

# Erratum to "Online Newton Step Algorithm for Exponentially Concave Functions"

September 5, 2015

## Abstract

We fix two typos in the statement of Theorem 4, and an error in Theorem 8. To be more clear, we rewrite the proof of the lower bound.

## 1 Statement of Theorem 4

$$|X_i^2| < R \rightarrow |X_i| \leq R$$

$$\sqrt{2R} \sqrt{\log \frac{2t+1}{\delta^2}} \rightarrow R \sqrt{\log \frac{2t+1}{\delta^2}}$$

## 2 Proof of the Lower Bound

We now show that for square loss, which is a special case of exponentially concave functions, the minimax risk is  $O(d/T)$ . As a result, the online Newton step algorithm achieves the almost optimal excess risk bound. The proof of the lower bound is built upon the distance-based Fano inequality (Duchi and Wainwright, 2013).

Let  $\mathcal{P}$  be a family of distributions on a sample space  $\mathcal{X}$ , and let  $\theta : \mathcal{P} \mapsto \Theta$  be a function mapping  $\mathcal{P}$  to some parameter space  $\Theta$ . Given a set of  $n$  samples  $X^n = \{X_1, \dots, X_n\}$  drawn i.i.d. from a distribution  $P \in \mathcal{P}$ , let  $\hat{\theta}(X^n)$  be a measurable function of  $X^n$ , which is an estimate of the unknown quantity  $\theta(P)$ . Then, the minimax risk for the family  $\mathcal{P}$  is given by

$$\mathfrak{M}_n(\theta(\mathcal{P}), \Phi \circ \rho) = \inf_{\hat{\theta}} \sup_{P \in \mathcal{P}} \mathbb{E}_P \left[ \Phi \left( \rho \left( \hat{\theta}(X^n), \theta(P) \right) \right) \right]$$

where  $\rho : \Theta \times \Theta \mapsto \mathbb{R}$  is a (semi)-metric on the parameter space, and  $\Phi : \mathbb{R}_+ \mapsto \mathbb{R}_+$  is a nondecreasing loss function. Our analysis is based on the following result from Duchi and Wainwright (2013).

**Lemma 1** (Corollary 2 of Duchi and Wainwright (2013)). Let  $\mathcal{V} \subseteq \Theta$  be a subset of the parameter space, and let  $\rho_{\mathcal{V}} : \mathcal{V} \times \mathcal{V} \mapsto \mathbb{R}$  be a (semi)-metric on  $\mathcal{V}$ . Suppose that  $\mathcal{V}$  is  $t$ -separated, i.e.,

$$\delta(t) := \sup \{ \rho_{\mathcal{V}}(\theta_{\mathbf{v}}, \theta_{\mathbf{w}}) \geq \delta \text{ or } \rho_{\mathcal{V}}(\theta_{\mathbf{v}}, \theta_{\mathbf{w}}) > t \}.$$

$P \in \mathcal{P}$

$$\mathfrak{M}_n(\theta(\mathcal{P}), \Phi \circ \rho) \geq \Phi \left( \frac{\delta(t)}{2} \right) \left( 1 - \frac{I(X^n; V) + \log 2}{\log |\mathcal{V}| - \log N_t^{\max}} \right), \quad \forall t$$

$$N_t^{\max} = \max_{\mathbf{v}} \text{card}\{\mathbf{v}' \in \mathcal{V} | \rho_{\mathbf{v}}(\mathbf{v}, \mathbf{v}') \leq t\}$$

In our case, we are interested the generalization error bound  $\mathcal{L}(\widehat{\mathbf{w}}) - \mathcal{L}(\mathbf{w}^*)$ . For square loss, the stochastic optimization problem is given by

$$\min_{\mathbf{w} \in \mathcal{W}} \mathcal{L}(\mathbf{w}) = \mathbb{E} \left[ (Y - X \mathbf{w})^2 \right]$$

where  $X$  is sampled from some underlying distribution  $P_X$ , and given  $X = \mathbf{x}$  the response  $Y$  is sampled from an Gaussian distribution  $\mathcal{N}(\mathbf{x} \mathbf{w}^*, 1)$ , where  $\mathbf{w}^* \in \mathbb{R}^d$  is the parameter vector. Furthermore, we assume  $\mathcal{W} = \mathbb{R}^d$ . Then, it is easy to verify that the excess risk of a solution  $\widehat{\mathbf{w}}$  is

$$\mathcal{L}(\widehat{\mathbf{w}}) - \mathcal{L}(\mathbf{w}^*) = \mathbb{E} \left[ (X \widehat{\mathbf{w}} - Y)^2 \right]$$

In addition, we have

$$I(V; (X, Y)^T) = TI(V; (X, Y))$$

and

$$\begin{aligned} I(V; (X, Y)) &= H(X, Y) - H(X, Y|V) \\ &= H(X) + H(Y|X) - H(X|V) - H(Y|X, V) = H(Y|X) - H(Y|X, V) \\ &\leq E \end{aligned}$$

## References

- John C. Duchi and Martin J. Wainwright. Distance-based and continuum fano inequalities with applications to statistical estimation. *ArXiv preprints*, arXiv:1311.2669, 2013.
- Shahar Mendelson, Alain Pajor, and Nicole Tomczak-Jaegermann. Uniform uncertainty principle for bernoulli and subgaussian ensembles. *Constructive Approximation*, 28(3):277–289, 2008.