(2024) 2024:30

RESEARCH

Open Access

The learning performance of the weak rescaled pure greedy algorithms



Qin Guo^{1*}, Xianghua Liu^{2†} and Peixin Ye^{2†}

*Correspondence: quogin 1985@163.com

¹School of Science, Shandong Jianzhu University, Jinan, 250101, China

Full list of author information is available at the end of the article [†]Equal contributors

Abstract

We investigate the regression problem in supervised learning by means of the weak rescaled pure greedy algorithm (WRPGA). We construct learning estimator by applying the WRPGA and deduce the tight upper bounds of the *K*-functional error estimate for the corresponding greedy learning algorithms in Hilbert spaces. Satisfactory learning rates are obtained under two prior assumptions on the regression function. The application of the WRPGA in supervised learning considerably reduces the computational cost while maintaining its powerful generalization capability when compared with other greedy learning algorithms.

Keywords: Least squares regression; Rescaled pure greedy algorithm; Weakness sequence; *K*-functional; Learning rate

1 Introduction

The applications of greedy algorithms to supervised learning have sparked great research interest because they have appealing generalization capability with lower computing burden than typical regularized methods, particularly in large-scale dictionary learning problem [1-6]. Big data sets for the most traditional learning algorithms frequently cause slow machine performance. To tackle this problem, many researchers [1-3, 7, 8] advocate greedy learning algorithms, which have greatly improved learning performance.

The approximation abilities of greedy-type algorithms for frames or more dictionaries \mathcal{D} were investigated in [7, 9–12], as well as various applications, see [3, 7, 13–19]. The pure greedy algorithm (PGA) can realize the best bilinear approximation, see [20, 21]. Although the PGA is outstanding at computing, the main problem is that it lacks optimal convergence properties for a general dictionary, and consequently the slower convergence rate than the best nonlinear approximation [11, 21–23] corrupts its learning performance. To improve the approximation rate, the orthogonal greedy algorithm (OGA), the relaxed greedy algorithm (RGA), the stepwise projection algorithm (SPA), and their weak versions have been proposed. It was shown that these greedy algorithms all achieved the optimal rate $\mathcal{O}(m^{-\frac{1}{2}})$ for approximating the elements in the class $\mathcal{A}_1(\mathcal{D})$, which will be defined in (14), where *m* is the iteration number, see [9, 11].

Both the OGA and the RGA have recently been employed successfully in machine learning [1-3, 7, 8]. For example, Barron et al. [7] established the optimal convergence rate

© The Author(s) 2024. **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit http://creativecommons.org/licenses/by/4.0/.



 $\mathcal{O}(n/\log n)^{-\frac{1}{2}}$, where *n* is the sample size. To reduce the OGA's computational load, Fang et al. [1] investigated the learning performance of the orthogonal super greedy algorithm (OSGA) and derived the almost same rate as the orthogonal greedy learning algorithm (OGLA). All these results demonstrate that each greedy learning algorithm has its advantages and disadvantages.

We study the applications of weak greedy algorithms to least squares regression in supervised learning. It is well known that the weak type are easier to implement than the usual greedy algorithms, see [12]. Specifically, the weak rescaled pure greedy algorithm (WRPGA), one fairly simple modification of the PGA, is the goal of our investigation, see [24, 25]. When compared to the OGA and the RGA, the WRPGA can also furthermore reduce the computational load. The best rate $\mathcal{O}(m^{-\frac{1}{2}})$ for functions in the basic sparse class has been proved [24]. Motivated by research results of [24], we proceed to use the same method employed for the RPGA in [24] to deduce the error bound of the *K*-functional estimate in the Hilbert space \mathcal{H} for the WRPGA. The WRPGA is a simple greedy algorithm with good approximation ability. Based on this, we propose the weak rescaled pure greedy learning algorithm (WRPGLA) for solving the kernel-based regression problems in supervised learning. Using the WRPGA's proven approximation result, we can derive that the WRPGLA has the almost same learning rate as the OGLA. Our results show that the WRPGLA further cuts down the computational complexity even more without reducing generalization capabilities.

The paper is organized as follows. In Sect. 2, we review least squares regression learning theory and the WRPGA. In Sect. 3, we propose the WRPGLA and state the main theorems on the error estimates. Section 4 is devoted to proofs of the main results. We present the convergence rates under two smoothness assumptions on the regression function f_{ρ} in the last section.

2 Preliminaries

Some preliminaries are presented in this section. Sections 2.1 and 2.2 provide a fast overview of least squares regression learning and the WRPGA, respectively.

2.1 Least squares regression

In this paper, the approximation problem is addressed in the following statistical learning context. Let *X* be a compact metric space and $Y = \mathbb{R}$. Let ρ be a Borel probability measure on $Z = X \times Y$. The generalization error for a function $f : X \to Y$ is defined by

$$\mathcal{E}(f) = \int_{Z} \left(f(x) - y \right)^2 d\rho, \tag{1}$$

which is minimized by the following regression function:

$$f_{\rho}(x) = \int_{Y} y \, d\rho(y|x),$$

where $\rho(\cdot|x)$ is the conditional distribution induced by ρ at $x \in X$. In regression learning, ρ is unknown, and what one can know is a set of samples $\mathbf{z} = \{z_i\}_{i=1}^n = \{(x_i, y_i)\}_{i=1}^n \in Z^n$ that are drawn independently and identically according to ρ . The goal of learning is to find a

good approximation f_z of f_ρ , which minimizes the empirical error

$$\mathcal{E}_{\mathbf{z}}(f) = \|y - f\|_n^2 \coloneqq \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2.$$
⁽²⁾

Denote the Hilbert space of the square integrable functions defined on *X* with respect to the measure ρ_X by $L^2_{\rho_X}(X)$, where ρ_X is the marginal measure of ρ on *X*. It is clear from the definition of $f_{\rho}(x)$ that for each $x \in X$, $\int_Y (f_{\rho}(x) - y) d\rho(y|x) = 0$. For any $f \in L^2_{\rho_X}(X)$, it holds that

$$\begin{split} \mathcal{E}(f) &= \int_{Z} \left(f(x) - f_{\rho}(x) + f_{\rho}(x) - y \right)^{2} d\rho \\ &= \int_{X} \left(f(x) - f_{\rho}(x) \right)^{2} d\rho_{X} + \int_{Z} \left(f_{\rho}(x) - y \right)^{2} d\rho \\ &+ 2 \int_{X} \left(f(x) - f_{\rho}(x) \right) d\rho_{X} \int_{Y} \left(f_{\rho}(x) - y \right) d\rho(y|x) \\ &= \int_{X} \left(f(x) - f_{\rho}(x) \right)^{2} d\rho_{X} + \mathcal{E}(f_{\rho}). \end{split}$$

Therefore,

$$\mathcal{E}(f) - \mathcal{E}(f_{\rho}) = \|f - f_{\rho}\|^2 \tag{3}$$

with the norm $\|\cdot\|$

$$\|f\| = \left(\int_{X} |f(x)|^{2} d\rho_{X}\right)^{\frac{1}{2}}.$$
(4)

The prediction accuracy of learning algorithms is measured by $E(||f_z - f_o||^2)$.

We will assume $|y| \le B$ for a positive real number $B < \infty$ almost surely. In this paper, we construct the learning estimator f_z by applying the WRPGA and estimate $E(||f_z - f_\rho||^2)$. So, in the following subsection, we recall this algorithm.

2.2 Weak rescaled pure greedy algorithm

We shall restrict our analysis to the situation in which approximation takes place in a real, separable Hilbert space \mathcal{H} with the inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ and the norm $\|\cdot\| := \|\cdot\|_{\mathcal{H}} = \langle \cdot, \cdot \rangle_{\mathcal{H}}^{\frac{1}{2}}$. Let $\mathcal{D} \subset \mathcal{H}$ be a given dictionary satisfying $\|g\| = 1$ for every $g \in \mathcal{D}$, $g \in \mathcal{D}$ implies $-g \in \mathcal{D}$ and $\overline{\text{Span}(\mathcal{D})} = \mathcal{H}$.

Petrova developed the rescaled pure greedy algorithm (RPGA) to enhance the PGA's convergence rate, which simply rescales f_m at the *m*th greedy step, see [24]. We begin by describing the weak rescaled pure greedy algorithm (WRPGA) also introduced by Petrova in [24].

WRPGA({ t_m }, D): *Step* 0: Let $f_0 := 0$. *Step* m ($m \ge 1$): (1) If $f = f_{m-1}$, then terminate the iterative process and define $f_k = f_{m-1} = f$ for $k \ge m$. (2) If $f \neq f_{m-1}$, then choose a direction $\varphi_m \in \mathcal{D}$ such that

$$\left| \langle f - f_{m-1}, \varphi_m \rangle \right| \ge t_m \sup_{\varphi \in \mathcal{D}} \left| \langle f - f_{m-1}, \varphi \rangle \right|, \tag{5}$$

where $\{t_m\}_{m=1}^{\infty}$ is a weakness sequence and $t_m \in (0, 1]$. Let

$$\lambda_m := \langle f - f_{m-1}, \varphi_m \rangle, \tag{6}$$

$$\hat{f}_m := f_{m-1} + \lambda_m \varphi_m, \tag{7}$$

$$s_m := \frac{\langle f, \hat{f}_m \rangle}{\|\hat{f}_m\|^2}.$$
(8)

The *m* step approximation f_m is defined as

$$f_m = s_m \hat{f}_m,\tag{9}$$

and proceed to Step m + 1.

Remark 1 When $t_m = 1$, this algorithm is the RPGA. Note that if the supremum is not attained, one can select $t_m < 1$ and proceed with the algorithm. In this case, it is easier to choose φ_m . If the output at the *m*th greedy step was \hat{f}_m rather than $f_m = s_m \hat{f}_m$, this would be the PGA. The WRPGA uses $s_m \hat{f}_m$, which is just suitable scaling of \hat{f}_m , and thus increases the rate to $\mathcal{O}(m^{-\frac{1}{2}})$ for functions in the closure of the convex hull of \mathcal{D} .

3 Weak rescaled pure greedy learning

We shall provide the WRPGLA for regression. From the definition of the WRPGA, computing $\sup_{\varphi \in \mathcal{D}} |\langle f - f_{m-1}, \varphi \rangle|$ may result in computation difficulty. Therefore we compute only over the truncation of the dictionary, which is a finite subset of \mathcal{D} . Let $\mathcal{D}_1 \subset \mathcal{D}_2 \subset$ $\cdots \subset \mathcal{D}$. Then \mathcal{D}_m is the truncation of \mathcal{D} with the cardinality $\#(\mathcal{D}_m) = m$. Here we assume that

$$m \le m(n) := \left| n^a \right|$$
 for some fixed $a \ge 1$. (10)

Then the WRPGLA is defined by the following simple processes.

WRPGLA:

Step 1: We apply the WRPGA for \mathcal{D}_m to the function $y(x_i) = y_i$ by utilizing the norm $\|\cdot\|_n$ associated with the empirical inner product, that is,

$$||f||_n := \left(\frac{1}{n}\sum_{i=1}^n |f(x_i)|^2\right)^{\frac{1}{2}}.$$

Step 2: The algorithms establish the approximation $f_{z,k} := f_k$ to the data at the *k*th greedy step. Then, we define our estimator as $f_z := Tf_{z,k^*}$, where $Tu := T_B \min\{B, |u|\} \operatorname{sgn}(u)$ and

$$k^* := \arg\min_{k>0} \left\{ \|y - Tf_{\mathbf{z},k}\|_n^2 + \kappa \frac{k \log n}{n} \right\},$$
(11)

where the constant $\kappa \ge \kappa_0 = 2568B^4(a+5)$, which will be discussed in proof of Theorem 1.

To discuss the approximation properties of WRPGLA, we introduce the class of functions

$$\mathcal{A}_{1}^{0}(\mathcal{D},M) := \left\{ f = \sum_{k \in \Lambda} c_{k}(f)\varphi_{k} : \varphi_{k} \in D, \#(\Lambda) < \infty, \sum_{k \in \Lambda} \left| c_{k}(f) \right| \le M \right\},$$
(12)

and

$$\mathcal{A}_1(\mathcal{D}, M) = \overline{\mathcal{A}_1^0(\mathcal{D}, M)}.$$
(13)

Then

$$\mathcal{A}_1(\mathcal{D}) = \bigcup_{M>0} \mathcal{A}_1(\mathcal{D}, M) \tag{14}$$

and

$$\|f\|_{\mathcal{A}_1(\mathcal{D})} \coloneqq \inf\{M : f \in \mathcal{A}_1(\mathcal{D}, M)\}.$$
(15)

We also use the following *K*-functional:

$$K(f,t) := K(f,t,\mathcal{H},\mathcal{A}_{1}(\mathcal{D})) := \inf_{h \in \mathcal{A}_{1}(\mathcal{D})} \{ \|f-h\|_{\mathcal{H}} + t \|h\|_{\mathcal{A}_{1}(\mathcal{D})} \}, \quad t > 0.$$
(16)

Since all the constants in this work depend at most on κ_0 , *B*, and *a*, we denote all of them by *C* for simplicity of notation. Now we take $\mathcal{H} = L^2_{\rho_X}(X)$ with the norm defined by (4).

Then, we provide our main results on the generalization error bounds for the WRPGLA.

Theorem 1 There exists κ_0 depending only on *B* and a such that if $\kappa \ge \kappa_0$, then for all k > 0and $h \in \text{Span}(\mathcal{D}_m)$, the learning estimator by applying the WRPGA satisfies

$$E(\|f_{\mathbf{z}} - f_{\rho}\|^{2}) \le 8 \frac{\|h\|_{\mathcal{A}_{1}(\mathcal{D}_{m})}^{2}}{\sum_{i=1}^{k} t_{i}^{2}} + 2\|f_{\rho} - h\|^{2} + C\frac{k\log n}{n}.$$
(17)

Furthermore, we have

$$E(\|f_{\mathbf{z}} - f_{\rho}\|) \le 2K\left(f_{\rho}, 2\left(\sum_{i=1}^{k} t_{i}^{2}\right)^{-\frac{1}{2}}\right) + C\frac{k\log n}{n}.$$
(18)

Applying Theorem 1 with $t_i = t_0$ for all $i \ge 1$ and $0 < t_0 \le 1$, we get the following theorem.

Theorem 2 Under the assumptions of Theorem 1, if $t_i = t_0$ for all $i \ge 1$ and $0 < t_0 \le 1$, then we have

$$E(\|f_{\mathbf{z}} - f_{\rho}\|^{2}) \le 8 \frac{\|h\|_{\mathcal{A}_{1}(\mathcal{D}_{m})}^{2}}{kt_{0}^{2}} + 2\|f_{\rho} - h\|^{2} + C\frac{k\log n}{n}.$$
(19)

Furthermore, we have

$$E(\|f_{\mathbf{z}} - f_{\rho}\|) \le 2K(f_{\rho}, 2k^{-\frac{1}{2}}t_{0}^{-1}) + C\frac{k\log n}{n}.$$
(20)

4 Proofs of the main results

To prove Theorem 1, we establish a lemma on the upper error bound for the WRPGA.

Lemma 4.1 If $f \in H$, $h \in A_1(D)$, then the output $(f_m)_{m \ge 0}$ of the WRPGA satisfies

$$e_m := \|f - f_m\| \le 2K \left(f, \left(\sum_{k=1}^m t_k^2 \right)^{-1/2} \right), \quad m = 0, 1, 2, \dots$$
(21)

Proof In terms of the definition of *K*-functional, we just need to prove that for $f \in \mathcal{H}$ and $h \in \mathcal{A}_1(\mathcal{D})$,

$$e_m^2 \le \|f - h\|^2 + \frac{4}{\sum_{k=1}^m t_k^2} \|h\|_{\mathcal{A}_1(\mathcal{D})}^2, \quad m = 1, 2, \dots$$
(22)

Since $\mathcal{A}_1^0(\mathcal{D}, M)$ is dense in $\mathcal{A}_1(\mathcal{D}, M)$, it suffices to prove (22) for functions *h* that are finite sums $\sum_j c_j \varphi_j$ with $\sum_j |c_j| \leq M$. We fix $\epsilon > 0$ and select a representation for $h = \sum_{\varphi \in \mathcal{D}} c_{\varphi} \varphi$, such that

$$\sum_{\varphi \in \mathcal{D}} |c_{\varphi}| < M + \epsilon.$$
(23)

Denote

$$a_m := e_m^2 - \|f - h\|^2, \quad m = 1, 2, \dots$$
 (24)

The nonincreasing of $\{e_m\}_{m=0}^{\infty}$ implies that $\{a_m\}_{m=0}^{\infty}$ is also a nonincreasing sequence.

Then we discuss these two cases separately.

Case 1: $a_0 := ||f||^2 - ||f - h||^2 \le 0$. Then, for every $m \ge 1$, we have $a_m \le 0$. Therefore inequality (22) holds true.

Case 2: $a_0 > 0$. Assume that $a_{m-1} > 0$, $m \ge 1$. Note that f_m is the orthogonal projection of f onto the linear space spanned by \hat{f}_m , it implies

$$\langle f - f_m, f_m \rangle = 0, \quad m \ge 0. \tag{25}$$

This together with the selection of φ_m implies

$$e_{m-1}^{2} = \langle f - f_{m-1}, f - f_{m-1} \rangle$$

= $\langle f - f_{m-1}, f \rangle$
= $\langle f - f_{m-1}, f - h \rangle + \langle f - f_{m-1}, h \rangle$
 $\leq e_{m-1} ||f - h|| + \sum_{\varphi \in \mathcal{D}} c_{\varphi} \langle f - f_{m-1}, \varphi \rangle$

$$\leq e_{m-1} \|f - h\| + t_m^{-1} |\langle f - f_{m-1}, \varphi_m \rangle| \sum_{\varphi \in \mathcal{D}} |c_\varphi|.$$

$$\tag{26}$$

By (23), we get

$$e_{m-1}^{2} \leq \frac{1}{2} \left(e_{m-1}^{2} + \|f - h\|^{2} \right) + t_{m}^{-1} \left| \langle f - f_{m-1}, \varphi_{m} \rangle \right| (M + \epsilon).$$
⁽²⁷⁾

Let $\epsilon \rightarrow 0$. Therefore

$$|\langle f - f_{m-1}, \varphi_m \rangle| \ge \frac{t_m (e_{m-1}^2 - ||f - h||^2)}{2M}.$$
 (28)

It has been proved in [24] that

$$e_m^2 \le e_{m-1}^2 - \langle f - f_{m-1}, \varphi_m \rangle^2, \quad m = 1, 2, \dots$$
 (29)

Then, using the assumption that $a_{m-1} > 0$, we have

$$e_m^2 \le e_{m-1}^2 - \frac{t_m^2 a_{m-1}^2}{4M^2}.$$
(30)

It yields

$$a_m \le a_{m-1} \left(1 - \frac{t_m^2 a_{m-1}}{4M^2} \right). \tag{31}$$

In particular, for m = 1, we have

$$a_1 \le a_0 \left(1 - \frac{t_1^2 a_0}{4M^2} \right). \tag{32}$$

Case 2.1: $0 < a_0 < \frac{4M^2}{t_1^2}$. Since $\psi(t) := t(1 - \frac{t_1^2 t}{4M^2})$ on $(0, \frac{4M^2}{t_1^2})$ has maximum $\frac{M^2}{t_1^2}$, it follows that

$$a_m \le \psi(a_0) \le \frac{M^2}{t_1^2} \le \frac{4M^2}{t_1^2}.$$

Therefore, either all $\{a_m\}_{m=0}^{\infty} \subset (0, \frac{4M^2}{t_1^2})$ and then satisfy (31), or we know that $a_{m^*} \leq 0$ for some $m^* \geq 1$. The analysis for $m \geq m^*$ is therefore the same as in Case 1. For the positive elements in $\{a_m\}_{m=0}^{\infty}$, by applying Lemma 2.2 from [24] with l = 1, $r_m = t_m^2$, $B = \frac{4M^2}{t_1^2}$, J = 0, and $r = 4M^2$, we obtain

$$a_m \le \frac{4M^2}{t_1^2 + \sum_{k=1}^m t_k^2} \le \frac{4M^2}{\sum_{k=1}^m t_k^2},\tag{33}$$

which gives inequality (22).

Case 2.2: $a_0 \ge \frac{4M^2}{t_1^2}$. It follows from (32) that $a_1 < 0$. That is, $e_1^2 < ||f - h||^2$, which yields (22) due to monotonicity. Lemma 4.1 is proved.

Now we prove Theorem 1.

Proof of Theorem 1 As shown in [5], $||f_z - f_\rho||^2$ can be decomposed as

$$\|f_{\mathbf{z}} - f_{\rho}\|^{2} \leq S_{1} + S_{2} + S_{3} + 2\left(\|y - f_{\mathbf{z}}\|_{n}^{2} + \kappa \frac{k^{*} \log n}{n} - \|y - Tf_{\mathbf{z},k}\|_{n}^{2} - \kappa \frac{k \log n}{n}\right),$$
(34)

where

$$S_{1} := \|f_{z} - f_{\rho}\|^{2} - 2\left(\|y - f_{z}\|_{n}^{2} - \|y - f_{\rho}\|_{n}^{2} + \kappa \frac{k^{*} \log n}{n}\right),$$

$$S_{2} := 2\left(\|y - f_{z,k}\|_{n}^{2} - \|y - h\|_{n}^{2}\right),$$

$$S_{3} := 2\left(\|y - h\|_{n}^{2} - \|y - f_{\rho}\|_{n}^{2} + \kappa \frac{k \log n}{n}\right),$$
(35)

and $h \in \text{Span}\{\mathcal{D}_m\}$.

We firstly estimate the bound of S_1 . To do this, we introduce Ω ,

$$\Omega = \left\{ \mathbf{z} : \mathbf{z} \in Z^{n}, \|f_{\mathbf{z}} - f_{\rho}\|^{2} \ge 2 \left(\|y - f_{\mathbf{z}}\|_{n}^{2} - \|y - f_{\rho}\|_{n}^{2} + \kappa \frac{k^{*} \log n}{n} \right) \right\}.$$
(36)

Let $\operatorname{Prob}(\Omega)$ be the probability that the sample point is a member of the set Ω . Then from $|y| \leq B$ and the definition of f_{ρ} and f_z , we have

$$E(\mathcal{S}_1) \le 6B^2 \operatorname{Prob}(\Omega). \tag{37}$$

For \mathcal{S}_2 , according to Lemma 4.1, we get

$$\|y - f_{\mathbf{z},k}\|_n^2 - \|y - h\|_n^2 \le 4 \frac{\|h\|_{\mathcal{A}_1^n}^2}{\sum_{k=1}^m t_k^2},$$
(38)

where

$$\mathcal{A}_{1}^{n}(\mathcal{D}) := \left\{ h : h = \sum_{i \in \Lambda} c_{i}^{n} \|g_{i}\|_{n} \frac{g_{i}}{\|g_{i}\|_{n}}, h \in \mathcal{A}_{1}(\mathcal{D}) \right\}$$
(39)

and

$$\|h\|_{\mathcal{A}_{1}^{n}(\mathcal{D})} := \inf_{h} \left\{ \sum_{i \in \Lambda} \left| c_{i}^{n} \right| \cdot \|g_{i}\|_{n}, h \in \mathcal{A}_{1}^{n}(\mathcal{D}) \right\}.$$

$$\tag{40}$$

It has been proved in Lemma 3.4 of [7] that

$$E\left(\|h\|_{\mathcal{A}_{1}^{n}}^{2}\right) \leq \|h\|_{\mathcal{A}_{1}}^{2},\tag{41}$$

which implies

$$E(S_2) \le 8 \frac{\|h\|_{\mathcal{A}_1}^2}{\sum_{k=1}^m t_k^2}.$$
(42)

For S_3 , from the property of mathematical expectation and (1), we have

$$E(\|y - h\|_{n}^{2} - \|y - f_{\rho}\|_{n}^{2}) = E(|y - h(x)|^{2}) - E(|y - f_{\rho}(x)|^{2})$$

= $\mathcal{E}(h) - \mathcal{E}(f_{\rho}).$ (43)

This together with (3) yields

$$E(S_3) = 2\|f_{\rho} - h\|^2 + 2\kappa \frac{k \log n}{n}.$$
(44)

Combining (37), (42), with (44), we obtain

$$E(\|f_{\mathbf{z}} - f_{\rho}\|^{2}) \le 6B^{2}\operatorname{Prob}(\Omega) + 8\frac{\|h\|_{\mathcal{A}_{1}}^{2}}{\sum_{k=1}^{m} t_{k}^{2}} + 2\|f_{\rho} - h\|^{2} + 2\kappa \frac{k\log n}{n}.$$
(45)

Next we bound $Prob(\Omega)$. To this end, we need the following known result in [10].

Lemma 4.2 Let \mathcal{F} be the class of functions $\mathcal{F} = \{|f| \leq B\}$ for some fixed constant B. For all n and $\alpha, \beta > 0$, we have

$$\operatorname{Prob}\left\{\exists f \in \mathcal{F} : \|f - f_{\rho}\|_{\rho_{X}}^{2} \geq 2\left(\|y - f\|_{n}^{2} - \|y - f_{\rho}\|_{n}^{2}\right) + \alpha + \beta\right\}$$
$$\leq 14 \sup_{x} \mathcal{N}\left(\frac{\beta}{40B}, \mathcal{F}, L_{1}(\vec{v}_{x})\right) \exp\left(-\frac{\alpha n}{2568B^{4}}\right), \tag{46}$$

where $\mathbf{x} = (x_1, \dots, x_n) \in X^n$ and $\mathcal{N}(t, \mathcal{F}, L_1(\vec{v}_x))$ is the covering number for the class \mathcal{F} by balls of radius t in $L_1(\vec{v}_x)$, with $\vec{v}_x := \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$ the empirical discrete measure.

We define $\mathcal{G}_{\Lambda} := \text{Span}\{g : g \in \Lambda \subset \mathcal{D}\}$ and $\mathcal{F}_k := \bigcup_{\Lambda \subset \mathcal{D}_m, \#(\Lambda) \leq k} \{Tf : f \in \mathcal{G}_{\Lambda}\}$. Consider the probability

$$p_{k} = \operatorname{Prob}\left\{\exists f \in \mathcal{F}_{k} : \|f - f_{\rho}\|^{2} \ge 2\left(\|y - f\|_{n}^{2} - \|y - f_{\rho}\|_{n}^{2} + \kappa \frac{k \log n}{n}\right)\right\}.$$

Applying Lemma 4.2 to \mathcal{F}_k with $\alpha = \kappa \frac{k \log n}{n}$, $\beta = \frac{1}{n}$, and $\kappa > 1$, we get

$$p_{k} \leq 14 \sup_{\mathbf{x}} \mathcal{N}\left(\frac{1}{40Bn}, \mathcal{F}_{k}, L_{1}(\vec{v}_{\mathbf{x}})\right) \exp\left(-\kappa \frac{k \log n}{2568B^{4}}\right)$$
$$= 14 \sup_{\mathbf{x}} \mathcal{N}\left(\frac{1}{40Bn}, \mathcal{F}_{k}, L_{1}(\vec{v}_{\mathbf{x}})\right) n^{-\frac{\kappa k}{2568B^{4}}}.$$
(47)

Lemma 3.3 of [7] provides the upper bound for $\mathcal{N}(t, \mathcal{F}_k, L_1(\vec{v}_x))$, which implies

$$p_k \le C n^{ak} n^{2(k+1)} n^{-\frac{\kappa k}{2568B^4}}.$$
(48)

Let $\kappa \ge \kappa_0 = 2568B^4(a + 5)$. Then the above inequality yields

$$p_k \le C n^{-3k+2} \le C n^{-2}. \tag{49}$$

So we have

$$\operatorname{Prob}(\Omega) \le \sum_{1 \le k \le \frac{Bn}{c}} p_k \le \frac{C}{n}.$$
(50)

By substituting the bound (50) of $Prob(\Omega)$ into (45), we get

$$E(\|f_{\mathbf{z}} - f_{\rho}\|^{2}) \le 8 \frac{\|h\|_{\mathcal{A}_{1}}^{2}}{\sum_{k=1}^{m} t_{k}^{2}} + 2\|f_{\rho} - h\|^{2} + C \frac{k \log n}{n}.$$
(51)

Next we derive the *K*-functional result of the upper bound (51). It is known from the property of variance that

$$E^{2}(\|f_{\mathbf{z}} - f_{\rho}\|) \leq E(\|f_{\mathbf{z}} - f_{\rho}\|)^{2}.$$
(52)

Combining (51) with (52), we have

$$E(\|f_{\mathbf{z}} - f_{\rho}\|) \leq \sqrt{8 \frac{\|h\|_{\mathcal{A}_{1}}^{2}}{\sum_{k=1}^{m} t_{k}^{2}} + 2 \|f_{\rho} - h\|^{2} + C \frac{k \log n}{n}}$$

$$\leq 2 \left(\frac{2 \|h\|_{\mathcal{A}_{1}}}{(\sum_{k=1}^{m} t_{k}^{2})^{1/2}} + \|f_{\rho} - h\|\right) + C \frac{k \log n}{n}$$

$$\leq 2 K \left(f_{\rho}, 2 \left(\sum_{k=1}^{m} t_{k}^{2} \right)^{-1/2} \right) + C \frac{k \log n}{n}.$$
(53)

This completes the proof of Theorem 1.

5 Convergence rate and universal consistency

In this section, we analyze Theorem 2 under two different prior assumptions on f_{ρ} . We begin with the definitions of $\mathcal{A}_1(\mathcal{D}_m)$, $\mathcal{A}_{1,r}$, and $\mathcal{B}_{p,r}$.

We define the space $\mathcal{A}_1(\mathcal{D}_m)$ to be the space $\text{Span}\{\mathcal{D}_m\}$ with the norm $\|\cdot\|_{\mathcal{A}_1(\mathcal{D}_m)}$ defined by (15). Note that now \mathcal{D} is replaced by \mathcal{D}_m .

For r > 0, we then introduce the space

$$\mathcal{A}_{1,r} = \{ f : \forall m, \exists h = h(m) \in \text{Span}\{\mathcal{D}_m\}, \|h\|_{\mathcal{A}_1(\mathcal{D}_m)} \le C, \|f - h\| \le Cm^{-r} \},$$
(54)

where $\|\cdot\|_{\mathcal{A}_{1,r}}$ is the minimum value of *C* such that (54) holds.

Furthermore, we present the following space:

$$\mathcal{B}_{p,r} := [\mathcal{H}, \mathcal{A}_{1,r}]_{\theta,\infty}, \quad 0 < \theta < 1, \tag{55}$$

with $\frac{1}{p} = \frac{1+\theta}{2}$. From the definition of interpolation spaces in [26], we know that $f \in [\mathcal{H}, \mathcal{A}_{1,r}]_{\theta,\infty}$ if and only if for any t > 0,

$$K(f, t, \mathcal{H}, \mathcal{A}_{1,r}) := \inf_{h \in \mathcal{A}_{1,r}} \left\{ \|f - h\|_{\mathcal{H}} + t \|h\|_{\mathcal{A}_{1,r}} \right\} \le Ct^{\theta}.$$
(56)

The minimum *C* such that (56) holds true is defined as the norm on $\mathcal{B}_{p,r}$.

Now we first consider $f_{\rho} \in \mathcal{A}_{1,r}$.

Corollary 5.1 Under the assumptions of Theorem 2, if $f_{\rho} \in A_{1,r}$ with $r > \frac{1}{2a}$, then we have

$$E(\|f_{\mathbf{z}} - f_{\rho}\|^{2}) \le C(1 + \|f_{\rho}\|_{\mathcal{A}_{1,r}})t_{0}^{-1}\left(\frac{n}{\log n}\right)^{-\frac{1}{2}}.$$
(57)

Proof From the definition of $A_{1,r}$, there exists $h := h(m) \in \text{Span}\{\mathcal{D}_m\}$ for every *m* that satisfies

$$\|h\|_{\mathcal{A}_1(\mathcal{D}_m)} \le M$$

and

$$\|f_{\rho}-h\|\leq Mm^{-r},$$

where $M := ||f_{\rho}||_{\mathcal{A}_{1,r}}$.

Theorem 2 thus implies

$$E(\|f_{\mathbf{z}} - f_{\rho}\|^2) \le C \min_{k>0} \left(\frac{M^2}{kt_0^2} + M^2 n^{-2ar} + \frac{k \log n}{n}\right).$$
(58)

Moreover, the mild restriction $2ar \ge 1$ with *a* arbitrarily large allows us to remove the term $M^2 n^{-2ar}$ in (58). To balance the errors in (58), we take $k := \left\lceil \frac{(M+1)^2}{t_0^2} \frac{n}{\log n} \right\rceil^{\frac{1}{2}}$. Then the desired result (57) can be obtained.

Next we consider $f_{\rho} \in \mathcal{B}_{p,r}$.

Corollary 5.2 Under the assumptions of Theorem 2, if $f_{\rho} \in \mathcal{B}_{p,r}$ with $r > \frac{1}{2a}$, then we have

$$E(\|f_{\mathbf{z}} - f_{\rho}\|^{2}) \le Ct_{0}^{-p} \left(1 + \|f_{\rho}\|_{\mathcal{B}_{p,r}}\right)^{p} \left(\frac{n}{\log n}\right)^{-1 + \frac{p}{2}}.$$
(59)

Proof By (56), if $f \in \mathcal{B}_{p,r}$, then for any t > 0, we can find a function $\tilde{f} \in \mathcal{A}_{1,r}$ that satisfies

$$\|\tilde{f}\|_{\mathcal{A}_{1,r}} \le \|f\|_{\mathcal{B}_{p,r}} t^{\theta-1} \tag{60}$$

and

$$\|f - \tilde{f}\| \le \|f\|_{\mathcal{B}_{p,r}} t^{\theta}.$$
(61)

For $\tilde{f} \in \mathcal{A}_{1,r}$, according to (54), there exists $h := h(m) \in \text{Span}\{\mathcal{D}_m\}$ for every *m* that satisfies

$$\|h\|_{\mathcal{A}_1(\mathcal{D}_m)} \le \|f\|_{\mathcal{A}_{1,r}} \tag{62}$$

and

$$\|\tilde{f} - h\| \le \|\tilde{f}\|_{\mathcal{A}_{1,r}} m^{-r}.$$
(63)

The relations (60), (62), and (63) imply

$$\|h\|_{\mathcal{A}_1(\mathcal{D}_m)} \le \|f\|_{\mathcal{B}_{P,r}} t^{\theta-1} \tag{64}$$

and

$$\|\tilde{f} - h\| \le \|f\|_{\mathcal{B}_{p,r}} t^{\theta - 1} m^{-r}.$$
(65)

Then combining (61) with (65), we obtain

$$\|f - h\| \le \|f\|_{\mathcal{B}_{p,r}} (t^{\theta} + t^{\theta - 1} m^{-r}).$$
(66)

From (64) and (66), there exists $h := h(m) \in \text{Span}\{\mathcal{D}_m\}$ for every *m* and t > 0 that satisfies

$$\|h\|_{\mathcal{A}_1(\mathcal{D}_m)} \le M t^{\theta-1}$$

and

$$\|f_{\rho}-h\|\leq M(t^{\theta}+t^{\theta-1}m^{-r}),$$

where $M = ||f_{\rho}||_{\mathcal{B}_{p,r}}$.

Therefore, Theorem 2 with $t = k^{-\frac{1}{2}}$ implies

$$E(\|f_{\mathbf{z}} - f_{\rho}\|^{2}) \le C \min_{k>0} \left(M^{2} t_{0}^{-2} k^{1-\frac{2}{p}} + M^{2} \left(k^{\frac{1}{2} - \frac{1}{p}} + k^{1-\frac{1}{p}} n^{-ar} \right)^{2} + \frac{k \log n}{n} \right).$$
(67)

The condition $2ar \ge 1$ also enables us to eliminate the term involving n^{-ar} . Then, by taking $k := \lceil \frac{(M+1)^2}{t_0^2} \frac{n}{\log n} \rceil^{\frac{p}{2}}$ in (67), we obtain the desired result (59).

Then we show the universal consistency of the WRPGLA.

Theorem 3 Under the assumptions of Theorem 2, if the dictionary \mathcal{D} is complete in $L^2_{\rho_X}(X)$, for any f_{ρ} , we have

$$\lim_{n \to +\infty} E\left(\|f_{\mathbf{z}} - f_{\rho}\|^2\right) = 0.$$
(68)

Proof Since \mathcal{D} is complete in $L^2_{\rho_X}(X)$, we can find $h \in \text{Span}\{\mathcal{D}_m\}$ satisfying $||f_{\rho} - h|| \le \varepsilon$, where $\varepsilon > 0$ and n is big enough. It follows from Theorem 2 that

$$E(\|f_{\mathbf{z}} - f_{\rho}\|^{2}) \le C \min_{k>0} \left(\frac{\|h\|_{\mathcal{A}_{1}(\mathcal{D}_{m})}^{2}}{kt_{0}^{2}} + \varepsilon^{2} + \frac{k \log n}{n} \right).$$
(69)

To balance the first and third error term, we choose $k := n^{\frac{1}{2}} t_0^{-1}$, which implies

$$E(\|f_{\mathbf{z}} - f_{\rho}\|^2) \le C(\varepsilon^2 + t_0^{-1} n^{-\frac{1}{2}} \log n).$$
(70)

 \Box

Thus, for *n* sufficiently large,

$$E(\|f_{\mathbf{z}} - f_{\rho}\|^2) \le 2C\varepsilon^2. \tag{71}$$

This completes the proof of Theorem 3.

Remark 3 It is known from [11] that the OGA and the RGA can achieve the optimal convergence rate $\mathcal{O}(m^{-\frac{1}{2}})$ on $\mathcal{A}_1(\mathcal{D})$. When $t_k = 1$, Lemma 4.1 shows that the WRPGA also attains the best rate. Meanwhile, we compare the WRPGLA with the OGLA and the relaxed greedy learning algorithm (RGLA). For $f_{\rho} \in \mathcal{A}_{1,r}$, we derive the same convergence rate $\mathcal{O}((n \log n)^{-1/2})$ of the WRPGLA as that of the OGLA and the RGLA in Ref. [7]. For $f_{\rho} \in \mathcal{B}_{p,r}$, when $p \to 1$, the rate $\mathcal{O}((n \log n)^{-1+\frac{p}{2}})$ of the WRPGLA can be arbitrarily close to $\mathcal{O}((n \log n)^{-1/2})$.

Moreover, from the viewpoint of the computational complexity, for the WRPGLA, the approximant f_k is constructed by solving a one-dimensional optimization problem since f_k is an orthogonal projection of f onto Span{ \hat{f}_k }. On the other hand, the OGLA is more expensive to implement since at each step, the algorithm requires the evaluation of orthogonal projection on a k-dimensional space, and the output is constructed by solving a k-dimensional optimization problem. And it is clear that the WRPGLA is simpler than the RGLA. Thus, the WRPGLA should essentially reduce the complexity and make the learning process accelerated.

In future research, it would be an interesting project to deduce the error bound of the WRPGLA in Banach spaces with modulus of smoothness $\rho(u) \leq \gamma u^q$, $1 < q \leq 2$ as [24, 27]. Furthermore, Guo and Ye [28, 29] derived the convergence rates of the moving least-squares learning algorithm for the weakly dependent and nonidentical samples. It remains open to explore the greedy learning algorithms in the non-i.i.d. and nonidentical sampling setting.

Funding

This research is supported by the National Science Foundation for Young Scientists of China (Grant No. 12001328), the National Natural Science Foundation of China (Grant No. 11671213), the Development Plan of Youth Innovation Team of University in Shandong Province (No. 2021KJ067), the Shandong Provincial Natural Science Foundation of China (No. ZR2022MF223), and the Qilu University of Technology (Shandong Academy of Sciences) Talent Research Project (No. 2023RCKY133).

Data availability

All data, models, and code generated or used during the study appear in the submitted article.

Declarations

Competing interests

The authors declare no competing interests.

Author contributions

All authors contributed substantially to this paper, participated in drafting and checking the manuscript. All authors read and approved the final manuscript.

Author details

¹School of Science, Shandong Jianzhu University, Jinan, 250101, China. ²School of Mathematical Sciences and LPMC, Nankai University, Tianjin, 300071, China.

Received: 2 September 2023 Accepted: 2 January 2024 Published online: 04 March 2024

References

- Fang, J., Lin, S.B., Xu, Z.B.: Learning and approximation capabilities of orthogonal super greedy algorithm. Knowl.-Based Syst. 95, 86–98 (2016)
- Chen, H., Li, L.Q., Pan, Z.B.: Learning rates of multi-kernel regression by orthogonal greedy algorithm. J. Stat. Plan. Inference 143, 276–282 (2013)
- Lin, S.B., Rong, Y.H., Sun, X.P., Xu, Z.B.: Learning capability of relaxed greedy algorithms. IEEE Trans. Neural Netw. Learn. Syst. 24(10), 1598–1608 (2013)
- Xu, L., Lin, S.B., Xu, Z.B.: Learning capability of the truncated greedy algorithm. Sci. China Inf. Sci. 59(5), 052103 (2016). https://doi.org/10.1007/s11432-016-5536-6
- Barron, A.R.: Universal approximation bounds for superposition of a sigmoidal function. IEEE Trans. Inf. Theory 39(3), 930–945 (1993)
- Guo, Q.: Distributed semi-supervised regression learning with coefficient regularization. Results Math. 77, 1–19 (2022)
 Barron, A.R., Cohen, A., Dahmen, W., DeVore, R.A.: Approximation and learning by greedy algorithms. Ann. Stat. 36(1),
- 64–94 (2008)
 8. Xu, L., Lin, S.B., Zeng, J.S., Liu, X., Xu, Z.B.: Greedy criterion in orthogonal greedy learning. IEEE Trans. Cybern. 48(3),
- AU, L., Lin, S.D., Zeng, J.S., Liu, A., Xu, Z.B.: Greedy criterion in orthogonal greedy learning. IEEE Trans. Cybern. 46(5), 955–966 (2018)
- 9. Jones, L.K.: A simple lemma on greedy approximation in Hilbert spaces and convergence rates for projection pursuit regression and neural network training. Ann. Stat. **20**(1), 608–613 (1992)
- Lee, W.S., Bartlett, P.L., Williamson, R.C.: Efficient agnostic learning of neural networks with bounded fan-in. IEEE Trans. Inf. Theory 42(6), 2118–2132 (1996)
- 11. DeVore, R.A., Temlyakov, V.N.: Some remarks on greedy algorithms. Adv. Comput. Math. 5, 173–187 (1996)
- 12. Temlyakov, V.N.: Greedy Approximation. Cambridge University Press, Cambridge (2011)
- 13. Dai, W., Milenkovic, O.: Subspace pursuit for compressive sensing signal recontruction. IEEE Trans. Inf. Theory 55(5), 2230–2249 (2009)
- 14. Kunis, S., Rauhut, H.: Random sampling of sparse trigonometric polynomials ii-orthogonal matching pursuit versus basis pursuit. Found. Comput. Math. **8**, 737–763 (2008)
- Donoho, D.L., Tsaig, Y., Drori, I., Starck, J.L.: Sparse solution of underdetermined systems of linear equations by stagewise orthogonal matching pursuit. IEEE Trans. Inf. Theory 58(2), 1094–1121 (2012)
- Tropp, J.A., Wright, S.: Computational methods for sparse solution of linear inverse problems. Proc. IEEE 98(6), 948–958 (2010)
- 17. Temlyakov, V.N., Zheltov, P.: On performance of greedy algorithms. J. Approx. Theory 163(9), 1134–1145 (2011)
- Donoho, D.L., Elad, M., Temlyakov, V.N.: On Lebesgue-type inequalities for greedy approximation. J. Approx. Theory 147(2), 185–195 (2007)
- Chen, H., Zhou, Y.C., Tang, Y.Y., Li, L.Q., Pan, Z.B.: Convergence rate of the semi-supervised greedy algorithm. Neural Netw. 44, 44–50 (2013)
- Schmidt, E.: Zur theorie der linearen und nicht linearen integralgleichungen zweite abhandlung. Math. Ann. 64, 161–174 (1907)
- 21. Temlyakov, V.N.: Greedy approximation. Acta Numer. 17, 235–409 (2008)
- 22. Livshitz, E.D., Temlyakov, V.N.: Two lower estimates in greedy approximation. Constr. Approx. 19, 509–523 (2003)
- 23. Livshits, E.D.: Lower bounds for the rate of convergence of greedy algorithms. Izv. Math. 73, 1197–1215 (2009)
- Petrova, G.: Rescaled pure greedy algorithm for Hilbert and Banach spaces. Appl. Comput. Harmon. Anal. 41, 852–866 (2016)
- 25. Jiang, B., Ye, P.X.: Efficiency of the weak rescaled pure greedy algorithm. Int. J. Wavelets Multiresolut. Inf. Process. 19(4), 2150001 (2021)
- 26. Bergh, J., Lofstrom, J.: Interpolation Spaces. Springer, Berlin (1976)
- 27. Temlyakov, V.N.: Greedy algorithms in Banach spaces. Adv. Comput. Math. 14, 277–292 (2001)
- Guo, Q., Ye, P.X.: Error analysis of the moving least-squares method with non-identical sampling. Int. J. Comput. Math. 96(4), 767–781 (2019)
- 29. Guo, Q., Ye, P.X.: Error analysis of the moving least-squares regression learning algorithm with β -mixing and non-identical sampling. Int. J. Comput. Math. **97**(8), 1586–1602 (2020)

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Submit your manuscript to a SpringerOpen^o journal and benefit from:

- Convenient online submission
- ► Rigorous peer review
- ► Open access: articles freely available online
- ► High visibility within the field
- Retaining the copyright to your article

Submit your next manuscript at > springeropen.com