



Generalization performance of support vector classifiers for density level detection



Hong Chen^a, Yicong Zhou^b, Yi Tang^c, Yuan Yan Tang^b, Zhibin Pan^{a,*}

^a College of Science, Huazhong Agricultural University, Wuhan 430070, China

^b Department of Computer and Information Science, University of Macau, Macau 999078, China

^c School of Mathematics and Computer Science, Yunnan University of Nationalities, Kunming 650031, China

ARTICLE INFO

Article history:

Received 29 October 2012

Received in revised form

11 January 2013

Accepted 23 March 2013

Communicated by D. Tao

Available online 25 April 2013

Keywords:

Learning rate

Density level detection

Rademacher average

Iterative technique

ABSTRACT

This paper investigates the generalization performance of support vector classifiers for density level detection (DLD) when the input term belongs to a separable Hilbert space. The estimate of learning rate for DLD problem is established by Rademacher average and iterative techniques, which is independent of the assumption of covering number used in the previous literature.

© 2013 Elsevier B.V. All rights reserved.

1. Introduction

A classification framework for density level detection (DLD) problem has been proposed in [11] and its error analysis has been well established in [10,6] based on the capacity assumption of covering numbers. The theoretical result is important to better understand the mathematical foundation of classification method for DLD. It is well known that the Rademacher complexity has been used successfully for mathematical analysis of machine learning algorithms, see e.g., [2,3,20]. In this paper, we consider establishing the generalization error analysis of the DLD problem by combining the Rademacher complexity with the iterative technique in [13,19,7].

Let us recall the background of the density level detection problem in Hilbert spaces (see [11,10,6]). Let $(H, \|\cdot\|)$ be a separable Hilbert space (possibly infinite dimensional) and let $X \subset H$ with $\|x\| \leq B$ for all $x \in X$. Let Q be an unknown data-generating distribution on X . One of the most common ways to define anomalies is by saying that anomalies are not concentrated. A reference distribution μ on X is introduced to describe the concentration of Q . Assume that Q has a density h with respect to μ , i.e. $dQ = h d\mu$. Given $\rho > 0$, the set $\{x : h(x) > \rho, x \in X\}$ is called ρ -level set of density h . To define anomalies in terms of the concentration one only has to fix a threshold $\rho > 0$ so that a sample $x \in X$ is considered to be anomalous whenever $h(x) \leq \rho$. The main task of the DLD problem is to find ρ -level set $\{x : h(x) > \rho, x \in X\}$. In this paper, we assume that

$\{x : h(x) = \rho, x \in X\}$ is a μ -zero set and hence it is also a Q -zero set (see e.g., [9,17]).

Let $S = \{x_i\}_{i=1}^k$ be a training set which is drawn independently from Q . Given S , a DLD algorithm learns a function $f_S : X \rightarrow \mathbb{R}$ such that the set $\{x : f_S(x) > 0\}$ is a good estimate of ρ -level set. For a measurable function $f : X \rightarrow \mathbb{R}$ the approximation performance is measured (see [11]) by

$$S_{\mu, h, \rho}(f) := \mu(\{f > 0\} \Delta \{h > \rho\}),$$

where Δ denotes the symmetric difference.

Unfortunately, there is no known method to estimate $S_{\mu, h, \rho}(f)$ from empirical data, and hence empirical comparison in terms of $S_{\mu, h, \rho}(f)$ is difficult. To overcome this difficulty, a novel performance measure has been proposed in [11] by interpreting the DLD problem as a binary classification problem. Let $Y = \{-1, 1\}$. The measure is defined as below.

Definition 1. Let Q and μ be probability measures on X and $s \in (0, 1)$. Then the probability measure $Q_{\ominus_s \mu}(A)$ on $X \times Y$ is defined by

$$Q_{\ominus_s \mu}(A) = s E_{x \sim Q} I_A(x, 1) + (1-s) E_{x \sim \mu} I_A(x, -1)$$

for all measurable subsets $A \subset X \times Y$. Here I_A denotes the indicator function of a set A .

From the definition we know that $P := Q_{\ominus_s \mu}$ can be associated with a binary classification problem where positive samples are drawn from Q and negative samples are drawn from μ .

* Corresponding author. Tel.: +86 15387173655.

E-mail address: zhibinpan2008@gmail.com (Z. Pan).

The misclassification risk for a measurable function $f : X \rightarrow \mathbb{R}$ and a distribution P on $Z = X \times Y$ is defined by

$$\mathcal{R}_P(f) = P(\{(x, y) : \text{sign}(f)(x) \neq y\}),$$

where $\text{sign } t = 1$ if $t > 0$ and $\text{sign } t = -1$ otherwise.

It is well known that the Bayes classifier $f_c = \text{sign}(2P(y = 1|\cdot) - 1)$ minimizes the misclassification risk $\mathcal{R}_P(f)$. Moreover, for $P = Q \otimes_s \mu$ and $s = 1/(1 + \rho)$, $f_c = I_{\{h > \rho\}} - I_{\{h \leq \rho\}}$.

As shown in [12], $S_{\mu, h, \rho}(f) \rightarrow 0$ if and only if $\mathcal{R}_P(f) \rightarrow \mathcal{R}_P(f_c)$. Thus, the problem of DLD can be transformed into finding a good function f such that $\mathcal{R}_P(f) \rightarrow \mathcal{R}_P(f_c)$. Based on the interpretation, a kernel-based method is introduced in [11] to realize DLD.

Recall that $K : X \times X \rightarrow \mathbb{R}$ is a Mercer kernel if it is continuous, symmetric, and positive semi-definite. The candidate reproducing kernel Hilbert space (RKHS) \mathcal{H}_K associated with a Mercer kernel K is defined as the closure of the linear span of the set of functions $\{K_x := K(x, \cdot) : x \in X\}$, equipped with the inner product $\langle \cdot, \cdot \rangle_K$ defined by $\langle K_x, K_y \rangle_K = K(x, y)$ (see [1]). The reproducing property is given by $\langle K_x, f \rangle_K = f(x)$ for all $x \in X$ and $f \in \mathcal{H}_K$.

For given positive labeled data $T^+ = \{x_i\}_{i=1}^n$ drawn independently from Q , the empirical quantity

$$\frac{1}{n(1 + \rho)} \sum_{i=1}^n (1 - f(x_i))_+ + \frac{\rho}{1 + \rho} E_{x' \sim \mu} (1 + f(x'))_+$$

is considered in [11]. As pointed out by Steinwart et al. in [11], although the measure μ is known, the expectation $E_{x' \sim \mu} (1 + f(x'))_+$ can be numerically computed through finite evaluation of f on $T^- = \{x'_j\}_{j=1}^m$. Here $T^- = \{x'_j\}_{j=1}^m$ are randomly drawn independently according to μ . The empirical risk of f is defined as

$$\mathcal{E}_T(f) = \frac{1}{n(1 + \rho)} \sum_{i=1}^n (1 - f(x_i))_+ + \frac{\rho}{m(1 + \rho)} \sum_{j=1}^m (1 + f(x'_j))_+.$$

The following regularized algorithm has been proposed in [11]

$$f_T = \arg \min_{f \in \mathcal{H}_K} \{\mathcal{E}_T(f) + \lambda \|f\|_K^2\}, \tag{1}$$

where $\lambda > 0$ is a regularization parameter.

Under the assumption on covering numbers, the convergence of (1) is well understood in [10,6]. In this paper, inspired by theoretical analysis in [3,5,7], we adopt the Rademacher average as the capacity measure of hypothesis space. Dimension-free bound of capacity can be derived in terms of the structural properties of Rademacher average. Without the covering number assumption, satisfactory learning rate is obtained by combining the Rademacher complexity with the iteration technique.

The rest of this paper is organized as follows. In Section 2, we introduce the necessary definitions and present the main result on learning rate. A detailed proof of the main result is provided in Section 3.

2. Error analysis

To establish the relationship between $S_{\mu, h, \rho}(f_T)$ and the excess risk $\mathcal{R}_P(f_T) - \mathcal{R}_P(f_c)$, we recall the following assumption [11,10].

Definition 2. Let μ be a distribution on X and let $h : X \rightarrow [0, \infty]$ be a measurable function with $\int h d\mu = 1$, i.e. h is a density with respect to μ . For $\rho > 0$ and $0 \leq q \leq \infty$, we say h has ρ -exponent q if there exists a constant $c > 0$ such that for all $t > 0$

$$\mu(\{|h - \rho| \leq t\}) \leq ct^q.$$

The assumption on h is closely related to the definition of Tsybakov noise in [18] for binary classification. If h has ρ -exponent $q \in (0, \infty]$, Theorem 10 [11] shows that there exists a constant $c > 0$

such that

$$S_{\mu, h, \rho}(\text{sign}(f)) \leq c(\mathcal{R}_P(f) - \mathcal{R}_P(f_c))^{q/(q+1)}. \tag{2}$$

According to $\mathcal{E}_T(f)$, we introduce the expected risk with a convex loss

$$\mathcal{E}(f) = \frac{1}{1 + \rho} E_{x \sim Q} (1 - f(x))_+ + \frac{\rho}{1 + \rho} E_{x' \sim \mu} (1 + f(x'))_+.$$

We know that for every measurable function $f : X \rightarrow \mathbb{R}$

$$\mathcal{R}_P(f) - \mathcal{R}_P(f_c) \leq \mathcal{E}(f) - \mathcal{E}(f_c) \tag{3}$$

according to Theorem 2.1 in [21] or Theorem 9.21 in [8].

Define the data independent regularization function

$$f_\lambda := \arg \min_{f \in \mathcal{H}_K} \{\mathcal{E}(f) + \lambda \|f\|_K^2\}. \tag{4}$$

From the definitions of f_T in (1) and f_λ in (4), we have

$$\mathcal{E}(f_T) - \mathcal{E}(f_c) \leq \mathcal{E}(f_T) - \mathcal{E}(f_c) + \lambda \|f_T\|_K^2 \leq \mathcal{S}(T, \lambda) + \mathcal{D}(\lambda), \tag{5}$$

where the sample error

$$\mathcal{S}(T, \lambda) = \{\mathcal{E}(f_T) - \mathcal{E}_T(f_T)\} + \{\mathcal{E}_T(f_\lambda) - \mathcal{E}(f_\lambda)\}$$

and the approximation error

$$\mathcal{D}(\lambda) = \mathcal{E}(f_\lambda) - \mathcal{E}(f_c) + \lambda \|f_\lambda\|_K^2.$$

The bounding technique for sample error $\mathcal{S}(T, \lambda)$ relies on complexity measure of hypothesis function space \mathcal{H}_K . To derive a dimension-free estimate, we introduce Rademacher complexity [2] as the measure of capacity.

Definition 3. Let ρ be a probability distribution on a set X and suppose that x_1, \dots, x_m are independent samples selected according to this distribution. Let \mathcal{F} be a class of real-valued functions defined on X . The empirical Rademacher average of \mathcal{F} is defined by

$$\hat{\mathcal{R}}_m(\mathcal{F}) = E_\sigma \left\{ \sup_{f \in \mathcal{F}} \left| \frac{1}{m} \sum_{i=1}^m \sigma_i f(x_i) \right| : x_1, \dots, x_m \right\},$$

where $\sigma_1, \dots, \sigma_m$ are independent uniform $\{\pm 1\}$ -valued random variables. The Rademacher complexity of \mathcal{F} is $\mathcal{R}_m(\mathcal{F}) = E \hat{\mathcal{R}}_m(\mathcal{F})$.

In this paper, we adopt the following condition for approximation error, which has been extensively used in the literature. See e.g., [4,19,8,20,6].

Definition 4. We say the target function f_c can be approximated with exponent $0 < \beta \leq 1$ in \mathcal{H}_K if there exists a constant $c_\beta \geq 1$, such that

$$\mathcal{D}(\lambda) \leq c_\beta \lambda^\beta, \quad \forall \lambda > 0. \tag{6}$$

It is now a position to present our main result on learning rate. The detailed proof will be given in the next section.

Theorem 1. Let $\rho > 0$. Let μ and Q be distributions on X such that Q has a density h with respect to μ . For $s = 1/(\rho + 1)$ we write $P = Q \otimes_s \mu$. Assume that h has ρ -exponent q and f_c can be approximated with exponent β in \mathcal{H}_K . Then, for any $0 < \delta < 1$, choosing $\lambda = (1/\sqrt{m} + 1/\sqrt{n})^{2/(\beta+1)}$, we have with confidence $1 - \delta$

$$S_{\mu, h, \rho}(f_T) \leq C \sqrt{\ln \left(\frac{k_0 + 1}{\delta} \right)} \left(\frac{1}{\sqrt{m}} + \frac{1}{\sqrt{n}} \right)^{q\beta/(\beta+1)(q+1) - q/2k_0(\beta+1)(q+1)},$$

where C is a constant independent of m, n, δ , and k_0 is a constant satisfying $(\sqrt{mn}/(\sqrt{m} + \sqrt{n}))^{2\beta/(\beta+1)2^{k_0+1}} = O(k)$.

From the result in Theorem 1, we know that the balance of samples is crucial to reach the fast learning rate. In particular, learning rate of f_T can be close to $O(n^{-q/(4q+4)})$ when $m = O(n)$ and $\beta \rightarrow 1$. It is worth noting that the presented convergence analysis is independent of the assumption on covering numbers in [6].

Now we give some comparisons on learning rates for any $\epsilon > 0$ and $m = O(n)$. It has been shown in [10] that if the density h has both ρ -exponent q and geometric ρ -exponent $\alpha \in (0, \infty)$, then the learning rates of f_T is $O(n^{-q\alpha/(1+q)(2\alpha+1)+\epsilon})$ for $\alpha < (q+2)/2q$ and $O(n^{-2q\alpha/(2\alpha(2+q)+3q+4)+\epsilon})$ otherwise. We can observe that our estimate of learning rate is faster than $O(n^{-q\alpha/(1+q)(2\alpha+1)+\epsilon})$ when $\alpha < (q+2)/2q$. In fact, by using the iterative technique, we can also improve the previous estimates on generalization error in [6].

Along the line of the present work, further research direction may establish the generalization estimate of the DLD problem with non-i.i.d samples [14,15] and with different analysis techniques [13,16].

3. Proof of Theorem 1

We introduce some properties of Rademacher complexity (see [2]) which are used in the sample error estimation.

Lemma 1. Let $\mathcal{G}, \mathcal{G}_1, \mathcal{G}_2$ be the classes of real functions. Then

- (1) $\mathcal{R}_m(|\mathcal{G}|) \leq \mathcal{R}_m(\mathcal{G})$ where $|\mathcal{G}| = \{f : f \in \mathcal{G}\}$.
- (2) $\mathcal{R}_m(\mathcal{G}_1 \oplus \mathcal{G}_2) \leq \mathcal{R}_m(\mathcal{G}_1) + \mathcal{R}_m(\mathcal{G}_2)$ where $\mathcal{G}_1 \oplus \mathcal{G}_2 = \{g_1 + g_2 : (g_1, g_2) \in \mathcal{G}_1 \times \mathcal{G}_2\}$.
- (3) If $\phi : \mathbb{R} \rightarrow \mathbb{R}$ is Lipschitz with constant L_ϕ and satisfies $\phi(0) = 0$, then $\mathcal{R}_m(|\phi \circ \mathcal{G}|) \leq 2L_\phi \mathcal{R}_m(\mathcal{G})$.

Now we give the estimate of Rademacher complexity for hypothesis function spaces in RKHS. The analysis technique used here is the same as Lemma 2.1 in [5]. We recall the key steps of proof for completeness.

Lemma 2. Define $\mathcal{F}_r = \{f \in \mathcal{H}_K : \|f\|_K \leq r\}$ and $\kappa = \sup_{x \in X} \sqrt{K(x, x)}$. Then, $\mathcal{R}_n(\mathcal{F}_r) \leq r\kappa / \sqrt{n}$.

Proof. Based on the reproducing property of $f \in \mathcal{F}_r$, we have

$$\begin{aligned} \mathcal{R}_n(\mathcal{F}_r) &= \mathbb{E} \mathbb{E}_\sigma \left(\sup_{f \in \mathcal{F}_r} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i \langle f, K_{x_i} \rangle \right| : x_1, \dots, x_n \right) \\ &\leq r \mathbb{E} \mathbb{E}_\sigma \left(\left\| \frac{1}{n} \sum_{i=1}^n \sigma_i K_{x_i} \right\| : x_1, \dots, x_n \right) \\ &= \frac{r}{n} \mathbb{E} \mathbb{E}_\sigma \left[\left(\sum_{i,j=1}^n \sigma_i \sigma_j K(x_i, x_j) \right)^{1/2} : x_1, \dots, x_n \right] \\ &\leq \frac{r}{n} \mathbb{E} \left(\sum_{i=1}^n K(x_i, x_i) \right)^{1/2} \\ &\leq \frac{r\kappa}{\sqrt{n}}. \quad \square \end{aligned}$$

Note that for all $f \in \mathcal{F}_r$,

$$\begin{aligned} |\mathcal{E}(f) - \mathcal{E}_T(f)| &\leq \frac{1}{1+\rho} |E_{X \sim Q}(1-f(x))_+ - \mathcal{E}_{T^+}(f)| \\ &\quad + \frac{\rho}{1+\rho} |E_{X' \sim \mu}(1+f(x'))_+ - \mathcal{E}_T(f)|, \end{aligned} \tag{7}$$

where $\mathcal{E}_{T^+}(f) = (1/n) \sum_{i=1}^n (1-f(x_i))_+$ and $\mathcal{E}_T(f) = (1/m) \sum_{j=1}^m (1+f(x'_j))_+$.

Now we turn to consider the two terms on the right side based on Rademacher average technique. The upper bound of sample error is presented as below.

Proposition 1. For any $f \in \mathcal{F}_r$, with probability at least $1-\delta$, there holds

$$\begin{aligned} |\mathcal{E}(f) - \mathcal{E}_T(f)| &\leq \left(2\kappa r + 4 + \kappa r \sqrt{2 \ln \frac{2}{\delta}} \right) \\ &\quad \left(\frac{1}{(1+\rho)\sqrt{m}} + \frac{\rho}{(1+\rho)\sqrt{n}} \right). \end{aligned}$$

Proof. For each $f \in \mathcal{F}_r$, we have $\|f\|_\infty \leq \kappa \|f\|_K \leq \kappa r$. Let \tilde{T}^+ be the same copy T^+ with k th sample replaced by sample \tilde{x}_k . Then

$$\begin{aligned} &\left| \sup_{f \in \mathcal{F}_r} |E_{X \sim Q}(1-f(x))_+ - \mathcal{E}_{T^+}(f)| - \sup_{f \in \mathcal{F}_r} |E_{X \sim Q}(1-f(x))_+ - \mathcal{E}_{\tilde{T}^+}(f)| \right| \\ &\leq \sup_{f \in \mathcal{F}_r} |\mathcal{E}_{T^+}(f) - \mathcal{E}_{\tilde{T}^+}(f)| = \frac{1}{n} \sup_{f \in \mathcal{F}_r} |(1-f(x_k))_+ - (1-f(\tilde{x}_k))_+| \\ &\leq \frac{1}{n} \sup_{f \in \mathcal{F}_r} |f(x_k) - f(\tilde{x}_k)| \leq \frac{2\kappa r}{n}. \end{aligned}$$

McDiarmid's inequality implies that with probability at least $1-\delta/2$

$$\begin{aligned} \sup_{f \in \mathcal{F}_r} |E_{X \sim Q}(1-f(x))_+ - \mathcal{E}_{T^+}(f)| &\leq \mathbb{E} \sup_{f \in \mathcal{F}_r} |E_{X \sim Q}(1-f(x))_+ - \mathcal{E}_{T^+}(f)| \\ &\quad + \kappa r \sqrt{\frac{2 \ln(2/\delta)}{n}}. \end{aligned} \tag{8}$$

Denote $\phi(f(x)) = (1-f(x))_+ - 1$. By the standard symmetrization arguments [2] and Lemma 1,

$$\begin{aligned} &\mathbb{E} \sup_{f \in \mathcal{F}_r} |E_{X \sim Q}(1-f(x))_+ - \mathcal{E}_{T^+}(f)| \\ &= \mathbb{E} \sup_{f \in \mathcal{F}_r} \left| E_{X \sim Q} \phi(f(x)) - \frac{1}{n} \sum_{i=1}^n \phi(f(x_i)) \right| \\ &\leq 2\mathbb{E} \sup_{f \in \mathcal{F}_r} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i \phi(f(x_i)) \right| \\ &\leq 2\mathbb{E} \sup_{f \in \mathcal{F}_r} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i f(x_i) \right| + 4\mathbb{E} \sup_{f \in \mathcal{F}_r} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i \right| \\ &\leq 2\mathcal{R}_n(\mathcal{F}_r) + \frac{4}{\sqrt{n}}. \end{aligned} \tag{9}$$

By combining (8) and (9), and Lemma 2, we have with probability at least $1-\delta/2$

$$\sup_{f \in \mathcal{F}_r} |E_{X \sim Q}(1-f(x))_+ - \mathcal{E}_{T^+}(f)| \leq (4 + 2\kappa r + \kappa r \sqrt{2 \ln(2/\delta)}) \frac{1}{\sqrt{n}}. \tag{10}$$

With the same fashion as above, we also have that

$$\sup_{f \in \mathcal{F}_r} |E_{X' \sim \mu}(1+f(x'))_+ - \mathcal{E}_T(f)| \leq (4 + 2\kappa r + \kappa r \sqrt{2 \ln(2/\delta)}) \frac{1}{\sqrt{m}}$$

holds with probability at least $1-\delta/2$. The desired result follows by combining (7) with (10). \square

For $R > 1$, we denote

$$\mathcal{W}(R) := \{T = (T^+, T^-) : T^+ \in X^n, T^- \in X^m, \|f_T\|_K \leq R\}.$$

Also, we denote

$$\Omega_T = \mathcal{E}(f_T) - \mathcal{E}(f_c) + \lambda \|f_T\|_K^2.$$

Lemma 3. Denote $V_R = (V_R^+, V_R^-)$, where $V_R^+ \in X^n, V_R^- \in X^m$. For all $t > 0$, there exists a set V_R with $P(V_R) \leq 4e^{-t}$ such that, for all $\mathcal{W}(R) \setminus V_R$

$$\Omega_T \leq \left(4\kappa \left(R + \sqrt{\frac{D(\lambda)}{\lambda}} \right) + 8 + 2\kappa \sqrt{t} \left(R + \sqrt{\frac{D(\lambda)}{\lambda}} \right) \right)$$

$$\left(\frac{1}{(1+\rho)\sqrt{m}} + \frac{\rho}{(1+\rho)\sqrt{n}}\right) + D(\lambda).$$

Proof. By the definitions of f_λ and $D(\lambda)$, we get $\|f_\lambda\|_K \leq \sqrt{D(\lambda)/\lambda}$. From Proposition 1, we have with confidence $1-\delta$

$$|\mathcal{E}(f_\lambda) - \mathcal{E}_T(f_\lambda)| \leq \left(2\kappa\sqrt{\frac{D(\lambda)}{\lambda}} + 4 + 2\kappa\sqrt{\frac{2D(\lambda)\ln\left(\frac{2}{\delta}\right)}{\lambda}}\right) \left(\frac{1}{(1+\rho)\sqrt{m}} + \frac{\rho}{(1+\rho)\sqrt{n}}\right).$$

For $T \in W(R)$, with confidence $1-\delta$

$$|\mathcal{E}(f_T) - \mathcal{E}_T(f_T)| \leq \left(2\kappa R + 4 + 2\kappa R\sqrt{2\ln\left(\frac{2}{\delta}\right)}\right) \left(\frac{1}{(1+\rho)\sqrt{m}} + \frac{\rho}{(1+\rho)\sqrt{n}}\right).$$

From (5), we have with confidence $1-\delta$

$$\Omega_T \leq \left(2\kappa\left(R + \sqrt{\frac{D(\lambda)}{\lambda}}\right) + 8 + 2\kappa\sqrt{t}\left(R + \sqrt{\frac{2D(\lambda)\ln\left(\frac{4}{\delta}\right)}{\lambda}}\right)\right) \left(\frac{1}{(1+\rho)\sqrt{m}} + \frac{\rho}{(1+\rho)\sqrt{n}}\right) + D(\lambda).$$

Setting $t = \ln(4/\delta)$, we have $\delta = 4e^{-t}$. Then, there exists a set V_R with $P(V_R) \leq 4e^{-t}$ such that, for all $\mathcal{W}(R) \setminus V_R$, the desired inequality in Lemma 3 holds true. \square

From the condition of Lemma 3, we also need an R such that $\mathcal{W}(R)$ contains all $T^+ \in X^n, T^- \in X^m$. By the definition of f_T , we get $\|f_T\|_K \leq 1/\sqrt{\lambda}$. Hence, $T \in \mathcal{W}(1/\sqrt{\lambda})$ for each T . In order to improve the estimate of learning rate, we shall consider an iteration technique used in [13, 19, 7]. We can also prove that $\|f\|_K$ can be bounded by $\sqrt{D(\lambda)/\lambda}$ with high confidence.

Now, we are in a position to prove the main result in Theorem 1.

Proof of Theorem 1. Let $t \geq 1$. Based on Lemma 3 and (6), we know that there exists a set V_R , with a probability at most e^{-t} , such that for every $T \in \mathcal{W}(R) \setminus V_R$,

$$\Omega_T \leq \tilde{c}_1 \sqrt{t} \left\{ \left(R + \lambda^{(\beta-1)/2}\right) \left(\frac{1}{\sqrt{m}} + \frac{1}{\sqrt{n}}\right) + \lambda^\beta \right\}, \quad (11)$$

where \tilde{c}_1 is a constant independent of m, n , and t . Choose $\lambda = (1/\sqrt{m} + 1/\sqrt{n})^{2/(\beta+1)}$. Then, we can easily check that $\lambda^{(\beta-1)/2}(1/\sqrt{m} + 1/\sqrt{n}) = \lambda^\beta$. Thus, (11) implies that

$$\Omega_T \leq \tilde{c}_1 \sqrt{t} R \left(\frac{1}{\sqrt{m}} + \frac{1}{\sqrt{n}}\right) + 2\tilde{c}_1 \sqrt{t} \lambda^\beta \leq \tilde{c}_1 \sqrt{t} R \lambda^{(\beta+1)/2} + 2\tilde{c}_1 \sqrt{t} \lambda^\beta. \quad (12)$$

Since $\|f_T\|_K \leq \sqrt{\Omega_T/\lambda}$, by using (12) iteratively, we can find a small ball \mathcal{F}_R that contains $f_{z,\lambda}$ with high confidence. Starting with $R = R^{(0)} = 1/\sqrt{\lambda}$, by (12) we know that each $T \in \mathcal{W}(R^{(0)}) \subseteq \mathcal{W}(R^{(1)}) \cup V_{R^{(0)}}$, where

$$R^{(1)} = \tilde{c}_1 \sqrt{t} \lambda^{(\beta-2)/4} + 2\tilde{c}_1 \sqrt{t} \lambda^{(\beta-1)/2}.$$

Based on (12), we iteratively derive

$$\mathcal{W}(R^{(0)}) \subseteq \mathcal{W}(R^{(1)}) \cup V_{R^{(0)}} \subseteq \dots \subseteq \mathcal{W}(R^{(k)}) \cup \left(\bigcup_{j=0}^{k-1} V_{R^{(j)}}\right),$$

where each $V_{R^{(j)}}$ has a probability at most e^{-t} and $R^{(k)}$ is given by

$$R^{(k)} = \tilde{c}_1 \sqrt{t} \lambda^{(\beta-1)/2 - (\beta/2)^{k+1}} + 2k\tilde{c}_1 \sqrt{t} \lambda^{(\beta-1)/2} \leq 2\tilde{c}_1 \sqrt{t} \lambda^{(\beta-1)/2} (\lambda^{-\beta/2^{k+1}} + k).$$

Choosing a constant k_0 such that $\lambda^{-\beta/2^{k+1}} = O(k)$, we get for $T \in \mathcal{W}(R^{(k_0)})$

$$\|f_T\|_K \leq 4\tilde{c}_1 \sqrt{t} \lambda^{((\beta-1)/2) - (\beta/2)^{k+1}}.$$

Together with (11) and taking $t = \ln(k_0 + 1/\delta)$, we have with a probability at most $1-\delta$

$$\mathcal{E}(f_T) - \mathcal{E}(f_c) \leq C \sqrt{\ln\left(\frac{k_0 + 1}{\delta}\right)} \left(\frac{1}{\sqrt{m}} + \frac{1}{\sqrt{n}}\right)^{(\beta/(\beta+1)) - (1/2^{k_0}(\beta+1))}.$$

Connecting this inequality with (2) and (3) derives the desired result. \square

Acknowledgments

The authors would like to thank the reviewers for their valued comments and suggestions. This work was supported in part by the National Natural Science Foundation of China under Grant Nos. 11001092, 11226304, 61105051, by the Fundamental Research Funds for the Central Universities (Program No. 2011PY130), and by the Macau Science and Technology Development Fund (FDCT) under Grant 017/2012/A1 and the Research Committee at University of Macau under grants MYRG113(Y1-L3)-FST12-ZYC, SRG010-FST11-TYY, MYRG187(Y1-L3)-FST11-TYY, MYRG205(Y1-L4)-FST11-TYY and MRG001/ZYC/2013/FST.

References

- [1] N. Aronszajn, Theory of reproducing kernels, Trans. Am. Math. Soc. 68 (1950) 337–404.
- [2] P.L. Bartlett, S. Mendelson, Rademacher and Gaussian complexities: risk bounds and structural results, J. Mach. Learn. Res. 3 (2002) 463–482.
- [3] G. Biau, L. Devroye, G. Lugosi, On the performance of clustering in Hilbert spaces, IEEE Trans. Inf. Theory 54 (2008) 781–790.
- [4] D.R. Chen, Q. Wu, Y. Ying, D.X. Zhou, Support vector machine soft margin classifiers: error analysis, J. Mach. Learn. Res. 5 (2004) 1143–1175.
- [5] D.R. Chen, H. Li, On the performance of regularized regression learning in Hilbert space, Neurocomputing 93 (2012) 41–47.
- [6] E.L. Cao, X. Xing, J.W. Zhao, Learning rates of support vector machine classifier for density level detection, Neurocomputing 82 (2012) 84–90.
- [7] H. Chen, L.Q. Li, Learning rates of multi-kernel regularized regression, J. Stat. Plann. Inference 140 (2010) 2562–2568.
- [8] F. Cucker, D.X. Zhou, Learning Theory: An Approximation Theory Viewpoint, Cambridge University Press, Cambridge, 2007.
- [9] W. Polonik, Measuring mass concentrations and estimating density contour clusters—an excess mass approach, Ann. Stat. 23 (1995) 855–881.
- [10] C. Scovel, D. Hush, I. Steinwart, Learning rates for density level detection, Anal. Appl. 3 (2005) 356–371.
- [11] I. Steinwart, D. Hush, C. Scovel, A classification framework for anomaly detection, J. Mach. Learn. Res. 6 (2005) 211–232.
- [12] I. Steinwart, C. Scovel, Fast rates for support vector machines using Gaussian kernels, Ann. Stat. 35 (2007) 575–607.
- [13] D. Tao, X. Li, X. Wu, S.J. Maybank, Geometric mean for subspace selection, IEEE Trans. Pattern Anal. Mach. Intell. 31 (2) (2009) 260–274.
- [14] C. Zhang, D. Tao, Generalization bound for infinitely divisible empirical process, J. Mach. Learn. Res. Proc. Track 15 (2011) 864–872.
- [15] C. Zhang, D. Tao, Generalization bounds of ERM-based learning processes for continuous-time Markov chains, IEEE Trans. Neural Network Learn. Syst. 23 (2012) 1872–1883.
- [16] C. Zhang, W. Bian, D. Tao, W. Lin, Discretized-Vapnik-Chervonenkis dimension for analyzing complexity of real function classes, IEEE Trans. Neural Network Learn. Syst. 23 (9) (2012) 1461–1472.
- [17] A.B. Tsybakov, On nonparameter estimation of density level sets, Ann. Stat. 25 (1997) 948–969.
- [18] A.B. Tsybakov, Optimal aggregation of classifiers in statistical learning, Ann. Stat. 32 (2004) 135–166.
- [19] Q. Wu, Y. Ying, D.X. Zhou, Learning rates of least square regularized regression, Found. Comput. Math. 6 (2006) 171–192.
- [20] Y.L. Xu, D.R. Chen, Learning rates of regularized regression for functional data, Int. J. Wavelets Multiresolution Inf. Process. 7 (2009) 839–850.
- [21] T. Zhang, Statistical behavior and consistency of classification methods based on convex risk minimization, Ann. Stat. 32 (2004) 56–85.



Hong Chen received the B.Sc. degree and the Ph.D. degree from Hubei University, Wuhan, China, in 2003 and 2009 respectively. He is an Associate Professor with Department of Mathematics and Informatics Sciences, College of Science, Huazhong Agricultural University, China. His current research interests include learning theory, pattern recognition, and approximation theory.



Yuan Yan Tang received the Ph.D. degree in computer science from Concordia University, Montreal, Canada. He is presently a Professor in the College of Computer at Chongqing University, a Chair Professor in the Department of Computer and Information Science at the University of Macau and a Honorary Professor at the Department of Computer Science at Hong Kong Baptist University. His current interests include wavelet theory and applications, pattern recognition, image processing, document processing, etc. He is the Founder and Editor-in-Chief of International Journal on Wavelets, Multiresolution, and Information Processing (IJWMIP). He is an IAPR Fellow and IEEE Fellow.



Yicong Zhou received his B.S. degree from Hunan University, Changsha, China, and his M.S. and Ph.D. degrees in electrical engineering from Tufts University, Massachusetts, USA, all degrees. Dr. Zhou is currently an Assistant Professor in the Department of Computer and Information Science at University of Macau, Macau, China. His ongoing research interests focus on multimedia security, image/signal processing, medical imaging and object recognition. He is a member of the IEEE and SPIE (International Society for Photo-Optical Instrumentations Engineers).



Zhibin Pan received M.Sc. degree in mathematics from Hubei University in 2004. He is with College of Science, Huazhong Agricultural University, China. He is now pursuing his Ph.D. degree from Huazhong University of Science and Technology, China. His research interests include machine learning and pattern recognition.



Yi Tang received the B.Sc. degree and the Ph.D. degree from Hubei University, Wuhan, China, in 2002 and 2010 respectively. He is with School of Mathematics and Computer Science, Yunnan University of Nationalities, Kunming, China. His current research interests include machine learning, pattern recognition, and computer vision.