



# Passive sampling in reproducing kernel Hilbert spaces using leverage scores<sup>☆</sup>

Pere Giménez-Febrer<sup>a,\*</sup>, Alba Pagès-Zamora<sup>a</sup>, Ignacio Santamaría<sup>b</sup>

<sup>a</sup>SPCOM Group, Universitat Politècnica de Catalunya-Barcelona Tech, Spain

<sup>b</sup>Dept. of Communications Engineering, Universidad de Cantabria, Spain

## ARTICLE INFO

### Article history:

Received 23 September 2021

Revised 13 March 2022

Accepted 23 April 2022

Available online 4 May 2022

### Keywords:

Kernel ridge regression

Leverage score

Nyström approximation

Passive sampling

Reproducing kernel Hilbert space

## ABSTRACT

This paper deals with the selection of the training dataset in kernel-based methods for function reconstruction, with a focus on kernel ridge regression. A functional analysis is performed which, in the absence of noise, links the optimal sampling distribution to the one minimizing the difference between the kernel matrix and its low-rank Nyström approximation. From this standpoint, a statistical passive sampling approach is derived which uses the leverage scores of the columns of the kernel matrix to design a sampling distribution that minimizes an upper bound of the risk function. The proposed approach constitutes a passive method, able to select the optimal subset of training samples using only information provided by the input set and the kernel, but without needing to know the values of the function to be approximated. Furthermore, the proposed approach is backed up by numerical tests on real datasets.

© 2022 Elsevier B.V. All rights reserved.

## 1. Introduction

Machine learning predicts labels by applying decision rules previously learned on a training dataset, which consists of a set of input samples and their labels. The training objective is to obtain a model that can generalize to previously unseen or unlabeled inputs. While in learning applications the focus is usually on achieving a small prediction error, which is heavily influenced by the choice of the training set, there is also the cost of label acquisition for training. This cost might be comparable or even superior to the prediction cost. For instance, in recommender systems, a request is sent to a user to rate an item; but this does not even guarantee a rating since the request might be ignored. In other areas such as sensor networks, medicine, or any scenario with human annotators, the time to label an input is also non-negligible. Hence, the design of optimal sample selection strategies is of paramount importance to improve algorithm performance both in terms of accuracy and cost. The existing sample selection strategies may be divided into active and passive, with the difference between them requiring labeled inputs or not.

<sup>☆</sup> This work has been funded by the Ministerio de Ciencia e Innovación (MICINN) of the Spanish Government and by the Agencia Estatal de Investigación (AEI/10.13039/501100011033) and ERDF funds (PID 2019-104958RB-C41/C43, RED2018-102668-T); and by the Catalan Government (2017 SGR 578).

\* Corresponding author.

E-mail address: [p.gimenez@upc.edu](mailto:p.gimenez@upc.edu) (P. Giménez-Febrer).

Active sampling schemes allow to choose the most promising set of inputs to label. They implement an iterative scheme that, given the availability of a small starting training set, executes two steps repeatedly: first, the function is learned or updated using the current training set and, second, a criterion is evaluated to decide which should be the next input to be labeled and then added to the training set. There exist a variety of criteria for active sampling, the most common being based on prediction uncertainty. For instance, in Gaussian processes [1,2], the predictive variance serves as an uncertainty measure.

While active sampling schemes generally provide good results, their performance may be degraded when the samples are too noisy or when online operation is not possible [3]. There are also batch versions of active sampling which acquire the labels in groups, but the issue with noisy samples still remains. Moreover, most active learning methods are designed for classification tasks [4,5], which often have built-in uncertainty measures, while the availability of methods for regression is more limited [6–8].

An alternative to active sampling is provided by the passive sampling schemes. In passive sampling, the set of inputs to be labeled is selected by observing the geometry of the input space only. For instance, the greedy sampling approach in [3] iteratively adds new input samples to the training set by choosing the one with the largest distance to the set. By relying only on the input space, passive sampling avoids iteratively recomputing the learned function for every additional input. The selected set of inputs can be labeled all at once when it is deemed complete. Furthermore, the impact of noisy samples is diminished.

This paper addresses passive sampling for function approximation in a reproducing kernel Hilbert space (RKHS), with a focus on kernel ridge regression (KRR) [9–11]. To do so, a functional analysis is conducted which connects the concept of optimal sampling to the Nyström approximation [12] - a low-rank approximation to the kernel matrix built from a subset of its columns. Building atop this finding, a sample selection approach is presented based on selecting the samples by means of the leverage scores [13] of the kernel matrix columns. The proposed method enables passive batch sampling, hence selecting the entire training set without needing any label; only the kernel is needed. Moreover, the use of kernels enables the direct application of passive sampling to any input space, e.g., non-metric spaces, provided that a kernel function exists for it.

The paper is organized as follows. Section 2 provides a theoretical analysis of the sampling procedure in RKHSs and proves that Nyström-based passive sampling is optimal, in the absolute error sense, for noiseless samples. Section 3 extends the analysis to KRR with noisy samples. Section 4 introduces the passive sampling approach for KRR based on the leverage scores of the kernel matrix. Finally, Section 5 presents numerical tests and Section 6 offers conclusions.

**Notation.** Boldface lower-case fonts denote column vectors, boldface upper-case fonts denote matrices, and calligraphic upper-case letters denote sets. The  $i$ th entry of a vector  $\mathbf{f}$  is denoted by  $\mathbf{f}(i)$ , the  $i$ th column of a matrix  $\mathbf{K}$  by  $\mathbf{K}(i)$ , and the entry at row  $i$  column  $j$  by  $\mathbf{K}_{ij}$ . The superscript  $T$  denotes the transpose,  $\dagger$  the pseudo inverse, and the hat  $\hat{\cdot}$  is used for estimates. The symbol  $\mathbf{I}$  stands for the identity matrix of appropriate size, specified by the context, and the trace operator is  $\text{Tr}(\cdot)$ .

## 2. Optimal sampling and reconstruction in RKHSs

Let  $f : \mathcal{X} \rightarrow \mathbb{R}$  be a function defined on an RKHS  $\mathcal{H}_N$ , and  $\bar{\mathcal{X}} = \{\bar{x}_i\}_{i=1}^s$  be a set of sampled points in  $\mathcal{X}$  where  $|\mathcal{X}| = N$ . Here, we focus on the problem of reconstructing  $f$  given a set of noiseless observations  $\{\bar{y}_i\}_{i=1}^s$  such that  $\bar{y}_i = f(\bar{x}_i) \forall i = 1, \dots, s$ .

The space  $\mathcal{H}_N$  is a Hilbert space of  $\mathbb{R}$ -valued functions on a non-empty set  $\mathcal{X}$  with inner product  $\langle \cdot, \cdot \rangle_{\mathcal{H}_N} : \mathcal{H}_N \times \mathcal{H}_N \rightarrow \mathbb{R}$  and induced norm  $\|f\|_{\mathcal{H}_N} = \langle f, f \rangle_{\mathcal{H}_N}$ . A function  $\kappa : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is a reproducing kernel of  $\mathcal{H}_N$ , and  $\mathcal{H}_N$  is a RKHS if

- $\forall x \in \mathcal{X}, \kappa(\cdot, x) \in \mathcal{H}_N$ ,
- $\forall x \in \mathcal{X}, \forall f \in \mathcal{H}_N, \langle f(\cdot), \kappa(\cdot, x) \rangle_{\mathcal{H}_N} = f(x)$ ,

where the latter is the reproducing property. In particular, the reproducing property implies that  $\langle \kappa(\cdot, x), \kappa(\cdot, x') \rangle_{\mathcal{H}_N} = \kappa(x, x')$ . In view of the Moore-Aronszajn theorem, one can construct the RKHS  $\mathcal{H}_N$  as the completion of the space of functions spanned by the set  $\{\kappa(\cdot, x_i)\}_{i=1}^N$ , i.e.,

$$\mathcal{H}_N := \left\{ f : f(x) = \sum_{i=1}^N \alpha_i \kappa(x, x_i), \alpha_i \in \mathbb{R} \right\} \quad (1)$$

with an inner product given by  $\langle f, f' \rangle_{\mathcal{H}_N} = \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha'_j \kappa(x_i, x_j)$ , where  $\{\alpha_i\}_{i=1}^N$  and  $\{\alpha'_j\}_{j=1}^N$  are the coefficients of  $f, f' \in \mathcal{H}_N$ , respectively. We refer to  $\kappa$  as the kernel function that spans  $\mathcal{H}_N$ , and to  $\mathbf{K}$  as its kernel matrix with  $\mathbf{K}_{ij} = \kappa(x_i, x_j)$ . Note that our RKHS definition (1) differs from the usual one with an infinite-dimensional RKHS since  $|\mathcal{X}| = N$  and, therefore,  $\mathcal{H}_N$  is a finite-dimensional space [9].

In order to motivate the passive sampling approach introduced in Section 4, this section builds on the functional analysis approach to the signal reconstruction process from [14] and derives results valid for any algorithm operating in a RKHS.

The recovery of  $f$  from the observation set can be thought of as the “inversion” of

$$\bar{\mathbf{f}} = \mathcal{A}f \quad (2)$$

where  $\bar{\mathbf{f}} = [f(\bar{x}_1), \dots, f(\bar{x}_s)]^T$ , and  $\mathcal{A} : \mathcal{H}_N \rightarrow \mathbb{R}^s$  is a sampling operator evaluating  $f$  at  $\{\bar{x}_i\}_{i=1}^s$ . Leveraging the reproducing property, the sampling operator is defined as

$$\mathcal{A}f = \sum_{n=1}^s \langle f, \kappa(\cdot, \bar{x}_n) \rangle_{\mathcal{H}_N} \mathbf{e}_n \quad (3)$$

where  $\mathbf{e}_n$  is the  $n$ th vector of the standard basis, i.e., the  $n$ th column of the  $s \times s$  identity matrix. Let  $\mathcal{H}_s \subseteq \mathcal{H}_N$  be the subspace of dimension  $s$  spanned by  $\{\kappa(\cdot, \bar{x}_n)\}_{n=1}^s$ . Then,  $f = f_s + f_s^\perp$ , where  $f_s$  is the projection of  $f$  onto  $\mathcal{H}_s$  and  $f_s^\perp$  is the projection onto the subspace orthogonal to  $\mathcal{H}_s$ , denoted by  $\mathcal{H}_s^\perp \subseteq \mathcal{H}_N$ , such that  $\langle f_s^\perp, \kappa(\cdot, \bar{x}) \rangle_{\mathcal{H}_N} = 0 \forall \bar{x} \in \bar{\mathcal{X}}$ . Thus, for  $\bar{x} \in \bar{\mathcal{X}}$ ,

$$\begin{aligned} f(\bar{x}) &= f_s(\bar{x}) + f_s^\perp(\bar{x}) = \langle f_s, \kappa(\cdot, \bar{x}) \rangle_{\mathcal{H}_N} + \langle f_s^\perp, \kappa(\cdot, \bar{x}) \rangle_{\mathcal{H}_N} \\ &= \langle f_s, \kappa(\cdot, \bar{x}) \rangle_{\mathcal{H}_N} = f_s(\bar{x}). \end{aligned} \quad (4)$$

Since  $f_s \in \mathcal{H}_s$ , it can be written as  $\sum_{n=1}^s \bar{\alpha}_n \kappa(\cdot, \bar{x}_n)$  for some real coefficients  $\{\bar{\alpha}_n\}_{n=1}^s$ . In turn,  $f(\bar{x}) = \sum_{n=1}^s \bar{\alpha}_n \kappa(\bar{x}, \bar{x}_n)$ . Therefore, since  $\bar{\mathbf{f}}(i) := f(\bar{x}_i)$ ,  $\mathcal{A}$  maps  $f$  into a new RKHS  $\bar{\mathcal{H}}_s \subseteq \mathbb{R}^s$  defined as

$$\bar{\mathcal{H}}_s := \left\{ \bar{\mathbf{f}} : \bar{\mathbf{f}}(i) = \sum_{n=1}^s \bar{\alpha}_n \kappa(\bar{x}_i, \bar{x}_n), \bar{\alpha}_n \in \mathbb{R} \right\} \quad (5)$$

with kernel matrix  $\bar{\mathbf{K}} \in \mathbb{R}^{s \times s}$  where  $\bar{\mathbf{K}}_{i,j} = \kappa(\bar{x}_i, \bar{x}_j)$  so that  $\bar{\mathbf{f}} = \bar{\mathbf{K}}\bar{\boldsymbol{\alpha}}$ , where  $\bar{\boldsymbol{\alpha}} = [\bar{\alpha}_1, \dots, \bar{\alpha}_s]^T$ , and  $\langle \bar{\mathbf{f}}, \bar{\mathbf{f}}' \rangle_{\bar{\mathcal{H}}_s} = \bar{\mathbf{f}}^T \bar{\mathbf{K}}^{-1} \bar{\mathbf{f}}'$  for  $\bar{\mathbf{f}}, \bar{\mathbf{f}}' \in \bar{\mathcal{H}}_s$  assuming that  $\bar{\mathbf{K}}$  is invertible. With  $\mathcal{A} : \mathcal{H}_N \rightarrow \bar{\mathcal{H}}_s$ , its adjoint operator is  $\mathcal{A}^* : \bar{\mathcal{H}}_s \rightarrow \mathcal{H}_N$ , and by definition it satisfies for any  $f'' \in \mathcal{H}_N$  and  $\bar{\mathbf{f}}' \in \bar{\mathcal{H}}_s$  that

$$\langle \mathcal{A}f'', \bar{\mathbf{f}}' \rangle_{\bar{\mathcal{H}}_s} = \langle f'', \mathcal{A}^* \bar{\mathbf{f}}' \rangle_{\mathcal{H}_N}. \quad (6)$$

Particularizing (6) for  $f'' = f$  and  $\bar{\mathbf{f}}' = \bar{\mathbf{f}}$ , with  $\bar{\mathbf{f}}$  in (2), we obtain

$$\begin{aligned} \langle \mathcal{A}f, \bar{\mathbf{f}} \rangle_{\bar{\mathcal{H}}_s} &= \left\langle \sum_{n=1}^s \langle f, \kappa(\cdot, \bar{x}_n) \rangle_{\mathcal{H}_N} \mathbf{e}_n, \bar{\mathbf{f}} \right\rangle_{\bar{\mathcal{H}}_s} \\ &= \sum_{n=1}^s \langle f, \kappa(\cdot, \bar{x}_n) \rangle_{\mathcal{H}_N} \langle \mathbf{e}_n, \bar{\mathbf{f}} \rangle_{\bar{\mathcal{H}}_s} \\ &= \langle f, \sum_{n=1}^s \langle \mathbf{e}_n, \bar{\mathbf{f}} \rangle_{\bar{\mathcal{H}}_s} \kappa(\cdot, \bar{x}_n) \rangle_{\mathcal{H}_N}. \end{aligned} \quad (7)$$

This leads to the definition

$$\mathcal{A}^* \bar{\mathbf{f}} = \sum_{n=1}^s \langle \mathbf{e}_n, \bar{\mathbf{f}} \rangle_{\bar{\mathcal{H}}_s} \kappa(\cdot, \bar{x}_n) = \sum_{n=1}^s \bar{\alpha}_n \kappa(\cdot, \bar{x}_n) \quad (8)$$

where  $\langle \mathbf{e}_n, \bar{\mathbf{f}} \rangle_{\bar{\mathcal{H}}_s} = \bar{\mathbf{f}}^T \bar{\mathbf{K}}^{-1} \mathbf{e}_n = \bar{\alpha}_n$ . The operator  $\mathcal{A}^*$  maps  $\bar{\mathbf{f}}$  back into  $\mathcal{H}_N$  using only the basis functions in  $\mathcal{H}_s$ , hence reversing the sampling in (2) and yielding an estimate for  $f$ . Moreover, the sampled values  $f(\bar{x}_i)_{i=1}^s$  are not modified by  $\mathcal{A}^*$ . Hence, since  $\mathcal{A}\mathcal{A}^* \bar{\mathbf{f}} = \bar{\mathbf{f}}$ , we have that  $\mathcal{A}\mathcal{A}^* = \mathcal{I}$ ; this can be seen by substituting (8) into (3).

After the sampling operator and its adjoint have been defined, the reconstruction of a function in  $\mathcal{H}_N$  from its samples can be achieved through the consecutive application of each operator, i.e.,  $\hat{f} = \mathcal{A}^* \mathcal{A}f$ . Hence,  $\mathcal{P} = \mathcal{A}^* \mathcal{A}$  is an orthogonal projector onto  $\mathcal{H}_s$  since it is a self-adjoint operator, i.e.,  $\langle \mathcal{P}f, f' \rangle_{\mathcal{H}_N} = \langle f, \mathcal{P}f' \rangle_{\mathcal{H}_N}$ , and it satisfies  $\mathcal{P}\mathcal{P} = \mathcal{P}$ . Then, using (8) the estimate  $\hat{f}$  is

$$\hat{f} = \mathcal{P}f = \mathcal{A}^* \bar{\mathbf{f}} = \sum_{n=1}^s \bar{\alpha}_n \kappa(\cdot, \bar{x}_n) = f_s. \quad (9)$$

Since an orthogonal projection onto a space yields an element with minimum distance to the original function, the reconstruction error at  $x \in \mathcal{X}$  is upper bounded by

$$\begin{aligned} |f(x) - \mathcal{P}f(x)| &= |\langle f, \kappa(\cdot, x) \rangle_{\mathcal{H}_N} - \langle \mathcal{P}f, \kappa(\cdot, x) \rangle_{\mathcal{H}_N}| \\ &\stackrel{(a)}{=} |\langle f, \kappa(\cdot, x) \rangle_{\mathcal{H}_N} - \langle f, \mathcal{P}\kappa(\cdot, x) \rangle_{\mathcal{H}_N}| \\ &= |\langle f, \kappa(\cdot, x) - \mathcal{P}\kappa(\cdot, x) \rangle_{\mathcal{H}_N}| \\ &\stackrel{(b)}{\leq} \|f\|_{\mathcal{H}_N} \|\kappa(\cdot, x) - \mathcal{P}\kappa(\cdot, x)\|_{\mathcal{H}_N} \end{aligned} \quad (10)$$

where (a) uses the self-adjoint property of  $\mathcal{P}$ , and (b) uses the Cauchy-Schwarz inequality. It can be seen that the error upper bound at  $x$  depends on the distance between the kernel function  $\kappa(\cdot, x)$  and its projection onto  $\mathcal{H}_S$ . Note that zero error is obtained at the sampled points since  $\mathcal{P}\kappa(\cdot, \bar{x}) = \kappa(\cdot, \bar{x}) \forall \bar{x} \in \bar{\mathcal{X}}$ . For  $x \notin \bar{\mathcal{X}}$ , the error can only be zero if  $\kappa(\cdot, x) \in \mathcal{H}_S$ ; this is the case when a kernel function is duplicated such that  $\kappa(\cdot, \bar{x}_i) = \kappa(\cdot, x_j)$  where  $\bar{x}_i \in \bar{\mathcal{X}}$  and  $x_j \notin \bar{\mathcal{X}}$ . This might happen, for instance, in a recommender system where two different users have identical tastes; the kernel function evaluated at each of the two users has identical value as well. Assuming no such duplicities exist,

$$|f(x) - \mathcal{P}f(x)| \leq \begin{cases} 0 & \text{if } x \in \bar{\mathcal{X}} \\ \|f\|_{\mathcal{H}_N} \|\kappa(\cdot, x) - \mathcal{P}\kappa(\cdot, x)\|_{\mathcal{H}_N} & \text{if } x \notin \bar{\mathcal{X}} \end{cases} \quad (11)$$

The absolute error over the function estimate is denoted by  $|f - \hat{f}| = \sum_{i=1}^N |f(x_i) - \hat{f}(x_i)|$ , and (11) shows that it can only be zero when  $f \in \mathcal{H}_S$  so that  $\mathcal{P}f = f$ . Therefore, in order to have zero error for any  $f \in \mathcal{H}_N$ , then  $\mathcal{H}_N = \mathcal{H}_S$  must be true.

The second norm in the last inequality of (10) provides the distance in  $\mathcal{H}_N$  between a kernel function and its closest approximation built with the functions in  $\mathcal{H}_S$ . Alternatively, this difference can be written in terms of  $\mathcal{A}$  as shown in the lemma below, which will be used later.

**Lemma 1.** *The distance between  $\kappa(\cdot, x)$  and its projection onto  $\mathcal{H}_S \forall x \in \mathcal{X}$  is equal to the difference between the norm of  $\kappa(\cdot, x)$  in  $\mathcal{H}_N$  and its sampled counterpart in  $\bar{\mathcal{H}}_S$ ,*

$$\|\kappa(\cdot, x) - \mathcal{P}\kappa(\cdot, x)\|_{\mathcal{H}_N}^2 = \|\kappa(\cdot, x)\|_{\mathcal{H}_N}^2 - \|\mathcal{A}\kappa(\cdot, x)\|_{\bar{\mathcal{H}}_S}^2 \quad (12)$$

**Proof.** Let us expand the second term in (10) as

$$\begin{aligned} &\|\kappa(\cdot, x) - \mathcal{P}\kappa(\cdot, x)\|_{\mathcal{H}_N}^2 \\ &= \langle \kappa(\cdot, x) - \mathcal{P}\kappa(\cdot, x), \kappa(\cdot, x) - \mathcal{P}\kappa(\cdot, x) \rangle_{\mathcal{H}_N} \\ &= \kappa(x, x) - 2\langle \kappa(\cdot, x), \mathcal{P}\kappa(\cdot, x) \rangle_{\mathcal{H}_N} + \langle \mathcal{P}\kappa(\cdot, x), \mathcal{P}\kappa(\cdot, x) \rangle_{\mathcal{H}_N} \\ &\stackrel{(a)}{=} \kappa(x, x) - 2\langle \kappa(\cdot, x), \mathcal{P}\kappa(\cdot, x) \rangle_{\mathcal{H}_N} + \langle \mathcal{A}\kappa(\cdot, x), \mathcal{A}\kappa(\cdot, x) \rangle_{\bar{\mathcal{H}}_S} \\ &\stackrel{(b)}{=} \kappa(x, x) - 2\langle \mathcal{A}\kappa(\cdot, x), \mathcal{A}\kappa(\cdot, x) \rangle_{\bar{\mathcal{H}}_S} + \langle \mathcal{A}\kappa(\cdot, x), \mathcal{A}\kappa(\cdot, x) \rangle_{\bar{\mathcal{H}}_S} \\ &= \kappa(x, x) - \langle \mathcal{A}\kappa(\cdot, x), \mathcal{A}\kappa(\cdot, x) \rangle_{\bar{\mathcal{H}}_S} \\ &= \|\kappa(\cdot, x)\|_{\mathcal{H}_N}^2 - \|\mathcal{A}\kappa(\cdot, x)\|_{\bar{\mathcal{H}}_S}^2 \end{aligned} \quad (13)$$

where, using  $\langle f, \mathcal{A}^* \bar{f} \rangle_{\mathcal{H}_N} = \langle \mathcal{A}f, \bar{f} \rangle_{\bar{\mathcal{H}}_S}$  and  $\mathcal{A}\mathcal{A}^* = \mathcal{I}$ , the identity  $\langle \mathcal{A}^* \mathcal{A}f, \mathcal{A}^* \mathcal{A}f \rangle_{\mathcal{H}_N} = \langle \mathcal{A}\mathcal{A}^* \mathcal{A}f, \mathcal{A}f \rangle_{\bar{\mathcal{H}}_S}$  has been applied in (a), and  $\langle f, \mathcal{A}^* \mathcal{A}f \rangle_{\mathcal{H}_N} = \langle \mathcal{A}f, \mathcal{A}f \rangle_{\bar{\mathcal{H}}_S}$  in (b).  $\square$

Before continuing with the analysis, let us first introduce the following definition:

**Nyström approximation.** Let  $\mathcal{B} \subseteq \mathbb{N}$  be a set indexing  $|\mathcal{B}| = r \ll N$  columns of  $\mathbf{K} \in \mathbb{R}^{N \times N}$ . Defining  $\mathbf{R} = [\mathbf{K}(i)]_{i \in \mathcal{B}}$  as the tall matrix formed by the indexed columns, its transpose  $\mathbf{R}^T$ , and  $\mathbf{C}$  as the submatrix indexed by  $\{(i, j)\}_{(i, j) \in \mathcal{B} \times \mathcal{B}}$ , the Nyström approximation to  $\mathbf{K}$  is

$$\mathbf{N} = \mathbf{R}\mathbf{C}^{\dagger}\mathbf{R}^T. \quad (14)$$

The Nyström approximation is a low-rank approximation that is used in place of the full-rank kernel matrix in kernel methods [15] to accelerate the inversion of the kernel matrix. While in this paper this matrix is not used for this purpose, it arises in the various analyses conducted from here on. Thus, although (10) addresses the per-sample error, Lemma 1 enables an upper bound formulated in terms of the Nyström approximation when considering the total reconstruction error across  $f$ . This bound constitutes the main statement in this section and is presented in the theorem below.

**Theorem 1.** *Let  $\mathbf{f} = [f(x_1), \dots, f(x_N)]^T$ , and  $\mathbf{S}$  be an  $s \times N$  binary sampling matrix with a non-zero element per row equal to 1 such that  $\mathbf{S}\mathbf{f} := \mathcal{A}f$ . Moreover, let  $\mathbf{K}$  be the kernel matrix of  $\mathcal{H}_N$ , and  $\mathbf{T} := \mathbf{K}\mathbf{S}^T \bar{\mathbf{K}}^{-1} \mathbf{S}\mathbf{K}$  with  $\bar{\mathbf{K}} = \mathbf{K}\mathbf{S}\mathbf{S}^T$ . Then, the absolute error across  $f$ , defined as  $|f - \mathcal{P}f| = \sum_{m=1}^N |f(x_m) - \mathcal{P}f(x_m)|$ , is bounded as*

$$|f - \mathcal{P}f| \leq \|f\|_{\mathcal{H}_N} \text{Tr}(\mathbf{K} - \mathbf{T}) \quad (15)$$

**Proof.** Using (10), (13), and  $\|\bar{f}\|_{\bar{\mathcal{H}}_S} = \bar{\mathbf{f}}^T \bar{\mathbf{K}}^{-1} \bar{\mathbf{f}}$ , it holds that

$$\begin{aligned} |f - \mathcal{P}f| &= \sum_{m=1}^N |f(x_m) - \mathcal{P}f(x_m)| \\ &\leq \|f\|_{\mathcal{H}_N} \sum_{m=1}^N \|\kappa(\cdot, x_m) - \mathcal{P}\kappa(\cdot, x_m)\|_{\mathcal{H}_N}^2 \\ &= \|f\|_{\mathcal{H}_N} \sum_{m=1}^N (\|\kappa(\cdot, x_m)\|_{\mathcal{H}_N}^2 - \|\mathcal{A}\kappa(\cdot, x_m)\|_{\bar{\mathcal{H}}_S}^2) \\ &= \|f\|_{\mathcal{H}_N} \sum_{m=1}^N \left( \kappa(x_m, x_m) - \sum_{i=1}^s \sum_{j=1}^s \kappa(\bar{x}_i, x_m) (\bar{\mathbf{K}}^{-1})_{i,j} \kappa(\bar{x}_j, x_m) \right) \end{aligned} \quad (16)$$

$$\begin{aligned} &= \|f\|_{\mathcal{H}_N} \text{Tr}(\mathbf{K} - \mathbf{K}\mathbf{S}^T \bar{\mathbf{K}}^{-1} \mathbf{S}\mathbf{K}) \\ &= \|f\|_{\mathcal{H}_N} \text{Tr}(\mathbf{K} - \mathbf{T}). \end{aligned} \quad (17)$$

$\square$

Matrix  $\mathbf{T}$  constitutes a Nyström approximation as defined in (14), with the specificity that its rank, given by  $s$ , and the selection of columns is tied to the sampling matrix  $\mathbf{S}$ . Again, note that full-rank kernel matrix is not being replaced with an approximation, the Nyström approximation appears through the analysis conducted in Theorem 1. The result in (15) shows that the upper bound on the reconstruction error of  $f$  is proportional to the difference between  $\mathbf{K}$  and  $\mathbf{T}$  on its diagonal entries. Moreover,  $\text{Tr}(\mathbf{K} - \mathbf{T}) \geq 0$  since in (12)  $\|\kappa(\cdot, x)\|_{\mathcal{H}_N}^2 \geq \|\mathcal{A}\kappa(\cdot, x)\|_{\bar{\mathcal{H}}_S}^2$ . Therefore, designing  $\mathcal{A}$  to maximize  $\sum_{n=1}^N \|\mathcal{A}\kappa(\cdot, x_n)\|_{\bar{\mathcal{H}}_S}^2 = \text{Tr}(\mathbf{T})$  minimizes the upper bound in Theorem 1.

### 3. Function sampling and reconstruction in kernel ridge regression

In KRR, the function to be recovered is a vector  $\mathbf{f} = [f(x_1), \dots, f(x_N)]^T \in \mathcal{H}_N \subseteq \mathbb{R}^N$ , where  $\mathbf{f} = \mathbf{K}\boldsymbol{\alpha}$ ,  $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_N]^T$ , and  $\langle \mathbf{f}, \mathbf{f}' \rangle_{\mathcal{H}_N} = \mathbf{f}^T \mathbf{K}^{-1} \mathbf{f}'$  with  $\mathbf{f}' \in \mathcal{H}_N$ . Given the  $s$ -dimensional vector of noisy observations  $\bar{\mathbf{y}} = \mathbf{S}\mathbf{y}$ , where  $\mathbf{S}$  is an  $s \times N$  binary sampling matrix with a single non-zero element per row,  $\mathbf{y} = [y_1, \dots, y_N]^T$ ,  $y_i = f(x_i) + w_i$  and  $w_i$  denotes noise, the problem is formulated as solving

$$\hat{\mathbf{f}} = \underset{\mathbf{f} \in \mathcal{H}_N}{\text{argmin}} \|\bar{\mathbf{y}} - \mathbf{S}\mathbf{f}\|_2^2 + \mu \|\mathbf{f}\|_{\mathcal{H}_N} \quad (18)$$

or, equivalently,

$$\hat{\boldsymbol{\alpha}} = \underset{\boldsymbol{\alpha} \in \mathbb{R}^N}{\operatorname{argmin}} \|\tilde{\mathbf{y}} - \mathbf{S}\mathbf{K}\boldsymbol{\alpha}\|_2^2 + \mu\boldsymbol{\alpha}^T\mathbf{K}\boldsymbol{\alpha} \quad (19)$$

where  $\mu > 0$  is a regularization parameter. The solution to (18) is

$$\hat{\mathbf{f}} = \mathbf{K}\mathbf{S}^T(\mathbf{S}\mathbf{K}\mathbf{S}^T + \mu\mathbf{I})^{-1}\mathbf{S}\mathbf{y}. \quad (20)$$

As shown in Section 2, the sampling operator in (2) takes the form  $\mathcal{A} := \mathbf{S}$ , where  $\mathbf{S}$  is an  $s \times N$  binary sampling matrix with a single non-zero element per row such that  $\mathbf{S}\mathbf{f} = \hat{\mathbf{f}}$ . Its adjoint, derived following (6), is  $\mathcal{A}^* := \mathbf{K}\mathbf{S}^T(\mathbf{S}\mathbf{K}\mathbf{S}^T)^{-1}$ , and the projection operator is  $\mathcal{P} = \mathbf{K}\mathbf{S}^T(\mathbf{S}\mathbf{K}\mathbf{S}^T)^{-1}\mathbf{S}$ . Thus, one observes that the solution (20) is not obtained by applying the projection operator  $\mathcal{P}$  to  $\mathbf{y}$  unless  $\mu = 0$ . However, since (20) lies on the column span of  $\mathbf{K}\mathbf{S}^T$ , the result is a non-orthogonal projection of  $\mathbf{y}$  onto  $\mathcal{H}_s = \operatorname{span}\{\mathbf{K}\mathbf{S}^T\}$ , i.e., the space spanned by the kernel functions evaluated at the sampled inputs. Still, a direct application of Theorem 1 is not possible unless  $\mu = 0$  and  $\mathbf{y}$  is noiseless.

To evaluate the performance of a specific sampling pattern, the risk function [16], which is equal to the mean squared error (MSE) except that it takes  $\mathbf{S}$  as a parameter, will be used. Thus, the risk is

$$\begin{aligned} R_\mu(\mathbf{S}) &= \mathbb{E}_{\mathbf{w}}\{\|\mathbf{f} - \hat{\mathbf{f}}\|_2^2\} \\ &= \|(\mathbf{I} - \mathbf{K}\mathbf{S}^T(\mathbf{S}\mathbf{K}\mathbf{S}^T + \mu\mathbf{I})^{-1}\mathbf{S})\mathbf{f}\|_2^2 \\ &\quad + \mathbb{E}_{\mathbf{w}}\left\{\|\mathbf{K}\mathbf{S}^T(\mathbf{S}\mathbf{K}\mathbf{S}^T + \mu\mathbf{I})^{-1}\mathbf{S}\mathbf{w}\|_2^2\right\} \end{aligned} \quad (21)$$

where  $\mathbf{w} = [w_1, \dots, w_N]^T$ . When  $\mu = 0$ , the risk becomes

$$\begin{aligned} R_0(\mathbf{S}) &= \|\mathbf{f} - \mathbf{K}\mathbf{S}^T(\mathbf{S}\mathbf{K}\mathbf{S}^T)^{-1}\mathbf{S}\mathbf{f}\|_2^2 + \mathbb{E}_{\mathbf{w}}\left\{\|\mathbf{K}\mathbf{S}^T(\mathbf{S}\mathbf{K}\mathbf{S}^T)^{-1}\mathbf{S}\mathbf{w}\|_2^2\right\} \\ &= \|\mathbf{f} - \mathcal{P}\mathbf{f}\|_2^2 + \mathbb{E}_{\mathbf{w}}\left\{\|\mathcal{P}\mathbf{w}\|_2^2\right\}. \end{aligned} \quad (22)$$

Hence, the first term in (22), i.e., the squared bias, takes a similar form as (15) in Theorem 1, whereas the second term or variance is the norm of the projection of the vectorized noise onto  $\mathcal{H}_s$ .

In a noiseless situation, the variance term in (22) would disappear and therefore the overall error would depend uniquely on the accuracy of the Nyström approximation. With noise, there needs to be a balance between bias and variance in order to make the risk small, i.e., sampling the most important points in  $\mathbf{f}$  while acquiring the least amount of noise. However, since the distribution of the noise is usually unknown, this strategy is not feasible. Assuming zero-mean Gaussian noise with variance  $\nu^2$ , then  $\mathbb{E}_{\mathbf{w}}\{\|\mathcal{P}\mathbf{w}\|_2^2\} = \nu^2 \operatorname{Tr}(\mathcal{P}) = \nu^2 s$  since  $\operatorname{Tr}(\mathcal{P}) = \operatorname{Tr}((\mathbf{S}\mathbf{K}\mathbf{S}^T)^{-1}\mathbf{S}\mathbf{K}\mathbf{S}^T) = s$ . Adapting the procedure in the proof of Theorem 1 to the 2-norm, we have that

$$R_0(\mathbf{S}) = \|\mathbf{f} - \mathcal{P}\mathbf{f}\|_2^2 + \nu^2 \operatorname{Tr}(\mathcal{P}) \leq \|\mathbf{f}\|_2^2 \operatorname{Tr}(\mathbf{K} - \mathbf{T}) + \nu^2 s$$

This expression shows that adding more samples reduces the approximation error to  $\mathbf{K}$ , and hence the bias, but it also increases the variance.

Since in KRR one typically uses  $\mu > 0$  to reduce the impact of the noise and risk of overfitting, the passive sampling strategy is derived in the next section by analyzing  $R_\mu(\mathbf{S})$  in (21).

#### 4. Passive sampling for kernel ridge regression

The previous sections have shown that the upper bound on the recovery error increases with the trace of the difference between the kernel matrix and its Nyström approximation. Since this approximation is built from a subset of columns of the kernel matrix, it is crucial to choose the columns that best approximate the spectrum of the full matrix to keep the approximation accurate. Indeed, this is known as the column subset selection problem [17], also related to matrix sketching [18–20]. In regards to the Nyström

approximation, while there exist deterministic [21] methods which order the columns according to a metric and then choose the top performers, most implementations opt for statistical approaches. In these, a sampling probability is assigned to each column and the required number of columns is sampled according to the resulting distribution. One example of such a technique is in [13,16], which measures column importance through the so-called leverage scores and provides theoretical guarantees on the approximation error. This section presents a passive sampling method based on kernel leverage scores, which is derived through an analysis of the estimation error in KRR.

##### 4.1. Risk function analysis

Rewriting the risk function as shown in Appendix A, assuming Gaussian noise with variance  $\nu^2$ , and using that  $\|\mathbf{A}\boldsymbol{\alpha}\|_2 \leq \|\mathbf{A}\|_2\|\boldsymbol{\alpha}\|_2$  for any matrix  $\mathbf{A}$  of appropriate size, the risk is upper bounded as

$$\begin{aligned} R_\mu(\mathbf{S}) &= \|(\mathbf{K} - \tilde{\mathbf{T}})\boldsymbol{\alpha}\|_2^2 + \mathbb{E}_{\mathbf{w}}\left\{\frac{1}{\mu^2}\|(\mathbf{K} - \tilde{\mathbf{T}})\mathbf{S}^T\tilde{\mathbf{w}}\|_2^2\right\} \\ &\leq \|\mathbf{K} - \tilde{\mathbf{T}}\|_2^2 \left(\|\boldsymbol{\alpha}\|_2^2 + \frac{\nu^2 s}{\mu^2}\right) \end{aligned} \quad (23)$$

where  $\tilde{\mathbf{T}} := \mathbf{K}\mathbf{S}^T(\mathbf{S}\mathbf{K}\mathbf{S}^T + \mu\mathbf{I})^{-1}\mathbf{S}\mathbf{K}$  is a regularized Nyström approximation.

Equation (23) shows that the spectral norm of the difference between the kernel matrix and its regularized Nyström approximation, which depends on  $\mu$ , also determines the risk bound for KRR. Therefore, it is well-grounded that one may pick the training inputs in  $\mathcal{X}$  by choosing the associated columns in  $\mathbf{K}$  that best build the regularized Nyström approximation. This section frames the sampling of the entries in  $\mathbf{y}$  as designing the sampling matrix  $\mathbf{S}$  that minimizes the error  $\mathbf{K} - \tilde{\mathbf{T}}$ . Through analyzing the risk function from a probabilistic standpoint,  $\mathbf{S}$  is cast as a weighted sampling matrix with weights set according to the leverage scores of the columns of  $\mathbf{K}$ . Since this process only involves the kernel matrix, it is a passive sampling approach. Once  $\mathbf{S}$  is obtained, the corresponding samples are obtained as  $\tilde{\mathbf{y}} = \mathbf{S}\mathbf{y}$  and KRR can then be applied to recover  $\mathbf{f}$ .

Let us first introduce the following lemma proven in Appendix B:

**Lemma 2.** *The difference  $\mathbf{K} - \tilde{\mathbf{T}}$  can be written as*

$$\mathbf{K} - \tilde{\mathbf{T}} = \mu\mathbf{Q}\boldsymbol{\Sigma}^{\frac{1}{2}}(\boldsymbol{\Sigma} + \mu\mathbf{I})^{-\frac{1}{2}}(\mathbf{I} - \mathbf{P})^{-1}(\boldsymbol{\Sigma} + \mu\mathbf{I})^{-\frac{1}{2}}\boldsymbol{\Sigma}^{\frac{1}{2}}\mathbf{Q}^T \quad (24)$$

where we have used the eigendecomposition  $\mathbf{K} = \mathbf{Q}\boldsymbol{\Sigma}\mathbf{Q}^T$  and

$$\mathbf{P} := \boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \mu\mathbf{I})^{-1} - (\boldsymbol{\Sigma} + \mu\mathbf{I})^{-\frac{1}{2}}\boldsymbol{\Sigma}^{\frac{1}{2}}\mathbf{Q}^T\mathbf{S}^T\mathbf{S}\mathbf{Q}\boldsymbol{\Sigma}^{\frac{1}{2}}(\boldsymbol{\Sigma} + \mu\mathbf{I})^{-\frac{1}{2}}. \quad (25)$$

Defining

$$\mathbf{V} = (\boldsymbol{\Sigma} + \mu\mathbf{I})^{-\frac{1}{2}}\boldsymbol{\Sigma}^{\frac{1}{2}}\mathbf{Q}^T \quad (26)$$

and knowing that the product of two diagonal matrices is commutative, substituting into (25) yields

$$\mathbf{P} = \mathbf{V}\mathbf{V}^T - \mathbf{V}\mathbf{S}^T\mathbf{S}\mathbf{V}^T \quad (27)$$

Therefore,  $\mathbf{P}$  is the difference between  $\mathbf{V}\mathbf{V}^T$  and an approximation built from a subset of the columns in  $\mathbf{V}$ , and it is the only element in (24) depending on the sampling distribution through  $\mathbf{S}$ . Thus, if  $\mathbf{S}$  is designed to minimize the norm of  $\mathbf{P}$  so that in (24)  $\|(\mathbf{I} - \mathbf{P})^{-1}\|_2$  is approximately minimized, a reduction on the difference  $\mathbf{K} - \tilde{\mathbf{T}}$  in (24) and also on the bound on  $R_\mu(\mathbf{S})$  in (23) will be achieved.

#### 4.2. Sampling based on leverage scores

Let  $\theta$  be an ordered  $s$ -tuple containing indices drawn with replacement from  $\{1, \dots, N\}$  with the probability distribution  $\Pr(i) = p_i$ . Here,  $p_i$  denotes the probability of sampling the  $i$ th column in  $\mathbf{K}$ , and  $\theta$  contains the indices of the sampled columns in non-descending order. Moreover, assume a weighted  $\mathbf{S} \in \mathbb{R}^{s \times N}$  with entries

$$\mathbf{S}_{i,j} = \begin{cases} \frac{1}{\sqrt{s p_j}} & \text{if } \theta(i) = j \\ 0 & \text{otherwise} \end{cases} \quad (28)$$

for  $i = 1, \dots, s$  and  $j = 1, \dots, N$ . Note that, since the indices in  $\theta$  are drawn with replacement, there might be repeated rows in  $\mathbf{S}$  sampling the same item. Moreover, a weighted  $\mathbf{S}$  is necessary in order to satisfy certain theoretical guarantees.

In [22] it is shown that the optimal probability distribution to minimize  $\mathbb{E}\{\|\mathbf{P}\|_F^2\}$  in (27) is

$$p_i = \frac{\|\mathbf{v}_i\|_2^2}{\|\mathbf{V}\|_F^2} \quad (29)$$

with  $\mathbf{v}_i$  denoting the  $i$ th column of  $\mathbf{V}$  in (26), and

$$\|\mathbf{v}_i\|_2^2 = \sum_{j=1}^N \frac{\sigma_j}{\sigma_j + \mu} \mathbf{q}_j^2(i) \quad (30)$$

where  $\sigma_i = (\boldsymbol{\Sigma})_{ii}$ , and  $\mathbf{q}_j$  is the  $j$ th column of  $\mathbf{Q}$ . Thus, the distribution obtained with (29) for  $i = 1, \dots, N$  can be expected to reduce  $\|\mathbf{K} - \tilde{\mathbf{T}}\|$  in the risk upper bound (23). Still, note that this distribution is obtained for a specific value of  $\mu$  since  $p_i$  depends on  $\mu$  through  $\mathbf{V}$ .

The KRR regularization parameter  $\mu$  is usually chosen experimentally via cross-validation to minimize the prediction error on a testing dataset. However, as it will be shown in the numerical results, it is possible to further reduce the prediction error by fine-tuning the sampling distribution. For this reason, we decouple the optimal distribution from the regularization parameter in KRR by redefining the  $i$ th leverage score as

$$l_i = \sum_{j=1}^N \frac{\sigma_j}{\sigma_j + \gamma} \mathbf{q}_j^2(i) \quad (31)$$

where  $\gamma > 0$  is now a tunable parameter different from the parameter  $\mu$ . Alternatively, (31) can be obtained as  $l_i = (\mathbf{K}(\mathbf{K} + \gamma\mathbf{I})^{-1})_{i,i}$ . Then, the chosen distribution is

$$p_i = \frac{l_i}{\sum_{j=1}^N l_j} \quad (32)$$

Thus, replacing  $\mu$  in (30) with  $\gamma$  to obtain (31) enables separate optimization of the regularization term in the regression problem and the probability distribution (32). We refer to the use of this distribution to select the training set as *leverage score sampling*.

The quantity  $l_i$  is the so-called regularized leverage score<sup>1</sup> of the  $i$ th column of  $\mathbf{K}$ , a concept that often appears in the context of random matrix approximations. For instance, [13,16] derive error bounds on the Nyström approximation when the columns are sampled according to their leverage scores. Similarly, [23] relies on the leverage scores to sample a subset of columns of the kernel matrix and reduce the computational cost of KRR. In [24], leverage scores are used to assess the approximation of the kernel matrix via random Fourier features, whereas [18] uses them to obtain accurate subspace embeddings. In contrast to these related works, our focus is not on obtaining the best possible kernel matrix approximation to reduce computational cost in KRR. We set out to

derive an optimal approach to select the samples in the training set, and the analysis of the risk function has led us to the leverage scores (31) as a good approach.

Since the  $i$ th leverage score is a weighted average of the  $i$ th row of the eigenvector matrix  $\mathbf{Q}$ , a high value indicates that the  $i$ th input stands out with respect to the other elements. Moreover,  $\gamma$  determines how much weight is given to the notable points, with  $\gamma = 0$  inducing a uniform probability distribution. Regarding computational cost, it is dominated by the inversion of an  $N \times N$  matrix, which could be reduced by replacing the matrix to be inverted with a low-rank approximation. This is for instance proposed in [16] where, given an  $N \times d$  binary sampling matrix  $\mathbf{D}$ , an approximate leverage score is obtained as

$$\tilde{l}_i = \frac{1}{\gamma} (\mathbf{K} - \mathbf{K}\mathbf{D}(\mathbf{D}^T\mathbf{K}\mathbf{D} + \gamma\mathbf{I})^{-1}\mathbf{S}^T\mathbf{K}). \quad (33)$$

If the columns of  $\mathbf{K}$  are sampled according to  $p_i$  in (32), the risk for a given  $\mathbf{S}$  can be bounded in relation to the risk of the full problem. To do so, let us first introduce the following theorem adapted from [13] to bound the spectral norm of  $\mathbf{P}$  in (27):

**Theorem 2.** [13] Let  $\mathbf{V} \in N \times N$ , and  $\mathbf{S}$  be an  $s \times N$  sparse weighted sampling matrix with  $\mathbf{S}_{i,j} = \frac{1}{\sqrt{s p_j}}$  if  $\theta(i) = j \forall i = 1, \dots, s, j \in \{1, \dots, N\}$ , where  $p_j$  denotes the probability of choosing the  $j$ th column in  $\mathbf{K}$ . Then,

$$\Pr(\|\mathbf{V}\mathbf{V}^T - \mathbf{V}\mathbf{S}^T\mathbf{S}\mathbf{V}^T\|_2 \geq t) \leq N \exp\left(\frac{-st^2/2}{\|\mathbf{V}\|_F^2 + t/3}\right). \quad (34)$$

Using Theorem 2,  $\mathbf{P} = \mathbf{V}\mathbf{V}^T - \mathbf{V}\mathbf{S}^T\mathbf{S}\mathbf{V}^T \leq t\mathbf{I}$  with probability greater than  $1 - \delta$  if  $s \geq \left(\frac{2d_{eff}}{t^2} + \frac{2}{3t}\right) \log \frac{N}{\delta}$ , where  $d_{eff}$  is the so-called effective dimension defined as  $d_{eff} = \sum_{i=1}^N l_i = \|\mathbf{V}\|_F^2$ , and  $\delta \in (0, 1]$ ; this condition on  $s$  is obtained by equating the right-hand side of (34) to  $\delta$ . Then, we have that  $(\mathbf{I} - \mathbf{P})^{-1} \leq \frac{1}{1-t}\mathbf{I}$  and, therefore, the eigenvalues of  $\mathbf{K} - \tilde{\mathbf{T}}$  in (24) are bounded as

$$\mathbf{K} - \tilde{\mathbf{T}} \leq \frac{\mu}{1-t}\mathbf{I} \quad (35)$$

with probability at least  $1 - \delta$ . This inequality is obtained by applying to (24) the property that the maximum eigenvalue of a matrix multiplication is smaller or equal to the multiplication of the maximum eigenvalue of each matrix [10]. Using (35), the risk is written in relation to the risk of the fully observed problem, i.e.,  $\mathbf{S} = \mathbf{I}$ , as Theorem 3 below shows with proof in Appendix C.

**Theorem 3.** Let  $\mathbf{f} = \mathbf{K}\boldsymbol{\alpha}$  be the function vector to be recovered,  $R_\mu(\cdot)$  the KRR risk function in (21), and  $\mathbf{S}$  a  $s \times N$  sparse weighted sampling matrix with  $\mathbf{S}_{i,j} = \frac{1}{\sqrt{s p_j}}$  if  $\theta(i) = j \forall i = 1, \dots, s, j \in \{1, \dots, N\}$ , where  $p_j$  denotes the probability of choosing the  $j$ th column in  $\mathbf{K}$ . Provided that  $s \geq \left(\frac{2d_{eff}}{t^2} + \frac{2}{3t}\right) \log \frac{N}{\delta}$ , it holds with probability greater than  $1 - \delta$  for  $\delta \in (0, 1]$  that the risk is upper bounded as

$$R_\mu(\mathbf{S}) \leq R_\mu(\mathbf{I}) + \frac{\mu t}{1-t} \|\boldsymbol{\alpha}\|_2^2 + \frac{\nu^2}{(1-t)^2} \text{Tr}(\mathbf{S}^T\mathbf{S} - \mathbf{I}). \quad (36)$$

Theorem 3 indicates that the excess risk when using  $s < N$  samples increases with the spectral norm of  $\mathbf{P}$  assumed smaller than  $t$ , with the second term in (36) being related to the excess bias, and the third to the excess variance; see (52) and (55) in Appendix C. As expected, increasing  $s$  reduces the bias through inducing a smaller  $t$ , whereas a larger number of samples may result in higher variance due to the additional noise.

While it does not have the same theoretical guarantees as the probabilistic approach, a greedy approach can also be derived using the leverage scores. Greedy methods are concerned with maximizing optimality in the stage at which they are executed, and

<sup>1</sup> Hereafter the words leverage score refer to the regularized leverage score.

do not consider the optimality of the final solution. Applied to the passive sampling paradigm we have developed, this means simply taking the samples with the highest associated leverage score. Defining  $\mathcal{L}$  as the set containing the  $s$  leverage scores with the highest value, the  $s$ -tuple of observation indices for *greedy leverage sampling* is denoted by

$$\theta_G = (i)_{i \in \mathcal{L}}. \tag{37}$$

**Remark 1.** Sampling based on leverage scores is applicable to the selection of inducing points in Gaussian processes. Assuming that  $f$  is a Gaussian process (GP) such that  $\Pr(f|\mathcal{X}) \sim \mathcal{N}(\mathbf{0}, \mathbf{K})$ , the distribution of  $f$  conditioned on the observations, inputs and variance  $v^2$  is

$$\Pr(f|\bar{\mathbf{m}}, \mathcal{X}, v^2) \sim \mathcal{N}[\mathbf{K}\mathbf{S}^T(\bar{\mathbf{K}} + v^2\mathbf{I})^{-1}\bar{\mathbf{y}}, \mathbf{K} - \mathbf{K}\mathbf{S}^T(\bar{\mathbf{K}} + v^2\mathbf{I})^{-1}\mathbf{K}\mathbf{S}]. \tag{38}$$

We observe in (38) that the mean is equal to a KRR with  $\mu = v^2$ , and that the covariance matrix is the difference between  $\mathbf{K}$  and its Nyström approximation. Hence, this suggests that an accurate approximation through the design of  $\mathbf{S}$  should yield a lower variance on the posterior distribution of  $f$ . Such a design can be produced by means of the leverage scores as presented in this section, which in turn induces a sampling distribution through  $\bar{\mathbf{y}}$ . Hence, to obtain either a leverage score sampling or greedy leverage sampling scheme for GPs, one just needs to calculate the leverage scores as in (31) for a specific  $\gamma$  and follow the indicated steps for each of the two methods.

### 5. Numerical tests

This section presents numerical tests comparing uniform sampling with  $p_i = \frac{1}{N}$  to leverage score sampling (32) and greedy leverage score sampling in (37) for different values of  $s$ , for both KRR and GP regression. Each combination of sampling method and algorithm is labeled as UKRR, LKRR, GKRR, UGP, LGP and GGP, where U, L and G stand for uniform, leverage and greedy, respectively. The optimal hyperparameters  $\mu$  and  $\gamma$  are found via grid search for each  $s$ . The search is performed randomly: intervals  $[\mu_{\min}, \mu_{\max}]$  and  $[\gamma_{\min}, \gamma_{\max}]$  are set for  $\mu$  and  $\gamma$  respectively, and 100 tuples of values within the intervals are selected at random. For the GP algorithm, the noise variance  $v^2$  is estimated through the log-likelihood. The  $s$  labeled samples form the training set, and the RMSE measures the estimation error as

$$\text{RMSE} = \sqrt{\frac{1}{N_r} \sum_{i=1}^{N_r} \|\hat{\mathbf{f}}_i - \mathbf{f}\|_F^2} \tag{39}$$

where  $N_r = 20$  is the number of realizations,  $\hat{\mathbf{f}}_i$  is the estimation at the  $i$ th realization, and a new training set is acquired at each realization. For KRR,  $\hat{\mathbf{f}}_i$  is given by (20), whereas for a GP it is given by the posterior mean in (38) with  $v$  obtained via maximum-likelihood estimation from the observed samples.

The Boston housing dataset<sup>2</sup> is the first to be evaluated. The input set  $\mathcal{X}$  is comprised of 506 feature vectors detailing the characteristics of 506 houses in Boston such as size or number of rooms. The variable to be predicted, i.e.,  $\mathbf{f}$ , is the price of each house. Here, the sampling algorithm selects a training set  $s$  of houses with their features, requests their prices, and estimates the price of the remaining ones. The used kernel is the Gaussian kernel applied on the feature vectors; its expression is  $\mathbf{K}_{i,j} = \exp(-\|x_i - x_j\|_2^2 / \xi)$ , where  $\xi > 0$ . Fig. 1 shows the RMSE for different values of  $s$  and the three sampling strategies. The figure shows that leverage score sampling attains a smaller error than uniform sampling for both

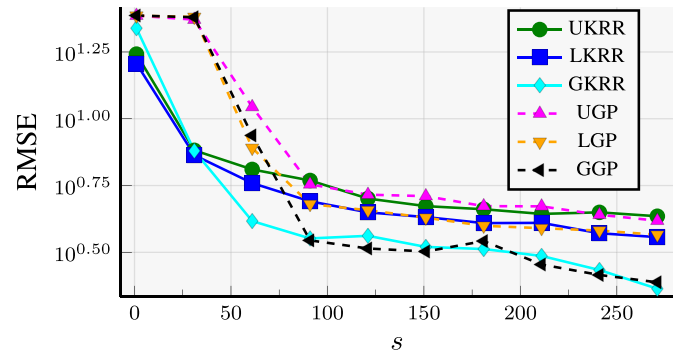


Fig. 1. NMSE vs.  $s$  for the Boston housing data.

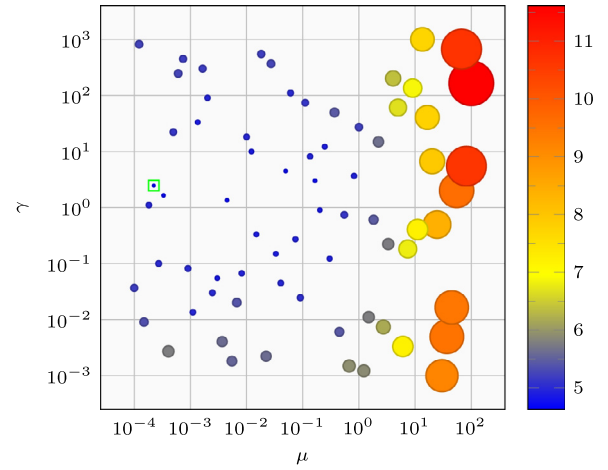


Fig. 2. Grid search for the Boston housing dataset for  $s = 91$ . The minimum is highlighted in a green square, centered at  $\mu^* = 2.2 \cdot 10^{-4}$ ,  $\gamma^* = 2.46$ . (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

KRR and GP regression. While the gains are small, they are in line with those obtained by other methods [3,13]. Nevertheless, for this particular dataset greedy sampling provides better performance for both algorithms.

Fig. 2 depicts the result of the grid search for the Boston housing dataset for the LKRR method with  $s = 91$ . We observe that, out of all the evaluated points, the combination of  $\mu$  and  $\gamma$  minimizing the error is  $\mu^* = 2.2 \cdot 10^{-4}$ ,  $\gamma^* = 2.46$ . This plot corroborates that decoupling the hyperparameter for the design of the sampling distribution from the one in the KRR regularization yields a smaller error.

To gain further insight into the impact of  $\gamma$ , Fig. 3 shows the sampling probability distribution for several values of  $\gamma$ . We observe that setting a very small value or one large enough results in a more uniform distribution, although for  $\gamma = 10$  the largest peaks still remain visible. On the other hand, an in-between value enhances the variability between probabilities. Thus, in the specific case in Fig. 2 with  $s = 91$ , the optimal  $\gamma$  is the one that maintains the largest sampling probabilities while promoting uniformity among the rest.

Fig. 4 shows the results for the Mushroom dataset. The vector to be recovered is of length 5,000 and the kernel is obtained from the Pearson correlation matrix of the features. Here, the feature vectors contain details about each mushroom such as shape or size, and the objective is to determine whether a mushroom is edible or poisonous. For KRR, both greedy and leverage score sampling show a much lower error than uniform sampling. On the other hand, in GP regression greedy sampling results in worse per-

<sup>2</sup> <https://www.cs.toronto.edu/~delve/data/boston/bostonDetail.html>.

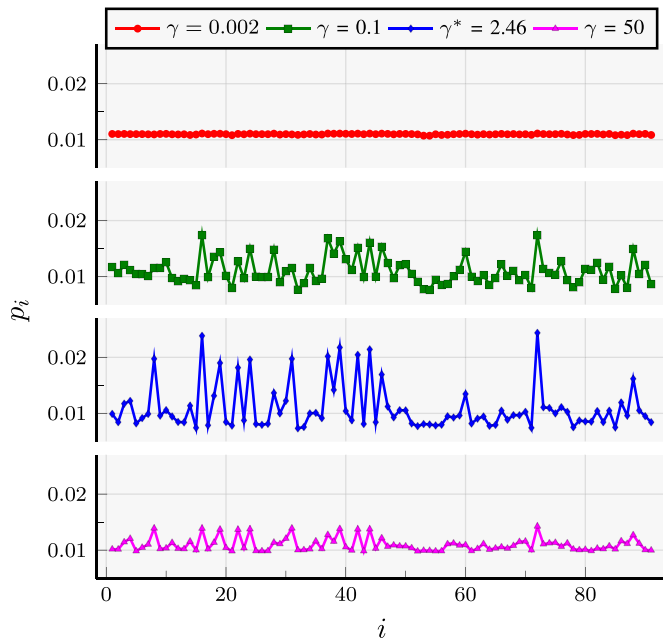


Fig. 3. Leverage scores and probabilities for the Boston housing dataset for different values of  $\gamma$ .

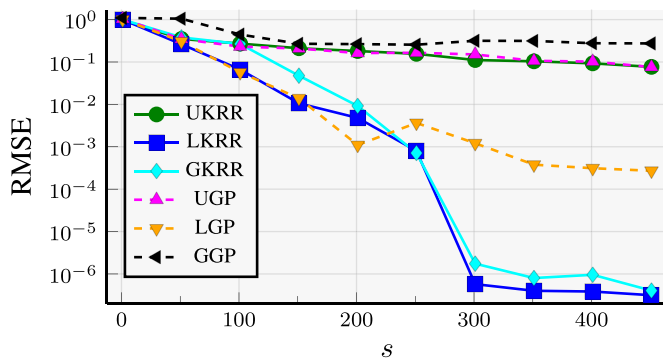


Fig. 4. RMSE vs.  $s$  for the Mushroom data.

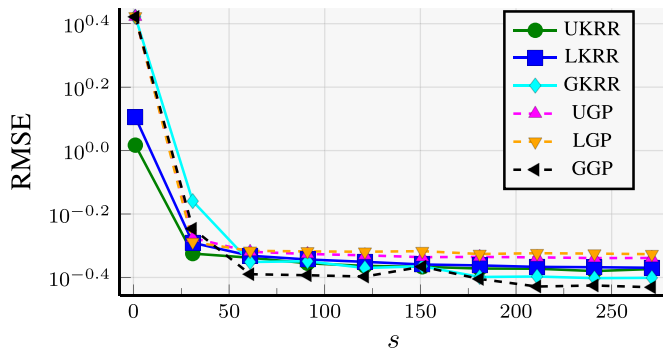


Fig. 5. RMSE vs.  $s$  for the Labor data.

formance when compared to uniform sampling. Note that, since this is a binary classification problem with classes  $[2, -1]$ , the NMSE evaluates the difference between the regression result and the actual numerical value of the class before applying any decision rule.

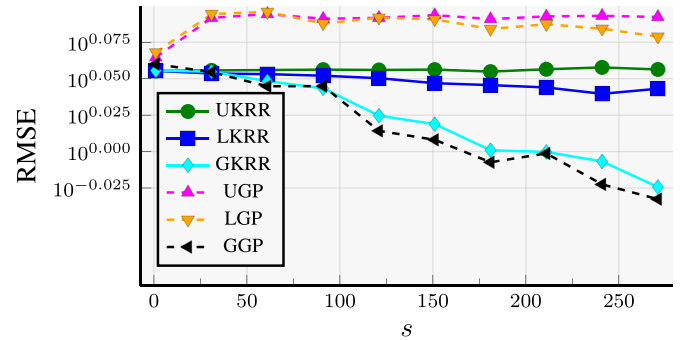


Fig. 6. RMSE vs.  $s$  for the Stocks data.

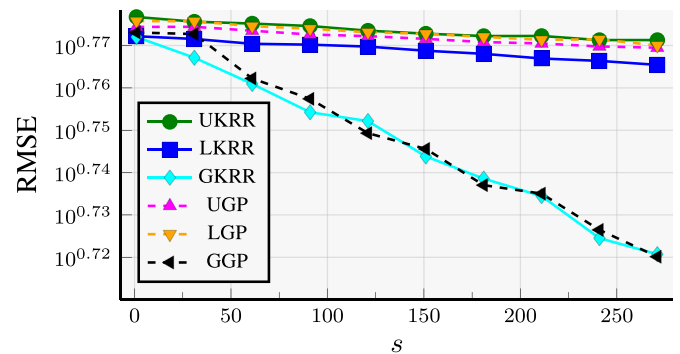


Fig. 7. RMSE vs.  $s$  for the Wine data.

Figs. 5, 6 and 7 show the test results for the Labor<sup>3</sup>, Stocks<sup>4</sup> and Wine<sup>5</sup> datasets, respectively, using the Gaussian kernel. Overall, we observe that the passive sampling approaches are able to reduce the error with respect to uniform sampling, with the reduction being more noticeable for greedy sampling in the Stocks and Wine datasets.

We can conclude that, for KRR, both the random leverage sampling and the deterministic greedy sampling result in smaller error than uniform sampling for the test data. However, greedy sampling is more vulnerable to a badly designed kernel and outliers as it always chooses the most distinct points; this can be mitigated in leverage sampling by setting a more uniform distribution. Moreover, it should be noted that the reported results show the average error over several realizations for leverage sampling whereas for greedy sampling they correspond to a single realization since the sampled set is deterministic.

### 5.1. Computational cost evaluation

In order to reduce the computational cost, one may resort to approximate leverage scores. Fig. 8 shows the time required to complete the passive sampling using the approximate scores in (33) plus the KRR (18) for the Mushroom data. We observe that the cost is reduced in proportion to the factor  $\frac{d}{N}$ . Moreover, since the cost to calculate the leverage scores is constant for different  $s$  as it only depends on  $\frac{d}{N}$ , the increase in time with  $s$  is due to the KRR step. Fig. 9 shows the resulting probability distribution for the different factor values and first 100 indices. Finally, Fig. 10 compares the error for KRR using exact leverage scores, and approximate leverage scores with  $\frac{d}{N} = 0.8$ . We observe that at this value

<sup>3</sup> <https://rdrr.io/rforge/Ecdat/man/Mroz.html>.

<sup>4</sup> <https://rdrr.io/cran/ISLR/man/Smarket.html>.

<sup>5</sup> <https://archive.ics.uci.edu/ml/datasets/wine>.

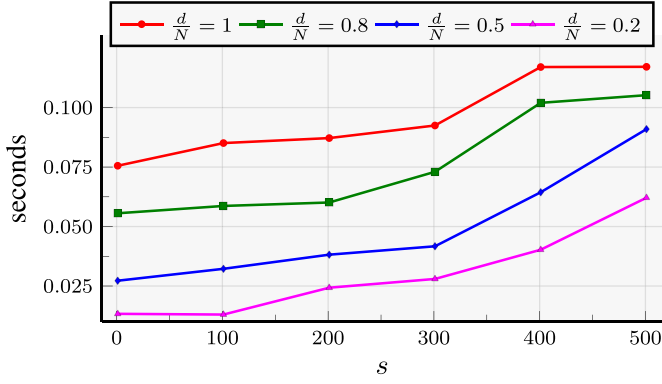


Fig. 8. Runtime for different values of  $\frac{d}{N}$  on the Mushroom data.

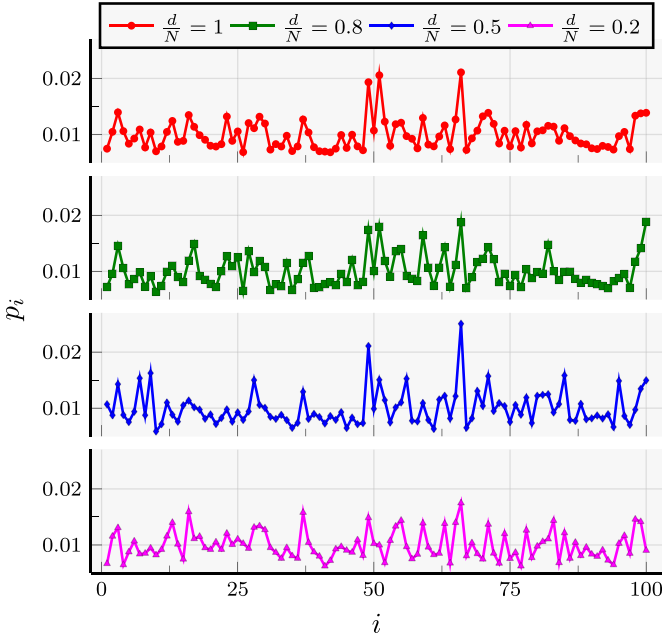


Fig. 9. Sampling probabilities of the first 100 points in the Mushroom data for different  $\frac{d}{N}$ .

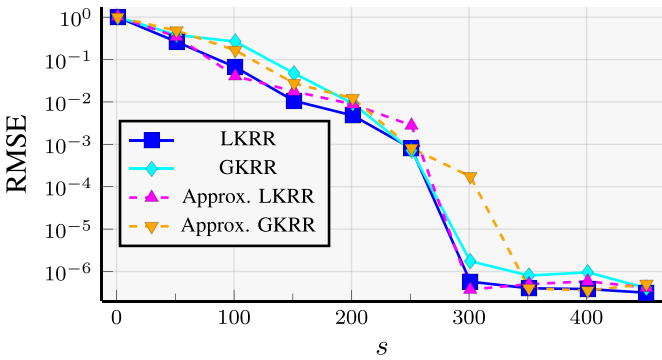


Fig. 10. RMSE vs.  $s$  for the Mushroom data using exact leverage scores, and approximate leverage scores with  $\frac{d}{N} = 0.8$ .

the performance is not significantly degraded, while the computational cost is reduced.

## 6. Conclusions

This paper has studied how the prior information embedded in the kernel functions spanning an RKHS can be used to deter-

mine which inputs are the most important to label in order to recover the complete function. First, through functional analysis, it has been shown that the recovery error of a sampled function in an RKHS is directly tied to the Nyström approximation to the kernel matrix. Hence, in the noiseless case, the error can be minimized by finding the best possible Nyström approximation and labeling the inputs corresponding to the chosen columns. This sampling approach is passive since it only involves knowledge of the input space and kernel matrix, which presents benefits over active sampling schemes which require online operation and are more vulnerable to noise. Given the theoretical background, the Nyström-based sampling approach has been applied to KRR and GP regression for the picking of the training set. In this, the weighted sampling matrix is designed according to the leverage scores of the kernel matrix, which measure the importance of each column in achieving a good Nyström approximation. Numerical tests have shown that the proposed approaches work well for the recovery of sampled vectors when compared to uniform sampling.

## Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## Acknowledgements

The authors would like to thank Prof. Daniel Romero from University of Adger, and Prof. Antonio G. Marqués from Universidad Rey Juan Carlos for their valuable insights.

## Appendix A. Risk function for KRR

For KRR, the MSE is

$$MSE := \mathbb{E}_{\mathbf{w}} \{ \|\mathbf{f} - \hat{\mathbf{f}}\|_2^2 \} = \mathbb{E}_{\mathbf{w}} \left\{ \left\| \mathbf{f} - \mathbf{K}\hat{\boldsymbol{\alpha}} \right\|_2^2 \right\}. \quad (39)$$

Plugging the estimator from (20) into (39) yields

$$\begin{aligned} MSE &= \mathbb{E}_{\mathbf{w}} \left\{ \left\| \mathbf{f} - \mathbf{K}\mathbf{S}^T (\mathbf{K}\mathbf{S}\mathbf{S}^T + \mu\mathbf{I})^{-1} (\mathbf{S}\mathbf{f} + \tilde{\mathbf{w}}) \right\|_2^2 \right\} \\ &= \left\| (\mathbf{I} - \mathbf{K}\mathbf{S}^T (\mathbf{K}\mathbf{S}\mathbf{S}^T + \mu\mathbf{I})^{-1} \mathbf{S}) \mathbf{f} \right\|_2^2 \\ &\quad + \mathbb{E}_{\mathbf{w}} \left\{ \left\| \mathbf{K}\mathbf{S}^T (\mathbf{K}\mathbf{S}\mathbf{S}^T + \mu\mathbf{I})^{-1} \tilde{\mathbf{w}} \right\|_2^2 \right\} \end{aligned} \quad (40)$$

where  $\tilde{\mathbf{w}} = \mathbf{S}\mathbf{w}$ , and we have used that  $\mathbb{E}\{\mathbf{w}\} = \mathbf{0}$ . Further, the first and second terms in (40) are the squared bias and variance of the KRR estimator, respectively. If we substitute  $\mathbf{f} = \mathbf{K}\boldsymbol{\alpha}$  into the first term of (40), we obtain the squared bias as

$$\begin{aligned} b^2 &= \left\| (\mathbf{I} - \mathbf{K}\mathbf{S}^T (\mathbf{K}\mathbf{S}\mathbf{S}^T + \mu\mathbf{I})^{-1} \mathbf{S}) \mathbf{K}\boldsymbol{\alpha} \right\|_2^2 \\ &= \left\| (\mathbf{K} - \mathbf{K}\mathbf{S}^T (\mathbf{K}\mathbf{S}\mathbf{S}^T + \mu\mathbf{I})^{-1} \mathbf{S}\mathbf{K}) \boldsymbol{\alpha} \right\|_2^2 \\ &= \left\| (\mathbf{K} - \tilde{\mathbf{T}}) \boldsymbol{\alpha} \right\|_2^2 \end{aligned} \quad (41)$$

where  $\tilde{\mathbf{T}} := \mathbf{K}\mathbf{S}^T (\mathbf{K}\mathbf{S}\mathbf{S}^T + \mu\mathbf{I})^{-1} \mathbf{S}\mathbf{K}$  is the regularized Nyström approximation of  $\mathbf{K}$ . On the other hand, the variance term is

$$\begin{aligned} var &= \mathbb{E}_{\mathbf{w}} \left\{ \left\| \mathbf{K}\mathbf{S}^T (\mathbf{K}\mathbf{S}\mathbf{S}^T + \mu\mathbf{I})^{-1} \tilde{\mathbf{w}} \right\|_2^2 \right\} \\ &= \mathbb{E}_{\mathbf{w}} \left\{ \frac{1}{\mu^2} \left\| \mathbf{K}\mathbf{S}^T (\mathbf{K}\mathbf{S}\mathbf{S}^T + \mu\mathbf{I})^{-1} (\mu\mathbf{I} + \mathbf{K}\mathbf{S}\mathbf{S}^T - \mathbf{K}\mathbf{S}\mathbf{S}^T) \tilde{\mathbf{w}} \right\|_2^2 \right\} \\ &= \mathbb{E}_{\mathbf{w}} \left\{ \frac{1}{\mu^2} \left\| \mathbf{K}\mathbf{S}^T - \mathbf{K}\mathbf{S}^T (\mathbf{K}\mathbf{S}\mathbf{S}^T + \mu\mathbf{I})^{-1} \mathbf{K}\mathbf{S}\mathbf{S}^T \tilde{\mathbf{w}} \right\|_2^2 \right\} \end{aligned}$$



$$\begin{aligned}
&= \mathbb{E}_{\mathbf{w}} \left\{ \frac{1}{\mu^2} \left\| (\mathbf{K} - \mathbf{K}\mathbf{S}^T (\mathbf{S}\mathbf{K}\mathbf{S}^T + \mu\mathbf{I})^{-1} \mathbf{S}\mathbf{K}) \mathbf{S}^T \tilde{\mathbf{w}} \right\|_2^2 \right\} \\
&= \mathbb{E}_{\mathbf{w}} \left\{ \frac{1}{\mu^2} \left\| (\mathbf{K} - \tilde{\mathbf{T}}) \mathbf{S}^T \tilde{\mathbf{w}} \right\|_2^2 \right\}. \quad (42)
\end{aligned}$$

Adding the two terms in (41) and (42), and assuming a fixed  $\mathbf{S}$ , we obtain the risk in (23).

### Appendix B. Proof of Lemma 2

With the eigendecomposition  $\mathbf{K} = \mathbf{Q}\Sigma\mathbf{Q}^T$ , we can write

$$\begin{aligned}
\mathbf{K} - \tilde{\mathbf{T}} &= \mathbf{K} - \mathbf{K}\mathbf{S}^T (\mathbf{S}\mathbf{K}\mathbf{S}^T + \mu\mathbf{I})^{-1} \mathbf{S}\mathbf{K} \\
&= \mathbf{Q}\Sigma^{\frac{1}{2}} \left[ \mathbf{I} - \Sigma^{\frac{1}{2}} \mathbf{Q}^T \mathbf{S}^T (\mathbf{S}\mathbf{Q}\Sigma^{\frac{1}{2}} \Sigma^{\frac{1}{2}} \mathbf{Q}^T \mathbf{S}^T + \mu\mathbf{I})^{-1} \mathbf{S}\mathbf{Q}\Sigma^{\frac{1}{2}} \right] \Sigma^{\frac{1}{2}} \mathbf{Q}^T. \quad (43)
\end{aligned}$$

Applying the Matrix Inversion Lemma to the matrix inside the square brackets of (43), we arrive at

$$\begin{aligned}
&\mathbf{I} - \Sigma^{\frac{1}{2}} \mathbf{Q}^T \mathbf{S}^T (\mathbf{S}\mathbf{Q}\Sigma^{\frac{1}{2}} \Sigma^{\frac{1}{2}} \mathbf{Q}^T \mathbf{S}^T + \mu\mathbf{I})^{-1} \mathbf{S}\mathbf{Q}\Sigma^{\frac{1}{2}} \\
&= \left( \mathbf{I} + \frac{1}{\mu} \Sigma^{\frac{1}{2}} \mathbf{Q}^T \mathbf{S}^T \mathbf{S}\mathbf{Q}\Sigma^{\frac{1}{2}} \right)^{-1}. \quad (44)
\end{aligned}$$

That in turn implies

$$\begin{aligned}
\mathbf{K} - \tilde{\mathbf{T}} &= \mu \mathbf{Q}\Sigma^{\frac{1}{2}} \left( \mu\mathbf{I} + \Sigma^{\frac{1}{2}} \mathbf{Q}^T \mathbf{S}^T \mathbf{S}\mathbf{Q}\Sigma^{\frac{1}{2}} \right)^{-1} \Sigma^{\frac{1}{2}} \mathbf{Q}^T \\
&= \mu \mathbf{Q}\Sigma^{\frac{1}{2}} \left( \Sigma + \mu\mathbf{I} - \Sigma + \Sigma^{\frac{1}{2}} \mathbf{Q}^T \mathbf{S}^T \mathbf{S}\mathbf{Q}\Sigma^{\frac{1}{2}} \right)^{-1} \Sigma^{\frac{1}{2}} \mathbf{Q}^T \\
&= \mu \mathbf{Q}\Sigma^{\frac{1}{2}} \left[ (\Sigma + \mu\mathbf{I})^{\frac{1}{2}} \left( \mathbf{I} - (\Sigma + \mu\mathbf{I})^{-\frac{1}{2}} \Sigma (\Sigma + \mu\mathbf{I})^{-\frac{1}{2}} \right. \right. \\
&\quad \left. \left. + (\Sigma + \mu\mathbf{I})^{-\frac{1}{2}} \Sigma^{\frac{1}{2}} \mathbf{Q}^T \mathbf{S}^T \mathbf{S}\mathbf{Q}\Sigma^{\frac{1}{2}} (\Sigma + \mu\mathbf{I})^{-\frac{1}{2}} \right) \right. \\
&\quad \left. (\Sigma + \mu\mathbf{I})^{\frac{1}{2}} \right]^{-1} \Sigma^{\frac{1}{2}} \mathbf{Q}^T \\
&= \mu \mathbf{Q}\Sigma^{\frac{1}{2}} (\Sigma + \mu\mathbf{I})^{-\frac{1}{2}} (\mathbf{I} - \mathbf{P})^{-1} (\Sigma + \mu\mathbf{I})^{-\frac{1}{2}} \Sigma^{\frac{1}{2}} \mathbf{Q}^T \quad (45)
\end{aligned}$$

with  $\mathbf{P}$  given in (25).

### Appendix C. Proof of Theorem 3

Let us first introduce the following lemma:

**Lemma 3.** Let  $\check{\mathbf{K}} = \mathbf{K}(\mathbf{K} + \mu\mathbf{I})^{-1}\mathbf{K}$ . Then,

$$\mathbf{K} - \check{\mathbf{K}} = \mu\mathbf{K}(\mathbf{K} + \mu\mathbf{I}) \quad (46)$$

**Proof.** Let

$$\mathbf{K} - \check{\mathbf{K}} = \mathbf{K} \left( \mathbf{I} - \frac{1}{\mu} \left( \frac{1}{\mu} \mathbf{K} + \mathbf{I} \right)^{-1} \mathbf{K} \right). \quad (47)$$

Using the Matrix Inversion Lemma,

$$\mathbf{A}^{-1} + \mathbf{A}^{-1} \mathbf{B} (\mathbf{D} - \mathbf{C} \mathbf{A}^{-1} \mathbf{B}) \mathbf{C} \mathbf{A} = (\mathbf{A} - \mathbf{B} \mathbf{D}^{-1} \mathbf{C})^{-1} \quad (48)$$

with  $\mathbf{A} = \mathbf{B} = \mathbf{D} = \mathbf{I}$  and  $\mathbf{C} = -\frac{1}{\mu} \mathbf{K}$ , we have that

$$\mathbf{K} \left( \mathbf{I} - \frac{1}{\mu} \left( \frac{1}{\mu} \mathbf{K} + \mathbf{I} \right)^{-1} \mathbf{K} \right) = \mathbf{K} \left( \frac{1}{\mu} \mathbf{K} + \mathbf{I} \right)^{-1} = \mu \mathbf{K} (\mathbf{K} + \mu\mathbf{I})^{-1}. \quad (49)$$

□

Knowing that  $R_{\mu}(\mathbf{S}) = R_{\mu}(\mathbf{I}) + (R_{\mu}(\mathbf{S}) - R_{\mu}(\mathbf{I}))$ , we first write the bias in terms of the one for the full problem, namely  $b_I$ , as

$$\begin{aligned}
b^2 &= \|(\mathbf{K} - \tilde{\mathbf{T}})\alpha\|_2^2 \leq \|(\mathbf{K} - \check{\mathbf{K}})\alpha\|_2^2 + \|(\check{\mathbf{K}} - \tilde{\mathbf{T}})\alpha\|_2^2 \\
&= b_I^2 + \|(\check{\mathbf{K}} - \tilde{\mathbf{T}})\alpha\|_2^2 \quad (50)
\end{aligned}$$

Provided that  $s \geq \left( \frac{2d_{\text{eff}}}{t^2} + \frac{2}{3t} \right) \log \frac{N}{\delta}$ , from (35) we know that  $\mathbf{K} - \tilde{\mathbf{T}} \leq \frac{\mu}{1-t} \mathbf{K}(\mathbf{K} + \mu\mathbf{I})^{-1}$  with probability  $1 - \delta$ . Hereafter assume that any equation involving  $t$  is satisfied with probability  $1 - \delta$ . Thus, to bound the last term in (50) we use Lemma 3 to show that

$$\begin{aligned}
\check{\mathbf{K}} - \tilde{\mathbf{T}} &= (\mathbf{K} - \tilde{\mathbf{T}}) - (\mathbf{K} - \check{\mathbf{K}}) \\
&\leq \frac{\mu}{1-t} \mathbf{K}(\mathbf{K} + \mu\mathbf{I})^{-1} - \mu \mathbf{K}(\mathbf{K} + \mu\mathbf{I})^{-1} \\
&= \frac{\mu t}{1-t} \mathbf{K}(\mathbf{K} + \mu\mathbf{I})^{-1} \leq \frac{\mu t}{1-t} \quad (51)
\end{aligned}$$

Applying the identity  $\|(\check{\mathbf{K}} - \tilde{\mathbf{T}})\alpha\|_2^2 \leq \lambda_{\max}(\check{\mathbf{K}} - \tilde{\mathbf{T}}) \|\alpha\|_2^2$  yields

$$b^2 = b_I^2 + \frac{\mu t}{1-t} \|\alpha\|_2^2. \quad (52)$$

Next, we proceed to calculate an upper bound to the difference between the variance of the sampled problem and that of the full problem, which is denoted as  $\text{var}_{\mathbf{K}}$ . Assuming Gaussian noise with variance  $v^2$  in (42), let  $\text{var} = \frac{v^2}{\mu^2} \text{Tr}((\mathbf{K} - \tilde{\mathbf{T}})^2 \mathbf{S}^T \mathbf{S})$  and  $\text{var}_I = \frac{v^2}{\mu^2} \text{Tr}((\mathbf{K} - \tilde{\mathbf{T}})^2)$ . Then,

$$\begin{aligned}
\text{var} - \text{var}_I &= \frac{v^2}{\mu^2} \text{Tr}((\mathbf{K} - \tilde{\mathbf{T}})^2 \mathbf{S}^T \mathbf{S}) - \frac{v^2}{\mu^2} \text{Tr}((\mathbf{K} - \tilde{\mathbf{T}})^2) \\
&= \frac{v^2}{\mu^2} \text{Tr}((\mathbf{K} - \tilde{\mathbf{T}})^2 (\mathbf{S}^T \mathbf{S} - \mathbf{I})) \quad (53)
\end{aligned}$$

Applying the property  $\mathbf{K} - \tilde{\mathbf{T}} \leq \frac{\mu}{1-t}$  from (35), this difference can be upper bounded as

$$\frac{v^2}{\mu^2} \text{Tr}((\mathbf{K} - \tilde{\mathbf{T}})^2 (\mathbf{S}^T \mathbf{S} - \mathbf{I})) \leq \frac{v^2}{(1-t)^2} \text{Tr}(\mathbf{S}^T \mathbf{S} - \mathbf{I}) \quad (54)$$

Thus,

$$\text{var} \leq \text{var}_I + \frac{v^2}{(1-t)^2} \text{Tr}(\mathbf{S}^T \mathbf{S} - \mathbf{I}). \quad (55)$$

Adding together (52) and (55) yields the upper bound in the theorem.

### References

- [1] C.E. Rasmussen, Gaussian processes in machine learning, in: Summer School on Machine Learning, Springer, 2003, pp. 63–71.
- [2] F. Pérez-Cruz, S. Van Vaerenbergh, J.J. Murillo-Fuentes, M. Lázaro-Gredilla, I. Santamaría, Gaussian processes for nonlinear signal processing: an overview of recent advances, *IEEE Signal Process. Mag.* 30 (4) (2013) 40–50.
- [3] H. Yu, S. Kim, Passive sampling for regression, in: 2010 IEEE International Conference on Data Mining, 2010, pp. 1151–1156, doi:10.1109/ICDM.2010.9.
- [4] Q. Gu, J. Han, Towards active learning on graphs: an error bound minimization approach, in: IEEE 12th International Conference on Data Mining, IEEE, 2012, pp. 882–887.
- [5] C. Orhan, O. Taştan, ALEVS: active learning by statistical leverage sampling, arXiv preprint arXiv:1507.04155(2015).
- [6] J. Goetz, A. Tewari, P. Zimmerman, Active learning for non-parametric regression using purely random trees, *Adv. Neural Inf. Process. Syst.* 31 (2018).
- [7] D. Wu, Pool-based sequential active learning for regression, *IEEE Trans. Neural Netw. Learn. Syst.* 30 (5) (2019) 1348–1359, doi:10.1109/TNNLS.2018.2868649.
- [8] D. Wu, C.T. Lin, J. Huang, Active learning for regression using greedy sampling, *Inf. Sci.* 474 (2019) 90–105.
- [9] D. Romero, M. Ma, G.B. Giannakis, Kernel-based reconstruction of graph signals, *IEEE Trans. Signal Process.* 65 (3) (2017) 764–778.
- [10] P. Giménez-Febrer, A. Pagès-Zamora, G.B. Giannakis, Matrix completion and extrapolation via kernel regression, *IEEE Trans. Signal Process.* 67 (19) (2019) 5004–5017, doi:10.1109/TSP.2019.2932875.
- [11] P. Giménez-Febrer, A. Pagès-Zamora, G.B. Giannakis, Generalization error bounds for kernel matrix completion and extrapolation, *IEEE Signal Process. Lett.* 27 (2020) 326–330, doi:10.1109/LSP.2020.2970306.
- [12] P. Drineas, M.W. Mahoney, On the Nyström method for approximating a Gram matrix for improved kernel-based learning, *J. Mach. Learn. Res.* 6 (2005) 2153–2175.
- [13] A. Alaoui, M.W. Mahoney, Fast randomized kernel ridge regression with statistical guarantees, in: Advances in Neural Information Processing Systems, vol. 28, 2015, pp. 775–783, Montreal, Canada.
- [14] A. Tanaka, H. Imai, M. Miyakoshi, Kernel-induced sampling theorem, *IEEE Trans. Signal Process.* 58 (7) (2010) 3569–3577.

- [15] C. Williams, M. Seeger, Using the Nyström method to speed up kernel machines, in: *Advances in Neural Information Processing Systems 13*, MIT Press, 2001, pp. 682–688.
- [16] C. Musco, C. Musco, Recursive sampling for the Nyström method, in: *Advances in Neural Information Processing Systems*, vol. 30, 2017, pp. 3833–3845.
- [17] A. Deshpande, L. Rademacher, Efficient volume sampling for row/column subset selection, in: *2010 IEEE 51st Annual Symposium on Foundations of Computer Science*, IEEE, 2010, pp. 329–338.
- [18] D.P. Woodruff, Sketching as a tool for numerical linear algebra, *Theor. Comput. Sci.* 10 (1–2) (2014) 1–157.
- [19] F. Gama, A.G. Marques, G. Mateos, A. Ribeiro, Rethinking sketching as sampling: a graph signal processing approach, *Signal Process.* 169 (2020) 107404.
- [20] S. Wang, A. Gittens, M.W. Mahoney, Sketched ridge regression: optimization perspective, statistical perspective, and model averaging, in: *International Conference on Machine Learning*, PMLR, 2017, pp. 3608–3616.
- [21] R. Patel, T. Goldstein, E. Dyer, A. Mirhoseini, R. Baraniuk, Deterministic column sampling for low-rank matrix approximation: Nyström vs. incomplete Cholesky decomposition, in: *Proceedings of the 2016 SIAM International Conference on Data Mining*, SIAM, 2016, pp. 594–602.
- [22] P. Drineas, R. Kannan, M.W. Mahoney, Fast Monte Carlo algorithms for matrices I: approximating matrix multiplication, *SIAM J. Comput.* 36 (1) (2006) 132–157.
- [23] A. Rudi, R. Camoriano, L. Rosasco, Less is more: Nyström computational regularization, *Advances in Neural Information Processing Systems* 28 (2015).
- [24] H. Avron, M. Kapralov, C. Musco, C. Musco, A. Velingker, A. Zandieh, Random Fourier features for kernel ridge regression: Approximation bounds and statistical guarantees, in: *International conference on machine learning*, PMLR, 2017, pp. 253–262.