



# Classification with non-i.i.d. sampling

Zheng-Chu Guo<sup>a,b</sup>, Lei Shi<sup>b,\*</sup>

<sup>a</sup> School of Mathematics and Computational Science, Sun Yat-sen University, Guangzhou 510275, PR China

<sup>b</sup> Department of Mathematics, City University of Hong Kong, Kowloon, Hong Kong, China

## ARTICLE INFO

### Article history:

Received 22 December 2010

Received in revised form 31 March 2011

Accepted 31 March 2011

### Keywords:

Learning theory

Regularized classification

$\beta$ -mixing sequence

Reproducing kernel Hilbert spaces

$\ell^2$ -empirical covering number

Capacity dependent error bounds

## ABSTRACT

We study learning algorithms for classification generated by regularization schemes in reproducing kernel Hilbert spaces associated with a general convex loss function in a non-i.i.d. process. Error analysis is studied and our main purpose is to provide an elaborate capacity dependent error bounds by applying concentration techniques involving the  $\ell^2$ -empirical covering numbers.

© 2011 Elsevier Ltd. All rights reserved.

## 1. Introduction

In this paper, we consider learning algorithms for classification with non-i.i.d. sampling processes.

In a binary classification problem, the input space is a compact subset  $X \subset \mathbb{R}^d$  and the outputs space  $Y = \{-1, 1\}$  represents two classes. Classification algorithms produce binary classifiers  $\mathcal{C} : X \rightarrow Y$ . Let  $\rho$  be a probability measure defined on  $Z := X \times Y$ . The prediction ability of a classifier  $\mathcal{C}$  is measured by the *misclassification error* which is defined as

$$\mathcal{R}(\mathcal{C}) = \text{Prob}_{(x,y) \in (Z,\rho)} \{\mathcal{C}(x) \neq y\} = \int_X \rho_x(y \neq \mathcal{C}(x)) d\rho_x. \quad (1.1)$$

Here  $\rho_x$  is the marginal distribution of  $\rho$  on  $X$  and  $\rho_x$  is the conditional distribution at  $x \in X$ . The best classifier that minimizes the misclassification error is the *Bayes rule* given by

$$f_c(x) = \begin{cases} 1, & \text{if } \rho_x(y = 1) \geq \rho_x(y = -1), \\ -1, & \text{if } \rho_x(y = 1) < \rho_x(y = -1). \end{cases} \quad (1.2)$$

Since  $\rho_x$  is unknown,  $f_c$  cannot be computed directly. The goal of classification algorithms is to find classifiers which approximate  $f_c$  from a finite sample  $\mathbf{z} = \{z_i = (x_i, y_i)\}_{i=1}^m \in Z^m$ . The classifiers considered here are induced by real-valued functions  $f : X \rightarrow \mathbb{R}$  as  $\mathcal{C} = \text{sgn}(f)$  which is defined by  $\text{sgn}(f)(x) = 1$  if  $f(x) \geq 0$  and  $\text{sgn}(f)(x) = -1$  otherwise. We define a loss function  $\phi : \mathbb{R} \rightarrow \mathbb{R}_+$  and use the error  $\phi(yf(x))$  to measure the difference between the output  $y$  and the prediction  $\text{sgn}(f)(x)$ .

**Definition 1.** A function  $\phi : \mathbb{R} \rightarrow \mathbb{R}_+$  is called a classifying loss function if it is convex, differentiable at 0 with  $\phi'(0) < 0$ , and the smallest zero of  $\phi$  is 1.

\* Corresponding author.

E-mail addresses: [gzhengchu@gmail.com](mailto:gzhengchu@gmail.com) (Z.-C. Guo), [leishi@cityu.edu.hk](mailto:leishi@cityu.edu.hk), [sl1983@mail.ustc.edu.cn](mailto:sl1983@mail.ustc.edu.cn) (L. Shi).