = R|z|^{-1/2}
\end{aligned} \tag{19}$$

To prove (12) we also need to bound σ_2, σ_3 in the norms $\|\cdot\|_{\mathcal{K}}$ and $\|\cdot\|_\rho$. We start with the decomposition

$$\sigma_2 = \sigma_{2,1} + \sigma_{2,2}, \tag{20}$$

where

$$\begin{aligned}
\sigma_{2,1} &= (\mathbf{I} - (\alpha_z \mathbf{I} + \mathbf{P}_{z^\nu} \mathbf{C}_z \mathbf{P}_{z^\nu}))^{-1} \mathbf{P}_{z^\nu} \mathbf{C}_z \mathbf{P}_{z^\nu} \varphi(\mathbf{P}_{z^\nu} \mathbf{C}_z \mathbf{P}_{z^\nu}) v^\dagger, \\
\sigma_{2,2} &= (\mathbf{I} - (\alpha_z \mathbf{I} + \mathbf{P}_{z^\nu} \mathbf{C}_z \mathbf{P}_{z^\nu}))^{-1} \mathbf{P}_{z^\nu} \mathbf{C}_z \mathbf{P}_{z^\nu} \sigma_{2,2,1}, \\
\sigma_{2,2,1} &= (\mathbf{P}_{z^\nu} \varphi(\mathbf{C}) - \mathbf{P}_{z^\nu} \varphi(\mathbf{C}) \mathbf{P}_{z^\nu} + \mathbf{P}_{z^\nu} \varphi(\mathbf{C}) \mathbf{P}_{z^\nu} \\
&\quad - \varphi(\mathbf{P}_{z^\nu} \mathbf{C} \mathbf{P}_{z^\nu}) + \varphi(\mathbf{P}_{z^\nu} \mathbf{C} \mathbf{P}_{z^\nu}) - \varphi(\mathbf{P}_{z^\nu} \mathbf{C}_z \mathbf{P}_{z^\nu})) v^\dagger.
\end{aligned}$$

From (14) it follows that

$$\|\sigma_{2,1}\|_{\mathcal{K}} \leq R \sup_t |(1 - (\alpha_z + t)^{-1}t)\varphi(t)| \leq Rh_{\varphi,0}\varphi(\alpha_z)$$

Moreover,

$$\begin{aligned}
\|\sigma_{2,1}\|_\rho &= \|\mathbf{C}^{1/2}\sigma_{2,1}\|_{\mathcal{K}} \\
&\leq \|\mathbf{C}_z^{1/2}\mathbf{P}_{z^\nu}\sigma_{2,1}\|_{\mathcal{K}} + \|(\mathbf{C}^{1/2} - \mathbf{C}_z^{1/2})\mathbf{P}_{z^\nu}\sigma_{2,1}\|_{\mathcal{K}},
\end{aligned}$$

and

$$\begin{aligned}
\|\mathbf{C}_z^{1/2}\mathbf{P}_{z^\nu}\sigma_{2,1}\|_{\mathcal{K}} &\leq \|(\mathbf{P}_{z^\nu} \mathbf{C}_z \mathbf{P}_{z^\nu})^{1/2}\sigma_{2,1}\|_{\mathcal{K}} \\
&\leq R \sup_t |(1 - (\alpha_z + t)^{-1}t)t^{1/2}\varphi(t)| \leq Rh_{\varphi, \frac{1}{2}}\alpha_z^{1/2}\varphi(\alpha_z).
\end{aligned}$$

Keeping in mind that $\psi(t) = \sqrt{t}$ is an operator monotone function, from (7), (15) and (17), we have

$$\|(\mathbf{C}^{1/2} - \mathbf{C}_z^{1/2})\mathbf{P}_{z^\nu}\sigma_{2,1}\|_{\mathcal{K}} \leq d_{1/2}\|\mathbf{C} - \mathbf{C}_z\|_{\mathcal{H}_{\mathcal{K}} \rightarrow \mathcal{H}_{\mathcal{K}}}^{1/2}\|\sigma_{2,1}\|_{\mathcal{K}} \leq d_{1/2}Rh_{\varphi,0}\alpha_z^{1/2}\varphi(\alpha_z).$$

All together this gives us the bound

$$\|\sigma_{2,1}\|_\rho = \mathcal{O}(\varphi(\alpha_{\mathbf{z}})\alpha_{\mathbf{z}}^{1/2}) = \mathcal{O}\left(\varphi(\Theta^{-1}(|\mathbf{z}|^{-1/2}))\sqrt{\Theta^{-1}(|\mathbf{z}|^{-1/2})}\right).$$

To estimate $\|\sigma_{2,2}\|_\rho$ we need to bound $\|\sigma_{2,2,1}\|_{\mathcal{K}}$. For this end, we use the following known estimate (see Proposition 3 [9])

$$\|\mathbf{P}_{\mathbf{z}^\nu}\varphi(\mathbf{C})\mathbf{P}_{\mathbf{z}^\nu} - \varphi(\mathbf{P}_{\mathbf{z}^\nu}\mathbf{C}\mathbf{P}_{\mathbf{z}^\nu})\|_{\mathcal{H}_{\mathcal{K}}\rightarrow\mathcal{H}_{\mathcal{K}}} \leq \bar{d}_\varphi\varphi(\|\mathbf{C}^{1/2}(\mathbf{I} - \mathbf{P}_{\mathbf{z}^\nu})\|_{\mathcal{H}_{\mathcal{K}}\rightarrow\mathcal{H}_{\mathcal{K}}}^2).$$

Moreover, (7), (8) and (15), (17) give us

$$\|\varphi(\mathbf{P}_{\mathbf{z}^\nu}\mathbf{C}\mathbf{P}_{\mathbf{z}^\nu}) - \varphi(\mathbf{P}_{\mathbf{z}^\nu}\mathbf{C}_{\mathbf{z}}\mathbf{P}_{\mathbf{z}^\nu})\|_{\mathcal{H}_{\mathcal{K}}\rightarrow\mathcal{H}_{\mathcal{K}}} \leq d_\varphi\varphi(\|\mathbf{C} - \mathbf{C}_{\mathbf{z}}\|_{\mathcal{H}_{\mathcal{K}}\rightarrow\mathcal{H}_{\mathcal{K}}}) \leq d_\varphi\varphi(\alpha_{\mathbf{z}}),$$

and

$$\begin{aligned} \|\mathbf{P}_{\mathbf{z}^\nu}\varphi(\mathbf{C}) - \mathbf{P}_{\mathbf{z}^\nu}\varphi(\mathbf{C})\mathbf{P}_{\mathbf{z}^\nu}\|_{\mathcal{H}_{\mathcal{K}}\rightarrow\mathcal{H}_{\mathcal{K}}} &\leq \|\varphi(\mathbf{C})(\mathbf{I} - \mathbf{P}_{\mathbf{z}^\nu})\|_{\mathcal{H}_{\mathcal{K}}\rightarrow\mathcal{H}_{\mathcal{K}}} \\ &= \|(\mathbf{I} - \mathbf{P}_{\mathbf{z}^\nu})\varphi(\mathbf{C})\|_{\mathcal{H}_{\mathcal{K}}\rightarrow\mathcal{H}_{\mathcal{K}}} \leq \varphi(\|\mathbf{C}^{1/2}(\mathbf{I} - \mathbf{P}_{\mathbf{z}^\nu})\|_{\mathcal{H}_{\mathcal{K}}\rightarrow\mathcal{H}_{\mathcal{K}}}^2) \leq \varphi(\alpha_{\mathbf{z}}). \end{aligned}$$

Therefore, $\|\sigma_{2,2,1}\|_{\mathcal{K}} \leq R(\bar{d}_\varphi + d_\varphi + 1)\varphi(\alpha_{\mathbf{z}})$, and

$$\|\sigma_{2,2}\|_{\mathcal{K}} \leq \|\sigma_{2,2,1}\|_{\mathcal{K}} \sup_t \left|1 - \frac{t}{\alpha_{\mathbf{z}} + t}\right| \leq \|\sigma_{2,2,1}\|_{\mathcal{K}} = \mathcal{O}(\varphi(\alpha_{\mathbf{z}})).$$

Then, using the same argument as for $\|\sigma_{2,2,1}\|_\rho$ we obtain

$$\begin{aligned} \|\sigma_{2,2}\|_\rho &= \mathcal{O}\left(\varphi(\Theta^{-1}(|\mathbf{z}|^{-1/2}))\sqrt{\Theta^{-1}(|\mathbf{z}|^{-1/2})}\right), \text{ and} \\ \|\sigma_2\|_\rho &= \mathcal{O}\left(\varphi(\Theta^{-1}(|\mathbf{z}|^{-1/2}))\sqrt{\Theta^{-1}(|\mathbf{z}|^{-1/2})}\right). \end{aligned}$$

Finally, we need to estimate $\|\sigma_3\|_\rho$. Observe that

$$\begin{aligned} \|\sigma_3\|_\rho &\leq \sup_t |(\alpha_{\mathbf{z}} + t)^{-1}| \|\mathbf{P}_{\mathbf{z}^\nu}\mathbf{C}_{\mathbf{z}}\mathbf{P}_{\mathbf{z}^\nu}f^\dagger - \mathbf{P}_{\mathbf{z}^\nu}\mathbf{S}_{\mathbf{z}}^*\mathbf{Y}\|_{\mathcal{K}} \\ &\leq \frac{1}{\alpha_{\mathbf{z}}} (\|\mathbf{P}_{\mathbf{z}^\nu}(\mathbf{C}_{\mathbf{z}}f^\dagger - \mathbf{S}_{\mathbf{z}}^*\mathbf{Y})\|_{\mathcal{K}} + \|\mathbf{P}_{\mathbf{z}^\nu}\mathbf{C}_{\mathbf{z}}f^\dagger - \mathbf{P}_{\mathbf{z}^\nu}\mathbf{C}_{\mathbf{z}}\mathbf{P}_{\mathbf{z}^\nu}f^\dagger\|_{\mathcal{K}}) \end{aligned}$$

Then using (8)–(10) we obtain

$$\|\mathbf{P}_{\mathbf{z}^\nu}(\mathbf{C}_{\mathbf{z}}f^\dagger - \mathbf{S}_{\mathbf{z}}^*\mathbf{Y})\|_{\mathcal{K}} \leq d_{2,\delta}|\mathbf{z}|^{-1/2},$$

$$\begin{aligned} & \|P_{\mathbf{z}^\nu} \mathbf{C}_z f^\dagger - P_{\mathbf{z}^\nu} \mathbf{C}_z P_{\mathbf{z}^\nu} f^\dagger\|_{\mathcal{K}} \leq \|P_{\mathbf{z}^\nu} (\mathbf{C}_z - \mathbf{C}) f^\dagger\|_{\mathcal{K}} + \|P_{\mathbf{z}^\nu} \mathbf{C} f^\dagger - P_{\mathbf{z}^\nu} \mathbf{C} P_{\mathbf{z}^\nu} f^\dagger\|_{\mathcal{K}} \\ & + \|P_{\mathbf{z}^\nu} (\mathbf{C} - \mathbf{C}_z) P_{\mathbf{z}^\nu} f^\dagger\|_{\mathcal{K}} \leq 2d_{1,\delta} \|f^\dagger\|_{\mathcal{K}} |\mathbf{z}|^{-1/2} + \|\mathbf{C}(\mathbf{I} - P_{\mathbf{z}^\nu})(\mathbf{I} - P_{\mathbf{z}^\nu})\varphi(\mathbf{C})v^\dagger\|_{\mathcal{K}} \\ & \leq d_{3,\delta} (|\mathbf{z}|^{-1/2} + \Delta_m \varphi(\Delta_m^2)) \leq d_{3,\delta} (|\mathbf{z}|^{-1/2} + \Theta_{1/2}(\Theta_{1/2}^{-1}(|\mathbf{z}|^{-1/2}))) = 2d_{3,\delta} |\mathbf{z}|^{-1/2}, \end{aligned}$$

that allows us to write

$$\begin{aligned} \|\sigma_3\|_{\mathcal{K}} &= \mathcal{O}(\alpha_{\mathbf{z}}^{-1} |\mathbf{z}|^{-1/2}) = \mathcal{O}(\alpha_{\mathbf{z}}^{-1} \Theta(\Theta^{-1}(|\mathbf{z}|^{-1/2}))) \\ &= \mathcal{O}([\Theta^{-1}(|\mathbf{z}|^{-1/2})]^{-1} \varphi(\Theta^{-1}(|\mathbf{z}|^{-1/2})) \Theta^{-1}(|\mathbf{z}|^{-1/2})) = \mathcal{O}(\varphi(\Theta^{-1}(|\mathbf{z}|^{-1/2}))). \end{aligned}$$

Using again the same argument as for $\|\sigma_{2,1}\|_{\rho}$ we obtain

$$\|\sigma_3\|_{\rho} = \mathcal{O}\left(\varphi(\Theta^{-1}(|\mathbf{z}|^{-1/2})) \sqrt{\Theta^{-1}(|\mathbf{z}|^{-1/2})}\right).$$

Summing up the above bounds for $\|\sigma_i\|$, $i = 1, 2, 3$, we prove the statement of the theorem. \square

3 Dealing with uncertainty in the sampling size $|\mathbf{z}^\nu|$

Theorem 2 contains a recipe for choosing the subsampling size $|\mathbf{z}^\nu|$ depending on the regularity of the target function and on the approximation power of the corresponding projection method. Both of them, especially the first, may not be exactly given in the form described above. Then several subsampling sizes $|\mathbf{z}^{\nu_1}|, |\mathbf{z}^{\nu_2}|, \dots, |\mathbf{z}^{\nu_l}|$ may be tried in Nyström method, provided that $|\mathbf{z}^{\nu_i}| = o(|\mathbf{z}|^{1/2})$, $i = 1, 2, \dots, l$. Of course, the number l of possible size candidates should not be too large to allow a calculation of all corresponding approximants $f_{\mathbf{z}, \mathbf{z}^{\nu_1}}^\alpha, f_{\mathbf{z}, \mathbf{z}^{\nu_2}}^\alpha, \dots, f_{\mathbf{z}, \mathbf{z}^{\nu_l}}^\alpha$ with a subquadratic complexity. Nevertheless, the question appears of how to select a good approximant among the calculated ones, or how to use all of them. This question is similar to the one in the regularization theory, where some strategy for aggregating all calculated regularized approximants has been discussed recently [4]. In [7] the strategy [4] has been adjusted in the context of learning and presented in several versions.

According to the simplest version, the intention is to approximate the vector $c^* = (c_1^*, c_2^*, \dots, c_l^*) \in \mathbb{R}^l$ solving the following minimization problem

$$\left\| f^\dagger - \sum_{i=1}^l c_i f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha \right\|_{\rho} \rightarrow \min. \quad (21)$$

Recall that $\|\cdot\|_\rho$ is the norm of the Hilbert space $L_2(X, \rho_X)$. Therefore, (21) is equivalent to the matrix problem

$$Gc = g^\dagger, \quad (22)$$

where G and g^\dagger are respectively a Gram matrix and a vector of inner products $\langle \cdot, \cdot \rangle_\rho$ in $L_2(X, \rho_X)$, i.e.

$$G = \left(\left\langle f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha, f_{\mathbf{z}, \mathbf{z}^{\nu_j}}^\alpha \right\rangle_\rho \right)_{i,j=1}^l, \quad g^\dagger = \left(\left\langle f^\dagger, f_{\mathbf{z}, \mathbf{z}^{\nu_j}}^\alpha \right\rangle_\rho \right)_{j=1}^l \quad (23)$$

Note that neither Gram matrix G nor the vector g^\dagger is accessible, since the target function f^\dagger is unknown and the marginal probability distribution ρ_X , which is involved in the definition of $\langle \cdot, \cdot \rangle_\rho$, is not assumed to be given.

On the other hand, $f^\dagger, f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha, i = 1, 2, \dots, l$, belong to the space \mathcal{H}_K . That is assumed to be continuously embedded into $L_2(X, \rho_X)$. Then, for example,

$$\begin{aligned} \left\langle f^\dagger, f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha \right\rangle_\rho &= \left\langle \mathbf{J}_K f^\dagger, \mathbf{J}_K f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha \right\rangle_\rho = \left\langle \mathbf{C} f^\dagger, f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha \right\rangle_K \\ &= \left\langle (\mathbf{C} - \mathbf{C}_z) f^\dagger, f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha \right\rangle_K + \left\langle \mathbf{C}_z f^\dagger - \mathbf{S}_z^* \mathbf{Y}, f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha \right\rangle_K + \left\langle \mathbf{S}_z^* \mathbf{Y}, f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha \right\rangle_K \end{aligned} \quad (24)$$

In view of (9) the first term of the last equality (24) can be estimated as follows:

$$\begin{aligned} \left| \left\langle (\mathbf{C} - \mathbf{C}_z) f^\dagger, f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha \right\rangle_K \right| &\leq \|\mathbf{C} - \mathbf{C}_z\|_{\mathcal{H}_K \rightarrow \mathcal{H}_K} \cdot \|f^\dagger\|_K \cdot \|f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha\|_K \\ &\leq d_{1,\delta} |\mathbf{z}|^{-1/2} \|f^\dagger\|_K \cdot \|f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha\|_K \end{aligned} \quad (25)$$

Moreover, the norm $\|f^\dagger\|_K$ does not depend on $|\mathbf{z}|, |\mathbf{z}^{\nu_i}|$, and the norm $\|f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha\|_K$ can be controlled. So, with a high probability it holds

$$\left| \left\langle (\mathbf{C} - \mathbf{C}_z) f^\dagger, f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha \right\rangle_K \right| = \mathcal{O} \left(|\mathbf{z}|^{-1/2} \right). \quad (26)$$

In the same way, with the use of (10) we have

$$\left| \left\langle \mathbf{C}_z f^\dagger - \mathbf{S}_z^* \mathbf{Y}, f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha \right\rangle_K \right| = \mathcal{O} \left(|\mathbf{z}|^{-1/2} \right). \quad (27)$$

As to the third term of the last equality (24), it can be directly calculated from the training data since

$$\left\langle \mathbf{S}_z^* \mathbf{Y}, f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha \right\rangle_K = \left\langle \mathbf{Y}, \mathbf{S}_z f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha \right\rangle_{\mathbb{R}^{|\mathbf{z}|}} = |\mathbf{z}|^{-1} \sum_{k=1}^{|\mathbf{z}|} y_k f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha(x_k) \quad (28)$$

Therefore, from (24)–(28) we have with high probability

$$\langle f^\dagger, f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha \rangle_\rho = |\mathbf{z}|^{-1} \sum_{k=1}^{|\mathbf{z}|} y_k f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha(x_k) + \mathcal{O}(|\mathbf{z}|^{-1/2}), i = 1, 2, \dots, l. \quad (29)$$

Similar reasoning gives us the relations

$$\langle f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha, f_{\mathbf{z}, \mathbf{z}^{\nu_j}}^\alpha \rangle_\rho = |\mathbf{z}|^{-1} \sum_{k=1}^{|\mathbf{z}|} f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha(x_k) f_{\mathbf{z}, \mathbf{z}^{\nu_j}}^\alpha(x_k) + \mathcal{O}(|\mathbf{z}|^{-1/2}), i, j = 1, 2, \dots, l. \quad (30)$$

In view of (29), (30) the matrix

$$\tilde{G} = \left(|\mathbf{z}|^{-1} \sum_{k=1}^{|\mathbf{z}|} f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha(x_k) f_{\mathbf{z}, \mathbf{z}^{\nu_j}}^\alpha(x_k) \right)_{i,j=1}^l$$

and the vector

$$\tilde{g} = \left(|\mathbf{z}|^{-1} \sum_{k=1}^{|\mathbf{z}|} y_k f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha(x_k) \right)_{i=1}^l$$

can be considered as approximations of G and g^\dagger respectively. Moreover, with probability at least $1 - \delta$

$$\|G - \tilde{G}\|_{\mathbb{R}^l} = \mathcal{O}\left(|\mathbf{z}|^{-1/2} \log \frac{1}{\delta}\right), \quad \|g^\dagger - \tilde{g}\|_{\mathbb{R}^l} = \mathcal{O}\left(|\mathbf{z}|^{-1/2} \log \frac{1}{\delta}\right).$$

With the matrix \tilde{G} in hand one can easily test whether or not \tilde{G}^{-1} exists. For sufficiently large $|\mathbf{z}|$ in case of positive test result a standard perturbation argument (see, e.g. [7] for details) implies the invertibility of G^{-1} , the existence of the vectors $c^* = G^{-1}g^\dagger$, $\tilde{c} = \tilde{G}^{-1}\tilde{g}$ and the bound

$$\|c^* - \tilde{c}\|_{\mathbb{R}^l} = \mathcal{O}\left(|\mathbf{z}|^{-1/2} \log \frac{1}{\delta}\right)$$

that holds with probability at least $1 - \delta$.

Consider now the function

$$f_{\mathbf{z}}^* = \sum_{i=1}^l c_i^* f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha,$$

that solves (21), and its approximation

$$\tilde{f}_{\mathbf{z}}^* = \sum_{i=1}^l \tilde{c}_i^* f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha,$$

where \tilde{c}_i , $i = 1, 2, \dots, l$, are the components of the vector $\tilde{c} = \tilde{G}^{-1} \tilde{g}$. Since $f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha$, $i = 1, 2, \dots, l$, are up to our choice, their norms can be controlled such that

$$\left\| f_{\mathbf{z}}^* - \tilde{f}_{\mathbf{z}} \right\| \leq l \max_i \left\| f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha \right\|_\rho \|c^* - \tilde{c}\|_{\mathbb{R}^l} = \mathcal{O} \left(|\mathbf{z}|^{-1/2} \log \frac{1}{\delta} \right).$$

This gives us the following statement

Theorem 3. *Assume that \tilde{G} is invertible and consider $\tilde{f}_{\mathbf{z}} = \sum_{i=1}^l \tilde{c}_i f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha$, $\tilde{c} = (\tilde{c}_i)_{i=1}^l = \tilde{G}^{-1} \tilde{g}$. Then under the conditions of Theorem 2 for sufficiently large $|\mathbf{z}|$ we have with probability at least $1 - \delta$*

$$\left\| f^\dagger - \tilde{f}_{\mathbf{z}} \right\|_\rho = \min_{c_i} \left\| f^\dagger - \sum_{i=1}^l c_i f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha \right\|_\rho + \mathcal{O} \left(|\mathbf{z}|^{-1/2} \log \frac{1}{\delta} \right),$$

where a coefficient implicit in \mathcal{O} -symbol may depend on the cardinality l of the family $\{f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha\}$ and on the distribution ρ , but does not depend on $|\mathbf{z}|$ and δ .

Note that in Theorem 3 the term $\mathcal{O} \left(|\mathbf{z}|^{-1/2} \log \frac{1}{\delta} \right)$ is negligible because, as we know from [3], $|\mathbf{z}|^{-1/2}$ is of higher order than the best guaranteed accuracy of a reconstruction of the target function $f^\dagger \in \mathcal{H}_K$ in $L_2(X, \rho_X)$ from a training set \mathbf{z} .

Thus, Theorem 3 tells us that the effectively constructed linear combination of the candidates $f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha$, $i = 1, 2, \dots, l$, is almost as accurate as the best linear aggregator of them.

In the next section we present some numerical experiments illustrating the performance of the aggregator $\tilde{f}_{\mathbf{z}}$.

4 Numerical experiments

For our first experiment we simulate data in the same way as in [17], where another strategy for learning with big data called divide and conquer algorithm

or distributed learning has been studied. Following that paper, we simulate training data sets $\mathbf{z} = \{(x_i, y_i)\}_{i=1}^{|\mathbf{z}|}$, $|\mathbf{z}| \in \{2^8, 2^9, \dots, 2^{13}\}$ from the regression model $y_i = f^\dagger(x_i) + \xi_i$, $i = 1, 2, \dots, |\mathbf{z}|$, where $f^\dagger(x) = \min\{x, 1-x\}$, the random samples x_i are uniformly distributed over $[0, 1]$, and the noise random variables ξ_i are normally distributed with zero mean and variance $\sigma^2 = 1/5$. This simulated problem can be seen as a supervised learning with $X = [0, 1]$ and $\rho_X = \text{Uni}[0, 1]$.

As in [17], all kernel ridge regression estimators appearing in this experiment are constructed in \mathcal{H}_K with $K(x, x') = 1 + \min\{x, x'\}$ and $\alpha = |\mathbf{z}|^{-2/3}$.

We perform plain Nyström subsampling and construct estimators $f_{\mathbf{z}, \mathbf{z}^{\nu_1}}^\alpha$, $f_{\mathbf{z}, \mathbf{z}^{\nu_2}}^\alpha$ with $|\mathbf{z}^{\nu_1}| = \lfloor |\mathbf{z}|^{4/10} \rfloor$ and $|\mathbf{z}^{\nu_2}| = \lfloor |\mathbf{z}|^{3/10} \rfloor$, such that the computational complexity of their construction is of order $o(|\mathbf{z}|^2)$, i.e. subquadratic. Then, as has been discussed in Theorem 3, we construct the aggregator $\tilde{f}_{\mathbf{z}} = \tilde{c}_1 f_{\mathbf{z}, \mathbf{z}^{\nu_1}}^\alpha + \tilde{c}_2 f_{\mathbf{z}, \mathbf{z}^{\nu_2}}^\alpha$.

The accuracy of $\tilde{f}_{\mathbf{z}}$ is compared with the one of divide and conquer algorithm [17]. That algorithm is based on splitting a large training set \mathbf{z} into p much smaller equal-sized subsets $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_p$, $|\mathbf{z}_i| = \lfloor |\mathbf{z}|/p \rfloor$, $i = 1, 2, \dots, p$; then, each data set \mathbf{z}_i is used as a training set for constructing the minimizer $f_{\mathbf{z}_i}^\alpha$ of (1), where \mathbf{z} is substituted for \mathbf{z}_i ; finally, the approximations $f_{\mathbf{z}_i}^\alpha$, $i = 1, 2, \dots, p$, are aggregated linearly with equal coefficients (averaged) into

$$f_{\mathbf{z}, p}^\alpha = p^{-1} \sum_{i=1}^p f_{\mathbf{z}_i}^\alpha.$$

In our experiment we compare the errors $\|f^\dagger - \tilde{f}_{\mathbf{z}}\|$, $\|f^\dagger - f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha\|$, $i = 1, 2$, and $\|f^\dagger - f_{\mathbf{z}, p}^\alpha\|$. As in [17] we consider $p = 1, 4, 16, 64$, and execute each simulation 20 times to obtain average values of the errors. In Figure 1 we plot these values versus the total number of samples $|\mathbf{z}|$, where the values corresponding to $\|f^\dagger - \tilde{f}_{\mathbf{z}}\|$, $\|f^\dagger - f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha\|$, and $\|f^\dagger - f_{\mathbf{z}, p}^\alpha\|$ are respectively depicted by dotted, dashed and solid lines.

Figure 1 shows that in the considered case the aggregated approximation $\tilde{f}_{\mathbf{z}}$ outperforms all others, including the baseline KRR-solution $f_{\mathbf{z}, 1}^\alpha$ constructed for the full sample \mathbf{z} . It is also interesting to note, that the Nyström approximation $f_{\mathbf{z}, \mathbf{z}^{\nu_2}}^\alpha$, $|\mathbf{z}^{\nu_2}| = \lfloor |\mathbf{z}|^{3/10} \rfloor$, performs poorly, but the aggregated approximation $\tilde{f}_{\mathbf{z}}$ automatically uses the best of available options.

In our second experiment we follow the paper [11], where the dataset `pumadyn32nh` and `cpuSmall` have been used for an empirical study of the

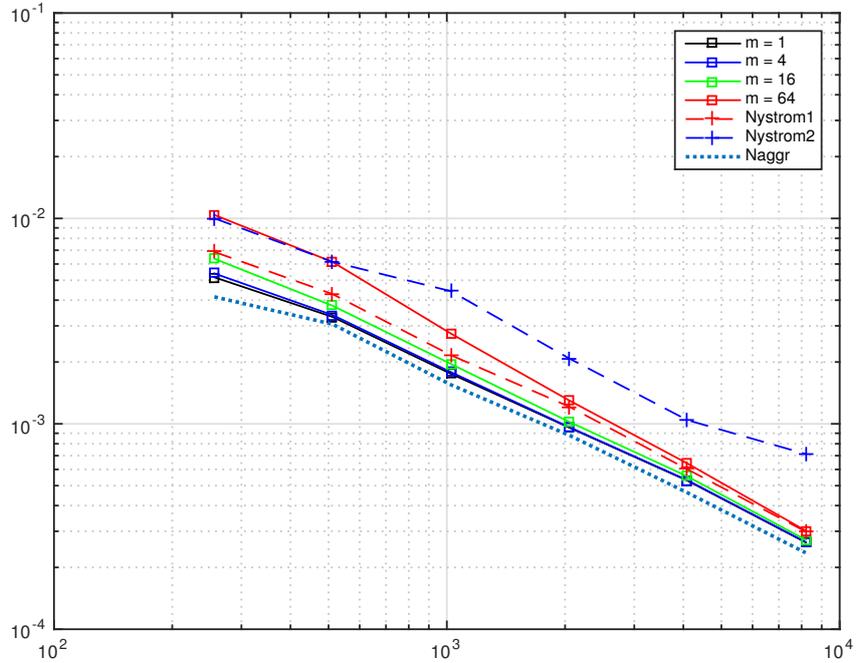


Figure 1: The mean square error between f^\dagger and the averaged estimate $f_{\mathbf{z},p}^\alpha$ for $p = 1, 4, 16, 64$ (solid), Nyström solutions $f_{\mathbf{z},\mathbf{z}^{\nu_1}}^\alpha$ ($|\mathbf{z}^{\nu_1}| = \lfloor |\mathbf{z}|^{4/10} \rfloor$) and $f_{\mathbf{z},\mathbf{z}^{\nu_2}}^\alpha$, $|\mathbf{z}^{\nu_2}| = \lfloor |\mathbf{z}|^{3/10} \rfloor$ (dashed) and aggregated solution $\tilde{f}_{\mathbf{z}}$ (dotted)

Nyström subsampling method. These datasets have been splitted in training and test sets and Gaussian kernels $\mathbf{K}(x, x') = \exp(-\|x - x'\|^2 / 2\sigma^2)$ have been used in construction of $f_{\mathbf{z},\mathbf{z}^\nu}^\alpha$. Moreover, 20% of the training points have been hold out for tuning such parameters as σ and α , and the performance of the selected models has been reported on the test sets.

In [11] the performance has been measured in particular by comparing the root-mean-square-errors (RMSE) of the approximations $f_{\mathbf{z},\mathbf{z}^{\nu_1}}^\alpha$, $f_{\mathbf{z},\mathbf{z}^{\nu_2}}^\alpha$ with large $|\mathbf{z}^{\nu_1}|$ and small $|\mathbf{z}^{\nu_2}|$.

It turns out that in the case of `cpuSmall` the effectiveness of the Nyström subsampling is not so high, since comparable values of RMSE of $f_{\mathbf{z},\mathbf{z}^{\nu_1}}^\alpha$, $f_{\mathbf{z},\mathbf{z}^{\nu_2}}^\alpha$ have been observed when both $|\mathbf{z}^{\nu_1}|$, $|\mathbf{z}^{\nu_2}|$, as well as $|\mathbf{z}|$, are of order of 10^3 .

At the same time, in the case of `pumadyn32nh` the same RMSE of 0.033

has been observed for $f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha$, $i = 1, 2$, with $|\mathbf{z}^{\nu_1}| = 1000$ and $|\mathbf{z}^{\nu_2}| = 62$.

Such different performances may hardly be explained by different capacities of the used hypothesis spaces \mathcal{H}_K , because in both considered cases they are generated by Gaussian kernels, and, moreover, the dimension of the input space X for `cpuSmall` is smaller than in case of `pumadyn32nh`.

In our Theorem 2 one may find a plausible explanation of the above mentioned behaviour of Nyström approximations. Namely, that is because of the regularities of the target functions corresponding to `pumadyn32nh` and `cpuSmall` are described by source condition (6) with functions φ tending to zero with essentially different rates. This is an example of how Theorem 2 can be used for interpreting empirical results and explaining limitations of the Nyström approach.

Now we use `pumadyn32nh` dataset for illustrating the performance of the aggregators $\tilde{f}_{\mathbf{z}}$. As in [11] we construct the approximants $f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha$, $i = 1, 2, 3$, in \mathcal{H}_K generated by the Gaussian kernel of width $\sigma = 2.66$, and we use $\alpha = 10^{-7}$, $|\mathbf{z}| = 4096$, $|\mathbf{z}^{\nu_1}| = 200$, $|\mathbf{z}^{\nu_2}| = 60$, $|\mathbf{z}^{\nu_3}| = 20$. Table 1 reports the performance of $f_{\mathbf{z}, \mathbf{z}^{\nu_i}}^\alpha$, $i = 1, 2, 3$, and $\tilde{\mathbf{z}}$.

Approximant	RMSE
$f_{\mathbf{z}, \mathbf{z}^{\nu_1}}^\alpha$	0.03381
$f_{\mathbf{z}, \mathbf{z}^{\nu_2}}^\alpha$	0.03325
$f_{\mathbf{z}, \mathbf{z}^{\nu_3}}^\alpha$	0.03442
Aggregator $\tilde{f}_{\mathbf{z}}$	0.03325

Table 1: Performance of Nyström approximants and their aggregator on a testing set of 4096 data points from `cpuSmall`

As can be seen from Table 1, the aggregation approach described in Section 3 again automatically uses the best of the available options and can be recommended as a reliable strategy to be implemented together with the Nyström subsampling when dealing with uncertainty in the subsampling size.

Acknowledgements

The authors affiliated with Johann Radon Institute for Computational and Applied Mathematics (RICAM) gratefully acknowledge the support of the the Austrian Science Fund (FWF), grant P25424.

References

- [1] Bach F 2013 Sharp analysis of low-rank kernel matrix approximations *JMLR W&CP* **30** 185–209
- [2] Bauer F, Pereverzev S and Rosasco L 2007 On regularization algorithms in learning theory *Journal of Complexity* **23** 52–72 ISSN 0885-064X
- [3] Caponnetto A and De Vito E 2007 Optimal rates for the regularized least-squares algorithm *Found. Comput. Math.* **7** 331–368
- [4] Chen J, Pereverzyev Jr S and Xu Y 2015 Aggregation of regularized solutions from multiple observation models *Inverse Problems* **31** 075005
- [5] De Vito E, Rosasco L, Caponnetto A, De Giovannini U and Odone F 2005 Learning from examples as an inverse problem *J. Mach. Learn. Res.* **6** 883–904
- [6] Kimeldorf G S and Wahba G 1970 A correspondence between bayesian estimation on stochastic processes and smoothing by splines *The Annals of Mathematical Statistics* **41** 495–502 ISSN 00034851
- [7] Kriukova G, Panasiuk O, Pereverzyev S V and Tkachenko P 2016 A linear functional strategy for regularized ranking *Neural Networks* **73** 26–35 ISSN 0893-6080
- [8] Mathé P and Hofmann B 2008 How general are general source conditions? *Inverse Problems* **24** 015009
- [9] Mathé P and Pereverzev S V 2003 Discretization strategy for linear ill-posed problems in variable Hilbert scales *Inverse Problems* **19** 1263–1277
- [10] Plato R and Vainikko G 1990 On the regularization of projection methods for solving ill-posed problems *Numerische Mathematik* **57** 63–79
- [11] Rudi A, Camoriano R and Rosasco L 2015 Less is more: Nyström computational regularization *Advances in Neural Information Processing Systems 28* ed Cortes C, Lawrence N, Lee D, Sugiyama M, Garnett R and Garnett R (Curran Associates, Inc.) pp 1648–1656 also arXiv:1507.04717 [stat.ML]

- [12] Schölkopf B and Smola A J 2001 *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond* (MIT Press) ISBN 0262194759
- [13] Shawe-Taylor J and Cristianini N 2004 *Kernel Methods for Pattern Analysis* (Cambridge University Press)
- [14] Smale S and Zhou D X 2007 Learning theory estimates via integral operators and their approximations *Constructive Approximation* **26** 153–172 ISSN 0176-4276
- [15] Smola A J and Schölkopf B 2000 Sparse greedy matrix approximation for machine learning *Proceedings of the Seventeenth International Conference on Machine Learning (ICML 2000), Stanford University, Stanford, CA, USA, June 29 - July 2, 2000* pp 911–918
- [16] Williams C and Seeger M 2001 Using the Nyström method to speed up kernel machines *Proceedings of the 14th Annual Conference on Neural Information Processing Systems* pp 682–688
- [17] Zhang Y, Duchi J C and Wainwright M J 2013 Divide and conquer kernel ridge regression: A distributed algorithm with minimax optimal rates *arXiv:1305.5029 [math.ST]* Also *JMLR W&CP* 30: 592–617