

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

Lijun Zhang

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

ZHANG LIJUN, LAMDA NJ, CHINA

Tianbao Yang

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

YANG TIANBAO, IO, USA

Rong Jin

Alibaba Group, Seattle, USA

JIN RONG, ALIBABA INC, CHINA

Abstract

Although there exist plenty of theoretical results on EM for supervised learning, current theoretical understandings of EM for a related problem, stochastic convex optimization (CO), are limited. In this work, we strengthen the results of EM for CO by exploiting smoothness and strong convexity conditions to prove the risk bounds. First, we establish an $(\frac{1}{n} + \sqrt{\frac{1}{n^3}})$ risk bound when θ

where $\mathcal{X} = \mathbb{R}^d$, $\mathcal{Y} = \mathbb{R}$ is a hypothesis class $(\mathcal{X}, \mathcal{Y})$, \mathcal{X} is an instance space, \mathbb{D} is a distribution on \mathcal{X} and $(\mathcal{Y}, \mathbb{R})$ is a loss function. In this paper we mainly focus on the convex version of online stochastic convex optimization (OCO) where both the domain \mathcal{X} and the expected function $\mathbb{E}[\ell(\cdot)]$ are convex.

Two classical approaches for solving stochastic optimization are stochastic approximation (SA) [Kushner and Yin, 1997] and the sample average approximation (SAA) method which is also referred to as empirical risk minimization (ERM) in the machine learning community [Bartlett and Mendelson, 1999].

While both SA and ERM have been extensively studied in recent years [Bartlett and Mendelson, 2002; Bartlett et al., 2005; Neirovs et al., 2009; Moulines and Bach, 2011; Hazan and Kale, 2009; Ahn et al., 2010; Agarwal et al., 2010; Bach and Moulines, 2011; Zhang et al., 2011; Mahdavi et al., 2011], most theoretical guarantees of ERM are restricted to supervised learning.

As pointed out in a seminal work of Havalivewartz et al. [2009], the success of ERM for supervised learning cannot be directly extended to stochastic optimization. Actually, Havalivewartz et al. [2009] have constructed an instance of CO that is learnable by SA but cannot be solved by ERM. Literature about ERM for stochastic optimization including CO are quite limited and we still lack a full understanding of the theory.

In ERM we are given functions ℓ_1, \dots, ℓ_n sampled independently from \mathbb{P} and a target minimize empirical objective function,

$$\min_{\mathbf{w} \in \mathcal{W}} \widehat{\ell}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{w})$$

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

$$\widehat{\text{excess}}(\widehat{\mathbf{w}}) = \widehat{\ell}(\widehat{\mathbf{w}}) - \min_{\mathbf{w} \in \mathcal{W}} \widehat{\ell}(\mathbf{w})$$

State of the art risk bounds of ERM include, an $\widetilde{O}(\sqrt{d})$ bound when the random function $\ell(\cdot)$ is Lipschitz continuous where d is the dimensionality of \mathbf{w} , an $\widetilde{O}(1)$ bound when $\ell(\cdot)$ is strongly convex [Havalivewartz et al., 2009], and an $\widetilde{O}(1)$ bound when $\ell(\cdot)$ is exponentially concave [Mehta et al., 2011]. From existing studies of ERM for supervised learning [Lepoint et al., 2011] we now 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 ERM for CO. This paper provides an affirmative answer to this question. Indeed, we propose a general approach for analyzing the excess risk of ERM which brings several improved risk bounds and new risk bounds as we

state our results we first introduce some notations. Let $\ell_* = \min_{\mathbf{w} \in \mathcal{W}} \ell(\mathbf{w})$ be the minimum risk, β be the modulus of strong convexity of $\ell(\cdot)$ and γ be the modulus of smoothness of $\ell(\cdot)$. Denote by $\kappa = \beta/\gamma$ the condition number of the problem. Our and previous results of ERM for CO are summarized in Table 1 where we also expect the assumptions on the random function $\ell(\cdot)$, the empirical function $\widehat{\ell}(\mathbf{w})$ and the expected function $\mathbb{E}[\ell(\cdot)]$. For our results of ERM for CO we assume the domain is bounded and the random function is nonnegative. The high light the significance of this work as follows.

¹we use the \widetilde{O} and $\widetilde{\Omega}$ notations to hide constant factors as well as polynomial factors in d and n .

able, utility of Excess $\mathbb{E}M$ Bounds of $\mathbb{E}M$ for COA 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, sc smooth, s exponentially concave, exp

	(\cdot)	$\widehat{(\cdot)}$	(\cdot)	$\mathbb{E}M$ Bounds
Haewhartz et al [9]	Lp			$\sim(\sqrt{\frac{d}{n}})$ $(\frac{1}{\lambda n})^*$
Mehta [10]	exp, Lp, b			$\sim(\frac{d}{\eta n})$
Hsiung	nn, c, s		Lp	$\sim(\frac{d}{n} + \sqrt{\frac{F_*}{n}})$
	nn, c, s		Lp, sc	$\sim(\frac{d}{n} + \frac{\kappa F_*}{n})$ $(\frac{1}{\lambda n^2} + \frac{\kappa F_*}{n})$ when $\kappa = \tilde{\Omega}(\cdot)$
	nn, s	c	sc	$\sim(\frac{\kappa d}{n} + \frac{\kappa F_*}{n}) = \sim(\frac{\kappa d}{n})$ $(\frac{1}{\lambda n^2} + \frac{\kappa F_*}{n})$ when $\kappa = \tilde{\Omega}(\cdot^2)$
	nn, s, g	c	sc	$(\frac{\kappa}{n} + \frac{\kappa F_*}{n}) = (\frac{\kappa}{n})$ $(\frac{1}{\lambda n^2} + \frac{\kappa F_*}{n})$ when $\kappa = \Omega(\cdot^2)$

When (\cdot) is both convex and smooth and (\cdot) is Lipschitz continuous we establish an $(\frac{d}{n} + \sqrt{\frac{F_*}{n}})$ risk bound cf Hsiung. In the optimistic case that κ is a $e^{-\kappa} = (\cdot^2)$, we obtain an $\sim(\cdot)$ risk bound which is analogous to the $\sim(1)$ optimistic rate of $\mathbb{E}M$ for supervised learning [11].

If (\cdot) is a so strongly convex we prove an $(\frac{d}{n} + \frac{\kappa F_*}{n})$ risk bound and prove that $(1[\cdot^2] + \frac{\kappa F_*}{n})$ when $\kappa = \tilde{\Omega}(\cdot)$ cf Hsiung. Thus for large κ is a $e^{-\kappa} = (1)$ we get an (\cdot^2) risk bound which to the best of our knowledge is the first (1^2) type of risk bound of $\mathbb{E}M$.

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

2. Related Work

In this section we give a brief introduction to previous work on $\mathbb{E}M$.

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 (ℓ) is Lipschitz continuous. [Haviv and Hertz](#) [9] demonstrate that \hat{w}_n converges to w^* uniformly over \mathcal{W} with an $\tilde{O}(\sqrt{\frac{1}{n}})$ error bound that holds with high probability, implying an $\tilde{O}(\sqrt{\frac{1}{n}})$ rate bound of ERM. They further establish an $(1 - \epsilon)$ rate bound of ERM that holds in expectation when (ℓ) is strongly convex and Lipschitz continuous. Stochastic optimization with exp-concave functions is studied recently [Koren and Levy](#) and [Mehta](#) [10]. [Mehta](#) [10] proves an $\tilde{O}(\frac{1}{\sqrt{n}})$ bound of ERM that holds with high probability when (ℓ) is exp-concave Lipschitz continuous and bounded. Lower bounds of ERM for stochastic optimization is investigated by [Fedan](#) [11] who exhibits a lower bound of $\Omega(\frac{1}{\sqrt{n}})$ sample complexity for uniform convergence that nearly matches the upper bound of [Haviv and Hertz](#) [9], and a lower bound of $\Omega(\frac{1}{\sqrt{n}})$ sample complexity of ERM which is matched by our $\tilde{O}(\frac{1}{\sqrt{n}} + \sqrt{\frac{1}{n}})$ bound when \mathcal{W} is a

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 $\hat{\ell}_n$ to ℓ over the hypothesis class. [Kochenc](#) [12]. In fact uniform convergence is a sufficient condition for learnability [Haviv, Hertz and Ben-David](#) [4] and in some special cases such as binary classification it is also a necessary condition [Agn](#) [9]. 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{W} has finite VC dimension the excess risk can be upper bounded by $\tilde{O}(\sqrt{\frac{VC(\mathcal{W})}{n}})$ where $VC(\mathcal{W})$ is the VC dimension of \mathcal{W} . If the loss (ℓ) is Lipschitz continuous with respect to its first argument we have a rate bound of $(1 - \epsilon)\sqrt{\frac{1}{n}} + \frac{1}{n} \ln(\frac{1}{\epsilon})$ where $\frac{1}{n} \ln(\frac{1}{\epsilon})$ is the Rademacher complexity of \mathcal{W} . The Rademacher complexity typically scales as $\frac{1}{\sqrt{n}}$ e.g. contains linear functions with bounded norm implying an $(1 - \epsilon)\sqrt{\frac{1}{n}}$ rate bound [Bartlett and Mendelson](#) [13]. There have been intensive efforts to derive rates faster than $(1 - \epsilon)\sqrt{\frac{1}{n}}$ under various conditions [Lee et al.](#) [9], [Panchenko](#) [14], [Bartlett et al.](#) [15], [Gonen and Haviv and Hertz](#) [16] such as smoothness [rebro et al.](#) [17] strong convexity [r dharan et al.](#) [9] to name a few amongst any peculiarity when the random function (ℓ) is nonnegative and smooth [rebro et al.](#) [17] have established a rate bound of $\tilde{O}(\frac{1}{n} + \frac{1}{\sqrt{n}})$ reducing to an $\tilde{O}(\frac{1}{\sqrt{n}})$ bound if $\frac{1}{n} \ln(\frac{1}{\epsilon}) = \frac{1}{\sqrt{n}}$ and $\epsilon = (1 - \epsilon)$. A generalized near-optimal bound is studied by [r dharan et al.](#) [9] and a rate bound of $(1 - \epsilon)\sqrt{\frac{1}{n}}$ is proved if the expected function (ℓ) is 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

the excess risk bounds for a regularized empirical 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 R , that is,

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

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

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

Assumption 3 The expected function $\bar{\ell}(\cdot)$ is L -Lipschitz continuous over \mathcal{W} , that is,

$$\|\bar{\ell}(\mathbf{w}) - \bar{\ell}(\mathbf{w}')\| \leq L \|\mathbf{w} - \mathbf{w}'\| \quad \forall \mathbf{w}, \mathbf{w}' \in \mathcal{W} \quad \cdot$$

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

- (a) The expected function $\bar{\ell}(\cdot)$ is convex over \mathcal{W} .
- (b) The expected function $\bar{\ell}(\cdot)$ is μ -strongly convex over \mathcal{W} , that is,

$$\bar{\ell}(\mathbf{w}) + \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}'\|^2 \leq \bar{\ell}(\mathbf{w}') \quad \forall \mathbf{w}, \mathbf{w}' \in \mathcal{W} \quad \nabla$$

- (c) The empirical function $\hat{\ell}(\cdot)$ is convex.
- (d) The random function $\ell(\cdot)$ is convex.

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

$$\|\nabla \ell(\mathbf{w}_*)\| \leq G \quad \mathbb{P} \quad 8$$

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 $\bar{\ell}(\cdot)$ implies the expected function $\bar{\ell}(\cdot)$ is smooth. By Jensen's inequality we have

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

and the empirical function $\hat{\ell}(\cdot)$ is also smooth. The condition number of $\bar{\ell}(\cdot)$ is defined as the ratio between L and $\mu = \frac{L}{\kappa} \geq 1$.

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

$$n(\epsilon, \delta) = 2 \left(\log \frac{2}{\epsilon} + \log \frac{6}{\delta} \right) \quad 9$$

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

$$\begin{aligned} & \mathbb{E} \left[\ell(\hat{\mathbf{w}}) - \ell(\mathbf{w}_*) \right] \\ & \leq \frac{16 \beta^2}{\beta} \left(\frac{\beta}{\beta} \right) + \frac{8 \beta \log(2 \beta)}{\beta} + 8 \sqrt{\frac{2 \beta^* \log(2 \beta)}{\beta}} + \left(8 \beta + \beta + \frac{4 \beta}{\beta} \left(\frac{\beta}{\beta} \right) \right) \end{aligned}$$

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

By choosing β large enough the last term in (10) that contains β becomes non dominating so by specifying we have the following corollary

Corollary 2 By setting $\beta = 1$ in Theorem 1, we have $\mathbb{E} \left[\ell(\hat{\mathbf{w}}) \right] = 2(\log(2 \beta) + \log(6 \beta))$, and with high probability

$$\mathbb{E} \left[\ell(\hat{\mathbf{w}}) - \ell(\mathbf{w}_*) \right] = \left(\frac{\log \beta}{\beta} + \sqrt{\frac{\beta^*}{\beta}} \right) = \tilde{\mathcal{O}} \left(\frac{1}{\beta} + \sqrt{\frac{\beta^*}{\beta}} \right)$$

Remark 2 The above corollary implies that under the smoothness and other common assumptions $\mathbb{E} \left[\ell(\hat{\mathbf{w}}) \right] \leq \tilde{\mathcal{O}} \left(\frac{1}{\beta} + \sqrt{\frac{\beta^*}{\beta}} \right)$ is a bound for CO then the naturalness assumption $\beta^* = \mathcal{O}(\beta^2)$ the rate is proved to $\tilde{\mathcal{O}} \left(\frac{1}{\beta} \right)$. Note that even under the smoothness assumption the near dependence on β is unavoidable. Fed and Leventheore 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

$$\begin{aligned} & \mathbb{E} \left[\ell(\hat{\mathbf{w}}) - \ell(\mathbf{w}_*) \right] \\ & \leq \frac{16 \beta^2}{\beta} \left(\frac{\beta}{\beta} \right) + \frac{8 \beta \log(2 \beta)}{\beta} + \frac{8 \beta^* \log(2 \beta)}{\beta} + \left(8 \beta + \beta + \frac{4 \beta}{\beta} \left(\frac{\beta}{\beta} \right) \right) \end{aligned}$$

Furthermore, if

$$\frac{4 \beta}{\beta} \left(\frac{\beta}{\beta} \right) = 4 \left(\frac{\beta}{\beta} \right)$$

we also have

$$\mathbb{E} \left[\ell(\hat{\mathbf{w}}) - \ell(\mathbf{w}_*) \right] \leq \frac{32 \beta^2 \log^2(2 \beta)}{2} + \frac{128 \beta^* \log(2 \beta)}{\beta} + \left(\frac{128 \beta^2}{\beta} + 16 \beta + 4 \beta^2 \right)$$

The above theorem can be simplified by choosing different values of

Corollary 4 By setting $\beta = 1$ in Theorem 3, we have $\mathbb{E} \left[\ell(\hat{\mathbf{w}}) \right] = 2(\log(2 \beta) + \log(6 \beta))$, and with high probability

$$\mathbb{E} \left[\ell(\hat{\mathbf{w}}) - \ell(\mathbf{w}_*) \right] = \left(\frac{\log \beta}{\beta} + \frac{\beta^*}{\beta} \right) = \tilde{\mathcal{O}} \left(\frac{1}{\beta} + \frac{\beta^*}{\beta} \right)$$

By setting $\beta = 1 - \beta^2$, we have $\mathbb{E} \left[\ell(\hat{\mathbf{w}}) \right] = 2(\log(2 \beta) + \log(6 \beta^2))$ and when $\beta = \Omega(\log \beta) = \tilde{\Omega}(\log \beta)$, with high probability

$$\mathbb{E} \left[\ell(\hat{\mathbf{w}}) - \ell(\mathbf{w}_*) \right] = \left(\frac{1}{2} + \frac{\beta^*}{\beta} \right)$$

Remark 3 The first part of Corollary 4 shows that $\mathbb{E} \| \sum_{k=1}^n X_k \|^2$ enjoys an $\tilde{O}(\dots + \dots)$ rate bound for stochastic optimization of strongly convex and smooth functions. In the literature the best comparable result is the (1.1) rate bound proved by [Haviv, Hertz, and Shalev-Shvartz \(2019\)](#) but with strong differences highlighted in [Haviv, Hertz, and Shalev-Shvartz \(2019\)](#) since the rate bound of [Haviv, Hertz, and Shalev-Shvartz \(2019\)](#) 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 3 is increased by a factor of $\sqrt{2}$. It is interesting to note that a similar phenomenon also happens in stochastic approximation recently a variance reduction technique called [Johnson and Zhang](#) or EMGD [Zhang et al](#) was proposed for stochastic optimization when both full gradients and stochastic gradients are available. In the analysis [3](#) 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 $(1 + \sqrt{2})$ risk bound for supervised learning. However a drawback of these theorems is that the lower bound of 3 depends on the dimensionality and thus cannot be applied to nonfinite dimensional cases e.g. kernel methods [Chopra and Karamchandani](#). In this section we exploit the structure of supervised learning to achieve the theory dimensionality independent. We focus on the generalized linear form of supervised learning.

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

where $\ell(\mathbf{w}; \mathbf{x}, y)$ is the loss of predicting y given \mathbf{x} when the true target is y and $\mathcal{R}(\mathbf{w})$ is a regularizer. Given training examples $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$ independent and identically distributed from \mathbb{D} the empirical risk is

$$\min_{\mathbf{w} \in \mathcal{W}} \hat{\mathcal{R}}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}; \mathbf{x}_i, y_i) + \mathcal{R}(\mathbf{w})$$

define

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

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

$$\|\nabla \ell(\mathbf{w}_*; \mathbf{x})\| \leq G \|\mathbf{x}\| \quad \mathbb{D}$$

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

$$\begin{aligned} \|\nabla \ell(\mathbf{w}; \mathbf{x}) - \nabla \ell(\mathbf{w}'; \mathbf{x})\| &= \|\nabla' \ell(\mathbf{w}; \mathbf{x}) - \nabla' \ell(\mathbf{w}'; \mathbf{x})\| \\ &\leq L \|\mathbf{w} - \mathbf{w}'\| \leq L \|\mathbf{w} - \mathbf{w}'\| \leq L \|\mathbf{w} - \mathbf{w}'\| \end{aligned}$$

By Jensen's inequality (1) is also L -smooth. Notice that L is the modulus of smoothness of (1) and is the modulus of strong convexity of (1) (a slight abuse of notation we define $L = 2$ and the condition number as the ratio between L and $\mu = \dots$). Finally we note that the regularizer (1) could be *non-smooth*.

We have the following excess risk bound of $E_{\mathcal{M}}$ for supervised learning

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

$$\begin{aligned} &= 4 \left(8 + \sqrt{2 \log \frac{2 \log_2(\cdot) + \log_2(2 \cdot)}{\epsilon}} \right) \\ &_* = (\mathbf{w}_*) = (\mathbf{w}_*) = (\mathbf{w}_*) \end{aligned} \quad \mathbf{4}$$

Under **Assumptions 4(b), 4(c), 6, 7, 8, 9, and 10** with probability at least $1 - 2\epsilon$, we have

$$(\widehat{\mathbf{w}}) \quad \ell(\widehat{\mathbf{w}}) - \ell(\mathbf{w}_*) \leq \max \left(\frac{\epsilon}{2} + \frac{4}{2^4} \frac{4 \cdot 2 \cdot 2 \cdot 2}{\epsilon} + \frac{4 \log(2 \cdot)}{\epsilon} + \frac{8 \cdot \log(2 \cdot)}{\epsilon} \right)$$

Furthermore, if

$$\frac{16 \cdot 2 \cdot 2}{2} = 16 \cdot 2 \cdot 2$$

with probability at least $1 - 2\epsilon$, we have

$$(\widehat{\mathbf{w}}) \quad \ell(\widehat{\mathbf{w}}) - \ell(\mathbf{w}_*) \leq \max \left(\frac{\epsilon}{2} + \frac{8}{2^4} \frac{2 \log^2(2 \cdot)}{\epsilon} + \frac{16 \cdot \log(2 \cdot)}{\epsilon} \right) \quad \mathbf{v}$$

Remark 8 The first part of Theorem 7 presents an (\cdot) risk bound similar to the (1.1) risk bound of [Rohrig et al. 2019](#). The second part is an $(1/\epsilon^2 + \dots)$ risk bound and in this case the lower bound of $\Omega(\epsilon^2)$ which is data independent. Thus Theorem 7 can be applied even when the data is non-stationary. Generally speaking the regularizer (1) is nonnegative and thus $\ell(\mathbf{w}_*) \leq \ell(\widehat{\mathbf{w}})$ so the second bound is even better than those in Theorem 7 and

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 (1) by a factor of ϵ . We investigate this question as a future work.

For brevity we treat C as a constant because t only has a *double* logarithmic dependence on n .

4. Analysis

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

4.1. The Key Idea

By the convexity of $\hat{f}(\cdot)$ and the optimality condition of \hat{w} (Boyd and Vandenberghe [4]), we have

(1) AND (1²) TYPE OF K BOND OF E M

Lemma 1 *Under Assumptions 2 and 4(d), with probability at least 1 - ϵ , $w_{024}(f_{083399}(b)0.0r.841(w4)0.510254(n)0$*

where the last step is due to

$$\|\widehat{\mathbf{w}} - \mathbf{w}_*\| \sqrt{\frac{(\cdot)(\widehat{\mathbf{w}} - \mathbf{w}_*)}{2}} \leq \frac{(\cdot)\widehat{\mathbf{w}} - \mathbf{w}_*^2}{2} + \frac{(\widehat{\mathbf{w}} - \mathbf{w}_*)}{2}$$

$$\|\widehat{\mathbf{w}} - \mathbf{w}_*\| \sqrt{\frac{(\cdot)}{2}} \leq \frac{(\cdot)\widehat{\mathbf{w}} - \mathbf{w}_*^2}{2} + \frac{(\cdot)}{2}$$

From 4 we get

$$\frac{1}{2}(\widehat{\mathbf{w}} - \mathbf{w}_*)$$

$$\leq \frac{2(\cdot)\widehat{\mathbf{w}} - \mathbf{w}_*^2}{2} + \frac{2 \log(2\cdot)\widehat{\mathbf{w}} - \mathbf{w}_*}{2} + \|\widehat{\mathbf{w}} - \mathbf{w}_*\| \sqrt{8\cdot \log(2\cdot)}$$

$$+ 2\|\widehat{\mathbf{w}} - \mathbf{w}_*\| + \frac{(\cdot)\widehat{\mathbf{w}} - \mathbf{w}_*}{2}$$

$$\leq \frac{8\cdot}{2} + \frac{4 \log(2\cdot)}{2} + 4\sqrt{2\cdot \log(2\cdot)} + \left(4\cdot + \frac{2\cdot}{2} + \frac{(\cdot)}{2}\right)$$

which implies

5. Conclusions and Future Work

In this paper we study the excess risk of $E_{\mathcal{M}}$ for CO. Our theoretical results show that it is possible to achieve $(1 - \epsilon)$ type of risk bounds under the smoothness and strong convexity conditions. The first part of the paper and the second part of the paper are exciting results that when s is large enough $E_{\mathcal{M}}$ has $(1 - \epsilon^2)$ type of risk bounds under the smoothness, strong convexity and strong convexity conditions.

In the context of CO there remain any open problems about $E_{\mathcal{M}}$.

Our current results are restricted to the Hilbert or Euclidean space because the smoothness and strong convexity are defined in terms of the L_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. $\tilde{O}(\cdot)$ or $\tilde{O}(1 - \epsilon)$ 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 ϵ in the second part of the paper. Ideally we expect that $\epsilon = \Omega(\cdot)$ is sufficient to derive an $(1 - \epsilon^2)$ risk bound.

4 The $(1 - \epsilon^2)$ 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.

Finally as far as we know there are no $(1 - \epsilon^2)$ 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 88 II 4 99 and the Collaborative Innovation Center of Novel Software Technology and Industrialization of Nanjing University.

References

Amit Agarwal, Peter L. Bartlett, Pradeep Ravuru, and Martin J. Wainwright. Information theoretic lower bounds on the oracle complexity of stochastic convex optimization. *IEEE Transactions on Information Theory*, 49, 2002.

François Bach and Ercole Moulines. Non strongly convex smooth stochastic approximation with convergence rate (1.1). In *Advances in Neural Information Processing Systems 26*, pages 49–57, 2013.

Peter L. Bartlett and Sahar Mendelson. Rademacher and gaussian complexities, bounds and structural results. *Journal of Machine Learning Research*, 4, 2003.

Peter L. Bartlett, Olivier Bousquet, and Sahar Mendelson. Local rademacher complexities. *The Annals of Statistics*, 4, 49, 2001.

Stephen Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.

Gautam Desai, Mehryar Mohr, and Arvind. Learning with deep cascades. In *Proceedings of the 26th International Conference on Algorithmic Learning Theory*, pages 4–19, 2015.

Yael Feder. Generalization of federated stochastic convex optimization. The dependence on the number of clients. *ArXiv e-prints*, arXiv:1404.4844, 2014.

Amnon Gonen and Elad Hazan. Average stability is invariant to data preconditioning. Preprint, 2014. *ArXiv e-prints*, arXiv:1404.4844, 2014.

Elad Hazan and Elad Hazan. Beyond the regret minimization barrier: an optimal algorithm for stochastic strong convex optimization. In *Proceedings of the 24th Annual Conference on Learning Theory*, pages 4–14, 2013.

Elad Hazan, Amit Agarwal, and Elad Hazan. Logarithmic regret algorithms for online convex optimization. *Machine Learning*, 9, 9, 9, 2013.

John D. Lee and Dong Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in Neural Information Processing Systems 26*, pages 49–57, 2013.

Radu Kojima. *Oracle Inequalities in Empirical Risk Minimization and Sparse Recovery Problems*. Springer, 2011.

Ofer Koren and Kfir Levy. Fast rates for exp-concave empirical risk minimization. In *Advances in Neural Information Processing Systems 28*, pages 49–57, 2015.

Harold J. Kushner and George Yin. *Stochastic Approximation and Recursive Algorithms and Applications*. Springer, 2003. ISBN 0-387-95314-2. <http://www.springer.com>

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*

Colin McDiarmid On the method of bounded differences In *Surveys in Combinatorics* pages 41–59

Nishant A. Mehta Fast rates with high probability in exponentially concave statistical learning *ArXiv e-prints* arXiv:1808.08899

Yong Ma and Yong Zhang Generalization error bounds for bayesian mixture algorithms *Journal of Machine Learning Research* 4:9–18

Eric Moulines and Francis Bach Non asymptotic analysis of stochastic approximation algorithms for machine learning In *Advances in Neural Information Processing Systems* 24 pages 44–52

Andrei Nemirovski, Anatoly Juditsky, G. Lan and Alexander Shapiro Robust stochastic approximation approach to stochastic programming *SIAM Journal on Optimization* 9(4):949–969

Yurii Nesterov *Introductory lectures on convex optimization: a basic course* volume 8 of *Applied optimization* Kluwer Academic Publishers 4

Dmitry Panchenko On extensions of an inequality of Vapnik and Chervonenkis *Electronic Communications in Probability*

Göran Persson *The volume of convex bodies and Banach space geometry* Cambridge Tracts in Mathematics No 94 Cambridge University Press 99

Yanyan Pan and Yoann Berthod One-bit compressed sensing by linear programming *Communications on Pure and Applied Mathematics* 62(9):1321–1347

Alexander A. Ahn, Ohad Hazan and Karthik Sridharan Managing gradient descent optimization for strongly convex stochastic optimization In *Proceedings of the 29th International Conference on Machine Learning* pages 449–457

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

John Haughey, Elad Hazan and Ben-David Elad *Understanding Machine Learning: From Theory to Algorithms* Cambridge University Press 4

Elad Hazan, Elad Hazan, Ohad Hazan, Nathan Srebro and Karthik Sridharan Stochastic convex optimization In *Proceedings of the 22nd Annual Conference on Learning Theory* 9

Alexander Shapiro, Darina Dentcheva and Andrzej Ruszczyński *Lectures on Stochastic Programming: Modeling and Theory* IAIM second edition 4

Teodoro S. G. and Dong Xuan Zhou Learning theory estimates via integral operators and the approximation *Constructive Approximation*

Nathanrebro Karthir dharan and Abubakar Optimal rates for learning with a smooth loss *ArXiv e-prints* arXiv: 1909.09091

Karthir dharan Haim Haevshwartz and Nathanrebro Fast rates for regularized objectives *In Advances in Neural Information Processing Systems 21* pages 4999

Alexandre Boussoffiyev Optimal aggregation of classifiers in statistical learning *The Annals of Statistics* 40(4) 1201-1234

Andrius Rapin *The Nature of Statistical Learning Theory* Springer second edition

Andrius Rapin *Statistical Learning Theory* Wiley InterScience 1998

Lun Zhang Mehrdad Mahdavi and Yong-Jin Linear convergence with condition number independent access of full gradients *In Advance in Neural Information Processing Systems 26* pages 989-998

Lun Zhang Anbao Yang Yong-Jin and Xiaofei He (log) projections for stochastic optimization of smooth and strongly convex functions *In Proceedings of the 30th International Conference on Machine Learning* 116-124

Appendix A. Proof of Lemma 1

we introduce Lemma of [Lai and Zhou](#)

Lemma 3 Let \mathcal{H} be a Hilbert space and let ξ be a random variable with values in \mathcal{H} . Assume $\|\xi\| \leq L$ almost surely. Denote $\mathbb{E}[\|\xi\|^2] = \sigma^2$. Let ξ_1, \dots, ξ_m be m independent draws of ξ . For any $0 < \delta < 1$, with confidence $1 - \delta$,

$$\left\| \frac{1}{m} \sum_{i=1}^m \langle \xi_i, \mathbb{E}[\xi] \rangle \right\| \leq \frac{2L \log(2/\delta)}{m} + \sqrt{\frac{2\sigma^2 \log(2/\delta)}{m}}$$

First consider a fixed w (1) since $i(\cdot)$ is smooth we have

$$\langle \nabla i(w) - \nabla i(w_*), w - w_* \rangle \leq \frac{L}{2} \|w - w_*\|^2$$

Because $i(\cdot)$ is both convex and smooth by (2) of [Nesterov](#) we have

$$\langle \nabla i(w) - \nabla i(w_*), w - w_* \rangle^2 \leq L \langle \nabla i(w) - \nabla i(w_*), w - w_* \rangle$$

taking expectation over both sides we have

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

where the last inequality follows from the optimality condition of w_* i.e.

$$\langle \nabla i(w_*) - \nabla i(w_*), w - w_* \rangle \leq 0$$

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

$$\begin{aligned} & \left\| \mathcal{L}(\mathbf{w}) - \mathcal{L}(\mathbf{w}_*) - [\mathcal{L}'(\widehat{\mathbf{w}}) - \mathcal{L}'(\widehat{\mathbf{w}}_*)] \right\| \\ &= \left\| \mathcal{L}(\mathbf{w}) - \mathcal{L}(\mathbf{w}_*) - \frac{1}{n} \sum_{i=1}^n [\ell_i(\mathbf{w}) - \ell_i(\mathbf{w}_*)] \right\| \\ &\leq \frac{2 \|\mathbf{w} - \mathbf{w}_*\| \log(2/\delta)}{\sqrt{n}} + \sqrt{\frac{2 \|\mathcal{L}'(\widehat{\mathbf{w}}) - \mathcal{L}'(\widehat{\mathbf{w}}_*)\| \log(2/\delta)}{n}} \end{aligned}$$

We obtain Lemma by taking the union bound over all $\mathbf{w} \in \mathcal{B}(\widehat{\mathbf{w}}_*)$. To this end we need an upper bound of the covering number $N(\cdot, \|\cdot\|, \epsilon)$.

Let \mathcal{B} be an unit ball of dimension d and $(\mathcal{B}, \|\cdot\|)$ be its metric space with finite cardinality. According to a standard volume comparison argument (Pisier 1989) we have

$$\log N(\mathcal{B}, \|\cdot\|, \epsilon) \leq \log \frac{3^d}{\epsilon^d}$$

Let $\mathcal{B}(\widehat{\mathbf{w}}_*)$ be a ball centered at origin with radius $\|\widehat{\mathbf{w}}_*\|$. Since we assume $\mathcal{B}(\widehat{\mathbf{w}}_*) \subset \mathcal{B}(\widehat{\mathbf{w}}_*)$ it follows that

$$\log N(\mathcal{B}(\widehat{\mathbf{w}}_*), \|\cdot\|, \epsilon) \leq \log \left| \mathcal{B}(\widehat{\mathbf{w}}_*) \left(\frac{\epsilon}{2\|\widehat{\mathbf{w}}_*\|} \right) \right| \leq \log \frac{6^d}{\epsilon^d}$$

where the first inequality is because the covering numbers are almost increasing by inclusion (Pisier and Vershynin 2009).

Appendix B. Proof of Lemma 2

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

$$\ell_i(\mathbf{w}_*)^2 \leq 4 \ell_i(\mathbf{w}_*)$$

and thus

$$E[\ell_i(\mathbf{w}_*)^2] \leq 4 E[\ell_i(\mathbf{w}_*)] = 4 \sigma_*^2$$

From Assumption 5 we have $\ell_i(\mathbf{w}_*) \leq \sigma_*^2$. Then according to Lemma with probability at least $1 - \delta$ we have

$$\left\| \mathcal{L}(\mathbf{w}_*) - \widehat{\mathcal{L}}(\mathbf{w}_*) \right\| = \left\| \mathcal{L}(\mathbf{w}_*) - \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{w}_*) \right\| \leq \frac{2 \sigma_*^2 \log(2/\delta)}{\sqrt{n}} + \sqrt{\frac{8 \sigma_*^2 \log(2/\delta)}{n}}$$

Appendix C. Proof of Theorem 3

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

$$\begin{aligned} & \left\| \mathcal{L}(\widehat{\mathbf{w}}) - \mathcal{L}(\mathbf{w}_*) + \frac{1}{2} \|\widehat{\mathbf{w}} - \mathbf{w}_*\|^2 \right\| \\ &\leq \left(\underbrace{\left\| \mathcal{L}(\widehat{\mathbf{w}}) - \mathcal{L}(\mathbf{w}_*) - [\mathcal{L}'(\widehat{\mathbf{w}}) - \mathcal{L}'(\widehat{\mathbf{w}}_*)] \right\|}_{:=A_1} + \underbrace{\left\| \mathcal{L}(\mathbf{w}_*) - \widehat{\mathcal{L}}(\mathbf{w}_*) \right\|}_{:=A_2} \right) \|\widehat{\mathbf{w}} - \mathbf{w}_*\| \end{aligned} \quad \mathfrak{g}$$

substitution and into 8 with probability at least $1 - 2^{-k}$ we have

$$\begin{aligned} & (\hat{w}) - (w_*) + \frac{1}{2} \hat{w} w_*^2 \\ & \leq \frac{(\cdot) \hat{w} w_*^2}{2} + \hat{w} w_* \sqrt{\frac{(\cdot)(\hat{w}) - (w_*)}{2}} \\ & + \frac{2 \log(2 \cdot) \hat{w} w_*}{2} + \hat{w} w_* \sqrt{\frac{8 \cdot \log(2 \cdot)}{2}} \\ & + 2 \hat{w} w_* + \hat{w} w_* \sqrt{\frac{(\cdot)}{2}} + \frac{(\cdot) \hat{w} w_*}{2} \end{aligned}$$

9

to prove we substitute and

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

into 9 and then obtain

$$\begin{aligned} & \frac{1}{2} (\hat{w}) - (w_*) \\ & \leq \frac{2 (\cdot) \hat{w} w_*^2}{2} + \frac{2 \log(2 \cdot) \hat{w} w_*}{2} + \frac{4 \cdot \log(2 \cdot)}{2} \\ & + 2 \hat{w} w_* + \frac{(\cdot) \hat{w} w_*}{2} \\ & \leq \frac{8 \cdot \log(2 \cdot)}{2} + \frac{4 \cdot \log(2 \cdot)}{2} + \frac{4 \cdot \log(2 \cdot)}{2} + \left(4 \cdot + \frac{1}{2} + \frac{2 (\cdot)}{2} \right) \end{aligned}$$

which gives

to prove we substitute

$$\begin{aligned} & \hat{w} w_* \sqrt{\frac{(\cdot)(\hat{w}) - (w_*)}{2}} < \frac{2 (\cdot)(\hat{w}) - (w_*)}{2} + \frac{1}{8} \hat{w} w_*^2 \\ & \frac{2 \log(2 \cdot) \hat{w} w_*}{2} < \frac{16 \cdot \log^2(2 \cdot)}{2} + \frac{1}{16} \hat{w} w_*^2 \\ & \hat{w} w_* \sqrt{\frac{8 \cdot \log(2 \cdot)}{2}} < \frac{64 \cdot \log(2 \cdot)}{32} + \frac{1}{32} \hat{w} w_*^2 \\ & 2 \hat{w} w_* < \frac{64 \cdot 2^2}{64} + \frac{1}{64} \hat{w} w_*^2 \\ & \hat{w} w_* \sqrt{\frac{(\cdot)}{2}} < \frac{32 (\cdot)}{128} + \frac{1}{128} \hat{w} w_*^2 \\ & \frac{(\cdot) \hat{w} w_*}{2} < \frac{32 \cdot 2^2 (\cdot)^2}{2} + \frac{1}{128} \hat{w} w_*^2 \end{aligned}$$

into (9) and then obtain

$$\begin{aligned}
 & \frac{(\widehat{\mathbf{w}}) - (\mathbf{w}_*)}{4} \|\widehat{\mathbf{w}} - \mathbf{w}_*\|^2 \\
 & \leq \frac{(\cdot) \|\widehat{\mathbf{w}} - \mathbf{w}_*\|^2}{4} + \frac{2(\cdot) \left(\frac{(\widehat{\mathbf{w}}) - (\mathbf{w}_*)}{4} \right) + \frac{16 \cdot \log^2(2 \cdot)}{2} + \frac{64 \cdot \log(2 \cdot)}{2}}{2} \\
 & \quad + \frac{64 \cdot \log^2(2 \cdot)}{2} + \frac{32(\cdot)}{2} + \frac{32 \cdot \log^2(2 \cdot)}{2} \\
 & \leq \frac{(\cdot) \|\widehat{\mathbf{w}} - \mathbf{w}_*\|^2}{4} + \frac{1}{2} \left(\frac{(\widehat{\mathbf{w}}) - (\mathbf{w}_*)}{4} \right) + \frac{16 \cdot \log^2(2 \cdot)}{2} + \frac{64 \cdot \log(2 \cdot)}{2} \\
 & \quad + \frac{64 \cdot \log^2(2 \cdot)}{2} + 8 \cdot \log^2(2 \cdot) + 2 \cdot \log^2(2 \cdot)
 \end{aligned}$$

which implies

Appendix D. Proof of Theorem 5

Without 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 $\mathbf{w} \in \mathcal{W}$, we have

$$\left\| \frac{(\widehat{\mathbf{w}}) - (\mathbf{w}_*)}{4} - \left[\frac{(\widehat{\mathbf{w}}) - (\mathbf{w}_*)}{4} \right] \right\| \leq \frac{(\cdot) \|\mathbf{w} - \mathbf{w}_*\|}{4} + \frac{1}{4} \|\mathbf{w} - \mathbf{w}_*\| \sqrt{\frac{(\cdot)}{4}}$$

where (\cdot) is defined in (9).

The above lemma is a direct consequence of Lemma 4 and the union 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

$$\begin{aligned}
 & \left\| \frac{(\widehat{\mathbf{w}}) - (\mathbf{w}_*)}{4} - \left[\frac{(\widehat{\mathbf{w}}) - (\mathbf{w}_*)}{4} \right] \right\| \\
 & \leq \frac{(\cdot) \|\widehat{\mathbf{w}} - \mathbf{w}_*\|}{4} + \frac{1}{4} \|\widehat{\mathbf{w}} - \mathbf{w}_*\| \sqrt{\frac{(\cdot)}{4}} + 2 \tag{4} \\
 & \leq \frac{(\cdot) \|\widehat{\mathbf{w}} - \mathbf{w}_*\|}{4} + \frac{1}{4} \|\widehat{\mathbf{w}} - \mathbf{w}_*\| \sqrt{\frac{(\cdot)}{4}} + \frac{(\cdot)}{4} + \frac{1}{4} \sqrt{\frac{(\cdot)}{4}} + 2
 \end{aligned}$$

Substituting (4) and (8) into (8) with probability at least $1 - 2\delta$, we have

$$\begin{aligned}
 & \frac{(\widehat{\mathbf{w}}) - (\mathbf{w}_*)}{2} \|\widehat{\mathbf{w}} - \mathbf{w}_*\|^2 \\
 & \leq \frac{(\cdot) \|\widehat{\mathbf{w}} - \mathbf{w}_*\|^2}{2} + \frac{1}{2} \|\widehat{\mathbf{w}} - \mathbf{w}_*\|^2 \sqrt{\frac{(\cdot)}{2}} \\
 & \quad + \frac{2 \cdot \log(2 \cdot) \|\widehat{\mathbf{w}} - \mathbf{w}_*\|}{2} + \frac{1}{2} \|\widehat{\mathbf{w}} - \mathbf{w}_*\| \sqrt{8 \cdot \log(2 \cdot)} \\
 & \quad + 2 \cdot \frac{1}{2} \|\widehat{\mathbf{w}} - \mathbf{w}_*\| + \frac{1}{2} \|\widehat{\mathbf{w}} - \mathbf{w}_*\| \sqrt{\frac{(\cdot)}{2}} + \frac{(\cdot)}{2} \|\widehat{\mathbf{w}} - \mathbf{w}_*\|
 \end{aligned}$$

8

to get (4) we substitute

$$\begin{aligned} \frac{1}{L} \widehat{w} w_*^2 \sqrt{\frac{(\cdot)}{(\cdot)}} &< \frac{2}{L} \widehat{w} w_*^2 + \frac{1}{4} \widehat{w} w_*^2 \\ \frac{1}{L} \widehat{w} w_*^2 \sqrt{\frac{8 \log(2 \cdot)}{(\cdot)}} &< \frac{8 \log(2 \cdot)}{L} + \frac{1}{4} \widehat{w} w_*^2 \end{aligned}$$

into (4) and then obtain

$$\begin{aligned} & \frac{(\widehat{w})}{L} (w_*) \\ & < \frac{(\cdot)}{L} \widehat{w} w_*^2 + \frac{2}{L} \widehat{w} w_*^2 + \frac{2 \log(2 \cdot)}{L} \widehat{w} w_*^2 + \frac{8 \log(2 \cdot)}{L} \\ & + 2 \frac{\widehat{w} w_*^2}{L} + \frac{\widehat{w} w_*^2}{L} \sqrt{\frac{(\cdot)}{(\cdot)}} + \frac{(\cdot)}{L} \widehat{w} w_*^2 \\ & \leq \frac{4 \cdot^2}{L} + \frac{4 \cdot^2}{L} + \frac{4 \log(2 \cdot)}{L} + \frac{8 \log(2 \cdot)}{L} \\ & + \left(4 \cdot + 2 \sqrt{\frac{(\cdot)}{(\cdot)}} + \frac{2 \cdot (\cdot)}{L} \right) \end{aligned}$$

which proves (4)

to get (4') we substitute

$$\begin{aligned} \frac{2 \log(2 \cdot)}{L} \widehat{w} w_*^2 &< \frac{8 \log^2(2 \cdot)}{L} + \frac{1}{8} \widehat{w} w_*^2 \\ \frac{1}{L} \widehat{w} w_*^2 \sqrt{\frac{8 \log(2 \cdot)}{(\cdot)}} &< \frac{32 \log(2 \cdot)}{L} + \frac{1}{16} \widehat{w} w_*^2 \\ \frac{2}{L} \widehat{w} w_*^2 &< \frac{32 \cdot^2}{L} + \frac{1}{32} \widehat{w} w_*^2 \\ \frac{1}{L} \widehat{w} w_*^2 \sqrt{\frac{(\cdot)}{(\cdot)}} &< \frac{16 \cdot^2}{L} + \frac{1}{64} \widehat{w} w_*^2 \\ \frac{(\cdot)}{L} \widehat{w} w_*^2 &< \frac{16 \cdot^2}{L} + \frac{1}{64} \widehat{w} w_*^2 \end{aligned}$$

into (4) and then obtain

$$\begin{aligned} & \frac{(\widehat{w})}{L} (w_*) + \frac{1}{4} \widehat{w} w_*^2 \\ & < \frac{(\cdot)}{L} \widehat{w} w_*^2 + \frac{1}{L} \widehat{w} w_*^2 \sqrt{\frac{(\cdot)}{(\cdot)}} + \frac{8 \log^2(2 \cdot)}{L} + \frac{32 \log(2 \cdot)}{L} \\ & + \left(\frac{32 \cdot^2}{L} + \frac{16 \cdot^2}{L} + \frac{16 \cdot^2}{L} \right) \\ & < \frac{2}{L} \widehat{w} w_*^2 + \frac{1}{5} \widehat{w} w_*^2 + \frac{8 \log^2(2 \cdot)}{L} + \frac{32 \log(2 \cdot)}{L} \\ & + \left(\frac{32 \cdot^2}{L} + \frac{16}{L} + \frac{16 \cdot^3}{625 \cdot^2} \right) \\ & \stackrel{\lambda/L \leq 16}{<} \frac{16}{25} \widehat{w} w_*^2 + \frac{8 \log^2(2 \cdot)}{L} + \frac{32 \log(2 \cdot)}{L} + \left(\frac{32 \cdot^2}{L} + \frac{416}{625} \right) \end{aligned}$$

By subtracting $\frac{1}{4} \|\hat{\mathbf{w}} - \mathbf{w}_*\|^2$ 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

$$\|\hat{\mathbf{w}} - \mathbf{w}_*\| \leq \frac{1}{2}$$

Since $\ell(\cdot)$ is smooth and $\ell(\cdot)$ is Lipschitz continuous we have

$$\begin{aligned} \ell(\hat{\mathbf{w}}) - \ell(\mathbf{w}_*) &= \ell(\hat{\mathbf{w}}) - \ell(\mathbf{w}_*) + \ell(\mathbf{w}_*) - \ell(\mathbf{w}_*) \\ &\leq \langle \hat{\mathbf{w}} - \mathbf{w}_*, \nabla \ell(\mathbf{w}_*) \rangle + \frac{1}{2} \|\hat{\mathbf{w}} - \mathbf{w}_*\|^2 + \|\hat{\mathbf{w}} - \mathbf{w}_*\| \\ &\leq \langle \hat{\mathbf{w}} - \mathbf{w}_*, \nabla \ell(\mathbf{w}_*) \rangle + \frac{1}{2} \|\hat{\mathbf{w}} - \mathbf{w}_*\|^2 + \|\hat{\mathbf{w}} - \mathbf{w}_*\| \leq \frac{1}{2} + \frac{1}{4} \end{aligned} \tag{4}$$

where the last step utilizes Jensen's inequality

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

Next we study the case

$$\frac{1}{2} < \|\hat{\mathbf{w}} - \mathbf{w}_*\| \tag{5}$$

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

$$\begin{aligned} &\ell(\hat{\mathbf{w}}) - \ell(\mathbf{w}_*) + \frac{1}{2} \|\hat{\mathbf{w}} - \mathbf{w}_*\|^2 \\ &\leq \ell(\hat{\mathbf{w}}) - \ell(\mathbf{w}_*) + \underbrace{\left[\ell(\hat{\mathbf{w}}) - \ell(\mathbf{w}_*) \right] \langle \hat{\mathbf{w}} - \mathbf{w}_*, \nabla \ell(\mathbf{w}_*) \rangle + \frac{1}{2} \|\hat{\mathbf{w}} - \mathbf{w}_*\|^2}_{:=B_1} \\ &= \ell(\hat{\mathbf{w}}) - \ell(\mathbf{w}_*) + \underbrace{\left[\ell(\hat{\mathbf{w}}) - \ell(\mathbf{w}_*) \right] \langle \hat{\mathbf{w}} - \mathbf{w}_*, \nabla \ell(\mathbf{w}_*) \rangle}_{:=B_2} + \frac{1}{2} \|\hat{\mathbf{w}} - \mathbf{w}_*\|^2 \\ &\leq \sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \|\hat{\mathbf{w}} - \mathbf{w}_*\|} \left\langle \nabla \ell(\mathbf{w}) - \nabla \ell(\mathbf{w}_*), \hat{\mathbf{w}} - \mathbf{w}_* \right\rangle \\ &\quad + \underbrace{\left\| \nabla \ell(\mathbf{w}_*) - \nabla \ell(\hat{\mathbf{w}}) \right\|}_{:=B_2} \|\hat{\mathbf{w}} - \mathbf{w}_*\| \end{aligned} \tag{4}$$

The first bound B_1 follows from the fact that the random variable $\langle \hat{\mathbf{w}} - \mathbf{w}_*, \nabla \ell(\mathbf{w}_*) \rangle$ lies in the range $(-1, 2]$ we develop the following lemma

Lemma 5 Under Assumptions [7](#) and [8](#), with probability at least $1 - \delta$, for all

$$\frac{1}{2} < \|\hat{\mathbf{w}} - \mathbf{w}_*\| < 2$$

the following bound holds:

$$\sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma} \left\langle \nabla \ell(\mathbf{w}) - \nabla \ell(\mathbf{w}_*), \hat{\mathbf{w}} - \mathbf{w}_* \right\rangle \leq \frac{4}{\gamma} \left(8 + \sqrt{2 \log \frac{1}{\delta}} \right)$$

where $\gamma = 2 \log_2 \left(\frac{1}{\delta} \right) + \log_2(2)$.

Based on the above we have with probability at least $1 - \epsilon$,

$$1 < \frac{4}{\lambda} \frac{\widehat{w} - w_*^2}{\lambda} \left(8 + \sqrt{2 \log \frac{1}{\epsilon}} \right) = \frac{\widehat{w} - w_*^2}{\lambda} \tag{4.4}$$

where λ is defined in

we then proceed to handle ϵ_2 which can be upper bounded in the same way as 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

$$\| \widehat{w} - w_* \|^2 < \frac{2 \log \frac{1}{\epsilon}}{\lambda} + \sqrt{\frac{8}{\lambda} \log \frac{1}{\epsilon}} \tag{4.5}$$

Substituting (4.4) and (4.5) into (4.3) with probability at least $1 - 2\epsilon$ we have

$$\begin{aligned} & \frac{(\widehat{w} - w_*)^2}{\lambda} + \frac{\widehat{w} - w_*^2}{2\lambda} \\ & < \frac{\widehat{w} - w_*^2}{\lambda} + \frac{2 \log \frac{1}{\epsilon}}{\lambda} \frac{\widehat{w} - w_*^2}{\lambda} + \frac{\widehat{w} - w_*^2}{\lambda} \sqrt{\frac{8}{\lambda} \log \frac{1}{\epsilon}} \end{aligned} \tag{4.6}$$

we substitute

$$\begin{aligned} \frac{\widehat{w} - w_*^2}{\lambda} & < \frac{2 \log \frac{1}{\epsilon}}{\lambda} \frac{\widehat{w} - w_*^2}{\lambda} + \frac{\widehat{w} - w_*^2}{4\lambda} \\ \frac{\widehat{w} - w_*^2}{\lambda} \sqrt{\frac{8}{\lambda} \log \frac{1}{\epsilon}} & < \frac{8 \log \frac{1}{\epsilon}}{\lambda} + \frac{\widehat{w} - w_*^2}{4\lambda} \end{aligned}$$

into (4.6) and then have

$$\begin{aligned} (\widehat{w} - w_*)^2 & < \frac{2 \log \frac{1}{\epsilon}}{\lambda} \frac{\widehat{w} - w_*^2}{\lambda} + \frac{2 \log \frac{1}{\epsilon}}{\lambda} \frac{\widehat{w} - w_*^2}{\lambda} + \frac{8 \log \frac{1}{\epsilon}}{\lambda} \\ & < \frac{4 \log \frac{1}{\epsilon}}{\lambda} + \frac{8 \log \frac{1}{\epsilon}}{\lambda} \end{aligned}$$

Combining the above inequality with (4.5) we obtain

to prove we substitute

$$\begin{aligned} \frac{2 \log \frac{1}{\epsilon}}{\lambda} \frac{\widehat{w} - w_*^2}{\lambda} & < \frac{8 \log^2 \frac{1}{\epsilon}}{\lambda} + \frac{\widehat{w} - w_*^2}{8\lambda} \\ \frac{\widehat{w} - w_*^2}{\lambda} \sqrt{\frac{8}{\lambda} \log \frac{1}{\epsilon}} & < \frac{16 \log \frac{1}{\epsilon}}{\lambda} + \frac{\widehat{w} - w_*^2}{8\lambda} \end{aligned}$$

into (4.6) and then have

$$\begin{aligned} & \frac{(\widehat{w} - w_*)^2}{\lambda} + \frac{\widehat{w} - w_*^2}{4\lambda} \\ & < \frac{\widehat{w} - w_*^2}{\lambda} + \frac{8 \log^2 \frac{1}{\epsilon}}{\lambda} + \frac{16 \log \frac{1}{\epsilon}}{\lambda} \\ & < \frac{\widehat{w} - w_*^2}{4\lambda} + \frac{8 \log^2 \frac{1}{\epsilon}}{\lambda} + \frac{16 \log \frac{1}{\epsilon}}{\lambda} \end{aligned}$$

Combining the above inequality with (4.5) we obtain

Appendix F. Proof of Lemma 5

First we partition the range $(1 - 2^{-k}, 2^{-k}]$ into $\lfloor 2 \log_2(1/\epsilon) + \log_2(2/\epsilon) \rfloor$ consecutive segments $\Delta_1, \Delta_2, \dots, \Delta_s$ such that

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

then we consider the case Δ_k for a fixed value of k . We have

$$\begin{aligned} & \sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma} \left\langle \mathbf{w}, \mathbf{w}_* \right\rangle [\hat{\mathbf{w}}, \hat{\mathbf{w}}_*] \\ & \leq \sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma_k^+} \left\langle \mathbf{w}, \mathbf{w}_* \right\rangle [\hat{\mathbf{w}}, \hat{\mathbf{w}}_*] \end{aligned} \tag{4.7}$$

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

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

$$\begin{aligned} & \sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma_k^+} \left\langle \mathbf{w}, \mathbf{w}_* \right\rangle [\hat{\mathbf{w}}, \hat{\mathbf{w}}_*] \\ & \leq \frac{(\gamma_k^+)^2}{\sqrt{n}} \left(8 + \sqrt{2 \log \frac{1}{\epsilon}} \right) \end{aligned} \tag{4.8}$$

since Δ_k we have

$$\gamma_k^+ = 2 \gamma_k^- < 2 \tag{4.9}$$

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

$$\begin{aligned} & \sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma} \left\langle \mathbf{w}, \mathbf{w}_* \right\rangle [\hat{\mathbf{w}}, \hat{\mathbf{w}}_*] \\ & \leq \frac{4}{\sqrt{n}} \left(8 + \sqrt{2 \log \frac{1}{\epsilon}} \right) \end{aligned} \tag{4.9}$$

we complete the proof by taking the union bound over segments

Appendix G. Proof of Lemma 7

to simplify the notation we define

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

$$\left(\frac{1}{n} \sum_{i=1}^n i(\mathbf{w}) \right) = \sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma_k^+} \left\langle \mathbf{w}, \mathbf{w}_* \right\rangle \frac{1}{n} \sum_{i=1}^n [i(\mathbf{w}) - i(\mathbf{w}_*)]$$

to upper bound $\left(\frac{1}{n} \sum_{i=1}^n i(\mathbf{w}) \right)$ we utilize the McDiarmid's inequality [McDiarmid 99](#)

Theorem 8 Let x_1, \dots, x_n be independent random variables taking values in a set A , and assume that $f: A^n \rightarrow \mathbb{R}$ satisfies

$$\sup_{x_1, \dots, x_n, x'_i \in A} \left| f(x_1, \dots, x_n) - f(x_1, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_n) \right| \leq \gamma_i$$

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

$$f(x_1, \dots, x_n) - E[f(x_1, \dots, x_n)] \leq \exp\left(-\frac{2\epsilon^2}{\sum_{i=1}^n \gamma_i^2}\right)$$

As pointed out in **Remark 7** Assumptions 7 and 8 imply the random function $f_i(\cdot)$ is smooth and thus

$$|f_i(\mathbf{w}) - f_i(\mathbf{w}_*)| \leq \|\mathbf{w} - \mathbf{w}_*\| \leq \left(\frac{\gamma_i}{k}\right)^2$$

As a result when a random function f_i changes the random variable $f_i(\cdot)$ can change by no more than $2\left(\frac{\gamma_i}{k}\right)^2$. To see this we have

$$\begin{aligned} & \left| f_i(\mathbf{w}) - f_i(\mathbf{w}_*) \right| \leq \left| f_i(\mathbf{w}) - f_i(\mathbf{w}_*) \right| \leq \frac{2}{k} \left(\frac{\gamma_i}{k}\right)^2 \\ & \leq \frac{1}{k} \sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma_k^+} \left\langle f'_i(\mathbf{w}) - f'_i(\mathbf{w}_*) \mid [f_i(\mathbf{w}) - f_i(\mathbf{w}_*)] \right\rangle \leq \frac{2}{k} \left(\frac{\gamma_i}{k}\right)^2 \end{aligned}$$

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

$$f(x_1, \dots, x_n) \leq E[f(x_1, \dots, x_n)] + \left(\frac{\gamma_i}{k}\right)^2 \sqrt{\frac{2}{\delta} \log \frac{1}{\delta}}$$

Let (x'_1, \dots, x'_n) be an independent copy of (x_1, \dots, x_n) and h_1, \dots, h_n be independent Rademacher variables with equal probability of being ± 1 . Using techniques of [Bartlett and Mendelson](#) we bound $E[f(x_1, \dots, x_n)]$ as follows.

$$\begin{aligned} & E_{h_1, \dots, h_n} \left[\sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma_k^+} \left\langle f(\mathbf{w}) - f(\mathbf{w}_*) \mid \frac{1}{n} \sum_{i=1}^n [f_i(\mathbf{w}) - f_i(\mathbf{w}_*)] \right\rangle \right] \\ &= \frac{1}{2} E_{h_1, \dots, h_n} \left[\sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma_k^+} \left\langle f(\mathbf{w}) - f(\mathbf{w}_*) \mid \sum_{i=1}^n [f_i(\mathbf{w}) - f_i(\mathbf{w}_*)] \right\rangle \right] \\ &\leq \frac{1}{2} E_{h_1, \dots, h_n, h'_1, \dots, h'_n} \left[\sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma_k^+} \left\langle f(\mathbf{w}) - f(\mathbf{w}_*) \mid \sum_{i=1}^n [f_i(\mathbf{w}) - f_i(\mathbf{w}_*)] \right\rangle \right] \\ &= \frac{1}{2} E_{h_1, \dots, h_n, h'_1, \dots, h'_n} \left[\sum_{i=1}^n \left\langle f'_i(\mathbf{w}) - f'_i(\mathbf{w}_*) \mid [f_i(\mathbf{w}) - f_i(\mathbf{w}_*)] \right\rangle \right] \end{aligned}$$

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

Note that $\sum_{i=1}^n \ell_i(\mathbf{w})$ is Lipschitz over $[\gamma_k^+, \gamma_k^-]$ and $\sum_{i=1}^n \ell_i(\mathbf{w}) + \sum_{i=1}^n \ell_i(\mathbf{w}')$ hence from the comparison theorem of Ledoux and Talagrand (1999) in particular Lemma 14 of Meir and Zhang (2009) we have

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

Similarly we have

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

Combining (1.1) and (1.2) we arrive at

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

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

$$\begin{aligned} & \left| \ell_i(\mathbf{w}) - \ell_i(\mathbf{w}') \right| = \frac{1}{\gamma} \left| \ell_i(\mathbf{w}; \mathbf{x}_i) - \ell_i(\mathbf{w}'; \mathbf{x}_i) \right| \\ & \leq \sqrt{\gamma} \left| \mathbf{w} \cdot \mathbf{x}_i - \mathbf{w}' \cdot \mathbf{x}_i \right| = \sqrt{\gamma} \left| \mathbf{x}_i \cdot \mathbf{w} - \mathbf{w}_* - \mathbf{x}_i \cdot \mathbf{w}' - \mathbf{w}_* \right| \end{aligned}$$

Applying the comparison theorem of Ledoux and Talagrand (1999) again we have

$$C_1 \leq \sqrt{\gamma} \mathbb{E} \left[\sup_{\mathbf{w}: \|\mathbf{w} - \mathbf{w}_*\| \leq \gamma_k^+} \sum_{i=1}^n \ell_i(\mathbf{x}_i; \mathbf{w} - \mathbf{w}_*) \right] = 2 \gamma_k^+ \sqrt{\gamma}$$

