



**Queensland University of Technology**  
Brisbane Australia

This is the author's version of a work that was submitted/accepted for publication in the following source:

Lee, Wee Sun, [Bartlett, Peter L.](#), & Williamson, Robert C. (2008) Correction to "The Importance of Convexity in Learning With Squared Loss" [Sep 98 1974-1980]. *IEEE Transactions on Information Theory*, 54(9), p. 4395.

This file was downloaded from: <http://eprints.qut.edu.au/43982/>

**© Copyright 2008 IEEE**

Personal use of this material is permitted. However, permission to reprint/republish this material for advertising or promotional purposes or for creating new collective works for resale or redistribution to servers or lists, or to reuse any copyrighted component of this work in other works must be obtained from the IEEE.

**Notice:** *Changes introduced as a result of publishing processes such as copy-editing and formatting may not be reflected in this document. For a definitive version of this work, please refer to the published source:*

<http://dx.doi.org/10.1109/TIT.2008.928242>

# Corrections

## Correction to “The Importance of Convexity in Learning With Squared Loss”

Wee Sun Lee, *Senior Member, IEEE*, Peter L. Bartlett, *Member, IEEE*, and Robert C. Williamson, *Member, IEEE*

**Abstract**—The paper “The Importance of Convexity in Learning with Squared Loss” gave a lower bound on the sample complexity of learning with quadratic loss using a nonconvex function class. The proof contains an error. We show that the lower bound is true under a stronger condition that holds for many cases of interest.

**Index Terms**—Agnostic learning, lower bound, sample complexity.

In [2, Theorem 2], it was claimed that if the closure of a function class  $F$  under the metric induced by some probability distribution is not convex, then the sample complexity for agnostically learning  $F$  with squared loss (using only hypotheses in  $F$ ) is  $\Omega(\ln(1/\delta)/\epsilon^2)$  where  $1 - \delta$  is the probability of success and  $\epsilon$  is the required accuracy. The proof of this result—in particular, the proof of Lemma 5—is incorrect. Thus, we only know that this theorem is true for cases where the function class  $F$  is finite dimensional. This weakens the result. However, the lower bound for the sample complexity for agnostic learning still holds for many cases of interest, including any case where the closure of the class of restrictions of functions to a finite subset of the input space  $\mathcal{X}$  is not convex. This is the case, for instance, for all of the examples mentioned in [2], including the set of linear combinations of a fixed number of linear threshold functions. (A counting argument, exploiting the finite pseudodimension of such a class, demonstrates this.)

The following is a corrected version of [2, Lemma 5]; it differs from that lemma by the addition of the words “finite dimensional.”

**Lemma 5’:** Suppose that  $P_{\mathcal{X}}$  is a probability distribution on  $\mathcal{X}$ ,  $H$  is the corresponding Hilbert space, and  $\mathcal{Y}'$  is a bounded interval in  $\mathbb{R}$ . Let  $H_{\mathcal{Y}'}$  denote the set of functions  $f$  in  $H$  with  $f(x) \in \mathcal{Y}'$  for all  $x \in \mathcal{X}$ . Let  $F$  be a totally bounded finite dimensional subset of  $H_{\mathcal{Y}'}$ . If  $\bar{F}$  is not convex, there is a bounded interval  $\mathcal{Y}$  in  $\mathbb{R}$  and functions  $c \in H_{\mathcal{Y}}$ , and  $f_1, f_2 \in \bar{F}$  satisfying  $\|f_1 - f_2\| \neq 0$ ,  $\|c - f_1\| = \|c - f_2\| > 0$ , and for all  $f \in \bar{F}$ ,  $\|c - f\| \geq \|c - f_1\|$ .

While it is true that for any closed, totally bounded nonconvex subset  $\bar{F}$  of  $H_{\mathcal{Y}'}$  there is a function  $c$  with two best approximations (see, e.g.,

Manuscript received April 23, 2008; revised May 15, 2008. Published August 27, 2008 (projected).

W. S. Lee is with the Department of Computer Science, National University of Singapore, Singapore 117590, Republic of Singapore.

P. L. Bartlett is with the Computer Science Division and Department of Statistics, University of California, Berkeley CA 94720-3860 USA.

R. C. Williamson is with National ICT Australia and the Research School of Information Sciences and Engineering, Australian National University, Canberra, ACT 2001, Australia.

Communicated by A. Krzyżak, Associate Editor for Pattern Recognition, Statistical Learning and Inference.

Digital Object Identifier 10.1109/TIT.2008.928242

Theorem 12.6 in [1]), we do not know if there is a uniformly bounded  $c$  with this property. However, if  $\bar{F}$  is finite dimensional, we can project a function  $c$  with two best approximations to the subspace spanned by functions in  $\bar{F}$ . This projection,  $c'$ , would still have two best approximations. As  $c'$  can be represented as a finite linear combination of functions in  $\bar{F}$  and every function in  $\bar{F}$  has bounded range,  $c'$  has bounded range as well, which proves Lemma 5’.

### ACKNOWLEDGMENT

Thanks to Shahar Mendelson for pointing out the error; see [3].

### REFERENCES

- [1] F. Deutsch, *Best Approximation in Inner Product Spaces*. New York: Springer-Verlag, 2001.
- [2] W. S. Lee, P. L. Bartlett, and R. C. Williamson, “The importance of convexity in learning with squared loss,” *IEEE Trans. Inf. Theory*, vol. 44, no. 5, pp. 1974–1980, Sep. 1998.
- [3] S. Mendelson, “Lower bounds for the empirical minimization algorithm,” *IEEE Trans. Inf. Theory*, to be published.

## Correction to “An Algebraic, Analytic, and Algorithmic Investigation on the Capacity and Capacity-Achieving Input Probability Distributions of Finite-Input–Finite-Output Discrete Memoryless Channels”

Xue-Bin Liang, *Member, IEEE*

The following misprints were introduced in the above paper [1]. In equation (13) of line 7 in the first column of page 1006, “ $\leq \frac{1}{2}$ ” should read as “ $< 1$ ”. In addition, in reference [32] on page 1023, the first author’s last name should read as “Shulman” and the page numbers “pp. 19–27” should be “pp. 1356–1362.”

### REFERENCES

- [1] X. B. Liang, “An algebraic, analytic, and algorithmic investigation on the capacity and capacity-achieving input probability distributions of finite-input–finite-output discrete memoryless channels,” *IEEE Trans. Inf. Theory*, vol. 54, no. 3, pp. 1003–1023, Mar. 2008.

Manuscript received February 25, 2008; revised March 17, 2008. Published August 27, 2008 (projected).

The author is with the Department of Electrical and Computer Engineering, Louisiana State University, Baton Rouge, LA 70803 USA (e-mail: xbliang@ee.lsu.edu).

Communicated by U. Mitra, Associate Editor At Large.

Digital Object Identifier 10.1109/TIT.2008.928236