

## Research Article

# On the Convergence Rate of Kernel-Based Sequential Greedy Regression

Xiaoyin Wang,<sup>1</sup> Xiaoyan Wei,<sup>2</sup> and Zhibin Pan<sup>1</sup>

<sup>1</sup> College of Sciences, Huazhong Agricultural University, Wuhan 430070, China

<sup>2</sup> Department of Statistics and Applied Mathematics, Hubei University of Economics, Wuhan 430205, China

Correspondence should be addressed to Zhibin Pan, zhibinpan2008@gmail.com

Received 13 October 2012; Accepted 27 November 2012

Academic Editor: Jean M. Combes

Copyright © 2012 Xiaoyin Wang et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

A kernel-based greedy algorithm is presented to realize efficient sparse learning with data-dependent basis functions. Upper bound of generalization error is obtained based on complexity measure of hypothesis space with covering numbers. A careful analysis shows the error has a satisfactory decay rate under mild conditions.

## 1. Introduction

Kernel methods have been extensively utilized in various learning tasks, and its generalization performance has been investigated from the viewpoint of approximation theory [1, 2]. Among these methods, a family of them can be considered as coefficient-based regularized framework in data-dependent hypothesis spaces; see, for example, [3–8]. For given samples  $\{(x_i, y_i)\}_{i=1}^n$ , the solution of these kernel methods has the following expression  $\sum_{i=1}^n \alpha_i K(x_i, \cdot)$ , where  $\alpha_i \in \mathcal{R}$  and  $K$  is a Mercer kernel. The aim of these coefficient-based algorithms is to search a set of coefficients  $\{\alpha_i\}$  with good predictive performance.

Inspired by greedy approximation methods in [9–12], we propose a sparse greedy algorithm for regression. The greedy approximation has two advantages over the regularization methods: one is that the sparsity is directly controlled by a greedy approximation algorithm, rather than by the regularization parameter; the other is that greedy approximation does not change the objective optimization function, while the regularized methods usually modify the objective function by including a sparse regularization term [13].

Before introducing the greedy algorithm, we recall some preliminary background for regression. Let the input space  $\mathcal{X} \subset \mathbb{R}^d$  be a compact subset and  $\mathcal{Y} = [-M, M]$  for some constant  $M > 0$ . In the regression model, the learner gets a sample set  $\mathbf{z} = \{(x_i, y_i)\}_{i=1}^n$ , where

$(x_i, y_i) \in \mathcal{Z} := \mathcal{X} \times \mathcal{Y}, 1 \leq i \leq n$ , are randomly independently drawn from an unknown distribution  $\rho$  on  $\mathcal{Z}$ . The goal of learning is to pick a function  $f : \mathcal{X} \rightarrow \mathcal{Y}$  with the expected error

$$\mathcal{E}(f) = \int_{\mathcal{Z}} (f(x) - y)^2 d\rho \quad (1.1)$$

as small as possible. Note that the regression function

$$f_\rho(x) = \int_{\mathcal{Y}} y d\rho(y | x), \quad x \in \mathcal{X}, \quad (1.2)$$

is the minimizer of  $\mathcal{E}(f)$ , where  $\rho(\cdot | x)$  is the conditional probability measure at  $x$  induced by  $\rho$ .

The empirical error is defined as

$$\mathcal{E}_z(f) = \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2. \quad (1.3)$$

We call a symmetric and positive semidefinite continuous function  $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  a Mercer kernel. The reproducing kernel Hilbert space (RKHS)  $\mathcal{H}_K$  is defined to be the closure of the linear span of the set of functions  $\{K_x := K(x, \cdot) : x \in \mathcal{X}\}$  with the inner product  $\langle \cdot, \cdot \rangle_K$  defined by  $\langle K_x, K_{x'} \rangle_K = K(x, x')$ . For all  $x \in \mathcal{X}$  and  $f \in \mathcal{H}_K$ , the reproducing property is given by  $\langle K_x, f \rangle_K = f(x)$ . We can see  $\kappa := \sup_{x \in \mathcal{X}} \sqrt{K(x, x)} < \infty$  because of the continuity of  $K$  and the compactness of  $\mathcal{X}$ .

Different from the coefficient-based regularized method [3–6], we use the idea of sequential greedy approximation to realize sparse learning in this paper. Denote  $\widehat{\mathcal{H}} = \{\widehat{h}_i\}_{i=1}^{2n}$ , where  $\widehat{h}_{2i-1} = K_{x_i}$  and  $\widehat{h}_{2i} = -K_{x_i}$ . The hypothesis space (depending on  $\mathbf{z}$ ) is defined as

$$CO_{2n}(\widehat{\mathcal{H}}) = \left\{ f : f(x) = \sum_{i=1}^{2n} \alpha_i \widehat{h}_i(x), \alpha_i \geq 0, \sum_{i=1}^{2n} \alpha_i \leq 1 \right\}. \quad (1.4)$$

For any hypothesis function space  $\mathcal{G}$ , we denote  $\beta\mathcal{G} = \{f : f = \beta g, g \in \mathcal{G}\}$ .

The definition of  $f_\rho$  tells us  $|f_\rho(x)| \leq M$ , so it is natural to restrict the approximating functions to  $[-M, M]$ . The projection operator has been used in error analysis of learning algorithms (see, e.g., [2, 14]).

*Definition 1.1.* The projection operator  $\pi = \pi_M$  is defined on the space of measurable functions  $f : \mathcal{X} \rightarrow \mathcal{R}$  as

$$\pi(f)(x) = \begin{cases} M, & \text{if } f(x) > M; \\ -M, & \text{if } f(x) < -M; \\ f(x), & \text{otherwise.} \end{cases} \quad (1.5)$$

The kernel-based greedy algorithm can be summarized as below. Let  $t$  be a stopping time and let  $\beta$  be a positive constant. Set  $\hat{f}_\beta^0 = 0$ . And then for  $\tau = 1, 2, \dots, t$ , define

$$\begin{aligned}\hat{h}_\tau, \hat{\alpha}_\tau, \hat{\beta}_\tau &= \arg \min_{h \in \mathcal{H}, 0 \leq \alpha \leq 1, 0 \leq \beta' \leq \beta} \mathcal{E}_z \left( (1 - \alpha) \hat{f}_\beta^{\tau-1} + \alpha \beta' h \right), \\ \hat{f}_\beta^\tau &= (1 - \hat{\alpha}_\tau) \hat{f}_\beta^{\tau-1} + \hat{\alpha}_\tau \hat{\beta}_\tau \hat{h}_\tau.\end{aligned}\quad (1.6)$$

Different from the regularized algorithms in [6, 12, 14–18], the above learning algorithm tries to realize efficient learning by greedy approximation. The study for its generalization performance can enrich the learning theory of kernel-based regression. In the remainder of this paper, we focus on establishing the convergence rate of  $\pi(\hat{f}_\beta^t)$  to the regression function  $f_\rho$  under choice of suitable parameters. The theoretical result is dependent on weaker conditions than the previous error analysis for kernel-based regularization framework in [4, 5].

## 2. Main Result

Define a data-free basis function set

$$\begin{aligned}\mathcal{H} &= \{h_i : h_{2i-1} = K_{u_i}, h_{2i} = -K_{u_i}, u_i \subset \mathcal{X}, i = 1, \dots, \infty\}, \\ CO(\mathcal{H}) &= \left\{ f : f(x) = \sum_{i=1}^{\infty} \alpha_i h_i(x), \alpha_i \geq 0, \sum_{i=1}^{\infty} \alpha_i \leq 1 \right\}.\end{aligned}\quad (2.1)$$

To investigate the approximation of  $\pi(\hat{f}_\beta^t)$  to  $f_\rho$ , we introduce a data-independent function

$$f_\beta^* = \arg \min_{f \in \beta CO(\mathcal{H})} \mathcal{E}(f). \quad (2.2)$$

Observe that

$$\begin{aligned}& \mathcal{E}(\pi(\hat{f}_\beta^t)) - \mathcal{E}(f_\rho) \\ & \leq \left\{ \mathcal{E}(\pi(\hat{f}_\beta^t)) - \mathcal{E}_z(\pi(\hat{f}_\beta^t)) + \mathcal{E}_z(f_\beta^*) - \mathcal{E}(f_\beta^*) \right\} + \left\{ \mathcal{E}_z(\pi(\hat{f}_\beta^t)) - \mathcal{E}_z(f_\beta^*) \right\} \\ & \quad + \left\{ \mathcal{E}(f_\beta^*) - \mathcal{E}(f_\rho) \right\}.\end{aligned}\quad (2.3)$$

Here, the three terms on the right-hand side are called as the sample error, the hypothesis error, and the approximation error, respectively.

To estimate the sample error, we usually need the complexity measure of hypothesis function space  $\mathcal{H}_K$ . For this reason, we introduce some definitions of covering numbers to measure the complexity.

*Definition 2.1.* Let  $(\mathcal{U}, d)$  be a pseudometric space and denote a subset  $S \subset \mathcal{U}$ . For every  $\epsilon > 0$ , the covering number  $\mathcal{N}(S, \epsilon, d)$  of  $S$  with respect to  $\epsilon, d$  is defined as the minimal number of balls of radius  $\epsilon$  whose union covers  $S$ , that is,

$$\mathcal{N}(S, \epsilon, d) = \min \left\{ l \in \mathbb{N} : S \subset \bigcup_{j=1}^l B(s_j, \epsilon) \text{ for some } \{s_j\}_{j=1}^l \subset \mathcal{U} \right\}, \quad (2.4)$$

where  $B(s_j, \epsilon) = \{s \in \mathcal{U} : d(s, s_j) \leq \epsilon\}$  is a ball in  $\mathcal{U}$ .

The empirical covering number with  $\ell_2$  metric is defined as below.

*Definition 2.2.* Let  $\mathcal{F}$  be a set of functions on  $\mathcal{X}$ ,  $\mathbf{u} = (u_i)_{i=1}^k$  and  $\mathcal{F}|_{\mathbf{u}} = \{(f(u_i))_{i=1}^k : f \in \mathcal{F}\} \subset \mathbb{R}^k$ . Set  $\mathcal{N}_{2,\mathbf{u}}(\mathcal{F}, \epsilon) = \mathcal{N}_{2,\mathbf{u}}(\mathcal{F}|_{\mathbf{u}}, \epsilon, d_2)$ . The  $\ell_2$  empirical covering number of  $\mathcal{F}$  is defined by

$$\mathcal{N}_2(\mathcal{F}, \epsilon) = \sup_{k \in \mathbb{N}} \sup_{\mathbf{u} \in \mathcal{X}^k} \mathcal{N}_{2,\mathbf{u}}(\mathcal{F}, \epsilon), \quad \epsilon > 0, \quad (2.5)$$

where  $\ell_2$  metric

$$d_2(\mathbf{a}, \mathbf{b}) = \left( \frac{1}{k} \sum_{i=1}^k |a_i - b_i|^2 \right)^{1/2}, \quad \forall \mathbf{a} = (a_i)_{i=1}^k \in \mathbb{R}^k, \mathbf{b} = (b_i)_{i=1}^k \in \mathbb{R}^k. \quad (2.6)$$

Denote  $\mathcal{B}_r$  as the ball of radius  $r$  with  $r > 0$ , where  $\mathcal{B}_r = \{f \in \mathcal{H}_K : \|f\|_{\mathcal{H}_K} \leq r\}$ . We need the following capacity assumption on  $\mathcal{H}_K$ , which has been used in [5, 6, 18].

*Assumption 2.3.* There exist an exponent  $p$ , with  $p \in (0, 2)$  and a constant  $c_{p,K} > 0$  such that

$$\log \mathcal{N}_2(\mathcal{B}_1, \epsilon) \leq c_{p,K} \epsilon^{-p}. \quad (2.7)$$

We now formulate the generalization error bounds for  $\pi(\hat{f}_\beta^t)$ . The result follows from Propositions 3.2–3.5 in the next section.

**Theorem 2.4.** *Under Assumption 2.3, for any  $0 < \delta < 1$ , the following inequality holds with confidence  $1 - \delta$*

$$\begin{aligned} \mathcal{E}(\pi(\hat{f}_\beta^t)) - \mathcal{E}(f_\rho) &\leq 4\left(\mathcal{E}(f_\beta^*) - \mathcal{E}(f_\rho)\right) + \frac{32\beta^2}{t} + \frac{4(3M + \kappa^2\beta)^2}{n} \log\left(\frac{2}{\delta}\right) \\ &\quad + 1280M^2(c_{p,K}(4M\kappa\beta)^p)^{2/(2+p)} \log\left(\frac{2}{\delta}\right) n^{-2/(2+p)}. \end{aligned} \quad (2.8)$$

From the result, we know there exists a constant  $C$  independent of  $n, t, \delta$  such that with confidence  $1 - \delta$

$$\mathcal{E}(\pi(\hat{f}_\beta^t)) - \mathcal{E}(f_\rho) \leq 4\left(\mathcal{E}(f_\beta^*) - \mathcal{E}(f_\rho)\right) + C \max\left\{\frac{\beta^2}{t}, \frac{\beta^2}{n}, \left(\frac{\beta^p}{n}\right)^{2/(2+p)}\right\} \log\left(\frac{2}{\delta}\right). \quad (2.9)$$

In particular, if  $f_\rho \in \tilde{\beta}CO(\mathcal{H})$  for some fixed constant  $\tilde{\beta}$  and  $t \geq n$ , we have  $\mathcal{E}(\pi(\hat{f}_\beta^t)) - \mathcal{E}(f_\rho) \rightarrow 0$  with decay rate  $O(n^{-2/(2+p)})$ . The learning rate is satisfactory as  $p \rightarrow 0$ .

Here, the estimate of the hypothesis error is simple and does not need the strict condition on  $\rho$  and  $\mathcal{K}$  in [3–5] for learning with data-dependent hypothesis spaces.

If there are some additional conditions on approximation error with the increasing of  $\beta$ , we can obtain the explicit learning rates with suitable parameter selection.

**Corollary 2.5.** *Assume that the RKHS  $\mathcal{H}_K$  satisfies (2.7) and  $\mathcal{E}(f_\beta^*) - \mathcal{E}(f_\rho) \leq c_\gamma \beta^{-\gamma}$  for some  $\gamma > 0$ . Choose  $\beta = n^{p/(4(2+p))}$ . For any  $0 < \delta < 1$  and  $t = n$ , one has*

$$\mathcal{E}(\pi(\hat{f}_\beta^t)) - \mathcal{E}(f_\rho) \leq C n^{-\min\{(4+p)/(4+2p), p\gamma/(8+4p)\}} \log\left(\frac{2}{\delta}\right) \quad (2.10)$$

with confidence  $1 - \delta$ . Here  $C$  is a constant independent of  $n, \delta$ .

Observe that the learning rate depends closely on the approximation condition between  $f_\rho$  and  $f_\beta^*$ . This means that only the target function can be well described by the functions from the hypothesis space, the learning algorithm can achieve good generalization performance. In fact, similar approximation assumption is extensively studied for error analysis in learning theory; see, for example, [1, 2, 4, 17].

From Corollary 2.5, when the kernel  $K \in C^\infty$ ,  $p > 0$  can be arbitrarily small, one can easily see that the learning rate is quite low. Future research direction may be furthered to improve the estimate by introducing some new analysis techniques.

### 3. Proof of Theorem 2.4

In this section, we provide the proof of Theorem 2.4 based on the upper bound estimates of sample error and hypothesis error. Denote

$$\begin{aligned} S_1 &= \left\{ \mathcal{E}_z(f_\beta^*) - \mathcal{E}_z(f_\rho) \right\} - \left\{ \mathcal{E}(f_\beta^*) - \mathcal{E}(f_\rho) \right\}, \\ S_2 &= \left\{ \mathcal{E}(\pi(\hat{f}_\beta^t)) - \mathcal{E}(f_\rho) \right\} - \left\{ \mathcal{E}_z(\pi(\hat{f}_\beta^t)) - \mathcal{E}_z(f_\rho) \right\}. \end{aligned} \quad (3.1)$$

We can observe that the sample error

$$\mathcal{E}(\pi(\hat{f}_\beta^t)) - \mathcal{E}_z(\pi(\hat{f}_\beta^t)) + \mathcal{E}_z(f_\beta^*) - \mathcal{E}(f_\beta^*) = S_1 + S_2. \quad (3.2)$$

Here  $S_1$  can be bounded by applying the following one-side Bernstein type probability inequality; see, for example, [1, 2, 14].

**Lemma 3.1.** *Let  $\xi$  be a random variable on a probability space  $Z$  with mean  $E\xi$  and variance  $\sigma^2(\xi) = \sigma^2$ . If  $|\xi(z) - E\xi| \leq B$  for almost all  $z \in Z$ , then for all  $\varepsilon > 0$ ,*

$$\text{Prob}_{z \in Z^n} \left\{ \frac{1}{n} \sum_{i=1}^n \xi(z_i) - E\xi \geq \varepsilon \right\} \leq \exp \left\{ -\frac{n\varepsilon^2}{2(\sigma^2 + B\varepsilon/3)} \right\}. \quad (3.3)$$

**Proposition 3.2.** For any  $0 < \delta < 1$ , with confidence  $1 - \delta$ , one has

$$S_1 \leq \mathcal{E}(f_\beta^*) - \mathcal{E}(f_\rho) + \frac{2(3M + \kappa^2\beta)^2}{n} \log\left(\frac{1}{\delta}\right). \quad (3.4)$$

*Proof.* Following the definition of  $S_1$ , we have  $S_1 = (1/n) \sum_{i=1}^n \xi(x_i) - E\xi$ , where random variable  $\xi(x) = (y - f_\beta^*(x))^2 - (y - f_\rho(x))^2$ .

From the definition of  $f_\beta^*$ , we know  $\|f_\beta^*\|_K \leq \kappa\beta$  and  $\|f_\beta^*\|_\infty \leq \kappa\|f_\beta^*\|_K \leq \kappa^2\beta$ . Then

$$|\xi(x)| = \left| (f_\beta^*(x) - f_\rho(x)) \left( (f_\beta^*(x) - y) + (f_\rho(x) - y) \right) \right| \leq (3M + \kappa^2\beta)^2 := c_1 \quad (3.5)$$

and  $|\xi - E\xi| \leq 2c_1$ . Moreover,

$$\begin{aligned} \sigma^2 &\leq E\xi^2 = \int_{\mathcal{X}} (f_\beta^*(x) - f_\rho(x))^2 \left( (f_\beta^*(x) - y) + (f_\rho(x) - y) \right)^2 d\rho \\ &\leq c_1 \left\{ \mathcal{E}(f_\beta^*) - \mathcal{E}(f_\rho) \right\}. \end{aligned} \quad (3.6)$$

Applying Lemma 3.1 with  $B = 2c_1$  and  $\sigma^2 = c_1 \{ \mathcal{E}(f_\beta^*) - \mathcal{E}(f_\rho) \}$ , we get

$$\frac{1}{n} \sum_{i=1}^n \xi(x_i) - E\xi \leq t \quad (3.7)$$

with confidence at least  $1 - \exp\{-nt^2/(2c_1(\mathcal{E}(f_\beta^*) - \mathcal{E}(f_\rho) + (2/3)t))\}$ . By setting  $-nt^2/(2c_1(\mathcal{E}(f_\beta^*) - \mathcal{E}(f_\rho) + (2/3)t)) = \log(\delta)$ , we derive the solution

$$\begin{aligned} t^* &= \frac{\left( (2c_1/3) \log(1/\delta) + \sqrt{((2c_1/3) \log(1/\delta))^2 + 2c_1 \log(1/\delta) (\mathcal{E}(f_\beta^*) - \mathcal{E}(f_\rho))} \right)}{n} \\ &\leq \frac{2c_1}{n} \log\left(\frac{1}{\delta}\right) + \mathcal{E}(f_\beta^*) - \mathcal{E}(f_\rho). \end{aligned} \quad (3.8)$$

Thus, with confidence  $1 - \delta$ , we have

$$\frac{1}{n} \sum_{i=1}^n \xi(x_i) - E\xi \leq \frac{2c_1}{n} \log\left(\frac{1}{\delta}\right) + \mathcal{E}(f_\beta^*) - \mathcal{E}(f_\rho). \quad (3.9)$$

This completes the proof.  $\square$

To establish the uniform upper bound of  $S_2$ , we introduce a concentration inequality established in [18].

**Lemma 3.3.** Assume that there are constants  $B, c > 0$  and  $\alpha \in [0, 1]$  such that  $\|f\|_\infty \leq B$  and  $Ef \leq c(Ef)^\alpha$  for every  $f \in \mathcal{F}$ . If for some  $a > 0$  and  $p \in (0, 2)$ ,

$$\log(\mathcal{N}_2(\mathcal{F}, \epsilon)) \leq a\epsilon^{-p}, \quad \forall \epsilon > 0, \quad (3.10)$$

then there exists a constant  $c'_p$  depending only on  $p$  such that for any  $t > 0$ , with probability at least  $1 - e^{-t}$ , there holds

$$Ef - \frac{1}{n} \sum_{i=1}^n f(z_i) \leq \frac{1}{2} \eta^{1-\alpha} (Ef)^\alpha + c'_p \eta + 2 \left( \frac{ct}{n} \right)^{1/(2-\alpha)} + \frac{18Bt}{n}, \quad \forall f \in \mathcal{F}, \quad (3.11)$$

where

$$\eta := \max \left\{ c^{(2-p)/(4-2\alpha+p\alpha)} \left( \frac{a}{n} \right)^{2/(4-2\alpha+p\alpha)}, B^{(2-p)/(2+p)} \left( \frac{a}{n} \right)^{2/(2+p)} \right\}. \quad (3.12)$$

**Proposition 3.4.** Under Assumption 2.3, for any  $0 < \delta < 1$ , one has with confidence at least  $1 - \delta$

$$S_2 \leq \frac{1}{2} \left\{ \mathcal{E}(\pi(\hat{f}_\beta^t)) - \mathcal{E}(f_\rho) \right\} + 640M^2 (c_{p,K} (4M\kappa\beta)^p)^{2/(2+p)} \log\left(\frac{1}{\delta}\right) n^{-2/(2+p)}. \quad (3.13)$$

*Proof.* From the definition of  $\hat{f}_\beta^t$ , we have  $\|\hat{f}_\beta^t\|_K \leq \kappa\beta$ . Denote

$$\mathcal{F}_{\kappa\beta} = \left\{ g(z) = (y - \pi(f)(x))^2 - (y - f_\rho(x))^2 : f \in \mathcal{B}_{\kappa\beta} \right\}. \quad (3.14)$$

We can see that  $Eg = \mathcal{E}(\pi(f)) - \mathcal{E}(f_\rho)$  and  $(1/n) \sum_{i=1}^n g(z_i) = \mathcal{E}_z(\pi(f)) - \mathcal{E}_z(f_\rho)$ . Since  $\|\pi(f)\|_\infty \leq M$  and  $|f_\rho(x)| \leq M$ , we have

$$\begin{aligned} |g(z)| &= |(\pi(f)(x) - f_\rho(x))((\pi(f)(x) - y) + (f_\rho(x) - y))| \leq 8M^2, \\ Eg^2 &= \int_{\mathcal{Z}} (\pi(f)(x) - f_\rho(x))^2 ((\pi(f)(x) - y) + (f_\rho(x) - y))^2 d\rho \leq 16M^2 Eg. \end{aligned} \quad (3.15)$$

For  $g_1, g_2 \in \mathcal{F}_{\kappa\beta}$ , we have

$$\begin{aligned} |g_1(z) - g_2(z)| &= |(y - \pi(f_1)(x))^2 - (y - \pi(f_2)(x))^2| \leq 4M |\pi(f_1)(x) - \pi(f_2)(x)| \\ &\leq 4M |f_1(x) - f_2(x)|. \end{aligned} \quad (3.16)$$

Then, from Assumption 2.3,

$$\mathcal{N}_{2,\mathcal{Z}}(\mathcal{F}_{\kappa\beta}, \epsilon) \leq \mathcal{N}_{2,\mathcal{X}}\left(\mathcal{B}_{\kappa\beta}, \frac{\epsilon}{4M}\right) \leq \mathcal{N}_{2,\mathcal{X}}\left(\mathcal{B}_1, \frac{\epsilon}{4M\kappa\beta}\right) \leq c_{p,K} (4M\kappa\beta)^p \epsilon^{-p}. \quad (3.17)$$

Applying Lemma 3.3 with  $B = c = 16M^2$  and  $a = c_{p,K}(4M\kappa\beta)^p$ , for any  $\delta \in (0, 1)$  and for all  $g \in \mathcal{F}_{\kappa\beta}$ ,

$$\begin{aligned} Eg - \frac{1}{n} \sum_{i=1}^n g(z_i) &\leq \frac{1}{2} Eg + c'_p (16M^2)^{(2-p)/(2+p)} \left( \frac{c_{p,K}(4M\kappa\beta)^p}{n} \right)^{2/(2+p)} + 320M^2 \frac{\log(1/\delta)}{n}, \\ &\leq \frac{1}{2} Eg + 640M^2 (c_{p,K}(4M\kappa\beta)^p)^{2/(2+p)} \log\left(\frac{1}{\delta}\right) n^{-2/(2+p)} \end{aligned} \quad (3.18)$$

holds with confidence  $1 - \delta$ . This completes the proof.  $\square$

Different from the previous studies related with regularized framework [3–5], we introduce the estimate of hypothesis error  $\mathcal{E}_z(\hat{f}_\beta^t) - \mathcal{E}_z(f_\beta^*)$  based on Theorem 4.2 in [11] for sequential greedy approximation.

**Proposition 3.5.** *For a fixed sample  $\mathbf{z}$ , one has*

$$\mathcal{E}_z(\hat{f}_\beta^t) - \mathcal{E}_z(f_\beta^*) \leq \frac{16\beta^2}{t}. \quad (3.19)$$

The desired result in Theorem 2.4 can be derived directly by combining Propositions 3.2–3.5.

## Acknowledgments

This work was supported partially by the National Natural Science Foundation of China under Grant no. 11001092, Humanities and Social Science Projects of the Ministry of Education of China (Program no. 11y3jc630197), and the Fundamental Research Funds for the Central Universities (Programs nos. 2011PY130, and 2011QC022).

## References

- [1] F. Cucker and S. Smale, "On the mathematical foundations of learning," *Bulletin of the American Mathematical Society*, vol. 39, no. 1, pp. 1–49, 2002.
- [2] F. Cucker and D. X. Zhou, *Learning Theory: An Approximation Theory Viewpoint*, Cambridge University Press, Cambridge, Mass, USA, 2007.
- [3] Q. Wu and D. X. Zhou, "Learning with sample dependent hypothesis spaces," *Computers & Mathematics with Applications*, vol. 56, no. 11, pp. 2896–2907, 2008.
- [4] Q. W. Xiao and D. X. Zhou, "Learning by nonsymmetric kernels with data dependent spaces and  $\ell^1$ -regularizer," *Taiwanese Journal of Mathematics*, vol. 14, no. 5, pp. 1821–1836, 2010.
- [5] L. Shi, Y. L. Feng, and D. X. Zhou, "Concentration estimates for learning with  $\ell^1$ -regularizer and data dependent hypothesis spaces," *Applied and Computational Harmonic Analysis*, vol. 31, no. 2, pp. 286–302, 2011.
- [6] Y. L. Feng and S. G. Lv, "Unified approach to coefficient-based regularized regression," *Computers & Mathematics with Applications*, vol. 62, no. 1, pp. 506–515, 2011.
- [7] S. G. Lv and J. D. Zhu, "Error bounds for  $l^p$ -norm multiple kernel learning with least square loss," *Abstract and Applied Analysis*, vol. 2012, Article ID 915920, 18 pages, 2012.



- [8] Y. K. Zhu and H. W. Sun, "Consistency analysis of spectral regularization algorithms," *Abstract and Applied Analysis*, vol. 2012, Article ID 436510, 16 pages, 2012.
- [9] A. R. Barron, A. Cohen, W. Dahmen, and R. A. DeVore, "Approximation and learning by greedy algorithms," *The Annals of Statistics*, vol. 36, no. 1, pp. 64–94, 2008.
- [10] S. Mannor, R. Meir, and T. Zhang, "Greedy algorithms for classification—consistency, convergence rates, and adaptivity," *Journal of Machine Learning Research*, vol. 4, no. 4, pp. 713–741, 2003.
- [11] T. Zhang, "Sequential greedy approximation for certain convex optimization problems," *IEEE Transactions on Information Theory*, vol. 49, no. 3, pp. 682–691, 2003.
- [12] H. Chen, L. Q. Li, and Z. B. Pan, "Learning rates of multi-kernel regression by orthogonal greedy algorithm," *Journal of Statistical Planning and Inference*, vol. 143, no. 2, pp. 276–282, 2013.
- [13] T. Zhang, "Approximation bounds for some sparse kernel regression algorithms," *Neural Computation*, vol. 14, no. 12, pp. 3013–3042, 2002.
- [14] D. R. Chen, Q. Wu, Y. Ying, and D. X. Zhou, "Support vector machine soft margin classifiers: error analysis," *Journal of Machine Learning Research*, vol. 5, pp. 1143–1175, 2004.
- [15] H. Chen, L. Q. Li, and J. T. Peng, "Error bounds of multi-graph regularized semi-supervised classification," *Information Sciences*, vol. 179, no. 12, pp. 1960–1969, 2009.
- [16] H. Chen, "On the convergence rate of a regularized ranking algorithm," *Journal of Approximation*, vol. 164, no. 12, pp. 1513–1519, 2012.
- [17] Z. C. Guo and D. X. Zhou, "Concentration estimates for learning with unbounded sampling," *Advances in Computational Mathematics*. In press.
- [18] Q. Wu, Y. Ying, and D. X. Zhou, "Multi-kernel regularized classifiers," *Journal of Complexity*, vol. 23, no. 1, pp. 108–134, 2007.