

Research Article

ERM Scheme for Quantile Regression

Dao-Hong Xiang

Department of Mathematics, Zhejiang Normal University, Jinhua, Zhejiang 321004, China

Correspondence should be addressed to Dao-Hong Xiang; daohongxiang@gmail.com

Received 30 November 2012; Accepted 21 February 2013

Academic Editor: Ding-Xuan Zhou

Copyright © 2013 Dao-Hong Xiang. 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.

This paper considers the ERM scheme for quantile regression. We conduct error analysis for this learning algorithm by means of a variance-expectation bound when a noise condition is satisfied for the underlying probability measure. The learning rates are derived by applying concentration techniques involving the ℓ^2 -empirical covering numbers.

1. Introduction

In this paper, we study empirical risk minimization scheme (ERM) for quantile regression. Let X be a compact metric space (input space) and $Y = \mathbb{R}$. Let ρ be a fixed but unknown probability distribution on $Z := X \times Y$ which describes the noise of sampling. The conditional quantile regression aims at producing functions to estimate quantile regression functions. With a prespecified quantile parameter $\tau \in (0, 1)$, a *quantile regression function* $f_{\tau, \rho}$ is defined by its value $f_{\tau, \rho}(x)$ to be a τ -quantile of $\rho(\cdot | x)$, that is, a value $t \in Y$ satisfying

$$\rho((-\infty, t] | x) \geq \tau, \quad \rho([t, \infty) | x) \geq 1 - \tau, \quad x \in X, \quad (1)$$

where $\rho(\cdot | x)$ is the conditional distribution of ρ at $x \in X$.

We consider a learning algorithm generated by ERM scheme associated with pinball loss and hypothesis space \mathcal{H} . The *pinball loss* $L_\tau : \mathbb{R} \rightarrow [0, \infty)$ is defined by

$$L_\tau(r) = \begin{cases} (\tau - 1)r, & \text{if } r \leq 0, \\ \tau r, & \text{if } r \geq 0. \end{cases} \quad (2)$$

The hypothesis space \mathcal{H} is a compact subset of $C(X)$. So there exists some $M > 0$ such that $\|f\|_{C(X)} \leq M$ for any $f \in \mathcal{H}$. We assume without loss of generality $\|f\|_{C(X)} \leq 1$ for any $f \in \mathcal{H}$.

The ERM scheme for quantile regression is defined with a sample $\mathbf{z} = \{(x_i, y_i)\}_{i=1}^m \in Z^m$ drawn independently from ρ as follows:

$$f_{\mathbf{z}} = \arg \min_{f \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m L_\tau(y_i - f(x_i)). \quad (3)$$

A family of kernel based learning algorithms for quantile regression has been widely studied in a large literature [1–4] and references therein. The form of the algorithms is a regularized scheme in a reproducing kernel Hilbert space \mathcal{H}_K (RKHS, see [5] for details) associated with a Mercer kernel K . Given a sample \mathbf{z} the kernel based regularized scheme for quantile regression is defined by

$$f_{\mathbf{z}, \lambda} = \arg \min_{f \in \mathcal{H}_K} \left\{ \frac{1}{m} \sum_{i=1}^m L_\tau(y_i - f(x_i)) + \lambda \|f\|_K^2 \right\}. \quad (4)$$

In [1, 3, 4], error analysis for general \mathcal{H}_K has been done. Learning with varying Gaussian kernel was studied in [2].

ERM scheme (3) is very different from kernel based regularized scheme (4). The output function $f_{\mathbf{z}}$ produced by the ERM scheme has a uniform bound, under our assumption, $\|f_{\mathbf{z}}\|_{C(X)} \leq 1$. However, we cannot expect it for $f_{\mathbf{z}, \lambda}$. It is easy to see that $\lambda \|f_{\mathbf{z}, \lambda}\|_K^2 \leq \sum_{i=1}^m |y_i|$ by choosing $f = 0$. It happens often that $\|f_{\mathbf{z}, \lambda}\|_K \rightarrow \infty$ as $\lambda \rightarrow 0$. The lack of a uniform bound for $f_{\mathbf{z}, \lambda}$ has a serious negative impact on the learning rates. So in the literature of kernel based regularized schemes for quantile regression, values of the output function $f_{\mathbf{z}, \lambda}$ are always projected onto the interval $[-1, 1]$, and error analysis is conducted for the projected function, not $f_{\mathbf{z}, \lambda}$ itself.

In this paper, we aim at establishing convergence and learning rates for the error $\|f_{\mathbf{z}} - f_{\tau, \rho}\|_{L^r_{\rho_X}}$ in the space $L^r_{\rho_X}$. Here $r > 0$ depends on the pair (ρ, τ) which will be decided in Section 2 and ρ_X is the marginal distribution of ρ on X . In the rest of this paper, we assume $Y = [-1, 1]$ which in turn leads to values of the target function $f_{\tau, \rho}$ lie in the same interval.

2. Noise Condition and Main Results

There has been a large literature in learning theory (see [6] and references therein) devoted to the least square regression. It aims at learning the regression function $f_\rho(x) = \int_Y y d\rho(y | x)$. The identity $\mathcal{E}_{\text{ls}}(f) - \mathcal{E}_{\text{ls}}(f_\rho) = \|f - f_\rho\|_{L^2_{\rho_X}}^2$ for the generalization error $\mathcal{E}_{\text{ls}}(f) = \int_Z (y - f(x))^2 d\rho$ leads to a variance-expectation bound with the form of $\mathbf{E}\xi^2 \leq 4\mathbf{E}\xi$, where $\xi = (y - f(x))^2 - (y - f_\rho(x))^2$ on (Z, ρ) . It plays an essential role in error analysis of kernel based regularized schemes.

However, this identity relation and expectation-variance bound fail in the setting of the quantile regression. The reason is that the pinball loss is lack of strong convexity. If we add some noise condition on distribution ρ named τ -quantile of p -average type q (see Definition 1), we can also get a similar identity relation which in turn enables us to have a variance-expectation bound stated in the following which is proved by Steinwart and Christman [1].

Definition 1. Let $p \in (0, \infty]$ and $q \in [1, \infty)$. A distribution ρ on $X \times [-1, 1]$ is said to have a τ -quantile of p -average type q if for ρ_X -almost every $x \in X$, there exist a τ -quantile $t^* \in \mathbb{R}$ and constants $\alpha_{\rho(\cdot|x)} \in (0, 2]$, $b_{\rho(\cdot|x)} > 0$ such that for each $s \in [0, \alpha_{\rho(\cdot|x)}]$,

$$\begin{aligned} \rho((t^* - s, t^*) | x) &\geq b_{\rho(\cdot|x)} s^{q-1}, \\ \rho((t^*, t^* + s) | x) &\geq b_{\rho(\cdot|x)} s^{q-1}, \end{aligned} \quad (5)$$

and that the function γ on X defined by $\gamma(x) = b_{\rho(\cdot|x)} \alpha_{\rho(\cdot|x)}^{q-1}$ satisfies $\gamma^{-1} \in L^p_{\rho_X}$.

We also need capacity of the hypothesis space \mathcal{H} for our learning rates. Here in this paper, we measure the capacity by empirical covering numbers.

Definition 2. Let (\mathfrak{M}, d) be a pseudometric space and S be a subset of \mathfrak{M} . For every $\varepsilon > 0$, the covering number $\mathcal{N}(S, \varepsilon, d)$ of S with respect to ε and d is defined as the minimal number of balls of radius ε whose union covers S , that is,

$$\begin{aligned} \mathcal{N}(S, \varepsilon, d) \\ = \min \left\{ \ell \in \mathbb{N} : S \subset \bigcup_{j=1}^{\ell} B(s_j, \varepsilon) \text{ for some } \{s_j\}_{j=1}^{\ell} \subset \mathfrak{M} \right\}, \end{aligned} \quad (6)$$

where $B(s_j, \varepsilon) = \{s \in \mathfrak{M} : d(s, s_j) \leq \varepsilon\}$ is a ball in \mathfrak{M} .

Definition 3. Let \mathcal{F} be a set of functions on X , $\mathbf{x} = (x_i)_{i=1}^k \subset X^k$ and $\mathcal{F}|_{\mathbf{x}} = \{(f(x_i))_{i=1}^k : f \in \mathcal{F}\} \subset \mathbb{R}^k$. Set $\mathcal{N}_{2, \mathbf{x}}(\mathcal{F}, \varepsilon) = \mathcal{N}(\mathcal{F}|_{\mathbf{x}}, \varepsilon, d_2)$. The ℓ^2 -empirical covering number of \mathcal{F} is defined by

$$\mathcal{N}_2(\mathcal{F}, \varepsilon) = \sup_{k \in \mathbb{N}} \sup_{\mathbf{x} \in X^k} \mathcal{N}_{2, \mathbf{x}}(\mathcal{F}, \varepsilon), \quad \varepsilon > 0. \quad (7)$$

Here d_2 is the normalized ℓ^2 -metric on the Euclidean space \mathbb{R}^k given by

$$d_2(\mathbf{a}, \mathbf{b}) = \left(\frac{1}{k} \sum_{i=1}^k |a_i - b_i|^2 \right)^{1/2} \quad \text{for } \mathbf{a} = (a_i)_{i=1}^k, \mathbf{b} = (b_i)_{i=1}^k \in \mathbb{R}^k. \quad (8)$$

Assumption. Assume that the empirical covering number of the hypothesis space \mathcal{H} is bounded for some $a > 0$ and $\iota \in (0, 2)$,

$$\log \mathcal{N}_2(\mathcal{H}, \varepsilon) \leq a\varepsilon^{-\iota}, \quad \forall \varepsilon > 0. \quad (9)$$

Theorem 4. Assume that ρ satisfies (5) with some $p \in (0, \infty]$ and $q \in [1, \infty)$. Denote $r = pq/(p+1)$. One further assumes that $f_{\tau, \rho}$ is uniquely defined. If $f_{\tau, \rho} \in \mathcal{H}$ and \mathcal{H} satisfies (9) with $\iota \in (0, 2)$, then for any $0 < \delta < 1$, with confidence $1 - \delta$, one has

$$\|f_z - f_{\tau, \rho}\|_{L^r_{\rho_X}} \leq \bar{C} m^{-\vartheta} \log \frac{2}{\delta}, \quad (10)$$

where

$$\vartheta = \frac{2(p+1)}{4q(p+1) - (2-\iota) \min\{2(p+1), pq\}} \quad (11)$$

and \bar{C} is a constant independent of m and δ .

Remark 5. In the ERM scheme, we can choose $f_{\tau, \rho} \in \mathcal{H}$ which in turn makes the approximation error described by (23) equal to zero. However, it is impossible for the kernel based regularized scheme because of the appearance of the penalty term $\lambda \|f\|_{\mathcal{K}}^2$.

If $q \leq 2$, all conditional distributions around the quantile behave similar to the uniform distribution. In this case $r = pq/(p+1) \leq 2$ and $\theta = \min\{2/q, p/(p+1)\} = p/(p+1)$ for all $p > 0$. Hence, $\vartheta = 2(p+1)/((2+\iota)pq + 4q)$. Furthermore, when p is large enough, the parameter r tends to q and the power index for the above learning rate arbitrarily approaches to $2/(2+\iota)q$ which shows that the learning rate power index for $\|f_z - f_{\tau, \rho}\|_{L^q_{\rho_X}}$ is arbitrarily close to $2/(2+\iota)$ independent of q . In particular, ι can be arbitrarily small when \mathcal{H} is smooth enough. In this case, the power index of the learning rates $2/(2+\iota)$ can be arbitrarily close to 1 which is the optimal learning rate for the least square regression.

Let us take some examples to demonstrate the above main result.

Example 6. Let \mathcal{H} be a unit ball of the sobolev space H^s with $s > 0$. Observe that the empirical covering number is bounded above by the uniform covering number defined in Definition 2. Hence we have (see [6, 7])

$$\log \mathcal{N}_2(\mathcal{H}, \varepsilon) \leq C_s \left(\frac{1}{\varepsilon} \right)^{n/s}, \quad (12)$$

where n is the dimension of the input space X and $C_s > 0$.

Under the same assumptions as Theorem 4, we get that by replacing ι by n/s , for any $0 < \delta < 1$, with confidence $1 - \delta$,

$$\|f_{\mathbf{z}} - f_{\tau, \rho}\|_{L_{\rho_X}^r} \leq \tilde{C} m^{-\vartheta} \log \frac{2}{\delta}, \quad (13)$$

where

$$\vartheta = \frac{2(p+1)}{4q(p+1) - (2-n/s) \min\{2(p+1), pq\}} \quad (14)$$

and \tilde{C} is a constant independent of m and δ .

We carry out the same discussions on the case of $q \leq 2$ and large enough p as Remark 5. Therefore the power index of the learning rates for $\|f_{\mathbf{z}} - f_{\tau, \rho}\|_{L_{\rho_X}^q}$ is arbitrarily close to $2/(2+n/s)$ independent of q . Furthermore, s can be arbitrarily large if the Sobolev space is smooth enough. In this special case, the learning rate power index arbitrarily approaches to 1.

Example 7. Let \mathcal{H} be a unit ball of the reproducing kernel Hilbert space \mathcal{H}_σ generated by a Gaussian kernel (see [5]). Reference [7] tells us that

$$\log \mathcal{N}(\mathcal{H}, \varepsilon) \leq C_{n, \sigma} \left(\log \frac{1}{\varepsilon} \right)^{n+1}, \quad \forall \varepsilon > 0, \quad (15)$$

where $C_{n, \sigma} > 0$ depends only on n and $\sigma > 0$. Obviously, the right-hand side of (15) is bounded by $C_{n, \sigma} (1/\varepsilon)^{n+1}$.

So from Theorem 4, we can get different learning rates with power index

$$\vartheta = \frac{2(p+1)}{4q(p+1) - (1-n) \min\{2(p+1), pq\}}. \quad (16)$$

If $q \leq 2$ and p is large enough, the power index of the learning rates for $\|f_{\mathbf{z}} - f_{\tau, \rho}\|_{L_{\rho_X}^q}$ is arbitrarily close to $2/(3+n)$ which is very slow if n is large. However, in most data sets the data are concentrated on a much lower dimensional manifold embedded in the high dimensional space. In this setting an analysis that replaces n by the intrinsic dimension of the manifold would be of great interest (see [8] and references therein).

3. Error Analysis

Define the noise-free error called *generalization error* associated with the pinball loss L_τ as

$$\mathcal{E}_\tau(f) = \int_Z L_\tau(y - f(x)) d\rho \quad \text{for } f: X \rightarrow [-1, 1]. \quad (17)$$

Then the measurable function $f_{\tau, \rho}$ is a minimizer of \mathcal{E}_τ . Obviously, $f_{\tau, \rho}(x) \in [-1, 1]$.

We need the following results from [1] for our error analysis.

Proposition 8. Let L_τ be the pinball loss. Assume that ρ satisfies (5) with some $p \in (0, \infty)$ and $q \in [1, \infty)$. Then for all $f: X \rightarrow [-1, 1]$ one has

$$\|f - f_{\tau, \rho}\|_{L_{\rho_X}^r} \leq 2^{1-1/q} q^{1/q} \|\gamma^{-1}\|_{L_{\rho_X}^p}^{1/q} (\mathcal{E}_\tau(f) - \mathcal{E}_\tau(f_{\tau, \rho}))^{1/q}. \quad (18)$$

Furthermore, with

$$\theta = \min\left\{\frac{2}{q}, \frac{p}{p+1}\right\} \in (0, 1], \quad (19)$$

$$C_\theta = 2^{2-\theta} q^\theta \|\gamma^{-1}\|_{L_{\rho_X}^p}^\theta > 0,$$

one has

$$\begin{aligned} & \mathbf{E} \left\{ \left(L_\tau(y - f(x)) - L_\tau(y - f_{\tau, \rho}(x)) \right)^2 \right\} \\ & \leq C_\theta (\mathcal{E}_\tau(f) - \mathcal{E}_\tau(f_{\tau, \rho}))^\theta, \quad \forall f: X \rightarrow Y. \end{aligned} \quad (20)$$

The above result implies that we can get convergence rates of $f_{\mathbf{z}}$ in the space $L_{\rho_X}^r$ by bounding the excess generalization error $\mathcal{E}_\tau(f_{\mathbf{z}}) - \mathcal{E}_\tau(f_{\tau, \rho})$.

To bound $\mathcal{E}_\tau(f_{\mathbf{z}}) - \mathcal{E}_\tau(f_{\tau, \rho})$, we need a standard error decomposition procedure [6] and a concentration inequality.

3.1. Error Decomposition. Define the empirical error associated with the pinball loss L_τ as

$$\mathcal{E}_{\mathbf{z}, \tau}(f) = \frac{1}{m} \sum_{i=1}^m L_\tau(y_i - f(x_i)) \quad \text{for } f: X \rightarrow [-1, 1]. \quad (21)$$

Define

$$f_{\mathcal{H}} = \arg \min_{f \in \mathcal{H}} \mathcal{E}_\tau(f) = \arg \min_{f \in \mathcal{H}} \left\{ \mathcal{E}_\tau(f) - \mathcal{E}_\tau(f_{\tau, \rho}) \right\}. \quad (22)$$

Lemma 9. Let L_τ be the pinball loss, $f_{\mathbf{z}}$ be defined by (3) and $f_{\mathcal{H}} \in \mathcal{H}$ by (22). Then

$$\begin{aligned} & \mathcal{E}_\tau(f_{\mathbf{z}}) - \mathcal{E}_\tau(f_{\tau, \rho}) \\ & \leq \left[\mathcal{E}_\tau(f_{\mathbf{z}}) - \mathcal{E}_\tau(f_{\tau, \rho}) \right] - \left[\mathcal{E}_{\mathbf{z}, \tau}(f_{\mathbf{z}}) - \mathcal{E}_{\mathbf{z}, \tau}(f_{\tau, \rho}) \right] \\ & \quad + \left[\mathcal{E}_{\mathbf{z}, \tau}(f_{\mathcal{H}}) - \mathcal{E}_{\mathbf{z}, \tau}(f_{\tau, \rho}) \right] - \left[\mathcal{E}_\tau(f_{\mathcal{H}}) - \mathcal{E}_\tau(f_{\tau, \rho}) \right] \\ & \quad + \mathcal{E}_\tau(f_{\mathcal{H}}) - \mathcal{E}_\tau(f_{\tau, \rho}). \end{aligned} \quad (23)$$

Proof. The excess generalization error can be written as

$$\begin{aligned} \mathcal{E}_\tau(f_{\mathbf{z}}) - \mathcal{E}_\tau(f_{\tau, \rho}) &= \left[\mathcal{E}_\tau(f_{\mathbf{z}}) - \mathcal{E}_{\mathbf{z}, \tau}(f_{\mathbf{z}}) \right] \\ & \quad + \left[\mathcal{E}_{\mathbf{z}, \tau}(f_{\mathbf{z}}) - \mathcal{E}_{\mathbf{z}, \tau}(f_{\mathcal{H}}) \right] \\ & \quad + \left[\mathcal{E}_{\mathbf{z}, \tau}(f_{\mathcal{H}}) - \mathcal{E}_\tau(f_{\mathcal{H}}) \right] \\ & \quad + \left[\mathcal{E}_\tau(f_{\mathcal{H}}) - \mathcal{E}_\tau(f_{\tau, \rho}) \right]. \end{aligned} \quad (24)$$

The definition of $f_{\mathbf{z}}$ implies that $\mathcal{E}_{\mathbf{z}, \tau}(f_{\mathbf{z}}) - \mathcal{E}_{\mathbf{z}, \tau}(f_{\mathcal{H}}) \leq 0$. Furthermore, by subtracting and adding $\mathcal{E}_\tau(f_{\tau, \rho})$ and $\mathcal{E}_{\mathbf{z}, \tau}(f_{\tau, \rho})$ in the first term and third term, we see that Lemma 9 holds true. \square

We call the term (23) *approximation error*. It has been studied in [9].

3.2. Concentration Inequality and Sample Error. Let us recall the one-sided Bernstein inequality as follows.

Lemma 10. *Let ξ be a random variable on a probability space Z with variance σ^2 satisfying $|\xi - \mathbb{E}(\xi)| \leq M_\xi$ for some constant M_ξ . Then for any $0 < \delta < 1$, with confidence $1 - \delta$, one has*

$$\frac{1}{m} \sum_{i=1}^m \xi(z_i) - \mathbb{E}(\xi) \leq \frac{2M_\xi \log(1/\delta)}{3m} + \sqrt{\frac{2\sigma^2 \log(1/\delta)}{m}}. \quad (25)$$

Proposition 11. *Let $f_{\mathcal{H}} \in \mathcal{H}$. Assume that ρ on $X \times [-1, 1]$ satisfies the variance bound (20) with index θ indicated in (19). For any $0 < \delta < 1$, with confidence $1 - \delta/2$, (3.6) can be bounded as*

$$\begin{aligned} & \left[\mathcal{E}_{z,\tau}(f_{\mathcal{H}}) - \mathcal{E}_{z,\tau}(f_{\tau,\rho}) \right] - \left[\mathcal{E}_\tau(f_{\mathcal{H}}) - \mathcal{E}_\tau(f_{\tau,\rho}) \right] \\ & \leq \frac{4 \log(2/\delta)}{3m} + \left(\frac{2C_\theta \log(2/\delta)}{m} \right)^{1/(2-\theta)} \\ & \quad + \mathcal{E}_\tau(f_{\mathcal{H}}) - \mathcal{E}_\tau(f_{\tau,\rho}). \end{aligned} \quad (26)$$

Proof. Let $\xi(z) = L_\tau(y - f_{\mathcal{H}}(x)) - L_\tau(y - f_{\tau,\rho}(x))$ which satisfies $|\xi| \leq 2$ and in turn $|\xi - \mathbb{E}(\xi)| \leq 2$. The variance bound (20) implies that

$$\mathbb{E}(\xi - \mathbb{E}(\xi))^2 \leq \mathbb{E}\xi^2 \leq C_\theta (\mathcal{E}_\tau(f_{\mathcal{H}}) - \mathcal{E}_\tau(f_{\tau,\rho}))^\theta. \quad (27)$$

Using (25) on the random variable $\xi(z) = L_\tau(y - f_{\mathcal{H}}(x)) - L_\tau(y - f_{\tau,\rho}(x))$, we can get the desired bound (28) with the help of Young's inequality. \square

Let us turn to estimate the sample error (3.5) involving the function f_z which runs over a set of functions since \mathbf{z} is a random sample itself. To estimate it, we use a concentration inequality below involving empirical covering numbers [10–12].

Lemma 12. *Let \mathcal{F} be a class of measurable functions on Z . Assume that there are constants $B, c > 0$ and $\alpha \in [0, 1]$ and $\|f\|_\infty \leq B$ and $\mathbb{E}f^2 \leq c(\mathbb{E}f)^\alpha$ for every $f \in \mathcal{F}$. If (7) holds, then there exists a constant c'_t depending only on ι such that for any $t > 0$, with probability at least $1 - e^{-t}$, there holds*

$$\begin{aligned} \mathbb{E}f - \frac{1}{m} \sum_{i=1}^m f(z_i) & \leq \frac{1}{2} \eta^{1-\alpha} (\mathbb{E}f)^\alpha \\ & \quad + c'_t \eta + 2 \left(\frac{ct}{m} \right)^{1/(2-\alpha)} + \frac{18Bt}{m}, \quad \forall f \in \mathcal{F}, \end{aligned} \quad (28)$$

where

$$\eta := \max \left\{ c^{(2-\iota)/(4-2\alpha+\iota)} \left(\frac{a}{m} \right)^{2/(4-2\alpha+\iota)}, \right. \\ \left. B^{(2-\iota)/(2+\iota)} \left(\frac{a}{m} \right)^{2/(2+\iota)} \right\}. \quad (29)$$

We apply Lemma 12 to a function set \mathcal{F} , where

$$\mathcal{F} = \{L_\tau(y - f(x)) - L_\tau(y - f_{\tau,\rho}(x)) : f \in \mathcal{H}\}. \quad (30)$$

Proposition 13. *Assume ρ on $X \times [-1, 1]$ satisfies the variance bound (20) with index θ indicated in (19). If \mathcal{H} satisfies (9) with $\iota \in (0, 2)$, then for any $0 < \delta < 1$, with confidence $1 - \delta/2$, one has*

$$\begin{aligned} & \left[\mathcal{E}_\tau(f) - \mathcal{E}_\tau(f_{\tau,\rho}) \right] \\ & \quad - \left[\mathcal{E}_{z,\tau}(f) - \mathcal{E}_{z,\tau}(f_{\tau,\rho}) \right] \\ & \leq \frac{1}{2} \left[\mathcal{E}_\tau(f) - \mathcal{E}_\tau(f_{\tau,\rho}) \right] \\ & \quad + C_{\iota,\theta} \log \frac{2}{\delta} \left(\frac{1}{m} \right)^{2/(4-2\theta+\iota)}, \quad \forall f \in \mathcal{H}, \end{aligned} \quad (31)$$

where

$$\begin{aligned} C_{\iota,\theta} & = \left(\frac{1}{2} + c'_t \right) \max \left\{ C_\theta^{(2-\iota)/(4-2\theta+\iota)} a^{2/(4-2\theta+\iota)}, \right. \\ & \quad \left. 2^{(2-\iota)/(2+\iota)} a^{2/(2+\iota)} \right\} + 2C_\theta^{1/(2-\theta)} + 36. \end{aligned} \quad (32)$$

Proof. Take $g \in \mathcal{F}$ with the form $g(z) = L_\tau(y - f(x)) - L_\tau(y - f_{\tau,\rho}(x))$ where $f \in \mathcal{H}$. Hence $\mathbb{E}g = \mathcal{E}_\tau(f) - \mathcal{E}_\tau(f_{\tau,\rho})$ and $(1/m) \sum_{i=1}^m g(z_i) = \mathcal{E}_{z,\tau}(f) - \mathcal{E}_{z,\tau}(f_{\tau,\rho})$.

The Lipschitz property of the pinball loss L_τ implies that

$$|g(z)| \leq |f(x) - f_{\tau,\rho}(x)| \leq 2. \quad (33)$$

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

$$\begin{aligned} & |g_1(z) - g_2(z)| \\ & = |L_\tau(y - f_1(x)) - L_\tau(y - f_2(x))| \\ & \leq |f_1(x) - f_2(x)|, \end{aligned} \quad (34)$$

where $f_1, f_2 \in \mathcal{H}$. It follows that

$$\mathcal{N}_{2,\mathbf{z}}(\mathcal{F}, \varepsilon) \leq \mathcal{N}_{2,\mathbf{x}}(\mathcal{H}, \varepsilon). \quad (35)$$

Hence

$$\log \mathcal{N}_2(\mathcal{F}, \varepsilon) \leq a\varepsilon^{-\iota}. \quad (36)$$

Applying Lemma 12 with $B = 2, \alpha = \theta$, and $c = C_\theta$, we know that for any $0 < \delta < 1$, with confidence $1 - \delta/2$, there holds

$$\begin{aligned} & \mathbb{E}f - \frac{1}{m} \sum_{i=1}^m g(z_i) \\ & \leq \frac{1}{2} \eta^{1-\theta} (\mathbb{E}g)^\theta + c'_t \eta + 2 \left(\frac{C_\theta \log(2/\delta)}{m} \right)^{1/(2-\theta)} \\ & \quad + \frac{36 \log(2/\delta)}{m} \\ & \leq \frac{1}{2} \mathbb{E}g + \left(\frac{1}{2} + c'_t \right) \eta + 2 \left(\frac{C_\theta \log(2/\delta)}{m} \right)^{1/(2-\theta)} \\ & \quad + \frac{36 \log(2/\delta)}{m}. \end{aligned} \quad (37)$$

Here

$$\eta \leq \max \left\{ C_\theta^{(2-i)/(4-2\theta+i\theta)} a^{2/(4-2\theta+i\theta)}, 2^{(2-i)/(2+i)} a^{2/(2+i)} \right\} \times \left(\frac{1}{m} \right)^{2/(4-2\theta+i\theta)}. \quad (38)$$

Note that

$$\left(\frac{1}{2} + c'_i \right) \eta + 2 \left(\frac{C_\theta \log(2/\delta)}{m} \right)^{1/(2-\theta)} + \frac{36 \log(2/\delta)}{m} \leq C_{i,\theta} \log \frac{2}{\delta} \left(\frac{1}{m} \right)^{2/(4-2\theta+i\theta)}, \quad (39)$$

where $C_{i,\theta}$ is indicated in (32). Then our desired bound holds true. \square

Proposition 14. Assume ρ on $X \times [-1, 1]$ satisfies the variance bound (20) with index θ indicated in (19). If \mathcal{H} satisfies (9) with $\iota \in (0, 2)$ Then for any $0 < \delta < 1$, with confidence $1 - \delta$, there holds

$$\begin{aligned} \mathcal{E}_\tau(f_z) - \mathcal{E}_\tau(f_{\tau,\rho}) &\leq 2 \left(\mathcal{E}_\tau(f_{\mathcal{H}}) - \mathcal{E}_\tau(f_{\tau,\rho}) \right) \\ &\quad + \frac{8 \log(2/\delta)}{3m} + 2 \left(\frac{2C_\theta \log(2/\delta)}{m} \right)^{1/(2-\theta)} \\ &\quad + 2C_{i,\theta} \log \frac{2}{\delta} \left(\frac{1}{m} \right)^{2/(4-2\theta+i\theta)}. \end{aligned} \quad (40)$$

The above bound follows directly from Propositions 11 and 13 with the fact that $f_z \in \mathcal{H}$.

3.3. Bounding the Total Error. Now we are in a position to present our general result on error analysis for algorithm (3).

Theorem 15. Assume that ρ satisfies (5) with some $p \in (0, \infty)$ and $q \in [1, \infty)$. Denote $\gamma = pq/(p+1)$. Further assume that \mathcal{H} satisfies (9) with $\iota \in (0, 2)$ and $f_{\tau,\rho}$ is uniquely defined. Then for any $0 < \delta < 1$, with confidence $1 - \delta$, one has

$$\begin{aligned} \|f_z - f_{\tau,\rho}\|_{L^r_{\rho_X}} &\leq \bar{C} \inf_{f \in \mathcal{H}} \left\{ \mathcal{E}_\tau(f) - \mathcal{E}_\tau(f_{\tau,\rho}) \right\}^{1/q} \\ &\quad + \bar{C} m^{-\vartheta} \log \frac{2}{\delta}, \end{aligned} \quad (41)$$

where

$$\vartheta = \frac{2(p+1)}{4q(p+1) - (2-i) \min\{2(p+1), pq\}} \quad (42)$$

and \bar{C} are constant independent of m and δ .

Proof. Combining (18), (19), and (40), with confidence $1 - \delta$, we have

$$\begin{aligned} \|f_z - f_{\tau,\rho}\|_{L^r_{\rho_X}} &\leq \bar{C} \left(\mathcal{E}_\tau(f_{\mathcal{H}}) - \mathcal{E}_\tau(f_{\tau,\rho}) \right)^{1/q} + \bar{C} \log \frac{2}{\delta} \left(\frac{1}{m} \right)^{2/q(4-2\theta+i\theta)} \\ &\leq \bar{C} \inf_{f \in \mathcal{H}} \left\{ \mathcal{E}_\tau(f) - \mathcal{E}_\tau(f_{\tau,\rho}) \right\}^{1/q} + \bar{C} m^{-\vartheta} \log \frac{2}{\delta}, \end{aligned} \quad (43)$$

where

$$\begin{aligned} \vartheta &= \frac{2(p+1)}{4q(p+1) - (2-i) \min\{2(p+1), pq\}}, \\ \bar{C} &= 2q^{1/q} \|\gamma^{-1}\|_{L^p_{\rho_X}}^{1/q} \left(\frac{4}{3} + (2C_\theta)^{1/(2-\theta)} + C_{i,\theta} \right)^{1/q}. \end{aligned} \quad (44)$$

\square

Proof of Theorem 1. The assumption $f_{\tau,\rho} \in \mathcal{H}$ implies that

$$\inf_{f \in \mathcal{H}} \left\{ \mathcal{E}_\tau(f) - \mathcal{E}_\tau(f_{\tau,\rho}) \right\} = 0. \quad (45)$$

Therefore, our desired result comes directly from Theorem 15.

4. Further Discussions

In this paper, we studied ERM algorithm (3) for quantile regression and provide convergence and learning rates. We showed some essential differences between ERM scheme and kernel based regularized scheme for quantile regression. We also point out the difficulty to deal with quantile regression: the lack of strong convexity of the pinball loss. To overcome this difficulty, some noise condition on ρ is proposed to enable us to get a variance-expectation bound similar to the one for the least square regression.

In our analysis we just consider $f \in \mathcal{H}$ and $\|f\|_{C(X)} \leq 1$. The case $\|f\|_{C(X)} \leq R$ for $R \geq 1$ would be interesting in the future work. The approximation error involving R can be estimated by the knowledge of interpolation space.

In our setting, the sample is drawn independently from the distribution ρ . However, in many practical problems, the i.i.d condition is a little demanding, so it would be interesting to investigate the ERM scheme for quantile regression with nonidentical distributions [13, 14] or dependent sampling [15].

Acknowledgment

This work described in this paper is supported by NSF of China under Grant 11001247 and 61170109.

References

- [1] I. Steinwart and A. Christman, *How SVMs Can Estimate Quantile and the Median*, vol. 20 of *Advances in Neural Information Processing Systems*, MIT Press, Cambridge, Mass, USA, 2008.

- [2] D. H. Xiang, "Conditional quantiles with varying Gaussians," *Advances in Computational Mathematics*, 2011.
- [3] D.-H. Xiang, T. Hu, and D.-X. Zhou, "Learning with varying insensitive loss," *Applied Mathematics Letters*, vol. 24, no. 12, pp. 2107–2109, 2011.
- [4] D.-H. Xiang, T. Hu, and D.-X. Zhou, "Approximation analysis of learning algorithms for support vector regression and quantile regression," *Journal of Applied Mathematics*, vol. 2012, Article ID 902139, 17 pages, 2012.
- [5] N. Aronszajn, "Theory of reproducing kernels," *Transactions of the American Mathematical Society*, vol. 68, pp. 337–404, 1950.
- [6] F. Cucker and D.-X. Zhou, *Learning Theory: An Approximation Theory Viewpoint*, vol. 24, Cambridge University Press, Cambridge, UK, 2007.
- [7] D.-X. Zhou, "The covering number in learning theory," *Journal of Complexity*, vol. 18, no. 3, pp. 739–767, 2002.
- [8] S. Mukherjee, Q. Wu, and D.-X. Zhou, "Learning gradients on manifolds," *Bernoulli*, vol. 16, no. 1, pp. 181–207, 2010.
- [9] S. Smale and D.-X. Zhou, "Estimating the approximation error in learning theory," *Analysis and Applications*, vol. 1, no. 1, pp. 17–41, 2003.
- [10] S. Mukherjee and Q. Wu, "Estimation of gradients and coordinate covariation in classification," *Journal of Machine Learning Research*, vol. 7, pp. 2481–2514, 2006.
- [11] Y. Yao, "On complexity issues of online learning algorithms," *Institute of Electrical and Electronics Engineers*, vol. 56, no. 12, pp. 6470–6481, 2010.
- [12] Y. Ying, "Convergence analysis of online algorithms," *Advances in Computational Mathematics*, vol. 27, no. 3, pp. 273–291, 2007.
- [13] T. Hu and D.-X. Zhou, "Online learning with samples drawn from non-identical distributions," *Journal of Machine Learning Research*, vol. 10, pp. 2873–2898, 2009.
- [14] S. Smale and D.-X. Zhou, "Online learning with Markov sampling," *Analysis and Applications*, vol. 7, no. 1, pp. 87–113, 2009.
- [15] Z.-C. Guo and L. Shi, "Classification with non-i.i.d. sampling," *Mathematical and Computer Modelling*, vol. 54, no. 5-6, pp. 1347–1364, 2011.