

# Empirical Risk Minimization for Stochastic Convex Optimization: $O(1/n)$ - and $O(1/n^2)$ -type of Risk Bounds

**Lijun Zhang**

ZHANGLJ@LAMDA.NJU.EDU.CN

*National Key Laboratory for Novel Software Technology, Nanjing University, Nanjing 210023, China*

**Tianbao Yang**

YANGTIANBAO@IOWA.EDU

*Department of Computer Science, the University of Iowa, Iowa City, IA 52242, USA*

**Rong Jin**

JINRONG@ALIBABA.COM

*Alibaba Group, Seattle, USA*

## Abstract

Although there exist plenty of theoretical results on ERM for supervised learning, current theoretical understandings of ERM for a related problem—stochastic convex optimization—CO—are limited. In this work, we strengthen the results of ERM for CO by exploiting smoothness and strong convexity conditions to prove the risk bounds  $\tilde{O}(d/n + \sqrt{F_*/n})$ .

where  $\mathbf{h} : \mathcal{X} \rightarrow \mathbb{R}$  is a hypothesis class  $(\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathbb{R}$  is an instance sampled from a distribution  $\mathbb{D}$  and  $(\cdot, \cdot) : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$  is certain loss. In this paper we mainly focus on the convex version of online stochastic convex optimization (CO) where both the domain  $\mathcal{W}$  and the expected function  $\mathbf{F}(\cdot)$  are convex.

Two classical approaches for solving stochastic optimization are stochastic approximation (A Kushner and Yin [10]) and the sample average approximation (AA) the latter of which is also referred to as empirical minimization (EM) in the machine learning community (e.g., [9]). Both A and EM have been extensively studied in recent years (Bartlett and Mendelson [2], Bartlett et al. [3], Neirovs et al. [11], Moulines and Bach [12], Hazan and Kale [4], Ahn et al. [5], Agarwal et al. [6], Bach and Moulines [7], Zhang et al. [8], Mahdavi et al. [13]). Most theoretical guarantees of EM are restricted to supervised learning. As pointed out in a seminal work of Havalwartz et al. [9] the success of EM for supervised learning cannot be directly extended to stochastic optimization. Actually Havalwartz et al. [9] have constructed an instance of CO that is learnable by A but cannot be solved by EM. Literature about EM for stochastic optimization including CO are quite limited and we still lack a full understanding of the theory.

In EM we are given  $n$  functions  $\mathbf{f}_1, \dots, \mathbf{f}_n$  sampled independently from  $\mathbb{P}$  and analyze an empirical objective function:

$$\min_{\mathbf{w} \in \mathcal{W}} \widehat{\mathbf{F}}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \mathbf{f}_i(\mathbf{w}).$$

Let  $\widehat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w} \in \mathcal{W}} \widehat{\mathbf{F}}(\mathbf{w})$  be an empirical minimizer. The performance of EM is measured in terms of the excess risk defined as

$$\mathbf{F}(\widehat{\mathbf{w}}) - \min_{\mathbf{w} \in \mathcal{W}} \mathbf{F}(\mathbf{w}).$$

State of the art risk bounds of EM include an  $\widetilde{\mathcal{O}}(\sqrt{\mathbf{d}/n})$  bound when the random function  $\mathbf{f}(\cdot)$  is Lipschitz continuous, where  $\mathbf{d}$  is the dimensionality of  $\mathbf{w}$ , an  $\mathcal{O}(1/n)$  bound when  $\mathbf{f}(\cdot)$  is strongly convex (Havalwartz et al. [9]), and an  $\widetilde{\mathcal{O}}(\mathbf{d}/n)$  bound when  $\mathbf{f}(\cdot)$  is exponentially concave (Mehta [14]). From existing studies of EM for supervised learning (e.g., [15]) we know that smoothness can be utilized to boost the risk bound. Thus it is natural to ask whether smoothness can also be exploited to prove the performance of EM for CO. This paper provides an affirmative answer to this question. Indeed we propose a general approach for analyzing the excess risk of EM which brings several improved risk bounds and new risk bounds as well.

To state our results we first introduce some notations. Let  $\mathbf{F}_* = \min_{\mathbf{w} \in \mathcal{W}} \mathbf{F}(\mathbf{w})$  be the minimum risk. Let  $\mathbf{L}$  be the modulus of strong convexity of  $\mathbf{F}(\cdot)$  and  $\mathbf{L}$  be the modulus of smoothness of  $\mathbf{f}(\cdot)$ . Denote by  $\mathbf{L} = \mathbf{L}/\mathbf{L}$  the condition number of the problem. Our and previous results of EM for CO are summarized in Table 1 where we abbreviate the assumptions on the random function  $\mathbf{f}(\cdot)$ , the empirical function  $\widehat{\mathbf{F}}(\mathbf{w})$  and the expected function  $\mathbf{F}(\cdot)$ . For our results of EM for CO we assume the domain is bounded and the random function is nonnegative. The high level significance of this work is as follows:

---

<sup>1</sup> We use the  $\widetilde{\mathcal{O}}$  and  $\widetilde{\Omega}$  notations to hide constant factors as well as polynomial factors in  $\mathbf{d}$  and  $n$ .

able to study of Excess s Bounds of E M for CO\_A bounds hold with high probability except the one derived by \* which holds in expectation. Abbreviations: bounded, b convex, c generalized near, g Lipschitz continuous, Lp nonnegative, nn strongly convex, s smooth, s exponentially concave, exp\_

	$f(\cdot)$	$\widehat{F}(\cdot)$	$F(\cdot)$	Bounds
Hartman et al. [9]	Lp			$\widetilde{O}(\sqrt{\frac{d}{n}})$
	Lp sc			$O(\frac{1}{\lambda n})^*$
Mehta [10]	exp Lp b			$\widetilde{O}(\frac{d}{\eta n})$
Hsiang et al.	nn c s		Lp	$\widetilde{O}(\frac{d}{n} + \sqrt{\frac{F_*}{n}})$
	nn c s		Lp sc	$\widetilde{O}(\frac{d}{n} + \frac{\kappa F_*}{n})$ $O(\frac{1}{\lambda n^2} + \frac{\kappa F_*}{n})$ when $n = \widetilde{\Omega}(d)$
	nn s	c	sc	$\widetilde{O}(\frac{\kappa d}{n} + \frac{\kappa F_*}{n}) = \widetilde{O}(\frac{\kappa d}{n})$ $O(\frac{1}{\lambda n^2} + \frac{\kappa F_*}{n})$ when $n = \widetilde{\Omega}(d^2)$
	nn s g c		sc	$O(\frac{\kappa}{n} + \frac{\kappa F_*}{n}) = O(\frac{\kappa}{n})$ $O(\frac{1}{\lambda n^2} + \frac{\kappa F_*}{n})$ when $n = \Omega(d^2)$

When  $f(\cdot)$  is both convex and smooth and  $F(\cdot)$  is Lipschitz continuous we establish an  $\widetilde{O}(d/n + \sqrt{F_*/n})$  risk bound cf. Heide [11]. In the optimal case that  $F_*$  is a  $\frac{1}{\lambda} F_* = O(d^2/n)$  we obtain an  $\widetilde{O}(d/n)$  risk bound which is analogous to the  $\widetilde{O}(1/n)$  optimal rate of E M for supervised learning [12]. If  $F(\cdot)$  is a so-called strongly convex we prove an  $\widetilde{O}(d/n + F_*/n)$  risk bound and prove that  $O(1/[n^2] + F_*/n)$  when  $n = \widetilde{\Omega}(d)$  cf. Heide [11]. Thus for large  $n$  and  $F_*$  is a  $\frac{1}{\lambda} F_* = O(1/n)$  we get an  $O(1/n^2)$  risk bound which to the best of our knowledge is the first  $O(1/n^2)$  type of risk bound of E M.

When convexity is not present in  $f(\cdot)$  as long as  $f(\cdot)$  is smooth  $\widehat{F}(\cdot)$  is convex and  $F(\cdot)$  is strongly convex we still obtain an improved risk bound of  $O(1/[n^2] + F_*/n)$  when  $n = \widetilde{\Omega}(d^2)$  which we further improve to  $O(1/n^2)$  if  $F_* = O(1/n)$  cf. Heide [11]. Finally we extend the  $O(1/[n^2] + F_*/n)$  risk bound to supervised learning with a generalized near for. Our analysis shows that in this case the lower bound of  $n$  can be replaced with  $\Omega(d^2)$  which is dimensionality independent cf. Heide [11]. Thus this result can be applied to nonlinear cases e.g. learning with neural networks.

## 2. Related Work

In this section we give a brief introduction to previous work on E M.

## 2.1. ERM for Stochastic Optimization

As we mentioned earlier there are few works devoted to ERM for stochastic optimization when  $\mathcal{W} \subseteq \mathbb{R}^d$  is bounded and  $\mathbf{f}(\cdot)$  is Lipschitz continuous. [Hartmann and Recht \(2019\)](#) demonstrate that  $\hat{\mathbf{F}}(\mathbf{w})$  converges to  $\mathbf{F}(\mathbf{w})$  uniformly over  $\mathcal{W}$  with an  $\tilde{\mathcal{O}}(\sqrt{d/n})$  error bound that holds with high probability by proving an  $\tilde{\mathcal{O}}(\sqrt{d/n})$   $\ell_2$  bound of ERM. They further establish an  $\mathcal{O}(1/n)$   $\ell_2$  bound of ERM that holds in expectation when  $\mathbf{f}(\cdot)$  is  $\mu$ -strongly convex and Lipschitz continuous. Stochastic optimization with exp-concave functions is studied recently [Koren and Levy \(2019\)](#) and [Mehta \(2019\)](#) proves an  $\tilde{\mathcal{O}}(d/n)$  bound of ERM that holds with high probability when  $\mathbf{f}(\cdot)$  is exp-concave, Lipschitz continuous and bounded. Lower bounds of ERM for stochastic optimization is investigated by [Fedan \(2019\)](#) who exhibits a lower bound of  $\Omega(d/n^2)$  sample complexity for uniform convergence that nearly matches the upper bound of [Hartmann and Recht \(2019\)](#), and a lower bound of  $\Omega(d/n)$  sample complexity of ERM which is matched by our  $\tilde{\mathcal{O}}(d/n + \sqrt{F_*/n})$  bound when  $\mathbf{F}_*$  is a  $\mu$ -

## 2.2. ERM for Supervised Learning

We note that there are extensive studies on ERM for supervised learning and hence the review here is non-exhaustive. In the context of supervised learning the performance of ERM is closely related to the uniform convergence of  $\hat{\mathbf{F}}(\cdot)$  to  $\mathbf{F}(\cdot)$  over the hypothesis class. [Kochenskiy \(2019\)](#). In fact uniform convergence is a sufficient condition for learnability [Hartmann and Recht \(2019\)](#) and in some special cases such as binary classification is also a necessary condition [Abernethy \(1998\)](#). The accuracy of uniform convergence as well as the quality of the empirical minimizer can be upper bounded in terms of the complexity of the hypothesis class including data-independent measures such as the VC dimension and data-dependent measures such as the Rademacher complexity.

Generally speaking when  $\mathcal{H}$  has finite VC dimension the excess risk can be upper bounded by  $\mathcal{O}(\sqrt{VC(\mathcal{H})/n})$  where  $VC(\cdot)$  is the VC dimension of  $\cdot$ . If the loss  $\ell(\cdot, \cdot)$  is Lipschitz continuous with respect to its first argument we have a  $\ell_2$  bound of  $\mathcal{O}(1/\sqrt{n} + \epsilon_n(\mathcal{H}))$  where  $\epsilon_n(\mathcal{H})$  is the Rademacher complexity of  $\mathcal{H}$ . The Rademacher complexity typically scales as  $\epsilon_n(\mathcal{H}) = \mathcal{O}(1/\sqrt{n})$  e.g. contains linear functions with bounded norm by proving an  $\mathcal{O}(1/\sqrt{n})$   $\ell_2$  bound [Bartlett and Mendelson \(2002\)](#). There have been intensive efforts to derive rates faster than  $\mathcal{O}(1/\sqrt{n})$  under various conditions [Lee et al. \(1999\)](#), [Panchenko \(2008\)](#), [Bartlett et al. \(2005\)](#), [Gonen and Hartmann \(2019\)](#) such as smoothness [Rebro et al. \(2019\)](#) strong convexity [Rdharan et al. \(2019\)](#) to name a few amongst any peculiarity when the random function  $\mathbf{f}(\cdot)$  is nonnegative and smooth [Rebro et al. \(2019\)](#) have established a  $\ell_2$  bound of  $\tilde{\mathcal{O}}(\frac{2}{n}(\mathbf{H}) + \epsilon_n(\mathbf{H}) \sqrt{F_*})$  reducing to an  $\tilde{\mathcal{O}}(1/n)$  bound if  $\epsilon_n(\cdot) = \mathcal{O}(1/\sqrt{n})$  and  $\mathbf{F}_* = \mathcal{O}(1/n)$ . A generalized near form of  $\tilde{\mathcal{O}}(1/n)$  is studied by [Rdharan et al. \(2019\)](#) and a  $\ell_2$  bound of  $\mathcal{O}(1/n)$  is proved if the expected function  $\mathbf{F}(\cdot)$  is  $\mu$ -strongly convex.

## 3. Faster Rates of ERM

We first introduce the assumptions used in our analysis then present theoretical results under different combinations of the and finally discuss a special case of supervised learning.

---

<sup>1</sup> the excess risk bounds for a regularized empirical risk minimizer.

### 3.1. Assumptions

In the following we use  $\|\cdot\|$  to denote the  $2$ -norm of vectors.

**Assumption 1** The domain  $\mathcal{W}$  is a convex subset of  $\mathbb{R}^d$ , and is bounded by  $\mathbf{R}$ , that is,

$$\|\mathbf{w}\| \leq \mathbf{R}, \forall \mathbf{w} \in \mathcal{W}. \quad 4$$

**Assumption 2** The random function  $\mathbf{f}(\cdot)$  is nonnegative, and  $\mathbf{L}$ -smooth over  $\mathcal{W}$ , that is,

$$\|\mathbf{f}(\mathbf{w}) - \mathbf{f}(\mathbf{w}')\| \leq \mathbf{L} \|\mathbf{w} - \mathbf{w}'\|, \forall \mathbf{w}, \mathbf{w}' \in \mathcal{W}, \mathbf{f} \in \mathbb{P}.$$

**Assumption 3** The expected function  $\mathbf{F}(\cdot)$  is  $\mathbf{G}$ -Lipschitz continuous over  $\mathcal{W}$ , that is,

$$\|\mathbf{F}(\mathbf{w}) - \mathbf{F}(\mathbf{w}')\| \leq \mathbf{G} \|\mathbf{w} - \mathbf{w}'\|, \forall \mathbf{w}, \mathbf{w}' \in \mathcal{W}.$$

**Assumption 4** We use different combinations of the following assumptions on convexity.

- (a) The expected function  $\mathbf{F}(\cdot)$  is convex over  $\mathcal{W}$ .
- (b) The expected function  $\mathbf{F}(\cdot)$  is  $\mu$ -strongly convex over  $\mathcal{W}$ , that is,

$$\mathbf{F}(\mathbf{w}) + \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}'\|^2 \leq \mathbf{F}(\mathbf{w}'), \forall \mathbf{w}, \mathbf{w}' \in \mathcal{W}. \quad \blacktriangledown$$

- (c) The empirical function  $\widehat{\mathbf{F}}(\cdot)$  is convex.
- (d) The random function  $\mathbf{f}(\cdot)$  is convex.

**Assumption 5** Let  $\mathbf{w}_* = \operatorname{argmin}_{\mathbf{w} \in \mathcal{W}} \mathbf{F}(\mathbf{w})$  be an optimal solution to (1). We assume the gradient of the random function at  $\mathbf{w}_*$  is upper bounded by  $\mathbf{M}$ , that is,

$$\|\nabla \mathbf{f}(\mathbf{w}_*)\| \leq \mathbf{M}, \forall \mathbf{f} \in \mathbb{P}.$$

**Remark 1** First note that **Assumption 4(a)** is implied by either **Assumption 4(b)** or **Assumption 4(d)** and **Assumption 4(c)** is implied by **Assumption 4(d)**. Second the smoothness assumption of  $\mathbf{f}(\cdot)$  implies the expected function  $\mathbf{F}(\cdot)$  is  $\mathbf{L}$ -smooth. By Jensen's inequality we have

$$\|\mathbf{F}(\mathbf{w}) - \mathbf{F}(\mathbf{w}')\| \leq \mathbb{E}_{\mathbf{f} \sim \mathbb{P}} \|\mathbf{f}(\mathbf{w}) - \mathbf{f}(\mathbf{w}')\| \leq \mathbf{L} \|\mathbf{w} - \mathbf{w}'\|, \forall \mathbf{w}, \mathbf{w}' \in \mathcal{W}.$$

Thirdly the empirical function  $\widehat{\mathbf{F}}(\cdot)$  is also  $\mathbf{L}$ -smooth. The condition number of  $\mathbf{F}(\cdot)$  is defined as the ratio between  $\mathbf{L}$  and  $\mu = \mathbf{L}/\kappa - 1$ .

### 3.2. Risk Bounds for SCO

We first present an excess risk bound under the smoothness condition.

**Theorem 1** For any  $0 < \epsilon < 1/2$ ,  $\delta > 0$ , define

$$\mathbf{C}(\epsilon, \delta) = 2 \left( \log \frac{2}{\epsilon} + d \log \frac{6\mathbf{R}}{\delta} \right). \quad 9$$

Under Assumptions 1, 2, 3, 4(d), and 5, with probability at least  $1 - 2^{-\beta}$ , we have

$$\mathbf{F}(\hat{\mathbf{w}}) - \mathbf{F}(\mathbf{w}_*) \leq \frac{16\mathbf{R}^2\mathbf{LC}(\beta, \gamma)}{n} + \frac{8\mathbf{R}\mathbf{M} \log(2/\delta)}{n} + 8\mathbf{R}\sqrt{\frac{2\mathbf{L}\mathbf{F}_* \log(2/\delta)}{n}} + \left(8\mathbf{R}\mathbf{L} + \mathbf{G} + \frac{4\mathbf{R}\mathbf{L}\mathbf{C}(\beta, \gamma)}{n}\right),$$

where  $\mathbf{F}_* = \mathbf{F}(\mathbf{w}_*)$  is the minimal risk.

By choosing  $\delta$  small enough the last term in (10) that contains  $\sqrt{\frac{2\mathbf{L}\mathbf{F}_* \log(2/\delta)}{n}}$  becomes non dominant and by specifying  $\delta$  we have the following corollary.

**Corollary 2** By setting  $\delta = 1/n$  in Theorem 1, we have  $\mathbf{C}(1/n, \beta) = 2(\log(2/\delta) + \mathbf{d} \log(6\mathbf{n}\mathbf{R}))$ , and with high probability

$$\mathbf{F}(\hat{\mathbf{w}}) - \mathbf{F}(\mathbf{w}_*) = \mathcal{O}\left(\frac{\mathbf{d} \log n}{n} + \sqrt{\frac{\mathbf{F}_*}{n}}\right) = \tilde{\mathcal{O}}\left(\frac{\mathbf{d}}{n} + \sqrt{\frac{\mathbf{F}_*}{n}}\right).$$

**Remark 2** The above corollary implies that under the smoothness and other common assumptions  $\tilde{\mathcal{O}}(\mathbf{d}/n + \sqrt{\mathbf{F}_*/n})$  is a tight bound for  $\mathbf{C}(\delta, \beta)$  when the noise is Gaussian. If  $\mathbf{F}_* = \mathcal{O}(\mathbf{d}^2/n)$  the rate is improved to  $\tilde{\mathcal{O}}(\mathbf{d}/n)$ . Note that even under the smoothness assumption the near dependence on  $\mathbf{d}$  is unavoidable. See Fed and Theorem 1.

We next present excess risk bounds under both the smoothness and strong convexity conditions.

**Theorem 3** Under Assumptions 1, 2, 3, 4(b), 4(d), and 5, with probability at least  $1 - 2^{-\beta}$ , we have

$$\mathbf{F}(\hat{\mathbf{w}}) - \mathbf{F}(\mathbf{w}_*) \leq \frac{16\mathbf{R}^2\mathbf{LC}(\beta, \gamma)}{n} + \frac{8\mathbf{R}\mathbf{M} \log(2/\delta)}{n} + \frac{8\mathbf{L}\mathbf{F}_* \log(2/\delta)}{n} + \left(8\mathbf{R}\mathbf{L} + \mathbf{G} + \frac{4\mathbf{R}\mathbf{L}\mathbf{C}(\beta, \gamma)}{n}\right).$$

Furthermore, if

$$n \geq \frac{4\mathbf{L}\mathbf{C}(\beta, \gamma)}{\epsilon},$$

we also have

$$\mathbf{F}(\hat{\mathbf{w}}) - \mathbf{F}(\mathbf{w}_*) \leq \frac{32\mathbf{M}^2 \log^2(2/\delta)}{n^2} + \frac{128\mathbf{L}\mathbf{F}_* \log(2/\delta)}{n} + \left(\frac{128\mathbf{L}^2}{n^2} + 16\mathbf{G} + 4\epsilon\right).$$

The above theorem can be simplified by choosing different values of  $\delta$ .

**Corollary 4** By setting  $\delta = 1/n$  in Theorem 3, we have  $\mathbf{C}(1/n, \beta) = 2(\log(2/\delta) + \mathbf{d} \log(6\mathbf{n}\mathbf{R}))$ , and with high probability

$$\mathbf{F}(\hat{\mathbf{w}}) - \mathbf{F}(\mathbf{w}_*) = \mathcal{O}\left(\frac{\mathbf{d} \log n}{n} + \frac{\mathbf{F}_*}{n}\right) = \tilde{\mathcal{O}}\left(\frac{\mathbf{d}}{n} + \frac{\mathbf{F}_*}{n}\right).$$

By setting  $\delta = 1/n^2$ , we have  $\mathbf{C}(1/n^2, \beta) = 2(\log(2/\delta) + \mathbf{d} \log(6\mathbf{n}^2\mathbf{R}))$  and when  $n = \Omega(\mathbf{d} \log n) = \tilde{\Omega}(\mathbf{d})$ , with high probability

$$\mathbf{F}(\hat{\mathbf{w}}) - \mathbf{F}(\mathbf{w}_*) = \mathcal{O}\left(\frac{1}{n^2} + \frac{\mathbf{F}_*}{n}\right).$$

**Remark 3** The first part of Corollary 4 shows that EM enjoys an  $\tilde{O}(d/n + F_*/n)$  risk bound for stochastic optimization of strongly convex and smooth functions. In the literature the best known parabola result is the  $O(1/n)$  risk bound proved by [Hartmann et al., 2019](#) but with strong differences highlighted in [Hartmann et al., 2019](#) since the risk bound of [Hartmann et al., 2019](#) is independent of the dimensionality.

**Remark 6** Comparing the second part of Corollaries 3 and 4 we can see that the risk bounds on the same order but the lower bound of  $n$  is increased by a factor of  $\frac{1}{\epsilon}$ . It is interesting to note that a similar phenomenon also happens in stochastic approximation recently a variance reduction technique called G Johnson and Zhang [1] or EMGD Zhang et al. [2] was proposed for stochastic optimization when both full gradients and stochastic gradients are available. In the analysis G assumes the stochastic function is convex while EMGD does not. From the theoretical results we observe that the nondifferentiable convexity leads to a difference of a factor in the sample complexity of stochastic gradients.

### 3.3. Risk Bounds for Supervised Learning

If the conditions of Theorem 3 or Theorem 4 are satisfied we can directly use them to establish an  $O(1/\sqrt{n}) + F_*/n$  risk bound for supervised learning. However a drawback of these theorems is that the lower bound of  $n$  depends on the dimensionality  $d$  and thus cannot be applied to nonlinear cases e.g. neural networks [3] and [4]. In this section we exploit the structure of supervised learning to derive the theory dimensionality independent. We focus on the generalized linear form of supervised learning:

$$\min_{\mathbf{w} \in \mathcal{W}} \mathbf{F}(\mathbf{w}) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{D}} [\ell(\mathbf{w}, \mathbf{x}, y)] + \mathbf{r}(\mathbf{w}),$$

where  $\ell(\mathbf{w}, \mathbf{x}, y)$  is the loss of predicting  $y$  given  $\mathbf{w}, \mathbf{x}$  when the true target is  $y$  and  $\mathbf{r}(\cdot)$  is a regularizer. Given  $n$  training examples  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$  independent and sampled from  $\mathbb{D}$  the empirical objective is

$$\min_{\mathbf{w} \in \mathcal{W}} \widehat{\mathbf{F}}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}, \mathbf{x}_i, y_i) + \mathbf{r}(\mathbf{w}).$$

define

$$\mathbf{H}(\mathbf{w}) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{D}} [\ell(\mathbf{w}, \mathbf{x}, y)] \text{ and } \widehat{\mathbf{H}}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}, \mathbf{x}_i, y_i)$$

to capture the stochastic component.

Besides 4(b) and 4(c) we introduce the following additional assumptions. We abuse the same notation  $\|\cdot\|$  to denote the norm induced by the inner product of a Hilbert space.

**Assumption 6** The domain  $\mathcal{H}$  is a convex subset of a Hilbert space  $\mathcal{H}$ , and is bounded by



**Assumption 10** Let  $\mathbf{w}_* = \operatorname{argmin}_{\mathbf{w} \in \mathcal{W}} \mathbf{F}(\mathbf{w})$  be an optimal solution to (17). We assume the gradient of the random function at  $\mathbf{w}_*$  is upper bounded by  $\mathbf{M}$ , that is,

$$\|\nabla_{\mathbf{w}} \ell(\mathbf{w}_*, \mathbf{x}, \mathbf{y})\| \leq \mathbf{M}, \forall (\mathbf{x}, \mathbf{y}) \in \mathcal{D}.$$

**Remark 7** The above assumption allows us to model any popular losses including squared loss such as regularized square loss and regularized logistic loss. **Assumptions 7 and 8** imply the random function  $\ell(\cdot, \mathbf{x}, \mathbf{y})$  is  $\mathbf{D}^2$ -smooth over  $\mathcal{W}$ . To see this for any  $\mathbf{w}, \mathbf{w}' \in \mathcal{W}$  we have

$$\begin{aligned} \|\nabla_{\mathbf{w}} \ell(\mathbf{w}, \mathbf{x}, \mathbf{y}) - \nabla_{\mathbf{w}'} \ell(\mathbf{w}', \mathbf{x}, \mathbf{y})\| &= \|\nabla_{\mathbf{w}} \ell(\mathbf{w}, \mathbf{x}, \mathbf{y}) - \nabla_{\mathbf{w}'} \ell(\mathbf{w}', \mathbf{x}, \mathbf{y})\| \\ &\leq \mathbf{D} \|\nabla_{\mathbf{w}} \ell(\mathbf{w}, \mathbf{x}, \mathbf{y}) - \nabla_{\mathbf{w}'} \ell(\mathbf{w}', \mathbf{x}, \mathbf{y})\| \leq \mathbf{D} \|\mathbf{w} - \mathbf{w}'\|. \end{aligned}$$

By Jensen's inequality  $\mathbf{H}(\cdot)$  is also  $\mathbf{D}^2$ -smooth. Notice that  $\mathbf{D}^2$  is the modulus of smoothness of  $\mathbf{H}(\cdot)$  and is the modulus of strong convexity of  $\mathbf{F}(\cdot)$ . This is a slight abuse of notation we define  $\mathbf{L} = \mathbf{D}^2$  and the condition number as the ratio between  $\mathbf{L}$  and  $\mu = \mathbf{L} - \mathbf{F}$ . Finally we note that the regularizer  $\mathbf{r}(\cdot)$  could be *non-smooth*.

We have the following general risk bound of ERM for supervised learning.

**Theorem 7** For any  $0 < \epsilon < 1/2$ , define

$$\begin{aligned} \mathbf{C} &= 4 \left( 8 + \sqrt{2 \log \frac{2 \log_2(n) + \log_2(2\mathbf{R})}{\epsilon}} \right), \\ \mathbf{H}_* &= \mathbf{H}(\mathbf{w}_*) = \mathbf{F}(\mathbf{w}_*) + \mathbf{r}(\mathbf{w}_*). \end{aligned} \quad \mathbf{4}$$

Under **Assumptions 4(b), 4(c), 6, 7, 8, 9, and 10** with probability at least  $1 - \epsilon$ , we have

$$\mathbf{F}(\hat{\mathbf{w}}) - \mathbf{F}(\mathbf{w}_*) \leq \max \left( \frac{\mathbf{M} + \mathbf{P}}{n^2} + \frac{\mathbf{L}}{2n^4}, \frac{4\mathbf{R}^2\mathbf{L}^2\mathbf{C}^2}{n} + \frac{4\mathbf{R}\mathbf{M} \log(2/\epsilon)}{n} + \frac{8\mathbf{L}\mathbf{H}_* \log(2/\epsilon)}{n} \right).$$

Furthermore, if

$$n \geq \frac{16\mathbf{L}^2\mathbf{C}^2}{2\epsilon} = 8\mathbf{L}^2\mathbf{C}^2/\epsilon,$$

with probability at least  $1 - \epsilon$ , we have

$$\mathbf{F}(\hat{\mathbf{w}}) - \mathbf{F}(\mathbf{w}_*) \leq \max \left( \frac{\mathbf{M} + \mathbf{P}}{n^2} + \frac{\mathbf{L}}{2n^4}, \frac{8\mathbf{M}^2 \log^2(2/\epsilon)}{n^2} + \frac{16\mathbf{L}\mathbf{H}_* \log(2/\epsilon)}{n} \right). \quad \mathbf{5}$$

**Remark 8** The first part of Theorem 7 presents an  $O(1/n)$  risk bound similar to the  $O(1/n)$  risk bound of [Dharan et al. 2009](#). The second part is an  $O(1/[n^2] + \mathbf{H}_*/n)$  risk bound and in this case the lower bound of  $n$  is  $\Omega(1/\epsilon^2)$  which is dimensionally independent. Thus Theorem 7 can be applied even when the dimensionality is infinite. Generally speaking the regularizer  $\mathbf{r}(\cdot)$  is nonnegative and thus  $\mathbf{H}_* \geq \mathbf{F}_*$ . So the second bound is even better than those in Theorems 7 and 8. Finally we note that Theorem 7 should be treated as a counterpart of Theorem 6 for supervised learning because both of them do not rely on the individual convexity. **Assumption 4(d)**. One may wonder whether it is possible to derive a counterpart of Theorem 7 that is whether it is possible to utilize the individual convexity to reduce the lower bound of  $n$  by a factor of  $1/\epsilon$ . We investigate this question as a future work.

For brevity we treat  $\mathbf{C}$  as a constant because it only has a double logarithmic dependence on  $n$ .

## 4. Analysis

We here present the key idea of our analysis and the proof of Theorem 4.1. The omitted ones can be found in appendices.

### 4.1. The Key Idea

By the convexity of  $\widehat{F}(\cdot)$  and the optimality condition of  $\widehat{w}$  (Boyd and Vandenberghe [14]) we have

$O(1/n)$  AND  $O(1/n^2)$  TYPE OF I K B O N D OF E M

**Lemma 1** *Under Assumptions 2 and 4(d), with probability at least 1*

where the last step is due to

$$\begin{aligned} & \|\widehat{\mathbf{w}} - \mathbf{w}_*\| \sqrt{\frac{\mathbf{LC}(\cdot, \cdot)(\mathbf{F}(\widehat{\mathbf{w}}) - \mathbf{F}(\mathbf{w}_*)))}{n}} = \frac{\mathbf{LC}(\cdot, \cdot)}{2n} \|\widehat{\mathbf{w}} - \mathbf{w}_*\|^2 + \frac{\mathbf{F}(\widehat{\mathbf{w}}) - \mathbf{F}(\mathbf{w}_*)}{2}, \\ & \|\widehat{\mathbf{w}} - \mathbf{w}_*\| \sqrt{\frac{\mathbf{LC}(\cdot, \cdot)\mathbf{G}}{n}} = \frac{\mathbf{LC}(\cdot, \cdot)}{2n} \|\widehat{\mathbf{w}} - \mathbf{w}_*\|^2 + \frac{\mathbf{G}}{2}. \end{aligned}$$

From [\(4\)](#) we get

$$\begin{aligned} & \frac{1}{2} (\mathbf{F}(\widehat{\mathbf{w}}) - \mathbf{F}(\mathbf{w}_*)) \\ & \frac{2\mathbf{LC}(\cdot, \cdot)}{n} \|\widehat{\mathbf{w}} - \mathbf{w}_*\|^2 + \frac{2\mathbf{M} \log(2/\delta)}{n} \|\widehat{\mathbf{w}} - \mathbf{w}_*\| + \|\widehat{\mathbf{w}} - \mathbf{w}_*\| \sqrt{\frac{8\mathbf{L}\mathbf{F}_* \log(2/\delta)}{n}} \\ & + 2\mathbf{L} \|\widehat{\mathbf{w}} - \mathbf{w}_*\| + \frac{\mathbf{G}}{2} + \frac{\mathbf{LC}(\cdot, \cdot)}{n} \|\widehat{\mathbf{w}} - \mathbf{w}_*\| \\ & \leq \frac{8\mathbf{R}^2\mathbf{LC}(\cdot, \cdot)}{n} + \frac{4\mathbf{R}\mathbf{M} \log(2/\delta)}{n} + 4\mathbf{R} \sqrt{\frac{2\mathbf{L}\mathbf{F}_* \log(2/\delta)}{n}} + \left( 4\mathbf{R}\mathbf{L} + \frac{\mathbf{G}}{2} + \frac{2\mathbf{R}\mathbf{LC}(\cdot, \cdot)}{n} \right), \end{aligned}$$

which implies  $\dots$

### 5. Conclusions and Future Work

In this paper we study the excess risk of EM for CO. Our theoretical results show that it is possible to achieve  $\mathcal{O}(1/n)$  type of risk bounds under the smoothness and standard assumptions  $\mathcal{E}$  or the smoothness and strong convexity conditions  $\mathcal{E}$ . The first part of Theorems 1 and 2. A more exciting result is that when  $n$  is large enough EM has  $\mathcal{O}(1/n^2)$  type of risk bounds under the smoothness strong convexity and standard assumptions  $\mathcal{E}$  the second part of Theorems 1 and 2.

In the context of CO there remain any open problems about EM.

Our current results are restricted to the Hilbert or Euclidean space because the smoothness and strong convexity are defined in terms of the  $\ell_2$  norm. We will extend our analysis to other geometries in the future.

As mentioned in Remark 3 under the strong convexity condition a dimensionally independent risk bound e.g.  $\widetilde{\mathcal{O}}(\cdot/n)$  or  $\widetilde{\mathcal{O}}(1/n)$  that holds with high probability is still missing.

As discussed in Remark 8 it is unclear whether the convexity of the loss can be exploited to prove the lower bound of  $n$  in the second part of Theorem 1. Ideally we expect that  $n = \Omega(\cdot)$  is sufficient to deliver an  $\mathcal{O}(1/[n^2] + \mathbf{H}_*/n)$  risk bound.

4. The  $\mathcal{O}(1/n^2)$  type of risk bounds require both the smoothness and strong convexity conditions. One may investigate whether strong convexity can be relaxed to other weaker conditions such as exponential concavity Hazan et al. [10].

Finally as far as we know there are no  $\mathcal{O}(1/n^2)$  type of risk bounds for stochastic approximation. We will try to establish such bounds for A.

### Acknowledgments

This work was partially supported by the NCFC Jiangsu F BK N F II 4 9 II 4 99 and the Collaborative Innovation Center of Novel Software Technology and Industrialization of Nanjing University.

## References

- A e h Agarwa Peter L Bart ett Pradeep av u ar and Mart n J\_ a nwr ght\_ Infor at on theoret c ower bounds on the orac e co p ex ty of stochast c convex opt zat on\_ *IEEE Transactions on Information Theory* . 49 0 -
- Franc s Bach and Er c Mou nes\_ Non strong y convex s ooth stochast c approx at on w th convergence rate  $O(1/n)$ \_ In *Advances in Neural Information Processing Systems 26* pages 49 -
- Peter L Bart ett and hahar Mende son\_ ad e acher and gauss an co p ex t es r s bounds and structural resu ts\_ *Journal of Machine Learning Research* 4 4 00 -
- Peter L Bart ett O v er Bousquet and hahar Mende son\_ Local rad e acher co p ex t es\_ *The Annals of Statistics* 4 49 00 -
- stephen Boyd and L even andenberghe\_ *Convex Optimization*\_ Ca br dge n vers ty Press 00 4 -
- G u a Desa vo Mehryar Mohr and ar yed\_ Learn ng w th deep cascades\_ In *Proceedings of the 26th International Conference on Algorithmic Learning Theory* pages 4 9 0 -
- ta y Fe d an\_ Genera zat on of er n stochast c convex opt zat on r he d ens on str es bac \_ *ArXiv e-prints* arX v: 0 0 44 4 0 -
- A on Gonen and ha ha ev hwhartz\_ Average stab ty s nvar ant to data precondition ng\_ p cat ons to exp concave e p r ca r s n zat on\_ *ArXiv e-prints* arX v: 0 0 4 0 -
- E ad Hazan and atyen Ka e\_ Beyond the regret n zat on barr er an opt a a gor th for stochast c strong y convex opt zat on\_ In *Proceedings of the 24th Annual Conference on Learning Theory* pages 4 4 0 -
- E ad Hazan A t Agarwa and atyen Ka e\_ Logarith c regret a gor th s for on ne convex opt zat on\_ *Machine Learning* 9 9 9 00 -
- e Johnson and ong Zhang\_ Acce erat ng stochast c grad ent descent us ng pred ct ve variance reduction\_ In *Advances in Neural Information Processing Systems 26* pages 4 -
- ad r Ko tch ns \_ *Oracle Inequalities in Empirical Risk Minimization and Sparse Recovery Problems*\_ pr nger 0 -
- o er Koren and K r Levy\_ Fast rates for exp concave e p r ca r s n zat on\_ In *Advances in Neural Information Processing Systems 28* pages 4 4 0 -
- Haro d J\_ Kushner and G\_ George Y n\_ *Stochastic Approximation and Recursive Algorithms and Applications*

Mehrdad Mahdavi, Lun Zhang and Yong Jin. Lower and upper bounds on the generalization of stochastic exponentally concave optimization. In *Proceedings of the 28th Conference on Learning Theory* –

Conrad McDiarmid. On the method of bounded differences. In *Surveys in Combinatorics* pages 499–

Nishant A. Mehta. Fast rates with high probability for exponentially concave statistical learning. *ArXiv e-prints* arXiv: –

Yong Zhang and Merrin. Generalization error bounds for Bayesian mixture algorithms. *Journal of Machine Learning Research* 4: 9 –

François Bach and Ercole Moulines. Non asymptotic analysis of stochastic approximation algorithms for machine learning. In *Advances in Neural Information Processing Systems* 24 pages 449 –

Andrei Nemirovski, Gennadiy Juditsky, and Alexander Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on Optimization* 94: 999–

Yurii Nesterov. *Introductory lectures on convex optimization: a basic course* volume of *Applied optimization*. Kluwer Academic Publishers –

Dmitry Panchenko. Extensions of an inequality of Vapnik and Chervonenkis. *Electronic Communications in Probability* –

Giles Pisier. *The volume of convex bodies and Banach space geometry*. Cambridge Tracts in Mathematics No. 94. Cambridge University Press 99–

Yaniv Plan and Omer Shram. One-bit compressed sensing by neural programming. *Communications on Pure and Applied Mathematics* –

Alexander A. Ohad and Karthikeyan Sridharan. Making gradient descent optimal for strongly convex stochastic optimization. In *Proceedings of the 29th International Conference on Machine Learning* pages 449 –

Bernhard Schölkopf and Alexander J. Smola. *Learning with kernels: support vector machines, regularization, optimization, and beyond*. MIT Press –

John D. Elfwortz and John D. Elfwortz. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press –

John D. Elfwortz, Ohad Shram, Nathan Srebro and Karthikeyan Sridharan. Stochastic convex optimization. In *Proceedings of the 22nd Annual Conference on Learning Theory* –

Alexander Shapiro, Darina Dentcheva and Andrzej Ruszczyński. *Lectures on Stochastic Programming: Modeling and Theory*. IAIM second edition –

Teo-Yee and Dong Xuan Zhou. Learning theory estimates via integral operators and the approximation. *Constructive Approximation* –

Nathanrebro, Karthikr dharan and Abuewar. Optimal rates for learning with a smooth loss. *ArXiv e-prints* arXiv:1909.00092.

Karthikr dharan, Haoheshwartz and Nathanrebro. Fast rates for regularized objectives. In *Advances in Neural Information Processing Systems 21*, pages 4009–4019.

Alexandre Boussoffier. Optimal aggregation of classifiers in statistical learning. *The Annals of Statistics*, 40(4):1469–1501, 2012.

Andrius J. Kulkarni. *The Nature of Statistical Learning Theory*, Springer, second edition, 2009.

Andrius J. Kulkarni. *Statistical Learning Theory*, Wiley-Interscience, 1999.

Lun Zhang, Mehrdad Mahdavi, and Tong Yin. Linear convergence with condition number independent access of full gradients. In *Advances in Neural Information Processing Systems 26*, pages 909–918, 2013.

Lun Zhang, Yanbao Yang, Tong Yin, and Xiaohe He.  $O(\log T)$  projections for stochastic optimization of smooth and strongly convex functions. In *Proceedings of the 30th International Conference on Machine Learning*, 2017.

## Appendix A. Proof of Lemma 1

We introduce Lemma of [Lai and Zhou](#) (2014).

**Lemma 3** Let  $\mathcal{H}$  be a Hilbert space and let  $\xi$  be a random variable with values in  $\mathcal{H}$ . Assume  $\|\xi\| \leq M$  almost surely. Denote  $\mathbb{E}[\xi] = \mathbb{E}[\xi]$ . Let  $\xi_1, \dots, \xi_m$  be  $m$  ( $m < \infty$ ) independent draws of  $\xi$ . For any  $0 < \epsilon < 1$ , with confidence  $1 - \epsilon$ ,

$$\left\| \frac{1}{m} \sum_{i=1}^m [\xi_i - \mathbb{E}[\xi]] \right\| \leq \frac{2M \log(2/\epsilon)}{m} + \sqrt{\frac{2 \mathbb{E}[\|\xi\|^2] \log(2/\epsilon)}{m}}.$$

First consider a fixed  $w \in \mathcal{H}$ . Since  $f_i(\cdot)$  is  $L$ -smooth we have

$$\|f_i(w) - f_i(w_*)\| \leq L \|w - w_*\|.$$

Because  $f_i(\cdot)$  is both convex and  $L$ -smooth by [Nesterov](#) (2004) we have

$$\|f_i(w) - f_i(w_*)\|^2 \leq L \langle f_i(w) - f_i(w_*), w - w_* \rangle.$$

Taking expectation over both sides we have

$$\mathbb{E} \left[ \|f_i(w) - f_i(w_*)\|^2 \right] \leq L \langle \mathbb{E}[f_i(w) - f_i(w_*)], w - w_* \rangle = L \langle \mathbb{E}[f_i(w) - f_i(w_*)], w - w_* \rangle$$

where the last inequality follows from the optimality condition of  $w_*$  [\(1\)](#)

$$\langle \mathbb{E}[f_i(w) - f_i(w_*)], w - w_* \rangle \leq 0, \forall w \in \mathcal{H}.$$

Following Lemma with probability at least  $1 - \delta$  we have

$$\begin{aligned} & \left\| \mathbf{F}(\mathbf{w}) - \mathbf{F}(\mathbf{w}_*) - [\hat{\mathbf{F}}(\mathbf{w}) - \hat{\mathbf{F}}(\mathbf{w}_*)] \right\| \\ &= \left\| \mathbf{F}(\mathbf{w}) - \mathbf{F}(\mathbf{w}_*) - \frac{1}{n} \sum_{i=1}^n [\mathbf{f}_i(\mathbf{w}) - \mathbf{f}_i(\mathbf{w}_*)] \right\| \\ & \leq \frac{2L \|\mathbf{w} - \mathbf{w}_*\| \log(2/\delta)}{n} + \sqrt{\frac{2L(\mathbf{F}(\mathbf{w}) - \mathbf{F}(\mathbf{w}_*)) \log(2/\delta)}{n}}. \end{aligned}$$

We obtain Lemma by taking the union bound over all  $\mathbf{w} \in (\mathcal{B}, \delta)$ . On this end we need an upper bound of the covering number  $(\mathcal{B}, \delta)$ .

Let  $\mathcal{B}$  be an unit ball of  $\mathbf{d}$  dimensions and  $(\mathcal{B}, \delta)$  be its covering number with respect to a cardinality. According to a standard volume comparison argument (Pisier 1999) we have

$$\log(\mathcal{B}, \delta) \leq \mathbf{d} \log \frac{3}{\delta}.$$

Let  $\mathcal{B}(\mathbf{R})$  be a ball centered at origin with radius  $\mathbf{R}$ . Since we assume  $\mathcal{B}(\mathbf{R})$  it follows that

$$\log(\mathcal{B}(\mathbf{R}), \frac{\mathbf{R}}{2}) \leq \log \left| \mathcal{B}(\mathbf{R}, \frac{\mathbf{R}}{2}) \right| \leq \mathbf{d} \log \frac{6\mathbf{R}}{\mathbf{R}}$$

where the first inequality is because the covering numbers are almost increasing by inclusion (Pan and Vershyn 2015).

### Appendix B. Proof of Lemma 2

To apply Lemma we need an upper bound of  $E[\|\mathbf{f}_i(\mathbf{w}_*)\|^2]$ . Since  $\mathbf{f}_i(\cdot)$  is  $L$ -smooth and nonnegative from Lemma 4 of (Broderick et al. 2015) we have

$$\|\mathbf{f}_i(\mathbf{w}_*)\|^2 \leq 4L\mathbf{f}_i(\mathbf{w}_*)$$

and thus

$$E[\|\mathbf{f}_i(\mathbf{w}_*)\|^2] \leq 4LE[\mathbf{f}_i(\mathbf{w}_*)] = 4L\mathbf{F}_*.$$

From Assumption 5 we have  $\mathbf{M} \leq \mathbf{M}$ . Then according to Lemma with probability at least  $1 - \delta$  we have

$$\left\| \mathbf{F}(\mathbf{w}_*) - \hat{\mathbf{F}}(\mathbf{w}_*) \right\| = \left\| \mathbf{F}(\mathbf{w}_*) - \frac{1}{n} \sum_{i=1}^n \mathbf{f}_i(\mathbf{w}_*) \right\| \leq \frac{2\mathbf{M} \log(2/\delta)}{n} + \sqrt{\frac{8L\mathbf{F}_* \log(2/\delta)}{n}}.$$

### Appendix C. Proof of Theorem 3

The proof follows the same logic as that of Theorem 1 under Assumption 4(b). It becomes

$$\begin{aligned} & \mathbf{F}(\hat{\mathbf{w}}) - \mathbf{F}(\mathbf{w}_*) + \frac{1}{2} \|\hat{\mathbf{w}} - \mathbf{w}_*\|^2 \\ & \left( \underbrace{\left\| \mathbf{F}(\hat{\mathbf{w}}) - \mathbf{F}(\mathbf{w}_*) - [\hat{\mathbf{F}}(\hat{\mathbf{w}}) - \hat{\mathbf{F}}(\mathbf{w}_*)] \right\|}_{:=A_1} + \underbrace{\left\| \mathbf{F}(\mathbf{w}_*) - \hat{\mathbf{F}}(\mathbf{w}_*) \right\|}_{:=A_2} \right) \|\hat{\mathbf{w}} - \mathbf{w}_*\|. \end{aligned}$$



ubst tut ng and nto w th probab ty at east 1 2 we have

$$\begin{aligned} F(\widehat{w}) - F(w_*) &+ \frac{1}{2} \|\widehat{w} - w_*\|^2 \\ &+ \frac{LC(\cdot, \cdot) \|\widehat{w} - w_*\|^2}{n} + \|\widehat{w} - w_*\| \sqrt{\frac{LC(\cdot, \cdot)(F(\widehat{w}) - F(w_*))}{n}} \\ &+ \frac{2M \log(2/\delta) \|\widehat{w} - w_*\|}{n} + \|\widehat{w} - w_*\| \sqrt{\frac{8LF_* \log(2/\delta)}{n}} \\ &+ 2L \|\widehat{w} - w_*\| + \|\widehat{w} - w_*\| \sqrt{\frac{LC(\cdot, \cdot)G}{n} + \frac{LC(\cdot, \cdot) \|\widehat{w} - w_*\|}{n}} \end{aligned}$$

9

o prove we subst tute and

$$\|\widehat{w} - w_*\| \sqrt{\frac{8LF_* \log(2/\delta)}{n}} + \frac{4LF_* \log(2/\delta)}{n} + \frac{1}{2} \|\widehat{w} - w_*\|^2$$

nto 9 and then obta n

$$\begin{aligned} &\frac{1}{2} (F(\widehat{w}) - F(w_*)) \\ &+ \frac{2LC(\cdot, \cdot) \|\widehat{w} - w_*\|^2}{n} + \frac{2M \log(2/\delta) \|\widehat{w} - w_*\|}{n} + \frac{4LF_* \log(2/\delta)}{n} \\ &+ 2L \|\widehat{w} - w_*\| + \frac{G}{2} + \frac{LC(\cdot, \cdot) \|\widehat{w} - w_*\|}{n} \\ &+ \frac{8R^2LC(\cdot, \cdot)}{n} + \frac{4RM \log(2/\delta)}{n} + \frac{4LF_* \log(2/\delta)}{n} + \left( 4RL + \frac{G}{2} + \frac{2RLC(\cdot, \cdot)}{n} \right). \end{aligned}$$

wh ch p es

o prove we subst tute

$$\begin{aligned} &\|\widehat{w} - w_*\| \sqrt{\frac{LC(\cdot, \cdot)(F(\widehat{w}) - F(w_*))}{n}} + \frac{2LC(\cdot, \cdot)(F(\widehat{w}) - F(w_*))}{n} + \frac{1}{8} \|\widehat{w} - w_*\|^2, \\ &\frac{2M \log(2/\delta) \|\widehat{w} - w_*\|}{n} + \frac{16M^2 \log^2(2/\delta)}{n^2} + \frac{1}{16} \|\widehat{w} - w_*\|^2, \\ &\|\widehat{w} - w_*\| \sqrt{\frac{8LF_* \log(2/\delta)}{n}} + \frac{64LF_* \log(2/\delta)}{n} + \frac{1}{32} \|\widehat{w} - w_*\|^2, \\ &2L \|\widehat{w} - w_*\| + \frac{64L^2}{64} + \frac{1}{64} \|\widehat{w} - w_*\|^2, \\ &\|\widehat{w} - w_*\| \sqrt{\frac{LC(\cdot, \cdot)G}{n}} + \frac{32LC(\cdot, \cdot)G}{n} + \frac{1}{128} \|\widehat{w} - w_*\|^2, \\ &\frac{LC(\cdot, \cdot) \|\widehat{w} - w_*\|}{n} + \frac{32L^2C^2(\cdot)}{n^2} + \frac{1}{128} \|\widehat{w} - w_*\|^2 \end{aligned}$$

into (9) and then obtain

$$\begin{aligned} & \mathbb{E} \left[ \mathbf{F}(\hat{\mathbf{w}}) - \mathbf{F}(\mathbf{w}_*) + \frac{1}{4} \|\hat{\mathbf{w}} - \mathbf{w}_*\|^2 \right] \\ & \leq \frac{\mathbf{LC}(\cdot, \cdot) \|\hat{\mathbf{w}} - \mathbf{w}_*\|^2}{n} + \frac{2\mathbf{LC}(\cdot, \cdot) (\mathbf{F}(\hat{\mathbf{w}}) - \mathbf{F}(\mathbf{w}_*))}{n} + \frac{16\mathbf{M}^2 \log^2(2/\delta)}{n^2} + \frac{64\mathbf{L}\mathbf{F}_* \log(2/\delta)}{n} \\ & \quad + \frac{64\mathbf{L}^2 \sigma^2}{n} + \frac{32\mathbf{LC}(\cdot, \cdot) \mathbf{G}}{n} + \frac{32\mathbf{L}^2 \mathbf{C}^2(\cdot, \cdot)}{n^2} \\ & \leq \frac{1}{4} \|\hat{\mathbf{w}} - \mathbf{w}_*\|^2 + \frac{1}{2} (\mathbf{F}(\hat{\mathbf{w}}) - \mathbf{F}(\mathbf{w}_*)) + \frac{16\mathbf{M}^2 \log^2(2/\delta)}{n^2} + \frac{64\mathbf{L}\mathbf{F}_* \log(2/\delta)}{n} \\ & \quad + \frac{64\mathbf{L}^2 \sigma^2}{n} + 8\mathbf{G} + 2\sigma^2 \end{aligned}$$

which implies

### Appendix D. Proof of Theorem 5

Throughout Assumption 4(d) Lemma 4 which is used in the proofs of Theorems 5 and 6 does not hold any more. Instead we will use the following version that only relies on the smoothness condition.

**Lemma 4** Under Assumption 2, with probability at least  $1 - \delta$ , for any  $\mathbf{w} \in \mathcal{W}(\delta, \epsilon)$ , we have

$$\left\| \mathbf{F}(\mathbf{w}) - \mathbf{F}(\mathbf{w}_*) - [\hat{\mathbf{F}}(\mathbf{w}) - \hat{\mathbf{F}}(\mathbf{w}_*)] \right\| \leq \frac{\mathbf{LC}(\cdot, \cdot) \|\mathbf{w} - \mathbf{w}_*\|}{n} + \mathbf{L} \|\mathbf{w} - \mathbf{w}_*\| \sqrt{\frac{\mathbf{C}(\cdot, \cdot)}{n}}$$

where  $\mathbf{C}(\cdot, \cdot)$  is defined in (9).

The above lemma is a direct consequence of Lemma 4 and the uniform bound.

The rest of the proof is similar to those of Theorems 5 and 6. We first derive a counterpart of Lemma 4. Combining with Lemma 4 with probability at least  $1 - \delta$  we have

$$\begin{aligned} & \left\| \mathbf{F}(\tilde{\mathbf{w}}) - \mathbf{F}(\mathbf{w}_*) - [\hat{\mathbf{F}}(\tilde{\mathbf{w}}) - \hat{\mathbf{F}}(\mathbf{w}_*)] \right\| \\ & \leq \frac{\mathbf{LC}(\cdot, \cdot) \|\tilde{\mathbf{w}} - \mathbf{w}_*\|}{n} + \mathbf{L} \|\tilde{\mathbf{w}} - \mathbf{w}_*\| \sqrt{\frac{\mathbf{C}(\cdot, \cdot)}{n}} + 2\mathbf{L} \\ & \leq \frac{\mathbf{LC}(\cdot, \cdot) \|\hat{\mathbf{w}} - \mathbf{w}_*\|}{n} + \mathbf{L} \|\hat{\mathbf{w}} - \mathbf{w}_*\| \sqrt{\frac{\mathbf{C}(\cdot, \cdot)}{n}} + \frac{\mathbf{LC}(\cdot, \cdot)}{n} + \mathbf{L} \sqrt{\frac{\mathbf{C}(\cdot, \cdot)}{n}} + 2\mathbf{L}. \end{aligned}$$

Substituting (10) and (11) into (8) with probability at least  $1 - 2\delta$  we have

$$\begin{aligned} & \mathbb{E} \left[ \mathbf{F}(\hat{\mathbf{w}}) - \mathbf{F}(\mathbf{w}_*) + \frac{1}{2} \|\hat{\mathbf{w}} - \mathbf{w}_*\|^2 \right] \\ & \leq \frac{\mathbf{LC}(\cdot, \cdot) \|\hat{\mathbf{w}} - \mathbf{w}_*\|^2}{n} + \mathbf{L} \|\hat{\mathbf{w}} - \mathbf{w}_*\|^2 \sqrt{\frac{\mathbf{C}(\cdot, \cdot)}{n}} \\ & \quad + \frac{2\mathbf{M} \log(2/\delta) \|\hat{\mathbf{w}} - \mathbf{w}_*\|}{n} + \|\hat{\mathbf{w}} - \mathbf{w}_*\| \sqrt{\frac{8\mathbf{L}\mathbf{F}_* \log(2/\delta)}{n}} \\ & \quad + 2\mathbf{L} \|\hat{\mathbf{w}} - \mathbf{w}_*\| + \mathbf{L} \|\hat{\mathbf{w}} - \mathbf{w}_*\| \sqrt{\frac{\mathbf{C}(\cdot, \cdot)}{n}} + \frac{\mathbf{LC}(\cdot, \cdot) \|\hat{\mathbf{w}} - \mathbf{w}_*\|}{n}. \end{aligned}$$

to get (4) we substitute

$$L(\hat{w}, w_*)^2 \sqrt{\frac{C(\cdot, \cdot)}{n}} + \frac{L^2 C(\cdot, \cdot)}{n} \frac{\hat{w} w_*^2}{n} + \frac{1}{4} \frac{\hat{w} w_*^2}{n},$$

$$\frac{1}{n} \frac{\hat{w} w_*^2}{n} \sqrt{\frac{8LF_* \log(2/\cdot)}{n}} + \frac{8LF_* \log(2/\cdot)}{n} + \frac{1}{4} \frac{\hat{w} w_*^2}{n}$$

into (4) and then obtain

$$F(\hat{w}) - F(w_*)$$

$$\frac{LC(\cdot, \cdot)}{n} \frac{\hat{w} w_*^2}{n} + \frac{L^2 C(\cdot, \cdot)}{n} \frac{\hat{w} w_*^2}{n} + \frac{2M \log(2/\cdot)}{n} \frac{\hat{w} w_*^2}{n} + \frac{8LF_* \log(2/\cdot)}{n}$$

$$+ 2L \frac{\hat{w} w_*^2}{n} + L \frac{\hat{w} w_*^2}{n} \sqrt{\frac{C(\cdot, \cdot)}{n}} + \frac{LC(\cdot, \cdot)}{n} \frac{\hat{w} w_*^2}{n}$$

$$(4) \frac{4R^2 LC(\cdot, \cdot)}{n} + \frac{4R^2 L^2 C(\cdot, \cdot)}{n} + \frac{4RM \log(2/\cdot)}{n} + \frac{8LF_* \log(2/\cdot)}{n}$$

$$+ \left( 4RL + 2RL \sqrt{\frac{C(\cdot, \cdot)}{n}} + \frac{2RLC(\cdot, \cdot)}{n} \right)$$

which proves (4) -

to get (5) we substitute

$$\frac{2M \log(2/\cdot)}{n} \frac{\hat{w} w_*^2}{n} + \frac{8M^2 \log^2(2/\cdot)}{n^2} + \frac{1}{8} \frac{\hat{w} w_*^2}{n},$$

$$\frac{1}{n} \frac{\hat{w} w_*^2}{n} \sqrt{\frac{8LF_* \log(2/\cdot)}{n}} + \frac{32LF_* \log(2/\cdot)}{n} + \frac{1}{16} \frac{\hat{w} w_*^2}{n},$$

$$\frac{2L}{n} \frac{\hat{w} w_*^2}{n} + \frac{32L^2}{n} \frac{\hat{w} w_*^2}{n} + \frac{1}{32} \frac{\hat{w} w_*^2}{n},$$

$$L \frac{\hat{w} w_*^2}{n} \sqrt{\frac{C(\cdot, \cdot)}{n}} + \frac{16L^2 C(\cdot, \cdot)}{n} \frac{\hat{w} w_*^2}{n} + \frac{1}{64} \frac{\hat{w} w_*^2}{n},$$

$$\frac{LC(\cdot, \cdot)}{n} \frac{\hat{w} w_*^2}{n} + \frac{16L^2 C^2(\cdot, \cdot)}{n^2} + \frac{1}{64} \frac{\hat{w} w_*^2}{n}$$

into (5) and then obtain

$$F(\hat{w}) - F(w_*) + \frac{1}{4} \frac{\hat{w} w_*^2}{n}$$

$$\frac{LC(\cdot, \cdot)}{n} \frac{\hat{w} w_*^2}{n} + L \frac{\hat{w} w_*^2}{n} \sqrt{\frac{C(\cdot, \cdot)}{n}} + \frac{8M^2 \log^2(2/\cdot)}{n^2} + \frac{32LF_* \log(2/\cdot)}{n}$$

$$+ \left( \frac{32L^2}{n} + \frac{16L^2 C(\cdot, \cdot)}{n} + \frac{16L^2 C^2(\cdot, \cdot)}{n^2} \right) \frac{\hat{w} w_*^2}{n}$$

$$(6) \frac{2}{25} \frac{\hat{w} w_*^2}{n} + \frac{1}{5} \frac{\hat{w} w_*^2}{n} + \frac{8M^2 \log^2(2/\cdot)}{n^2} + \frac{32LF_* \log(2/\cdot)}{n}$$

$$+ \left( \frac{32L^2}{n} + \frac{16}{25} + \frac{16^3}{625L^2} \right) \frac{\hat{w} w_*^2}{n}$$

$$\lambda/L \leq 16 \frac{1}{25} \frac{\hat{w} w_*^2}{n} + \frac{8M^2 \log^2(2/\cdot)}{n^2} + \frac{32LF_* \log(2/\cdot)}{n} + \left( \frac{32L^2}{n} + \frac{416}{625} \right) \frac{\hat{w} w_*^2}{n}$$

By subtracting  $\frac{1}{2}\|\widehat{\mathbf{w}} - \mathbf{w}_*\|^2/4$  from both sides we complete the proof of [\(9\)](#).

### Appendix E. Proof of Theorem 7

We consider two cases. In the first case we assume that

$$\frac{1}{n^2} < \frac{1}{2R}.$$

Since  $\mathbf{H}(\cdot)$  is  $L$ -smooth and  $\mathbf{r}(\cdot)$  is  $P$ -Lipschitz continuous we have

$$\begin{aligned} \mathbf{F}(\widehat{\mathbf{w}}) - \mathbf{F}(\mathbf{w}_*) &= \mathbf{H}(\widehat{\mathbf{w}}) + \mathbf{r}(\widehat{\mathbf{w}}) - \mathbf{H}(\mathbf{w}_*) - \mathbf{r}(\mathbf{w}_*) \\ &\leq \frac{L}{2}\|\widehat{\mathbf{w}} - \mathbf{w}_*\|^2 + P\|\widehat{\mathbf{w}} - \mathbf{w}_*\| \\ &\leq \frac{L}{2}\|\widehat{\mathbf{w}} - \mathbf{w}_*\|^2 + P\|\widehat{\mathbf{w}} - \mathbf{w}_*\| \leq \frac{M+P}{n^2} + \frac{L}{2n^4} \end{aligned} \quad (4)$$

where the last step utilizes Jensen's inequality

$$\|\mathbf{H}(\mathbf{w}_*)\| = \left\| \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathbb{D}} [\mathbf{H}(\mathbf{w}_*, \mathbf{x}, \mathbf{y})] \right\| \leq \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathbb{D}} [\|\mathbf{H}(\mathbf{w}_*, \mathbf{x}, \mathbf{y})\|] \leq M.$$

Next we study the case

$$\frac{1}{n^2} < \frac{1}{2R}.$$

From [\(9\)](#) we have

$$\begin{aligned} \mathbf{F}(\widehat{\mathbf{w}}) - \mathbf{F}(\mathbf{w}_*) &+ \frac{1}{2}\|\widehat{\mathbf{w}} - \mathbf{w}_*\|^2 \\ &\leq \mathbf{F}(\widehat{\mathbf{w}}) - \mathbf{F}(\mathbf{w}_*) + \frac{1}{2}\|\widehat{\mathbf{w}} - \mathbf{w}_*\|^2 \\ &= \mathbf{H}(\widehat{\mathbf{w}}) - \mathbf{H}(\mathbf{w}_*) + \frac{1}{2}\|\widehat{\mathbf{w}} - \mathbf{w}_*\|^2 \\ &\leq \underbrace{\sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \|\widehat{\mathbf{w}} - \mathbf{w}_*\|} \left\langle \mathbf{H}(\mathbf{w}) - \mathbf{H}(\mathbf{w}_*) + \frac{1}{2}\|\mathbf{w} - \mathbf{w}_*\|^2, \mathbf{w} - \mathbf{w}_* \right\rangle}_{:=B_1} \\ &\quad + \underbrace{\left\| \mathbf{H}(\mathbf{w}_*) - \widehat{\mathbf{H}}(\mathbf{w}_*) \right\|}_{:=B_2} \|\widehat{\mathbf{w}} - \mathbf{w}_*\|. \end{aligned} \quad (4)$$

The first bound  $B_1$  follows from the fact that the random variable  $\|\widehat{\mathbf{w}} - \mathbf{w}_*\|$  lies in the range  $(1/n^2, 2R]$  we develop the following lemma.

**Lemma 5** Under Assumptions 7 and 8, with probability at least  $1 - \delta$ , for all

$$\frac{1}{n^2} < \frac{1}{2R}$$

the following bound holds:

$$\sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma} \left\langle \mathbf{H}(\mathbf{w}) - \mathbf{H}(\mathbf{w}_*) + \frac{1}{2}\|\mathbf{w} - \mathbf{w}_*\|^2, \mathbf{w} - \mathbf{w}_* \right\rangle \leq \frac{4L}{n} \left( 8 + \sqrt{2 \log \frac{\mathbf{s}}{\delta}} \right)$$

where  $\mathbf{s} = 2 \log_2(n) + \log_2(2R)$ .

Based on the above lemma we have with probability at least 1

$$\mathbf{B}_1 \quad \frac{4L}{n} \frac{\widehat{w} w_*^2}{n} \left( 8 + \sqrt{2 \log \frac{s}{\epsilon}} \right) = \frac{LC}{n} \frac{\widehat{w} w_*^2}{n} \quad (44)$$

where  $\mathbf{C}$  is defined in (43)

we then proceed to handle  $\mathbf{B}_2$  which can be upper bounded in the same way as  $\mathbf{A}_2$ . In particular we have the following lemma

**Lemma 6** Under Assumptions 7, 8, and 10, with probability at least  $1 - \epsilon$ , we have

$$\left\| \mathbf{H}(\mathbf{w}_*) - \widehat{\mathbf{H}}(\mathbf{w}_*) \right\| \leq \frac{2M \log(2/\epsilon)}{n} + \sqrt{\frac{8LH_* \log(2/\epsilon)}{n}}. \quad (45)$$

Substituting (44) and (45) into (42) with probability at least  $1 - 2\epsilon$  we have

$$\begin{aligned} \mathbf{F}(\widehat{\mathbf{w}}) &= \mathbf{F}(\mathbf{w}_*) + \frac{2L}{n} \frac{\widehat{w} w_*^2}{n} \\ &\leq \frac{LC}{n} \frac{\widehat{w} w_*^2}{n} + \frac{2M \log(2/\epsilon)}{n} \frac{\widehat{w} w_*^2}{n} + \frac{\widehat{w} w_*^2}{n} \sqrt{\frac{8LH_* \log(2/\epsilon)}{n}}. \end{aligned} \quad (46)$$

we substitute

$$\begin{aligned} &\frac{LC}{n} \frac{\widehat{w} w_*^2}{n} \leq \frac{L^2 C^2}{n} \frac{\widehat{w} w_*^2}{n} + \frac{\widehat{w} w_*^2}{4n}, \\ &\frac{\widehat{w} w_*^2}{n} \sqrt{\frac{8LH_* \log(2/\epsilon)}{n}} \leq \frac{8LH_* \log(2/\epsilon)}{n} + \frac{\widehat{w} w_*^2}{4n} \end{aligned}$$

into (46) and then have

$$\begin{aligned} \mathbf{F}(\widehat{\mathbf{w}}) &= \mathbf{F}(\mathbf{w}_*) + \frac{L^2 C^2}{n} \frac{\widehat{w} w_*^2}{n} + \frac{2M \log(2/\epsilon)}{n} \frac{\widehat{w} w_*^2}{n} + \frac{8LH_* \log(2/\epsilon)}{n} \\ &\leq \frac{4R^2 L^2 C^2}{n} + \frac{4RM \log(2/\epsilon)}{n} + \frac{8LH_* \log(2/\epsilon)}{n}. \end{aligned}$$

Combining the above inequality with (41) we obtain

to prove (40) we substitute

$$\begin{aligned} &\frac{2M \log(2/\epsilon)}{n} \frac{\widehat{w} w_*^2}{n} \leq \frac{8M^2 \log^2(2/\epsilon)}{n^2} + \frac{\widehat{w} w_*^2}{8n}, \\ &\frac{\widehat{w} w_*^2}{n} \sqrt{\frac{8LH_* \log(2/\epsilon)}{n}} \leq \frac{16LH_* \log(2/\epsilon)}{n} + \frac{\widehat{w} w_*^2}{8n} \end{aligned}$$

into (46) and then have

$$\begin{aligned} \mathbf{F}(\widehat{\mathbf{w}}) &= \mathbf{F}(\mathbf{w}_*) + \frac{\widehat{w} w_*^2}{4n} \\ &\leq \frac{LC}{n} \frac{\widehat{w} w_*^2}{n} + \frac{8M^2 \log^2(2/\epsilon)}{n^2} + \frac{16LH_* \log(2/\epsilon)}{n} \\ &\quad + \frac{\widehat{w} w_*^2}{4n} + \frac{8M^2 \log^2(2/\epsilon)}{n^2} + \frac{16LH_* \log(2/\epsilon)}{n}. \end{aligned}$$

Combining the above inequality with (41) we obtain

**Appendix F. Proof of Lemma 5**

First we partition the range  $(1/n^2, 2R]$  into  $s = \lceil 2 \log_2(n) + \log_2(2R) \rceil$  consecutive segments  $\Delta_1, \Delta_2, \dots, \Delta_s$  such that

$$\Delta_k = \left( \underbrace{\frac{2^{k-1}}{n^2}}_{:=\gamma_k^-}, \underbrace{2^k}_{:=\gamma_k^+} \right), \mathbf{k} = 1, \dots, s.$$

then we consider the case  $\Delta_k$  for a fixed value of  $\mathbf{k}$ . We have

$$\begin{aligned} & \sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma} \left\langle \mathbf{H}(\mathbf{w}) - \mathbf{H}(\mathbf{w}_*) - [\hat{\mathbf{H}}(\mathbf{w}) - \hat{\mathbf{H}}(\mathbf{w}_*)], \mathbf{w} - \mathbf{w}_* \right\rangle \\ & \sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma_k^+} \left\langle \mathbf{H}(\mathbf{w}) - \mathbf{H}(\mathbf{w}_*) - [\hat{\mathbf{H}}(\mathbf{w}) - \hat{\mathbf{H}}(\mathbf{w}_*)], \mathbf{w} - \mathbf{w}_* \right\rangle. \end{aligned}$$

Based on the McDiarmid's inequality [McDiarmid 99](#) and the Rademacher complexity [Bartlett and Mendelson 00](#) we have the following lemma to upper bound the last term.

**Lemma 7** *Under Assumptions 7 and 8, with probability at least  $1 - \delta$ , we have*

$$\begin{aligned} & \sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma_k^+} \left\langle \mathbf{H}(\mathbf{w}) - \mathbf{H}(\mathbf{w}_*) - [\hat{\mathbf{H}}(\mathbf{w}) - \hat{\mathbf{H}}(\mathbf{w}_*)], \mathbf{w} - \mathbf{w}_* \right\rangle \\ & \frac{L}{n} \left( \frac{\gamma_k^+}{k} \right)^2 \left( 8 + \sqrt{2 \log \frac{1}{\delta}} \right). \end{aligned}$$

Since  $\Delta_k$  we have

$$\gamma_k^+ = 2 \gamma_k^- \cdot 2.$$

thus with probability at least  $1 - \delta$  we have

$$\begin{aligned} & \sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma} \left\langle \mathbf{H}(\mathbf{w}) - \mathbf{H}(\mathbf{w}_*) - [\hat{\mathbf{H}}(\mathbf{w}) - \hat{\mathbf{H}}(\mathbf{w}_*)], \mathbf{w} - \mathbf{w}_* \right\rangle \\ & \frac{4L}{n} \left( \frac{\gamma}{k} \right)^2 \left( 8 + \sqrt{2 \log \frac{1}{\delta}} \right). \end{aligned}$$

we complete the proof by taking the union bound over  $s$  segments.

**Appendix G. Proof of Lemma 7**

To simplify the notation we define

$$\mathbf{h}_i(\mathbf{w}) = (\langle \mathbf{w}, \mathbf{x}_i \rangle, \mathbf{y}_i), \mathbf{i} = 1, \dots, n,$$

$$l(\mathbf{h}_1, \dots, \mathbf{h}_n) = \sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma_k^+} \left\langle \mathbf{H}(\mathbf{w}) - \mathbf{H}(\mathbf{w}_*) - \frac{1}{n} \sum_{i=1}^n [\mathbf{h}_i(\mathbf{w}) - \mathbf{h}_i(\mathbf{w}_*)], \mathbf{w} - \mathbf{w}_* \right\rangle.$$

To upper bound  $l(\mathbf{h}_1, \dots, \mathbf{h}_n)$  we utilize the McDiarmid's inequality [McDiarmid 99](#).

**Theorem 8** Let  $\mathbf{X}_1, \dots, \mathbf{X}_n$  be independent random variables taking values in a set  $\mathbf{A}$ , and assume that  $\mathbf{H} : \mathbf{A}^n \rightarrow \mathbb{R}$  satisfies

$$\sup_{x_1, \dots, x_n, x'_i \in A} |\mathbf{H}(x_1, \dots, x_n) - \mathbf{H}(x_1, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_n)| \leq c_i$$

for every  $1 \leq i \leq n$ . Then, for every  $t > 0$ ,

$$|\mathbf{P} \cdot \mathbf{H}(\mathbf{X}_1, \dots, \mathbf{X}_n) - \mathbb{E}[\mathbf{H}(\mathbf{X}_1, \dots, \mathbf{X}_n)]| \leq t \exp\left(-\frac{2t^2}{\sum_{i=1}^n c_i^2}\right).$$

As pointed out in Remark 7 Assumptions 7 and 8 imply the random function on  $\mathbf{h}_i(\cdot)$  is  $L$  smooth and thus

$$|\mathbf{h}_i(\mathbf{w}) - \mathbf{h}_i(\mathbf{w}_*)| \leq L \|\mathbf{w} - \mathbf{w}_*\| \leq L \left(\frac{\gamma_k^+}{k}\right)^2.$$

As a result when a random function on  $\mathbf{h}_i$  changes the random variable  $\mathbf{l}(\mathbf{h}_1, \dots, \mathbf{h}_n)$  can change by no more than  $2L \left(\frac{\gamma_k^+}{k}\right)^2/n$ . To see this we have

$$\begin{aligned} & |\mathbf{l}(\mathbf{h}_1, \dots, \mathbf{h}_n) - \mathbf{l}(\mathbf{h}_1, \dots, \mathbf{h}_{i-1}, \mathbf{h}'_i, \mathbf{h}_{i+1}, \dots, \mathbf{h}_n)| \\ & \leq \frac{1}{n} \sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma_k^+} \left| \langle \mathbf{h}'_i(\mathbf{w}) - \mathbf{h}'_i(\mathbf{w}_*), [\mathbf{h}_i(\mathbf{w}) - \mathbf{h}_i(\mathbf{w}_*)], \mathbf{w} - \mathbf{w}_* \rangle \right| \leq \frac{2}{n} L \left(\frac{\gamma_k^+}{k}\right)^2. \end{aligned}$$

McDermid's inequality implies that with probability at least  $1 - \epsilon$

$$|\mathbf{l}(\mathbf{h}_1, \dots, \mathbf{h}_n) - \mathbb{E}[\mathbf{l}(\mathbf{h}_1, \dots, \mathbf{h}_n)]| \leq L \left(\frac{\gamma_k^+}{k}\right)^2 \sqrt{\frac{2}{n} \log \frac{1}{\epsilon}}. \quad \bullet$$

Let  $(\mathbf{h}'_1, \dots, \mathbf{h}'_n)$  be an independent copy of  $(\mathbf{h}_1, \dots, \mathbf{h}_n)$  and  $\mathbf{1}, \dots, \mathbf{n}$  be  $n$ -dimensional indicator variables with equal probability of being 1 using techniques of indicator complexities [Bartlett and Mendelson](#)  $\bullet\bullet$  we bound  $\mathbb{E}[\mathbf{l}(\mathbf{h}_1, \dots, \mathbf{h}_n)]$  as follows:

$$\begin{aligned} & \mathbb{E}_{h_1, \dots, h_n} \left[ \sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma_k^+} \left\langle \mathbf{H}(\mathbf{w}) - \mathbf{H}(\mathbf{w}_*), \frac{1}{n} \sum_{i=1}^n [\mathbf{h}_i(\mathbf{w}) - \mathbf{h}_i(\mathbf{w}_*)], \mathbf{w} - \mathbf{w}_* \right\rangle \right] \\ & = \frac{1}{n} \mathbb{E}_{h_1, \dots, h_n} \left[ \sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma_k^+} \left\langle \mathbf{H}(\mathbf{w}) - \mathbf{H}(\mathbf{w}_*), \sum_{i=1}^n [\mathbf{h}_i(\mathbf{w}) - \mathbf{h}_i(\mathbf{w}_*)], \mathbf{w} - \mathbf{w}_* \right\rangle \right] \\ & = \frac{1}{n} \mathbb{E}_{h_1, \dots, h_n, h'_1, \dots, h'_n} \left[ \sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma_k^+} \left\langle \mathbf{H}(\mathbf{w}) - \mathbf{H}(\mathbf{w}_*), \sum_{i=1}^n [\mathbf{h}'_i(\mathbf{w}) - \mathbf{h}'_i(\mathbf{w}_*)], \mathbf{w} - \mathbf{w}_* \right\rangle \right] \\ & = \frac{1}{n} \mathbb{E}_{h_1, \dots, h_n, h'_1, \dots, h'_n} \left[ \sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma_k^+} \left\langle \mathbf{H}(\mathbf{w}) - \mathbf{H}(\mathbf{w}_*), \sum_{i=1}^n [\mathbf{h}_i(\mathbf{w}) - \mathbf{h}_i(\mathbf{w}_*)], \mathbf{w} - \mathbf{w}_* \right\rangle \right] \end{aligned}$$

$$= \frac{1}{n} \mathbb{E}_{h_1, \dots, h_n, h'_1, \dots, h'_n, \epsilon_1, \dots, \epsilon_n} \left[ \sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma_k^+} \right]$$
$$\sum_i^n$$



Note that  $\mathbf{x}^2$  is 2a Lipschitz over  $[\mathbf{a}, \mathbf{a}]$  and  $\mathbf{p}_i(\mathbf{w}) + \mathbf{q}_i(\mathbf{w}) \in [2\gamma_k^+ \mathbf{D}^-, 2\gamma_k^+ \mathbf{D}^-]$ , then from the comparison theorem of adachi copex t es [Ledoux and a agrand 99](#) in part. cu ar Lemma of [Me r and Zhang 00](#) we have

$$\begin{aligned} & \mathbb{E} \left[ \sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma_k^+} \sum_{i=1}^n i (\mathbf{p}_i(\mathbf{w}) + \mathbf{q}_i(\mathbf{w}))^2 \right] \\ & \leq 4\gamma_k^+ \mathbf{D} \sqrt{\mathbb{E} \left[ \sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma_k^+} \sum_{i=1}^n i (\mathbf{p}_i(\mathbf{w}) + \mathbf{q}_i(\mathbf{w})) \right]} \\ & \leq 4\gamma_k^+ \mathbf{D} \sqrt{\left( \mathbb{E} \left[ \sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma_k^+} \sum_{i=1}^n i \mathbf{p}_i(\mathbf{w}) \right] + \mathbb{E} \left[ \sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma_k^+} \sum_{i=1}^n i \mathbf{q}_i(\mathbf{w}) \right] \right)}. \end{aligned}$$

any we have

$$\begin{aligned} & \mathbb{E} \left[ \sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma_k^+} \sum_{i=1}^n i (\mathbf{p}_i(\mathbf{w}) - \mathbf{q}_i(\mathbf{w}))^2 \right] \\ & \leq 4\gamma_k^+ \mathbf{D} \sqrt{\left( \mathbb{E} \left[ \sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma_k^+} \sum_{i=1}^n i \mathbf{p}_i(\mathbf{w}) \right] + \mathbb{E} \left[ \sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma_k^+} \sum_{i=1}^n i \mathbf{q}_i(\mathbf{w}) \right] \right)}. \end{aligned}$$

Combining and 4 we arrive at

$$\begin{aligned} & \mathbb{E} \left[ \sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma_k^+} \sum_{i=1}^n i \|\mathbf{h}_i(\mathbf{w}) - \mathbf{h}_i(\mathbf{w}_*)\| \|\mathbf{w} - \mathbf{w}_*\| \right] \\ & \leq 2\gamma_k^+ \mathbf{D} \sqrt{\left( \underbrace{\mathbb{E} \left[ \sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma_k^+} \sum_{i=1}^n i \mathbf{p}_i(\mathbf{w}) \right]}_{:=C_1} + \underbrace{\mathbb{E} \left[ \sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma_k^+} \sum_{i=1}^n i \mathbf{q}_i(\mathbf{w}) \right]}_{:=C_2} \right)}. \end{aligned}$$

we proceed to upper bound  $C_1$ . From our definition of  $\mathbf{p}_i(\mathbf{w})$  we have

$$\begin{aligned} \|\mathbf{p}_i(\mathbf{w}) - \mathbf{p}_i(\mathbf{w}')\| &= \frac{1}{\sqrt{2}} \left| \langle \mathbf{w}, \mathbf{x}_i \rangle - \langle \mathbf{w}', \mathbf{x}_i \rangle \right| \\ &= \frac{1}{\sqrt{2}} \|\mathbf{x}_i\| \|\mathbf{w} - \mathbf{w}'\|. \end{aligned}$$

Applying the comparison theorem of adachi copex t es again we have

$$C_1 \leq \sqrt{\mathbb{E} \left[ \sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma_k^+} \sum_{i=1}^n i \|\mathbf{x}_i\| \|\mathbf{w} - \mathbf{w}_*\| \right]} = C_2.$$

