# Stability and Sharper Risk Bounds with Convergence Rate $O(1/n^2)$

Anonymous Author(s) Affiliation Address email

## Abstract

The sharpest known high probability excess risk bounds are up to O(1/n) for 1 2 empirical risk minimization and projected gradient descent via algorithmic stability 3 [Klochkov and Zhivotovskiy, 2021]. In this paper, we show that high probability excess risk bounds of order up to  $O(1/n^2)$  are possible. We discuss how high prob-4 ability excess risk bounds reach  $O(1/n^2)$  under strongly convexity, smoothness 5 and Lipschitz continuity assumptions for empirical risk minimization, projected 6 gradient descent and stochastic gradient descent. Besides, to the best of our knowl-7 edge, our high probability results on the generalization gap measured by gradients 8 9 for nonconvex problems are also the sharpest.

## 10 1 Introduction

Algorithmic stability is a fundamental concept in learning theory [Bousquet and Elisseeff, 2002], 11 which can be traced back to the foundational works of Vapnik and Chervonenkis [1974] and has 12 a deep connection with learnability [Rakhlin et al., 2005, Shalev-Shwartz et al., 2010, Shalev-13 Shwartz and Ben-David, 2014]. It is not difficult for only providing in-expectation error bounds via 14 stability arguments. However, high probability bounds are beneficial to understand the robustness of 15 optimization algorithms [Bousquet et al., 2020, Klochkov and Zhivotovskiy, 2021] and are much 16 more challenging [Feldman and Vondrak, 2019, Bousquet et al., 2020, Lv et al., 2021]. In this paper, 17 our goal is to improve the high probability risk bounds via algorithmic stability. 18

19 Let us start with some standard notations. We have a set of independent and identically distributed 20 observations  $S = \{z_1, \ldots, z_n\}$  sampled from a probability measure  $\rho$  defined on a sample space 21  $\mathcal{Z} := \mathcal{X} \times \mathcal{Y}$ . Based on the training set S, our goal is to build a model  $h : \mathcal{X} \mapsto \mathcal{Y}$  for prediction, 22 where the model is determined by parameter  $\mathbf{w}$  from parameter space  $\mathcal{W} \subset \mathbb{R}^d$ . The performance of 23 a model  $\mathbf{w}$  on an example z can be quantified by a loss function  $f(\mathbf{w}; z)$ , where  $f : \mathcal{W} \times \mathcal{Z} \mapsto \mathbb{R}_+$ . 24 Then the population risk and the empirical risk of  $\mathbf{w} \in \mathcal{W}$ , respectively as

$$F(\mathbf{w}) := \mathbb{E}_z \left[ f(\mathbf{w}; z) \right], \quad F_S(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^n f(\mathbf{w}; z_i),$$

where  $\mathbb{E}_z$  denotes the expectation w.r.t. z.

Let  $\mathbf{w}^* \in \arg\min_{\mathbf{w} \in \mathcal{W}} F(\mathbf{w})$  be the model with the minimal population risk in  $\mathcal{W}$  and  $\mathbf{w}^*(S) \in \arg\min_{\mathbf{w} \in \mathcal{W}} F_S(\mathbf{w})$  be the model with the minimal empirical risk w.r.t. dataset S. Let A(S) be the

output of a (possibly randomized) algorithm A on the dataset S. Let  $\|\cdot\|_2$  denote the Euclidean norm

and  $\nabla q(\mathbf{w})$  denote a subgradient of q at w.

<sup>30</sup> Traditional generalization analysis aims to bound the generalization error  $F(A(S)) - F_S(A(S))$  w.r.t

the algorithm A and the dataset S. Based on the technique developed by Feldman and Vondrak [2018,

Submitted to 38th Conference on Neural Information Processing Systems (NeurIPS 2024). Do not distribute.

<sup>32</sup> 2019], Bousquet et al. [2020] provide the sharpest high probability bounds of  $O(L/\sqrt{n})$ , where the

loss function  $f(\cdot, \cdot)$  is bounded by M. No matter how stable the algorithm is, the high probability

generalization bound will not be smaller than  $O(L/\sqrt{n})$ . This is sampling error term scaling as

 $O(1/\sqrt{n})$  that controls the generalization error [Klochkov and Zhivotovskiy, 2021].

A frequently used alternative to generalization bounds, that can avoid the sampling error, are the excess risk bounds. The excess risk of algorithm A w.r.t. the dataset S is  $F(A(S)) - F(\mathbf{w}^*)$ , which

is more essential because it considers both generalization error and optimization error. Recently,
 Klochkov and Zhivotovskiy [2021] provided the best high probability excess risk bounds of order

- 40 up to  $O(\log n/n)$  for empirical risk minimization (ERM) and projected gradient descent (PGD)
- 41 algorithms via algorithmic stability.

<sup>42</sup> On the other hand, Zhang et al. [2017], Li and Liu [2021], Xu and Zeevi [2024] derived high <sup>43</sup> probability excess risk bounds with  $O(1/n^2)$  for ERM and stochastic gradient descent (SGD) via <sup>44</sup> uniform convergence when the sample number satisfies  $n = \Omega(d)$ , which implied that the rate <sup>45</sup>  $O(1/n^2)$  is possible. However, the results obtained by the uniform convergence technique are related <sup>46</sup> to the dimension d, which is unacceptable in high-dimensional learning problems. Since stability <sup>47</sup> analysis can yield dimension-free bounds, we naturally have the following question:

48 Can algorithmic stability provide high probability excess risk bounds with the rate beyond O(1/n)?

<sup>49</sup> The main results of this paper answers this question positively. We provides the first high probability <sup>50</sup> bounds that are dimension-free with the rate  $O(1/n^2)$  for ERM, PGD and SGD. Our framework can <sup>51</sup> also be used to solve other stable algorithms.

To this end, we develop the generalization gap measured by gradients. Our bounds under *nonconvex* 52 setting are tighter than existing works based on both algorithmic stability [Fan and Lei, 2024] and 53 uniform convergence [Xu and Zeevi, 2024]. This is why we can achieve dimension-free excess risk 54 bounds of order  $O(1/n^2)$ . In fact, in nonconvex problems, optimization algorithms can only find 55 a local minimizer and we can only obtain optimization error bounds for  $\|\nabla F_S(A(S))\|_2$  [Ghadimi 56 and Lan, 2013]. Therefore, it is important to study the generalization behavior of A(S) measured by 57 gradients. Under Polyak-Lojasiewicz condition, we also obtain sharper results for both generalization 58 59 bounds of gradients and excess risk bounds. Our route to excess risk bounds can also be applied to various stable algorithms and complex learning scenarios. In this paper, we take ERM, PGD, 60 and SGD as examples to explore the stability of stochastic convex optimization algorithms with 61 strongly convex losses. We provide tighter high probability dimension-free excess risk bounds of 62 up to  $O(1/n^2)$  comapring with existing works based on both algorithmic stability [Klochkov and 63 Zhivotovskiy, 2021, Fan and Lei, 2024] and uniform convergence [Zhang et al., 2017, Li and Liu, 64 2021, Xu and Zeevi, 2024]. 65

Besides, to obtain tighter bounds, we obtain a tighter *p*-moment bound for sums of vector-valued
 functions by introducing the optimal Marcinkiewicz-Zygmund's inequality for random variables
 taking values in a Hilbert space in the proof, which has more potential applications in vector-valued
 functional data.

This paper is organized as follows. The related work is reviewed in Section 2. In Section 3, we present our main results for stability and generalization. We give applications to ERM, PGD and SGD in Section 4. The conclusion is given in Section 5. All the proofs and additional lemmata are deferred to the Appendix.

# 74 2 Related Work

Algorithmic stability. Algorithmic stability is a classical approach in generalization analysis, which 75 can be traced back to the foundational works of [Vapnik and Chervonenkis, 1974]. It gave the 76 generalization bound by analyzing the sensitivity of a particular learning algorithm when changing 77 one data point in the dataset. Modern framework of stability analysis was established by Bousquet 78 and Elisseeff [2002], where they presented an important concept called uniform stability. Since 79 then, a lot of works based on uniform stability have emerged. On one hand, generalization bounds 80 with algorithmic stability have been significantly improved by Feldman and Vondrak [2018, 2019], 81 Bousquet et al. [2020], Klochkov and Zhivotovskiy [2021]. On the other hand, different algorithmic 82 stability measures such as uniform argument stability [Liu et al., 2017, Bassily et al., 2020], on 83 average stability [Shalev-Shwartz et al., 2010, Kuzborskij and Lampert, 2018], hypothesis stability 84

[Bousquet and Elisseeff, 2002, Charles and Papailiopoulos, 2018], hypothesis set stability [Foster 85 et al., 2019], pointwise uniform stability [Fan and Lei, 2024], PAC-Bayesian stability [Li et al., 2020], 86 locally elastic stability [Deng et al., 2021], collective stability [London et al., 2016] and uniform 87 stability in gradients [Lei, 2023, Fan and Lei, 2024] have been developed. Most of them provided the 88 connection on stability and generalization in expectation. Bousquet and Elisseeff [2002], Elisseeff 89 et al. [2005], Feldman and Vondrak [2018, 2019], Bousquet et al. [2020], Klochkov and Zhivotovskiy 90 [2021], Fan and Lei [2024] considered high probability bounds. However, only Fan and Lei [2024] 91 developed vector-valued bounds (eg: generalization bounds of gradients), which can be the order at 92 most  $O(M/\sqrt{n})$  and remain improvement. 93 Uniform convergence. Uniform convergence is another popular approach in statistical learning 94 theory to study generalization bounds [Fisher, 1922, Vapnik, 1999, Van der Vaart, 2000]. The main 95 idea is to bound the generalization gap by its supremum over the whole (or a subset) of the hypothesis 96 space via some space complexity measures, such as VC dimension, covering number and Rademacher 97 complexity. For finite-dimensional problem, Kleywegt et al. [2002] provided that the generalization 98

error is  $O\left(\sqrt{d/n}\right)$  depended on the sample size n and the dimension of parameters d in high

probability. In nonconvex settings, Mei et al. [2018] showed that the empirical of generalization error

is  $O(\sqrt{d/n})$ . Xu and Zeevi [2024] developed a novel "uniform localized convergence" framework using generic chaining for the minimization problem and provided the localized generalization bounds

in gradients  $O\left(\max\left\{\|\mathbf{w} - \mathbf{w}^*\|_2, \frac{1}{n}\right\}\left(\sqrt{\frac{d}{n}} + \frac{d}{n}\right)\right)$ , which is the optimal result when we only consider the order of *n*. However, uniform convergence results are related to the dimension *d*, which

consider the order of n. However, uniform convergence results are related to the dimension d, which is unacceptable in high-dimensional learning problems.

#### **106 3 Stability and Generalization**

To derive sharper generalization bounds of gradients, we need to develop a novel concentration inequality which provide *p*-moment bound for sums of vector-valued functions. For a real-valued random variable *Y*, the  $L_p$ -norm of *Y* is defined by  $||Y||_p := (\mathbb{E}[|Y|^p])^{\frac{1}{p}}$ . Similarly, let  $|| \cdot ||$  denote the norm in a Hilbert space  $\mathcal{H}$ . Then for a random variable *X* taking values in a Hilbert space, the  $L_p$ -norm of *X* is defined by  $|||\mathbf{X}|||_p := (\mathbb{E}[|\mathbf{X}||^p])^{\frac{1}{p}}$ .

#### 112 3.1 A Moment Bound for Sums of Vector-valued Functions

Here we present our sharper moment bound for sums of vector-valued functions of n independent variables.

**Theorem 1.** Let  $\mathbf{Z} = (Z_1, ..., Z_n)$  be a vector of independent random variables each taking values in  $\mathcal{Z}$ , and let  $\mathbf{g}_1, ..., \mathbf{g}_n$  be some functions:  $\mathbf{g}_i : \mathcal{Z}^n \mapsto \mathcal{H}$  such that the following holds for any  $i \in [n]$ :

- 118  $\|\mathbb{E}[\mathbf{g}_i(\mathbf{Z})|Z_i]\| \leq M \ a.s.,$
- 119  $\mathbb{E}\left[\mathbf{g}_{i}(\mathbf{Z})|Z_{[n]\setminus\{i\}}\right] = 0 \text{ a.s.,}$
- $\mathbf{g}_i$  satisfies the bounded difference property with  $\beta$ , namely, for any i = 1, ..., n, the following inequality holds

$$\sup_{z_1,\dots,z_n,z'_j} \|\mathbf{g}_i(z_1,\dots,z_{j-1},z_j,z_{j+1},\dots,z_n) - \mathbf{g}_i(z_1,\dots,z_{j-1},z'_j,z_{j+1},\dots,z_n)\| \le \beta.$$
(1)

122 Then, for any  $p \ge 2$ , we have

$$\left\| \left\| \sum_{i=1}^{n} \mathbf{g}_{i} \right\| \right\|_{p} \leq 2(\sqrt{2p}+1)\sqrt{n}M + 4 \times 2^{\frac{1}{2p}} \left( \sqrt{\frac{p}{e}} \right) (\sqrt{2p}+1)n\beta \left\lceil \log_{2} n \right\rceil.$$

**Remark 1.** The proof is motivated by Bousquet et al. [2020]. Under the same assumptions, Fan and Lei [2024] also established the following inequality<sup>1</sup>

$$\left\| \left\| \sum_{i=1}^{n} \mathbf{g}_{i} \right\| \right\|_{p} \leq 2(\sqrt{2}+1)\sqrt{np}M + 4(\sqrt{2}+1)np\beta \left\lceil \log_{2} n \right\rceil.$$
(2)

It is easy to verify that our result is tighter than result provided by Fan and Lei [2024] for both the first and second term. Comparing Theorem 1 with (2), the larger p is, the tighter our result is relative to (2). In the worst case, when p = 2, the constant of our first term is 0.879 times tighter than (2), and the constant of our second term is 0.634 times tighter than (2). This is because we derive the optimal Marcinkiewicz-Zygmund's inequality for random variables taking values in a Hilbert space in the proof.

The improvement of this concentration inequality is meaningful. On one hand, we derive the optimal Marcinkiewicz-Zygmund's inequality for random variables taking values in a Hilbert space. On the other hand, in Section 3.2, we will carefully construct vector-valued functions which satisfies all the assumptions in Theorem 1 and ensures M = 0 at the same time. Under this condition, we can eliminate the first term. When we use Theorem 1 instead of (2) in the whole proofs, at least 0.634 times tighter bound can be obtained strictly.

#### 137 3.2 Sharper Generalization Bounds in Gradients

Let  $S = \{z_1, \ldots, z_n\}$  be a set of independent random variables each taking values in Z and  $S' = \{z'_1, \ldots, z'_n\}$  be its independent copy. For any  $i \in [n]$ , define  $S^{(i)} = \{z_i, \ldots, z_{i-1}, z'_i, z_{i+1}, \ldots, z_n\}$  be a dataset replacing the *i*-th sample in S with another i.i.d. sample  $z'_i$ . We introduce some basic definitions here and we want to emphasize that our main Theorem 2 and Theorem 3 do not need smoothness assumption and PL condition.

143 **Definition 1.** Let  $g : W \mapsto \mathbb{R}$ . Let  $\gamma, \mu < 0$ .

• We say g is 
$$\gamma$$
-smooth if

1

$$\|\nabla g(\mathbf{w}) - \nabla g(\mathbf{w}')\|_2 \le \gamma \|\mathbf{w} - \mathbf{w}'\|_2, \quad \forall \mathbf{w}, \mathbf{w}' \in \mathcal{W}.$$

145 • Let  $g * = \min_{\mathbf{w} \in \mathcal{W}} g(\mathbf{w})$ . We say g satisfies the Polyak-Lojasiewicz (PL) condition with 146 parameter  $\mu > 0$  on  $\mathcal{W}$  if

$$g(\mathbf{w}) - g * \leq \frac{1}{2\mu} \|\nabla g(\mathbf{w})\|_2^2, \quad \forall \mathbf{w} \in \mathcal{W}.$$

147 Then we define uniform stability in gradients.

**Definition 2** (Uniform Stability in Gradients). Let A be a randomized algorithm. We say A is  $\beta$ -uniformly-stable in gradients if for all neighboring datasets  $S, S^{(i)}$ , we have

$$\sup_{z} \left\| \nabla f(A(S); z) - \nabla f(A(S^{(i)}); z) \right\|_{2} \le \beta.$$
(3)

Remark 2. Gradient-based stability was firstly introduced by Lei [2023], Fan and Lei [2024] to 150 describe the generalization performance for nonconvex problems. In nonconvex problems, we can 151 only find a local minimizer by optimization algorithms which may be far away from the global 152 minimizer. Thus the convergence does not make much sense in function values. Instead, the 153 convergence of  $\|\nabla F_S(A(S))\|_2$  was often studied in the optimization community [Ghadimi and Lan, 154 2013]. Since the population risk of gradients  $\|\nabla F(A(S))\|_2$  can be decomposed as the convergence 155 of  $\|\nabla F_S(A(S))\|_2$  and the generalization gap  $\|\nabla F(A(S)) - \nabla F_S(A(S))\|_2$ , the generalization 156 analysis of  $\|\nabla F(A(S)) - \nabla F_S(A(S))\|_2$  is important, which can be achieved by uniform stability 157 in gradients. 158

**Theorem 2** (Generalization via Stability in Gradients). Assume for any S and any z, 160  $\|\nabla f(A(S);z)\|_2 \leq M$ . If A is  $\beta$ -uniformly-stable in gradients, then for any  $\delta \in (0,1)$ , the

<sup>&</sup>lt;sup>1</sup>They assume  $n = 2^k, k \in \mathbb{N}$ . Here we give the version of their result with general n.

following inequality holds with probability at least  $1 - \delta$ 

$$\begin{aligned} \|\nabla F(A(S)) - \nabla F_S(A(S))\|_2 \\ \leq 2\beta + \frac{4M\left(1 + e\sqrt{2\log\left(e/\delta\right)}\right)}{\sqrt{n}} + 8 \times 2^{\frac{1}{4}}(\sqrt{2}+1)\sqrt{e}\beta \left\lceil \log_2 n \right\rceil \log\left(e/\delta\right). \end{aligned}$$

**Remark 3.** Theorem 2 is a direct application via Theorem 1 where we denote  $\mathbf{g}_i(S) = \mathbb{E}_{z'_i} \left[ \mathbb{E}_Z \left[ \nabla f(A(S^{(i)}), Z) \right] - \nabla f(A(S^{(i)}), z_i) \right]$  and find that  $\mathbf{g}_i(S)$  satisfies all the assumptions in Theorem 1. As a comparison, Fan and Lei [2024] also developed high probability bounds under same assumptions, but our bounds are sharper since our moment inequality for sums of vector-valued functions are tighter as we have discussed in Remark 1. Next, we derive sharper generalization bound of gradients under same assumptions.

**Theorem 3** (Sharper Generalization via Stability in Gradients). Assume for any S and any z,  $\|\nabla f(A(S); z)\|_2 \leq M$ . If A is  $\beta$ -uniformly-stable in gradients, then for any  $\delta \in (0, 1)$ , the following inequality holds with probability at least  $1 - \delta$ 

$$\begin{aligned} \|\nabla F(A(S)) - \nabla F_S(A(S))\|_2 \\ \leq & \sqrt{\frac{4\mathbb{E}_Z \left[ \|\nabla f(A(S); Z)\|_2^2 \right] \log \frac{6}{\delta}}{n}} + \sqrt{\frac{\left(\frac{1}{2}\beta^2 + 32n\beta^2 \log(3/\delta)\right) \log \frac{6}{\delta}}{n}} + \frac{M \log \frac{6}{\delta}}{n} \\ & + 16 \times 2^{\frac{3}{4}} \sqrt{e}\beta \left\lceil \log_2 n \right\rceil \log (3e/\delta) + 32\sqrt{e}\beta \left\lceil \log_2 n \right\rceil \sqrt{\log 3e/\delta}. \end{aligned}$$

**Remark 4.** Note that the factor in Theorem 2 before  $1/\sqrt{n}$  is  $O\left(M\sqrt{\log(e/\delta)}\right)$ , which depends on the bound of  $\|\nabla f(\cdot, \cdot)\|_2$ . However, the factor in Theorem 3 before  $1/\sqrt{n}$  is  $O\left(\sqrt{\mathbb{E}_Z\left[\|\nabla f(A(S); Z)\|_2^2\right]\log 1/\delta} + \beta \log(1/\delta)\right)$ , not involving the possibly large term M. As is known, optimization algorithms often provide parameters approaching the optimal solution, which make the term  $\mathbb{E}_Z[\|\nabla f(A(S); Z)\|_2^2]$  much more smaller than M. We will give further reasonable results under more assumptions such as smoothness in Lemma 1 and Lemma 2.

On the other hand, best high probability bounds based on uniform convergence [Xu and Zeevi, 2024] is

$$\|\nabla F(A(S)) - \nabla F_S(A(S))\|_2$$
  
$$\lesssim \sqrt{\frac{\mathbb{E}_Z \left[\nabla \|f(\mathbf{w}^*; Z)\|_2^2\right] \log(1/\delta)}{n}} + \frac{\log(1/\delta)}{n} + \max\left\{\|\mathbf{w} - \mathbf{w}^*\|_2, \frac{1}{n}\right\} \left(\sqrt{\frac{d}{n}} + \frac{d}{n}\right), \quad (4)$$

which is the optimal result when we only consider the order of n. However, uniform convergence

results are related to the dimension *d*, which is unacceptable in high-dimensional learning problems.

Note that (4) requires an additional smoothness-type assumption. As a comparison, when f is  $\gamma$ -smoothness, our result in Theorem 3 can be easily derived as

$$\begin{aligned} \|\nabla F(A(S)) - \nabla F_S(A(S))\|_2 \\ \lesssim \beta \log n \log(1/\delta) + \frac{\log(1/\delta)}{n} + \sqrt{\frac{\mathbb{E}_Z \left[\nabla \|f(\mathbf{w}^*; Z)\|_2^2\right] \log(1/\delta)}{n}} + \|A(S) - \mathbf{w}^*\| \sqrt{\frac{\log(1/\delta)}{n}}. \end{aligned}$$

This result implies that when the uniformly stable in gradients parameter  $\beta$  is smaller than  $1/\sqrt{n}$ , our bound is tighter than (4) and is dimension independent. Note that Theorem 3 holds in nonconvex problems, to the best of our knowledge, this is the sharpest upper bound in both uniform convergence and algorithmic stability analysis.

Here we give the proof sketch of Theorem 3, which is motivated by the analysis in Klochkov and Zhivotovskiy [2021]. The key idea is to build vector functions  $\mathbf{q}_i(S) = \mathbf{h}_i(S) - \mathbb{E}_{S\{z_i\}}[\mathbf{h}_i(S)]$ where we define  $\mathbf{h}_i(S) = \mathbb{E}_{z'_i} \left[ \mathbb{E}_Z \left[ \nabla f(A(S^{(i)}), Z) \right] - \nabla f(A(S^{(i)}), z_i) \right]$ . These functions satisfy all the assumptions in Theorem 1 and ensure the factor M in Theorem 1 to 0. Then the term  $O(1/\sqrt{n})$ can be eliminated. **Lemma 1.** Let assumptions in Theorem 3 hold. Suppose the function f is  $\gamma$ -smooth and the population risk F satisfies the PL condition with parameter  $\mu$ . Then for any  $\delta \in (0, 1)$ , when  $n \ge \frac{16\gamma^2 \log \frac{6}{\delta}}{\mu^2}$ , with probability at least  $1 - \delta$ , we have

$$\begin{aligned} \|\nabla F(A(S)) - \nabla F_S(A(S))\|_2 \\ \leq \|\nabla F_S(A(S))\|_2 + 4\sqrt{\frac{\mathbb{E}_Z\left[\|\nabla f(\mathbf{w}^*;Z)\|_2^2\right]\log\frac{6}{\delta}}{n}} + 2\sqrt{\frac{\left(\frac{1}{2}\beta^2 + 32n\beta^2\log(3/\delta)\right)\log\frac{6}{\delta}}{n}} \\ + \frac{2M\log\frac{6}{\delta}}{n} + 32 \times 2^{\frac{3}{4}}\sqrt{e}\beta\left[\log_2 n\right]\log\left(3e/\delta\right) + 64\sqrt{e}\beta\left[\log_2 n\right]\sqrt{\log 3e/\delta}. \end{aligned}$$

Remark 5. The following inequality can be easily derived using triangle inequality and Cauchy Bunyakovsky-Schwarz inequality:

$$F(A(S)) - F(\mathbf{w}^*) \lesssim \|\nabla F_S(A(S))\|_2 + \frac{F(\mathbf{w}^*)\log(1/\delta)}{n} + \frac{\log^2(1/\delta)}{n^2} + \beta^2\log^2 n\log^2(1/\delta).$$
(5)

Above inequality implies that excess risk can be bound by the optimization gradient error 197  $\|\nabla F_S(A(S))\|_2$  and uniform stability in gradients  $\beta$ . Note that the assumption  $F(\mathbf{w}^*) = O(1/n)$ 198 is common and can be found in Srebro et al. [2010], Lei and Ying [2020], Liu et al. [2018], Zhang 199 et al. [2017], Zhang and Zhou [2019]. This is natural since  $F(\mathbf{w}^*)$  is the minimal population risk. 200 On the other hand, we can derive that under  $\mu$ -strongly convex and  $\gamma$ -smooth assumptions for the 201 objective function f, uniform stability in gradients can be bounded of order O(1/n) for ERM and 202 PGD. Thus high probability excess risk can be bounded of order up to  $O(1/n^2)$  under these common 203 assumptions via algorithmic stability. Comparing with current best related work [Klochkov and 204 Zhivotovskiy, 2021], they are insensitive to the stability parameter being smaller than O(1/n) and 205 their best rates can only up to O(1/n). Although we involve extra smoothness and PL condition 206 assumptions, these assumptions are also common in optimization community and our work can fully 207 208 utilize these assumptions.

Besides, we discuss uniform stability in gradients for common algorithms such as ERM, PGD, and
SGD in Section 4. Our results can be easily extended to other stable algorithms. Due to smoothness's
property to link the uniform stability in gradients with uniform argument stability, many works
[Bassily et al., 2020, Feldman and Vondrak, 2019, Hardt et al., 2016] exploring uniform argument
stability can also use our framework.

Finally, the population risk of gradients  $\|\nabla F(A(S))\|_2$  can be gracefully bounded by the empirical risk of gradients  $\|\nabla F_S(A(S))\|_2$  under strong growth condition (SGC), that connects the rates at which the stochastic gradients shrink relative to the full gradient Vaswani et al. [2019].

217 **Definition 3** (Strong Growth Condition). We say SGC holds if

$$\mathbb{E}_{Z}\left[\|\nabla f(\mathbf{w}; Z)\|_{2}^{2}\right] \leq \rho \|\nabla F(\mathbf{w})\|_{2}^{2}.$$

**Remark 6.** There has been some related work that takes SGC into assumption Solodov [1998], Vaswani et al. [2019], Lei [2023]. Vaswani et al. [2019] has proved that the squared-hinge loss with linearly separable data and finite support features satisfies the SGC. Note that we only suppose this condition holds in Lemma 2.

Lemma 2 (SGC case). Let assumptions in Theorem 3 hold and suppose SGC holds. Then for any  $\delta > 0$ , with probability at least  $1 - \delta$ , we have

$$\|\nabla F(A(S))\| \lesssim (1+\eta) \|\nabla F_S(A(S))\| + \frac{1+\eta}{\eta} \left(\frac{M}{n}\log\frac{6}{\delta} + \beta\log n\log\frac{1}{\delta}\right)$$

**Remark 7.** Lemma 2 build a connection between the population gradient error and the empirical gradient error under Lipschitz, nonconvex, nonsmooth and SGC case and elucidate that the population gradient error can be bounded of up to O(1/n) under nonconvex problems.

## 227 4 Application

In this section, we analysis stochastic convex optimization with strongly convex losses. The most common setting is where at each round, the learner gets information on f through a stochastic gradient oracle [Rakhlin et al., 2012]. To derive uniform stability in gradients for algorithms, we firstly introduce the strongly convex assumption.

#### **Definition 4.** We say g is $\mu$ -strongly convex if

$$g(\mathbf{w}) \geq g(\mathbf{w}') + \langle \mathbf{w} - \mathbf{w}', \nabla g(\mathbf{w}') \rangle + \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}'\|_2^2, \quad \forall \mathbf{w}, \mathbf{w}' \in \mathcal{W}.$$

#### 233 4.1 Empirical Risk Minimizer

Empirical risk minimizer is one of the classical approaches for solving stochastic optimization (also referred to as sample average approximation (SAA)) in machine learning community. The following lemma shows the uniform stability in gradient for ERM under  $\mu$ -strongly convexity and  $\gamma$ -smoothness assumptions.

**Lemma 3** (Stability of ERM). Suppose the objective function f is  $\mu$ -strongly-convex and  $\gamma$ -smooth. For any  $\mathbf{w} \in \mathcal{W}$  and any z, suppose that  $\|\nabla f(\mathbf{w}; z) \leq M\|$ . Let  $\hat{\mathbf{w}}^*(S^{(i)})$  be the ERM of  $F_{S^{(i)}}(\mathbf{w})$ that denotes the empirical risk on the samples  $S^{(i)} = \{z_1, ..., z'_i, ..., z_n\}$  and  $\hat{\mathbf{w}}^*(S)$  be the ERM of  $F_S(\mathbf{w})$  on the samples  $S = \{z_1, ..., z_i, ..., z_n\}$ . For any  $S^{(i)}$  and S, there holds the following uniform stability bound of ERM:

$$\forall z \in \mathcal{Z}, \quad \left\| \nabla f(\hat{\mathbf{w}}^*(S^{(i)}); z) - \nabla f(\hat{\mathbf{w}}^*(S); z) \right\|_2 \le \frac{4M\gamma}{n\mu}.$$

Then, we present the application of our main sharper Theorem 3. In the strongly convex and smooth case, we provide a up to  $O(1/n^2)$  high probability excess risk guarantee valid for any algorithms depending on the optimal population error  $F(\mathbf{w}^*)$ .

**Theorem 4.** Let assumptions in Theorem 3 and Lemma 3 hold. Suppose the function f is nonnegative.

247 Then for any  $\delta \in (0,1)$ , when  $n \geq \frac{16\gamma^2 \log \frac{6}{\delta}}{\mu^2}$ , with probability at least  $1 - \delta$ , we have

$$F(\hat{\mathbf{w}}) - F(\mathbf{w}^*) \lesssim \frac{F(\mathbf{w}^*)\log(1/\delta)}{n} + \frac{\log^2 n \log^2(1/\delta)}{n^2}.$$

Furthermore, assume  $F(\mathbf{w}^*) = O(\frac{1}{n})$ , we have

$$F(\hat{\mathbf{w}}) - F(\mathbf{w}^*) \lesssim \frac{\log^2 n \log^2(1/\delta)}{n^2}.$$

**Remark 8.** Theorem 4 shows that when the objective function f is  $\mu$ -strongly convex,  $\gamma$ -smooth 249 and nonnegative, high probability risk bounds can even up to  $O(1/n^2)$  for ERM. The most related 250 work to ours is Zhang et al. [2017]. They also obtain the  $O(1/n^2)$ -type bounds for ERM by uniform 251 convergence of gradients approach. However, they need the sample number  $n = \Omega(\gamma d/\mu)$ , which 252 is related to the dimension d. Our risk bounds are dimension independent and only require the 253 sample number  $n = \Omega(\gamma^2/\mu^2)$ . Comparing with Klochkov and Zhivotovskiy [2021], we add two 254 assumptions, smoothness and  $F(\mathbf{w}^*) = O(1/n)$ , but our bounds also tighter, from O(1/n) to 255  $O(1/n^2).$ 256

#### 257 4.2 Projected Gradient Descent

Note that when the objective function f is strongly convex and smooth, the optimization error can be ignored. However, the generalization analysis framework proposed by Klochkov and Zhivotovskiy [2021] does not use smoothness assumption, which only derive high probability excess risk bound of order O(1/n) after  $T = O(\log n)$  steps under strongly convex and smooth assumptions. In this subsection, we provide sharper risk bound under the same iteration steps, which is because our generalization analysis also fully utilized the smooth assumptions. Here we give the definition of PGD.

**Definition 5** (Projected Gradient Descent). Let  $\mathbf{w}_1 = o \in \mathbb{R}^d$  be an initial point and  $\{\eta_t\}_t$  be a sequence of positive step sizes. PGD updates parameters by

$$\mathbf{w}_{t+1} = \Pi_{\mathcal{W}} \left( \mathbf{w}_t - \eta_t \nabla F_S \left( \mathbf{w}_t \right) \right),$$

where  $\nabla F_S(\mathbf{w}_t)$  denotes a subgradient of F w.r.t.  $\mathbf{w}_t$  and  $\Pi_W$  is the projection operator onto W.

**Lemma 4** (Stability of Gradient Descent). Suppose the objective function f is  $\mu$ -strongly-convex and  $\gamma$ -smooth. For any  $\mathbf{w} \in W$  and any z, suppose that  $\|\nabla f(\mathbf{w}; z)\|_2 \leq M$ . Let  $\mathbf{w}_t^i$  be the output of  $F_{S^{(i)}}(\mathbf{w})$  on t-th iteration on the samples  $S^{(i)} = \{z_1, ..., z'_i, ..., z_n\}$  in running PGD, and  $\mathbf{w}_t$  be the output of  $F_S(\mathbf{w})$  on t-th iteration on the samples  $S = \{z_1, ..., z'_i, ..., z_n\}$  in running PGD. Let the constant step size  $\eta_t = 1/\gamma$ . For any  $S^{(i)}$  and S, there holds the following uniform stability bound of PGD:

$$\forall z \in \mathcal{Z}, \quad \left\| \nabla f(\hat{\mathbf{w}}^*(S^{(i)}); z) - \nabla f(\hat{\mathbf{w}}^*(S); z) \right\|_2 \le \frac{4M\gamma}{n\mu}.$$

**Remark 9.** The derivations of Feldman and Vondrak [2019] in Section 4.1.2 (See also Hardt et al. [2016] in Section 3.4) imply that if the objective function f is  $\gamma$ -smooth in addition to  $\mu$ -strongly convexity and M-Lipschitz property, then PGD with the constant step size  $\eta = 1/\gamma$  is  $\left(\frac{2M}{n\mu}\right)$ uniformly argument stable for any number of steps, which means that PGD is  $\left(\frac{2M\gamma}{n\mu}\right)$ -uniformlystable in gradients regardless of iteration steps.

**Theorem 5.** Let assumptions in Theorem 3 and Lemma 3 hold. Suppose the function f is nonnegative. Let  $\{\mathbf{w}_t\}_t$  be the sequence produced by PGD with  $\eta_t = 1/\gamma$ . Then for any  $\delta \in (0,1)$ , when  $n \ge \frac{16\gamma^2 \log \frac{6}{\delta}}{\mu^2}$ , with probability at least  $1 - \delta$ , we have

$$F(\mathbf{w}) - F(\mathbf{w}^*) \lesssim \left(1 - \frac{\mu}{\gamma}\right)^{2T} + \frac{F(\mathbf{w}^*)\log\left(1/\delta\right)}{n} + \frac{\log^2 n \log^2(1/\delta)}{n^2}$$

*Furthermore, assume*  $F(\mathbf{w}^*) = O(\frac{1}{n})$  *and let*  $T \asymp \log n$ *, we have* 

$$F(\hat{\mathbf{w}}) - F(\mathbf{w}^*) \lesssim \frac{\log^2 n \log^2(1/\delta)}{n^2}$$

**Remark 10.** Theorem 5 shows that under the same assumptions as Klochkov and Zhivotovskiy [2021], our bound is  $O\left(\frac{F(\mathbf{w}^*)\log(1/\delta)}{n} + \frac{\log^2 n \log^2(1/\delta)}{n^2}\right)$ . Comparing with their bound  $O\left(\frac{\log n \log(1/\delta)}{n}\right)$ , we are sharper because  $F(\mathbf{w}^*)$  is the minimal population risk, which is a common assumption towards sharper risk bounds Srebro et al. [2010], Lei and Ying [2020], Liu et al. [2018], Zhang et al. [2017], Zhang and Zhou [2019].

#### 288 4.3 Stochastic Gradient Descent

Stochastic gradient descent optimization algorithm has been widely used in machine learning due to its simplicity in implementation, low memory requirement and low computational complexity per iteration, as well as good practical behavior. Here we give the definition of standard SGD.

**Definition 6** (Stochastic Gradient Descent). Let  $\mathbf{w}_1 = o \in \mathbb{R}^d$  be an initial point and  $\{\eta_t\}_t$  be a sequence of positive step sizes. SGD updates parameters by

$$\mathbf{w}_{t+1} = \Pi_{\mathcal{W}} \left( \mathbf{w}_t - \eta_t \nabla f \left( \mathbf{w}_t; z_{i_t} \right) \right),$$

where  $\nabla f(\mathbf{w}_t; z_{i_t})$  denotes a subgradient of f w.r.t.  $\mathbf{w}_t$  and  $i_t$  is independently drawn from the uniform distribution over  $[n] := \{1, 2, ..., n\}$ .

**Lemma 5** (Stability of SGD). Suppose the objective function f is  $\mu$ -strongly-convex and  $\gamma$ -smooth. For any  $\mathbf{w} \in W$  and any z, suppose that  $\|\nabla f(\mathbf{w}; z)\|_2 \leq M$ . Let  $\mathbf{w}_t^i$  be the output of  $F_{S^{(i)}}(\mathbf{w})$  on t-th iteration on the samples  $S^{(i)} = \{z_1, ..., z'_i, ..., z_n\}$  in running PGD and and  $\mathbf{w}_t$  be the output of  $F_S(\mathbf{w})$  on t-th iteration on the samples  $S = \{z_1, ..., z_i, ..., z_n\}$  in running SGD. For any  $S^{(i)}$  and S, there holds the following uniform stability bound of SGD:

$$\left\|\nabla f(\mathbf{w}_t; z) - \nabla f(\mathbf{w}_t^i; z)\right\|_2 \le 2\gamma \sqrt{\frac{2\epsilon_{opt}(\mathbf{w}_t)}{\mu}} + \frac{4M\gamma}{n\mu}, \quad \forall z \in \mathcal{Z},$$

301 where  $\epsilon_{opt}(\mathbf{w}_t) = F_S(\mathbf{w}_t) - F_S(\hat{\mathbf{w}}^*(S))$  and  $\hat{\mathbf{w}}^*(S)$  is the ERM of  $F_S(\mathbf{w})$ .

<sup>302</sup> Next, we introduce a necessary assumption in stochastic optimization theory.

**Assumption 1.** Assume the existence of  $\sigma > 0$  satisfying

$$\mathbb{E}_{i_t}[\|\nabla f(\mathbf{w}_t; z_{i_t}) - \nabla F_S(\mathbf{w}_t)\|_2^2] \le \sigma^2, \quad \forall t \in \mathbb{N},$$
(6)

where  $\mathbb{E}_{i_t}$  denotes the expectation w.r.t.  $i_t$ .

**Remark 11.** Assumption 1 is a standard assumption from the stochastic optimization theory [Nemirovski et al., 2009, Ghadimi and Lan, 2013, Ghadimi et al., 2016, Kuzborskij and Lampert, 2018, Zhou et al., 2018, Bottou et al., 2018, Lei and Tang, 2021], which essentially bounds the variance of the stochastic gradients for dataset *S*.

**Theorem 6.** Let assumptions in Theorem 3 and Lemma 5 hold. Suppose Assumption 1 holds and the

sum function f is nonnegative. Let  $\{\mathbf{w}_t\}_t$  be the sequence produced by SGD with  $\eta_t = \eta_1 t^{-\theta}, \theta \in (0, 1)$ and  $\eta_1 < \frac{1}{2\gamma}$ . Then for any  $\delta \in (0, 1)$ , when  $n \ge \frac{16\gamma^2 \log \frac{6}{\delta}}{\mu^2}$ , with probability at least  $1 - \delta$ , we have

$$\begin{pmatrix} \sum_{t=1}^{T} \eta_t \end{pmatrix}^{-1} \sum_{t=1}^{T} \eta_t \|\nabla F(\mathbf{w}_t)\|_2^2 \\ = \begin{cases} O\left(\frac{\log^2 n \log^3(1/\delta)}{T^{-\theta}}\right) + O\left(\frac{\log^2 n \log^2(1/\delta)}{n^2} + \frac{F(\mathbf{w}^*) \log^2(1/\delta)}{n}\right), & \text{if } \theta < 1/2 \\ O\left(\frac{\log^2 n \log^3(1/\delta)}{T^{-\frac{1}{2}}}\right) + O\left(\frac{\log^2 n \log^2(1/\delta)}{n^2} + \frac{F(\mathbf{w}^*) \log^2(1/\delta)}{n}\right), & \text{if } \theta = 1/2 \\ O\left(\frac{\log^2 n \log^3(1/\delta)}{T^{\theta-1}}\right) + O\left(\frac{\log^2 n \log^2(1/\delta)}{n^2} + \frac{F(\mathbf{w}^*) \log^2(1/\delta)}{n}\right), & \text{if } \theta > 1/2 \end{cases}$$

**Remark 12.** When  $\theta < 1/2$ , we take  $T \simeq n^{2/\theta}$ . When  $\theta = 1/2$ , we take  $T \simeq n^4$  and when  $\theta > 1/2$ , we set  $T \simeq n^{2/(1-\theta)}$ . Then according to Theorem 6, the population risk of gradient is bounded by  $O\left(\frac{\log^2 n \log^3(1/\delta)}{n^2} + \frac{F(\mathbf{w}^*) \log^2(1/\delta)}{n}\right)$ . To the best of our knowledge, this is the first high probability population gradient bound  $\|\nabla F(\mathbf{w}_t)\|_2$  for SGD via algorithmic stability.

**Theorem 7.** Let Assumptions in Theorem 3 and Lemma 5 hold. Suppose Assumption 1 holds and the function f is nonnegative. Let  $\{\mathbf{w}_t\}_t$  be the sequence produced by SGD with  $\eta_t = \frac{2}{\mu(t+t_0)}$  such that  $t_0 \ge \max\left\{\frac{4\gamma}{\mu}, 1\right\}$ . Then for any  $\delta > 0$ , when  $n \ge \frac{16\gamma^2 \log \frac{6}{\delta}}{\mu^2}$  and  $T \asymp n^2$ , with probability at least  $1 - \delta$ , we have

$$F(\mathbf{w}_{T+1}) - F(\mathbf{w}^*) = O\left(\frac{\log^4 n \log^5(1/\delta)}{n^2} + \frac{F(\mathbf{w}^*) \log(1/\delta)}{n}\right).$$

S20 Furthermore, assume  $F(\mathbf{w}^*) = O(\frac{1}{n})$ , we have

$$F(\mathbf{w}_{T+1}) - F(\mathbf{w}^*) = O\left(\frac{\log^4 n \log^5(1/\delta)}{n^2}\right).$$

**Remark 13.** Theorem 7 implies that high probability risk bounds for SGD optimization algorithm can be up to  $O(1/n^2)$  and the rate is dimension-free in high-dimensional learning problems. We compare Theorem 7 with most related work. For algorithmic stability, high probability risk bounds in Fan and Lei [2024] is up to O(1/n) when choosing optimal iterate number T for SGD optimization algorithm. To the best of knowledge, we are faster than all the existing bounds. The best high probability risk bounds of order  $O(1/n^2)$  are given by Li and Liu [2021] via uniform convergence, which require the sample number  $n = \Omega(\gamma d/\mu)$  depending on dimension d.

## 328 5 Conclusion

In this paper, we improve a *p*-moment concentration inequality for sums of vector-valued functions. By carefully constructing functions, we apply this moment concentration to derive sharper generalization bounds in gradients in nonconvex problems, which can further be used to obtain sharper high probability excess risk bounds for stable optimization algorithms. In application, we study three common algorithms: ERM, PGD, SGD. To the best of our knowledge, we provide the sharpest high probability dimension independent  $O(1/n^2)$ -type for these algorithms.

## 335 References

- R. Bassily, V. Feldman, C. Guzmán, and K. Talwar. Stability of stochastic gradient descent on nonsmooth convex
   losses. In *Proceedings of the 34th International Conference on Neural Information Processing Systems* (*NeurIPS*), volume 33, pages 4381–4391, 2020.
- L. Bottou, F. E. Curtis, and J. Nocedal. Optimization methods for large-scale machine learning. *SIAM review*, 60(2):223–311, 2018.
- O. Bousquet and A. Elisseeff. Stability and generalization. *The Journal of Machine Learning Research*, 2: 499–526, 2002.
- O. Bousquet, Y. Klochkov, and N. Zhivotovskiy. Sharper bounds for uniformly stable algorithms. In *Conference* on *Learning Theory*, pages 610–626. PMLR, 2020.
- Z. Charles and D. Papailiopoulos. Stability and generalization of learning algorithms that converge to global
   optima. In *International conference on machine learning*, pages 745–754. PMLR, 2018.
- 247 P. J. Davis. Gamma function and related functions. Handbook of mathematical functions, 256, 1972.
- V. De la Pena and E. Giné. *Decoupling: from dependence to independence*. Springer Science & Business Media,
   2012.
- Z. Deng, H. He, and W. Su. Toward better generalization bounds with locally elastic stability. In *International Conference on Machine Learning*, pages 2590–2600. PMLR, 2021.
- A. Elisseeff, T. Evgeniou, M. Pontil, and L. P. Kaelbing. Stability of randomized learning algorithms. *Journal of Machine Learning Research*, 6(1), 2005.
- J. Fan and Y. Lei. High-probability generalization bounds for pointwise uniformly stable algorithms. *Applied and Computational Harmonic Analysis*, 70:101632, 2024.
- V. Feldman and J. Vondrak. Generalization bounds for uniformly stable algorithms. *Advances in Neural Information Processing Systems*, 31, 2018.
- V. Feldman and J. Vondrak. High probability generalization bounds for uniformly stable algorithms with nearly
   optimal rate. In *Conference on Learning Theory*, pages 1270–1279. PMLR, 2019.
- R. A. Fisher. On the mathematical foundations of theoretical statistics. *Philosophical transactions of the Royal Society of London. Series A, containing papers of a mathematical or physical character*, 222(594-604):
   309–368, 1922.
- D. J. Foster, S. Greenberg, S. Kale, H. Luo, M. Mohri, and K. Sridharan. Hypothesis set stability and generalization. *Advances in Neural Information Processing Systems*, 32, 2019.
- S. Ghadimi and G. Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM journal on optimization*, 23(4):2341–2368, 2013.
- S. Ghadimi, G. Lan, and H. Zhang. Mini-batch stochastic approximation methods for nonconvex stochastic composite optimization. *Mathematical Programming*, 155(1):267–305, 2016.
- M. Hardt, B. Recht, and Y. Singer. Train faster, generalize better: Stability of stochastic gradient descent. In *International conference on machine learning*, pages 1225–1234. PMLR, 2016.
- H. Karimi, J. Nutini, and M. Schmidt. Linear convergence of gradient and proximal-gradient methods under the
   polyak-łojasiewicz condition. In *ECML*, pages 795–811. Springer, 2016.
- A. J. Kleywegt, A. Shapiro, and T. Homem-de Mello. The sample average approximation method for stochastic discrete optimization. *SIAM Journal on optimization*, 12(2):479–502, 2002.
- Y. Klochkov and N. Zhivotovskiy. Stability and deviation optimal risk bounds with convergence rate o(1/n). Advances in Neural Information Processing Systems, 34:5065–5076, 2021.
- I. Kuzborskij and C. Lampert. Data-dependent stability of stochastic gradient descent. In *Proceedings of the 35th International Conference on Machine Learning (ICML)*, pages 2815–2824. PMLR, 2018.
- R. Latała and K. Oleszkiewicz. On the best constant in the khinchin-kahane inequality. *Studia Mathematica*, 109(1):101–104, 1994.
- Y. Lei. Stability and generalization of stochastic optimization with nonconvex and nonsmooth problems. In *The Thirty Sixth Annual Conference on Learning Theory*, pages 191–227. PMLR, 2023.
- Y. Lei and K. Tang. Learning rates for stochastic gradient descent with nonconvex objectives. *IEEE Transactions* on Pattern Analysis and Machine Intelligence, 43(12):4505–4511, 2021.
- Y. Lei and Y. Ying. Fine-grained analysis of stability and generalization for stochastic gradient descent. In International Conference on Machine Learning, pages 5809–5819. PMLR, 2020.
- J. Li, X. Luo, and M. Qiao. On generalization error bounds of noisy gradient methods for non-convex learning.
   In *International Conference on Learning Representations*, 2020.
- S. Li and Y. Liu. Improved learning rates for stochastic optimization: Two theoretical viewpoints. *arXiv preprint* arXiv:2107.08686, 2021.
- M. Liu, X. Zhang, L. Zhang, R. Jin, and T. Yang. Fast rates of erm and stochastic approximation: Adaptive to
- error bound conditions. Advances in Neural Information Processing Systems, 31, 2018.

- T. Liu, G. Lugosi, G. Neu, and D. Tao. Algorithmic stability and hypothesis complexity. In *International Conference on Machine Learning*, pages 2159–2167. PMLR, 2017.
- B. London, B. Huang, and L. Getoor. Stability and generalization in structured prediction. *The Journal of Machine Learning Research*, 17(1):7808–7859, 2016.
- X. Luo and D. Zhang. Khintchine inequality on normed spaces and the application to banach-mazur distance.
   *arXiv preprint arXiv:2005.03728*, 2020.
- s. Lv, J. Wang, J. Liu, and Y. Liu. Improved learning rates of a functional lasso-type svm with sparse multi-kernel
   representation. *Advances in Neural Information Processing Systems*, 34:21467–21479, 2021.
- S. Mei, Y. Bai, and A. Montanari. The landscape of empirical risk for nonconvex losses. *The Annals of Statistics*, 46(6A):2747–2774, 2018.
- A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on optimization*, 19(4):1574–1609, 2009.
- I. Pinelis. Optimum bounds for the distributions of martingales in banach spaces. *The Annals of Probability*,
   pages 1679–1706, 1994.
- A. Rakhlin, S. Mukherjee, and T. Poggio. Stability results in learning theory. *Analysis and Applications*, 3(04):
   397–417, 2005.
- A. Rakhlin, O. Shamir, and K. Sridharan. Making gradient descent optimal for strongly convex stochastic
   optimization. In *Proceedings of the 29th International Coference on International Conference on Machine Learning*, pages 1571–1578, 2012.
- 412 O. Rivasplata, E. Parrado-Hernández, J. S. Shawe-Taylor, S. Sun, and C. Szepesvári. Pac-bayes bounds for stable 413 algorithms with instance-dependent priors. *Advances in Neural Information Processing Systems*, 31, 2018.
- S. Shalev-Shwartz and S. Ben-David. Understanding machine learning: From theory to algorithms. Cambridge university press, 2014.
- S. Shalev-Shwartz, O. Shamir, N. Srebro, and K. Sridharan. Learnability, stability and uniform convergence.
   *The Journal of Machine Learning Research*, 11:2635–2670, 2010.
- S. Smale and D.-X. Zhou. Learning theory estimates via integral operators and their approximations. *Constructive approximation*, 26(2):153–172, 2007.
- 420 M. V. Solodov. Incremental gradient algorithms with stepsizes bounded away from zero. *Computational* 421 *Optimization and Applications*, 11:23–35, 1998.
- N. Srebro, K. Sridharan, and A. Tewari. Optimistic rates for learning with a smooth loss. *arXiv preprint arXiv:1009.3896*, 2010.
- 424 A. W. Van der Vaart. Asymptotic statistics, volume 3. Cambridge university press, 2000.
- 425 V. Vapnik and A. Chervonenkis. Theory of Pattern Recognition. 1974.
- V. N. Vapnik. An overview of statistical learning theory. *IEEE transactions on neural networks*, 10(5):988–999,
   1999.
- S. Vaswani, F. Bach, and M. Schmidt. Fast and faster convergence of sgd for over-parameterized models and an
   accelerated perceptron. In *The 22nd international conference on artificial intelligence and statistics*, pages
   1195–1204. PMLR, 2019.
- R. Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47.
   Cambridge university press, 2018.
- Y. Xu and A. Zeevi. Towards optimal problem dependent generalization error bounds in statistical learning
   theory. *Mathematics of Operations Research*, 2024.
- L. Zhang and Z.-H. Zhou. Stochastic approximation of smooth and strongly convex functions: Beyond the o(1/t)
   convergence rate. In *Conference on Learning Theory*, pages 3160–3179. PMLR, 2019.
- L. Zhang, T. Yang, and R. Jin. Empirical risk minimization for stochastic convex optimization: O(1/n)-and o(1/n\*\*2)-type of risk bounds. In *Conference on Learning Theory*, pages 1954–1979. PMLR, 2017.
- Y. Zhou, Y. Liang, and H. Zhang. Generalization error bounds with probabilistic guarantee for sgd in nonconvex
   optimization. *arXiv preprint arXiv:1802.06903*, 2018.

## 441 A Additional definitions and lemmata

Lemma 6 (Equivalence of tails and moments for random vectors [Bassily et al., 2020]). Let X be a
 random variable with

$$||X||_p \le \sqrt{p}a + pb$$

444 for some  $a, b \ge 0$  and for any  $p \ge 2$ . Then for any  $\delta \in (0, 1)$  we have, with probability at least  $1 - \delta$ ,

$$|X| \le e\left(a\sqrt{\log\left(\frac{e}{\delta}\right)} + b\log\frac{e}{\delta}\right)$$

Lemma 7 (Vector Bernstein's inequality [Pinelis, 1994, Smale and Zhou, 2007]). Let  $\{X_i\}_{i=1}^n$  be a sequence of i.i.d. random variables taking values in a real separable Hilbert space. Assume that  $\mathbb{E}[X_i] = \mu$ ,  $\mathbb{E}[||X_i - \mu||^2] = \sigma^2$ , and  $||X_i|| \le M$ ,  $\forall 1 \le i \le n$ , then for all  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$  we have

$$\left\|\frac{1}{n}\sum_{i=1}^{n}X_{i}-\mu\right\| \leq \sqrt{\frac{2\sigma^{2}\log(\frac{2}{\delta})}{n}}+\frac{M\log\frac{2}{\delta}}{n}.$$

**Definition 7** (Weakly self-Bounded Function). Assume that a, b > 0. A function  $f : \mathbb{Z}^n \mapsto [0, +\infty)$ is said to be (a, b)-weakly self-bounded if there exist functions  $f_i : \mathbb{Z}^{n-1} \mapsto [0, +\infty)$  that satisfies for all  $\mathbb{Z}^n \in \mathbb{Z}^n$ ,

$$\sum_{i=1}^{n} (f_i(Z^n) - f(Z^n))^2 \le af(Z^n) + b.$$

**Lemma 8** ([Klochkov and Zhivotovskiy, 2021]). Suppose that  $z_1, \ldots, z_n$  are independent random variables and the function  $f: \mathbb{Z}^n \mapsto [0, +\infty)$  is (a, b)-weakly self-bounded and the corresponding

function  $f_i$  satisfy  $f_i(Z^n) \ge f(Z^n)$  for  $\forall i \in [n]$  and any  $Z^n \in \mathbb{Z}^n$ . Then, for any t > 0,

$$Pr(\mathbb{E}f(z_1,\ldots,z_n) \ge f(z_1,\ldots,z_n) + t) \le \exp\left(-\frac{t^2}{2a\mathbb{E}f(z_1,\ldots,z_n) + 2b}\right)$$

**Definition 8.** A Rademacher random variable is a Bernoulli variable that takes values  $\pm 1$  with probability  $\frac{1}{2}$  each.

## 457 **B Proofs of Section 3.1**

The proof of Theorem 1 is motivated by Bousquet et al. [2020], which need the Marcinkiewicz-Zygmund's inequality for random variables taking values in a Hilbert space and the McDiarmid's inequality for vector-valued functions.

Firstly, we derive the optimal constants in the Marcinkiewicz-Zygmund's inequality for random variables taking values in a Hilbert space.

**Lemma 9** (Marcinkiewicz-Zygmund's Inequality for Random Variables Taking Values in a Hilbert Space). Let  $\mathbf{X}_1, \ldots, \mathbf{X}_n$  be random variables taking values in a Hilbert space with  $\mathbb{E}[\mathbf{X}_i] = 0$  for all  $i \in [n]$ . Then for  $p \ge 2$  we have

$$\left\| \left\| \sum_{i=1}^{n} \mathbf{X}_{i} \right\| \right\|_{p} \leq 2 \cdot 2^{\frac{1}{2p}} \sqrt{\frac{np}{e}} \left( \frac{1}{n} \sum_{i=1}^{n} \left\| \left\| \mathbf{X}_{i} \right\| \right\|_{p}^{p} \right)^{\frac{1}{p}}.$$

Remark 14. Comparing with Marcinkiewicz-Zygmund's inequality given by Fan and Lei [2024],
 we provide best constants. Next, we give the proof of Lemma 9.

468 The Marcinkiewicz-Zygmund's inequality can be proved by using its connection to Khintchine-

Kahane's inequality. Thus, we introduce the best constants in Khintchine-Kahane's inequality for

<sup>470</sup> random variables taking values from a Hilbert space here.

**Lemma 10** (Best constants in Khintchine-Kahane's inequality in Hilbert space [Latała and Oleszkiewicz, 1994, Luo and Zhang, 2020]). For all  $p \in [2, \infty)$  and for all choices of Hilbert space  $\mathcal{H}$ , finite sets of vectors  $\mathbf{X}_i, \ldots, \mathbf{X}_n \in \mathcal{X} \in \mathcal{H}$ , and independent Rademacher variables  $r_1, \ldots, r_n$ ,

$$\left[\mathbb{E}\left\|\sum_{i=1}^{n} r_{i} \mathbf{X}_{i}\right\|^{p}\right]^{\frac{1}{p}} \leq C_{p} \cdot \left[\sum_{i=1}^{n} \|\mathbf{X}_{i}\|^{2}\right]^{\frac{1}{2}},$$

475 where  $C_p = 2^{\frac{1}{2}} \left\{ \frac{\Gamma(\frac{p+1}{2})}{\sqrt{\pi}} \right\}^{\frac{1}{p}}$ .

476 Proof of Lemma 9. The symmetrization argument goes as follows: Let  $(r_1, \ldots, r_n)$  be i.i.d. with

<sup>477</sup>  $\mathbb{P}(r_i = 1) = \mathbb{P}(r_i = -1) = 1/2$  and besides such that  $r_1, \ldots, r_n$  and  $(\mathbf{X}_1, \ldots, \mathbf{X}_n)$  are independent. <sup>478</sup> Then by independence and symmetry, according to Lemma 1.2.6 of De la Pena and Giné [2012], <sup>479</sup> conditioning on  $(\mathbf{X}_1, \ldots, \mathbf{X}_n)$  yields

$$\mathbb{E}\left[\left\|\sum_{i=1}^{n} \mathbf{X}_{i}\right\|^{p}\right] = 2^{p} \mathbb{E}\left[\left\|\sum_{i=1}^{n} r_{i} \mathbf{X}_{i}\right\|^{p}\right] \leq 2^{p} \mathbb{E}\left[\mathbb{E}\left[\left\|\sum_{i=1}^{n} r_{i} \mathbf{X}_{i}\right\|^{p} \middle| \mathbf{X}_{1}, \dots, \mathbf{X}_{n}\right]\right].$$
 (7)

480 As for the conditional expectation in (7), notice that by independence

$$\mathbb{E}\left[\left\|\sum_{i=1}^{n} r_{i} \mathbf{X}_{i}\right\|^{p} \left| \mathbf{X}_{1} = \mathbf{x}_{1}, \dots, \mathbf{X}_{n} = \mathbf{x}_{n}\right] = \mathbb{E}\left[\left\|\sum_{i=1}^{n} r_{i} \mathbf{X}_{i}\right\|^{p}\right]$$
(8)

According to Lemma 10, for  $v_n$ -almost every  $\mathbf{x}_1, \ldots, \mathbf{x}_n \in \mathbb{R}^n$ , where  $v_n := \mathbb{P} \circ (\mathbf{X}_1, \ldots, \mathbf{X}_n)^{-1}$ denotes the distribution of  $(\mathbf{X}_1, \ldots, \mathbf{X}_n)$ , we have

$$\left[\mathbb{E}\left\|\sum_{i=1}^{n} r_i \mathbf{x}_i\right\|^p\right] \le C \cdot \left[\sum_{i=1}^{n} \|\mathbf{x}_i\|^2\right]^{\frac{p}{2}},\tag{9}$$

where  $C = 2^{\frac{p}{2}} \frac{\Gamma(\frac{p+1}{2})}{\sqrt{\pi}}$  and C is optimal. This means that for any constant C' such that

$$\left[\mathbb{E}\left\|\sum_{i=1}^{n} r_{i} \mathbf{x}_{i}\right\|^{p}\right] \leq C' \cdot \left[\sum_{i=1}^{n} \left\|\mathbf{x}_{i}\right\|^{2}\right]^{\frac{p}{2}},$$
(10)

for all  $n \in \mathbb{N}$  and for each collection of vectors  $\mathbf{x}_1, \ldots, \mathbf{x}_n$ , it follows that  $C' \geq C$ .

485 From (8) and (9), we can infer that

$$\mathbb{E}\left[\left\|\sum_{i=1}^{n} r_{i} \mathbf{X}_{i}\right\|^{p} \middle| \mathbf{X}_{1} = \mathbf{x}_{1}, \dots, \mathbf{X}_{n} = \mathbf{x}_{n}\right] \leq C \cdot \left[\sum_{i=1}^{n} \|\mathbf{X}_{i}\|^{2}\right]^{\frac{p}{2}}.$$

486 Taking expectations in the above inequalities and (7) yield that

$$\mathbb{E}\left[\left\|\sum_{i=1}^{n} \mathbf{X}_{i}\right\|^{p}\right] \leq C \cdot \mathbb{E}\left[\sum_{i=1}^{n} \left\|\mathbf{X}_{i}\right\|^{2}\right]^{\frac{p}{2}}.$$
(11)

- 487 To see optimality let the above statement hold for some constants C' in place of C. Then if we choose
- 488  $\mathbf{X}_i := \mathbf{x}_i r_i, 1 \le i \le n$  with arbitrary reals vectors  $\mathbf{x}_1, \ldots, \mathbf{x}_n$ , it follows that

$$\mathbb{E}\left[\left\|\sum_{i=1}^{n} r_{i} \mathbf{x}_{i}\right\|^{p}\right] \leq C' \cdot \mathbb{E}\left[\sum_{i=1}^{n} \left\|\mathbf{x}_{i}\right\|^{2}\right]^{\frac{p}{2}},$$

whence we can conclude from (10) that  $C' \ge C$ . Thus we obtain that C' = C.

490 Notice that by Holder's inequality

$$\left[\sum_{i=1}^{n} \|\mathbf{X}_{i}\|^{2}\right]^{\frac{p}{2}} \leq n^{p/2-1} \sum_{i=1}^{n} \|\mathbf{X}_{i}\|^{p}.$$
(12)

491 Plugging (12) into (11), we have

$$\mathbb{E}\left[\left\|\sum_{i=1}^{n} \mathbf{X}_{i}\right\|^{p}\right] \leq C \cdot 2^{p} n^{p/2-1} \cdot \mathbb{E}\left[\sum_{i=1}^{n} \left\|\mathbf{X}_{i}\right\|^{p}\right],$$

492 where  $C = 2^{\frac{p}{2}} \frac{\Gamma\left(\frac{p+1}{2}\right)}{\sqrt{\pi}}$  is a constant.

- <sup>493</sup> Next, we use the following form of Stirling's formula for the Gamma-function, which follows from
- (6.1.5), (6.1.15) and (6.1.38) in Davis [1972] to bound the constant C. For every x > 0, there exists a  $\mu(x) \in (0, 1/(12x))$  such that

$$\Gamma(x) = \sqrt{2\pi} x^{x-1/2} e^{-x} e^{\mu(x)}.$$

496 Thus

$$C = 2^{\frac{p}{2}} \frac{\Gamma\left(\frac{p+1}{2}\right)}{\sqrt{\pi}} = g(p)\sqrt{2}e^{-p/2}p^{p/2},$$

497 with  $g(p) = \left(1 + \frac{1}{p}\right)^{p/2} e^{v(p) - 1/2}$ , where 0 < v(p) < 1/(6(p+1)). By Taylor's formula we have 498 that

$$\log(1+x) = \sum_{m=1}^{\infty} \frac{1}{m} (-1)^{m-1} x^m, \quad \forall x \in (-1,1],$$

499 and that for every  $k \in \mathbb{N}_0$ 

$$\sum_{m=1}^{2k} \frac{1}{m} (-1)^{m-1} x^m \le \log(1+x) \le \sum_{m=1}^{2k+1} \frac{1}{m} (-1)^{m-1} x^m, \forall x \ge 0.$$

500 Therefor we obtain with k = 1 that

$$\log g(p) = \frac{p}{2}\log(1+\frac{1}{p}) + v(p) - \frac{1}{2} \le -\frac{1}{4p} + \frac{1}{6p^2} + \frac{1}{6(p+1)} \le -\frac{1}{18p},$$

<sup>501</sup> where the last equality follows from elementary calculus. Similarly,

$$\log g(p) = \frac{p}{2}\log(1+\frac{1}{p}) + v(p) - \frac{1}{2} \ge -\frac{1}{4p} + v(p) \ge -\frac{1}{4p},$$

502 Thus, we have

$$e^{-\frac{1}{4p}}\sqrt{2}e^{-p/2}p^{p/2} < C < e^{-\frac{1}{18p}}\sqrt{2}e^{-p/2}p^{p/2}$$

- which implies that C is strictly smaller than  $\sqrt{2}e^{-p/2}p^{p/2}$  for all  $p \ge 2$ .
- Since  $C = \frac{1}{g(p)}\sqrt{2}e^{-p/2}p^{p/2}$  and  $g(p) \ge e^{-\frac{1}{4p}}$ , we can obtain that the relative error between C and  $\sqrt{2}e^{-p/2}p^{p/2}$  is equal to

$$\frac{1}{g(p)} - 1 \le e^{-\frac{1}{4p}} - 1 \le \frac{1}{4p}e^{\frac{1}{4p}}$$

- using Mean Value Theorem. This implies that the corresponding relative errors between C and  $\sqrt{2}e^{-p/2}p^{p/2}$  converge to zero as p tends to infinity.
- 508 The proof is complete.

509

- Then we introduce the McDiarmid's inequality for vector-valued functions. We firstly consider 510 real-valued functions, which follows from the standard tail-bound of McDiarmid's inequality and 511
- Proposition 2.5.2 in Vershynin [2018]. 512

**Lemma 11** (McDiarmid's Inequality for real-valued functions). Let  $Z_i, \ldots, Z_n$  be independent random variables, and  $f : \mathbb{Z}^n \mapsto \mathbb{R}$  such that the following inequality holds for any 513 514  $z_i,\ldots,z_{i-1},z_{i+1},\ldots,z_n$ 515

$$\sup_{z_i, z'_i} |f(z_1, \dots, z_{i-1}, z_i, z_{i+1}, \dots, z_n) - f(z_1, \dots, z_{i-1}, z'_i, z_{i+1}, \dots, z_n)| \le \beta,$$

Then for any p > 1 we have 516

$$\|f(Z_1,\ldots,Z_n) - \mathbb{E}f(Z_1,\ldots,Z_n)\|_p \le \sqrt{2pn\beta}.$$

To derive the McDiarmid's inequality for vector-valued functions, we need the expected distance 517 between  $\mathbf{f}(Z_1, \ldots, Z_n)$  and its expectation. 518

Lemma 12 ([Rivasplata et al., 2018]). Let  $Z_i, \ldots, Z_n$  be independent random variables, and 519  $\mathbf{f}: \mathcal{Z}^n \mapsto \mathcal{H}$  is a function into a Hilbert space  $\mathcal{H}$  such that the following inequality holds for 520 any  $z_i, ..., z_{i-1}, z_{i+1}, ..., z_n$ 521

$$\sup_{z_i, z'_i} \|\mathbf{f}(z_1, \dots, z_{i-1}, z_i, z_{i+1}, \dots, z_n) - \mathbf{f}(z_1, \dots, z_{i-1}, z'_i, z_{i+1}, \dots, z_n)\| \le \beta,$$

Then we have 522

$$\mathbb{E}\left[\left\|\mathbf{f}(Z_1,\ldots,Z_n)-\mathbb{E}\mathbf{f}(Z_1,\ldots,Z_n)\right\|\right] \leq \sqrt{n}\beta.$$

Now, we can easily derive the *p*-norm McDiarmid's inequality for vector-valued functions which 523 refines from Fan and Lei [2024] with better constants. 524

**Lemma 13** (McDiarmid's inequality for vector-valued functions). Let  $Z_i, \ldots, Z_n$  be independent 525

- random variables, and  $\mathbf{f}: \mathbb{Z}^n \mapsto \mathcal{H}$  is a function into a Hilbert space  $\mathcal{H}$  such that the following 526
- inequality holds for any  $z_i, \ldots, z_{i-1}, z_{i+1}, \ldots, z_n$ 527

$$\sup_{z_i, z'_i} \|\mathbf{f}(z_1, \dots, z_{i-1}, z_i, z_{i+1}, \dots, z_n) - \mathbf{f}(z_1, \dots, z_{i-1}, z'_i, z_{i+1}, \dots, z_n)\| \le \beta,$$
(13)

Then for any p > 1 we have 528

$$\left\|\left\|\mathbf{f}(Z_1,\ldots,Z_n)-\mathbb{E}\mathbf{f}(Z_1,\ldots,Z_n)\right\|\right\|_p \leq (\sqrt{2p}+1)\sqrt{n\beta}.$$

*Proof of Lemma 13.* Define a real-valued function  $h : \mathbb{Z}^n \mapsto \mathbb{R}$  as

$$h(z_1,\ldots,z_n) = \|\mathbf{f}(z_1,\ldots,z_n) - \mathbb{E}[\mathbf{f}(Z_1,\ldots,Z_n)]\|.$$

- We notice that this function satisfies the increment condition. For any i and  $z_1, \ldots, z_{i-1}, z_{i+1}, \ldots, z_n$ , 529 530 we have
  - $\sup |h(z_1, \dots, z_{i-1}, z_i, z_{i+1}, \dots, z_n) h(z_1, \dots, z_{i-1}, z'_i, z_{i+1}, \dots, z_n)|$  $z_i, z'_i$  $= \sup |||\mathbf{f}(z_1, \dots, z_n) - \mathbb{E}[\mathbf{f}(Z_1, \dots, Z_n)]|| - ||\mathbf{f}(z_1, \dots, z_{i-1}, z'_i, z_{i+1}, \dots, z_n) - \mathbb{E}[\mathbf{f}(Z_1, \dots, Z_n)]|||$  $z_i, z'_i$  $\leq \sup \left\| \left\| \mathbf{f}(z_1,\ldots,z_n) - \mathbf{f}(z_1,\ldots,z_{i-1},z'_i,z_{i+1},\ldots,z_n) \right\| \leq \beta.$  $z_i, z'_i$

Therefore, we can apply Lemma 11 to the real-valued function h and derive the following inequality 531

$$\|h(Z_1,\ldots,Z_n)-\mathbb{E}[h(Z_1,\ldots,Z_n)]\|_p \le \sqrt{2pn\beta}.$$

According to Lemma 12, we know the following inequality  $\mathbb{E}[h(Z_1, \ldots, Z_n)] \leq \sqrt{n\beta}$ . Combing the 532 above two inequalities together and we can derive the following inequality 533

$$\begin{aligned} &\|\|\mathbf{f}(Z_{1},\ldots,Z_{n}) - \mathbb{E}\mathbf{f}(Z_{1},\ldots,Z_{n})\|\|_{p} \\ \leq &\|h(Z_{1},\ldots,Z_{n}) - \mathbb{E}[h(Z_{1},\ldots,Z_{n})]\|_{p} + \|\mathbb{E}[h(Z_{1},\ldots,Z_{n})]\|_{p} \\ \leq &(\sqrt{2p}+1)\sqrt{n}\beta. \end{aligned}$$

\_ \....

The proof is complete. 534

535

- Proof of Theorem 1. For  $\mathbf{g}(Z_1,\ldots,Z_n)$  and  $A \subset [n]$ , we write  $\|\|\mathbf{g}\|\|_p(Z_A) = (\mathbb{E}[\|f\|^p Z_A])^{\frac{1}{p}}$ . 536
- Without loss of generality, we suppose that  $n = 2^k$ . Otherwise, we can add extra functions equal to 537 zero, increasing the number of therms by at most two times.
- 538
- Consider a sequence of partitions  $\mathcal{P}_0, \ldots, \mathcal{P}_k$  with  $\mathcal{P}_0 = \{\{i\} : i \in [n]\}, \mathcal{P}_k$  with  $\mathcal{P}_n = \{[n]\}, and$ 539 to get  $\mathcal{P}_l$  from  $\mathcal{P}_{l+1}$  we split each subset in  $\mathcal{P}_{l+1}$  into two equal parts. We have 540

$$\mathcal{P}_0 = \{\{1\}, \dots, \{2^k\}\}, \quad \mathcal{P}_1 = \{\{1, 2\}, \{3, 4\}, \dots, \{2^k - 1, 2^k\}\}, \quad \mathcal{P}_k = \{\{1, \dots, 2^k\}\}.$$

We have  $|\mathcal{P}_l| = 2^{k-l}$  and  $|P| = 2^l$  for each  $P \in \mathcal{P}_l$ . For each  $i \in [n]$  and  $l = 0, \ldots, k$ , denote by  $P^l(i) \in \mathcal{P}_l$  the only set from  $\mathcal{P}_l$  that contains *i*. In particular,  $P^0(i) = \{i\}$  and  $P^K(i) = [n]$ . 541 542

For each  $i \in [n]$  and every l = 0, ..., k consider the random variables 543

$$\mathbf{g}_{i}^{l} = \mathbf{g}_{i}^{l}(Z_{i}, Z_{[n] \setminus P^{l}(i)}) = \mathbb{E}[\mathbf{g}_{i} | Z_{i}, Z_{[n] \setminus P^{l}(i)}],$$

i.e. conditioned on  $z_i$  and all the variables that are not in the same set as  $Z_i$  in the partition  $\mathcal{P}_l$ . In 544 particular,  $\mathbf{g}_i^0 = \mathbf{g}_i$  and  $\mathbf{g}_i^k = \mathbb{E}[\mathbf{g}_i | Z_i]$ . We can write a telescopic sum for each  $i \in [n]$ , 545

$$\mathbf{g}_i - \mathbb{E}[\mathbf{g}_i | Z_i] = \sum_{l=1}^{k-1} \mathbf{g}_i^l - \mathbf{g}_i^{l+1}$$

Then, by the triangle inequality 546

$$\left\| \left\| \sum_{i=1}^{n} \mathbf{g}_{i} \right\| \right\|_{p} \leq \left\| \left\| \sum_{i=1}^{n} \mathbb{E}[\mathbf{g}_{i} | Z_{i}] \right\| \right\|_{p} + \sum_{l=0}^{k-1} \left\| \left\| \sum_{i=1}^{n} \mathbf{g}_{i}^{l} - \mathbf{g}_{i}^{l+1} \right\| \right\|_{p}.$$
 (14)

- $\mathbf{f}(Z_1,\ldots,Z_n) = \sum_{i=1}^n \mathbb{E}[\mathbf{g}_i|Z_i] \cong M$ , we can check that the vector-valued function  $\mathbf{f}(Z_1,\ldots,Z_n) = \sum_{i=1}^n \mathbb{E}[\mathbf{g}_i|Z_i]$  satisfies (13) with  $\beta = 2M$ , and  $\mathbb{E}[\mathbb{E}[\mathbf{g}_i|Z_i]] = 0$ , applying Lemma 13 with  $\beta = 2M$ , we have To bound the first term, since  $\|\mathbb{E}[\mathbf{g}_i|Z_i]\| \leq M$ , we can check that the vector-valued function 547 548
- 549

$$\left\| \left\| \sum_{i=1}^{n} \mathbb{E}[\mathbf{g}_{i} | Z_{i}] \right\| \right\|_{p} \le 2(\sqrt{2p} + 1)\sqrt{n}M.$$
(15)

Then we start to bound the second term of the right hand side of (14). Observe that 550

$$\mathbf{g}_i^{l+1}(Z_i, Z_{[n]\setminus P^{l+1}(i)}) = \mathbb{E}\left[\mathbf{g}_i^l(Z_i, Z_{[n]\setminus P^l(i)}) \middle| Z_i, Z_{[n]\setminus P^{l+1}(i)}\right],$$

where the expectation is taken with respect to the variables  $Z_j, j \in P^{l+1}(i) \setminus P^l(i)$ . Changing any 551  $Z_i$  would change  $\mathbf{g}_i^l$  by  $\beta$ . Therefore, we apply Lemma 13 with  $\mathbf{f} = \mathbf{g}_i^l$  where there are  $2^l$  random 552

variables and obtain a uniform bound 553

$$\|\|\mathbf{g}_{i}^{l}-\mathbf{g}_{i}^{l+1}\|\|_{p}\left(Z_{i}, Z_{[n]\setminus P^{l+1}(i)}\right) \leq (\sqrt{2p}+1)\sqrt{2^{l}}\beta, \quad \forall p \geq 2,$$

Taking integration over  $(Z_i, Z_{[n] \setminus P^{l+1}(i)})$ , we have  $\left\| \left\| \mathbf{g}_i^l - \mathbf{g}_i^{l+1} \right\| \right\|_p \le (\sqrt{2p} + 1)\sqrt{2^l}\beta$  as well. 554

Next, we turn to the sum  $\sum_{i \in P^l} \mathbf{g}_i^l - \mathbf{g}_i^{l+1}$  for any  $P^l \in \mathcal{P}_l$ . Since  $\mathbf{g}_i^l - \mathbf{g}_i^{l+1}$  for  $i \in P^l$  depends only on  $Z_i, Z_{[n] \setminus P^l}$ , the terms are independent and zero mean conditioned on  $Z_{[n] \setminus P^l}$ . Applying 555 556 Lemma 9, we have for any  $p \ge 2$ , 557

$$\left\| \left\| \sum_{i \in P^{l}} \mathbf{g}_{i}^{l} - \mathbf{g}_{i}^{l+1} \right\| \right\|_{p}^{p} (Z_{[n] \setminus P^{l}}) \leq \left( 2 \cdot 2^{\frac{1}{2p}} \sqrt{\frac{2^{l}p}{e}} \right)^{p} \frac{1}{2^{l}} \sum_{i \in P^{l}} \left\| \left\| \mathbf{g}_{i}^{l} - \mathbf{g}_{i}^{l+1} \right\| \right\|_{p}^{p} (Z_{[n] \setminus P^{l}}).$$

Integrating with respect to  $(Z_{[n]\setminus P^l})$  and using  $\left\| \left\| \mathbf{g}_i^l - \mathbf{g}_i^{l+1} \right\| \right\|_p \le (\sqrt{2p} + 1)\sqrt{2^l}\beta$ , we have 558

$$\begin{split} \left\| \left\| \sum_{i \in P^l} g_i^l - g_i^{l+1} \right\| \right\|_p &\leq \left( 2 \cdot 2^{\frac{1}{2p}} \sqrt{\frac{2^l p}{e}} \right) \frac{1}{2^l} \times 2^l (\sqrt{2p} + 1) \sqrt{2^l} \beta \\ &= 2^{1 + \frac{1}{2p}} \left( \sqrt{\frac{p}{e}} \right) (\sqrt{2p} + 1) 2^l \beta. \end{split}$$

Then using triangle inequality over all sets  $P^l \in \mathcal{P}_l$ , we have

$$\begin{split} \left\| \left\| \sum_{i \in [n]} \mathbf{g}_i^l - \mathbf{g}_i^{l+1} \right\| \right\|_p &\leq \sum_{P^l \in \mathcal{P}_l} \left\| \left\| \sum_{i \in P^l} \mathbf{g}_i^l - \mathbf{g}_i^{l+1} \right\| \right\|_p \\ &\leq 2^{k-l} \times 2^{1+\frac{1}{2p}} \left( \sqrt{\frac{p}{e}} \right) (\sqrt{2p} + 1) 2^l \beta \\ &\leq 2^{1+\frac{1}{2p}} \left( \sqrt{\frac{p}{e}} \right) (\sqrt{2p} + 1) 2^k \beta. \end{split}$$

Recall that  $2^k \le n$  due to the possible extension of the sample. Then we have

$$\sum_{i=0}^{k-1} \left\| \left\| \sum_{i=1}^{n} \mathbf{g}_{i}^{l} - \mathbf{g}_{i}^{i+1} \right\| \right\|_{p} \leq 2^{2+\frac{1}{2p}} \left( \sqrt{\frac{p}{e}} \right) (\sqrt{2p} + 1) n\beta \left\lceil \log_{2} n \right\rceil.$$

<sup>561</sup> We can plug the above bound together with (15) into (14), to derive the following inequality

$$\left\| \left\| \sum_{i=1}^{n} \mathbf{g}_{i} \right\| \right\|_{p} \leq 2(\sqrt{2p}+1)\sqrt{n}M + 2^{2+\frac{1}{2p}} \left(\sqrt{\frac{p}{e}}\right)(\sqrt{2p}+1)n\beta \left\lceil \log_{2} n \right\rceil.$$

562 The proof is completed.

563

## 564 C Proofs of Section 3

Proof of Theorem 2. Let  $S = \{z_1, \ldots, z_n\}$  be a set of independent random variables each taking values in  $\mathcal{Z}$  and  $S' = \{z'_1, \ldots, z'_n\}$  be its independent copy. For any  $i \in [n]$ , define  $S^{(i)} = \{z_i, \ldots, z_{i-1}, z'_i, z_{i+1}, \ldots, z_n\}$  be a dataset replacing the *i*-th sample in S with another i.i.d. sample  $z'_i$ . Then we can firstly write the following decomposition

$$\begin{split} n\nabla F(A(S)) &- n\nabla F_S(A(S)) \\ = \sum_{i=1}^n \mathbb{E}_Z \left[ \nabla f(A(S); Z) \right] - \mathbb{E}_{z'_i} \left[ \nabla f(A(S^{(i)}), Z) \right] \right] \\ &+ \sum_{i=1}^n \mathbb{E}_{z'_i} \left[ \mathbb{E}_Z \left[ \nabla f(A(S^{(i)}), Z) \right] - \nabla f(A(S^{(i)}), z_i) \right] \\ &+ \sum_{i=1}^n \mathbb{E}_{z'_i} \left[ \nabla f(A(S^{(i)}), z_i) \right] - \sum_{i=1}^n \nabla f(A(S), z_i). \end{split}$$

We denote that  $\mathbf{g}_i(S) = \mathbb{E}_{z'_i} \left[ \mathbb{E}_Z \left[ \nabla f(A(S^{(i)}), Z) \right] - \nabla f(A(S^{(i)}), z_i) \right]$ , thus we have  $\| n \nabla F(A(S)) - n \nabla F_C(A(S)) \|$ 

$$\|n\nabla F(A(S)) - n\nabla F_{S}(A(S))\|_{2}$$

$$= \left\|\sum_{i=1}^{n} \mathbb{E}_{Z} \left[\nabla f(A(S); Z)\right] - \mathbb{E}_{z_{i}'} \left[\nabla f(A(S^{(i)}), Z)\right]\right]$$

$$+ \sum_{i=1}^{n} \mathbb{E}_{z_{i}'} \left[\mathbb{E}_{Z} \left[\nabla f(A(S^{(i)}), Z)\right] - \nabla f(A(S^{(i)}), z_{i})\right]$$

$$+ \sum_{i=1}^{n} \mathbb{E}_{z_{i}'} \left[\nabla f(A(S^{(i)}), z_{i})\right] - \sum_{i=1}^{n} \nabla f(A(S), z_{i})\right\|_{2}$$

$$\leq 2n\beta + \left\|\sum_{i=1}^{n} \mathbf{g}_{i}(S)\right\|_{2},$$
(16)

- <sup>570</sup> where the inequality holds from the definition of uniform stability in gradients.
- According to our assumptions, we get  $\|\mathbf{g}_i(S)\|_2 \leq 2M$  and

$$\mathbb{E}_{z_i}[\mathbf{g}_i(S)] = \mathbb{E}_{z_i}\mathbb{E}_{z'_i}\left[\mathbb{E}_Z\left[\nabla f(A(S^{(i)}); Z)\right] - \nabla f(A(S^{(i)}); z_i)\right]$$
$$= \mathbb{E}_{z'_i}\left[\mathbb{E}_Z\left[\nabla f(A(S^{(i)}); Z)\right] - \mathbb{E}_{z_i}\left[\nabla f(A(S^{(i)}); z_i)\right]\right] = 0,$$

where this equality holds from the fact that  $z_i$  and Z follow from the same distribution. For any  $i \in [n]$ , any  $j \neq i$  and any  $z''_j$ , we have

$$\begin{aligned} & \left\| \mathbf{g}_{i}(z_{1},\ldots,z_{j-1},z_{j},z_{j+1},\ldots,z_{n}) - \mathbf{g}_{i}(z_{1},\ldots,z_{j-1},z_{j}'',z_{j+1},\ldots,z_{n}) \right\|_{2} \\ & \leq \left\| \mathbb{E}_{z_{i}'} \left[ \mathbb{E}_{Z} \left[ \nabla f(A(S^{(i)});Z) \right] - \nabla f(A(S^{(i)});z_{i}) \right] - \mathbb{E}_{z_{i}'} \left[ \mathbb{E}_{Z} \left[ \nabla f(A(S^{(i)}_{j});Z) \right] - \nabla f(A(S^{(i)}_{j});Z) \right] \right\|_{2} \\ & \leq \left\| \mathbb{E}_{z_{i}'} \left[ \mathbb{E}_{Z} \left[ \nabla f(A(S^{(i)});Z) - \nabla f(A(S^{(i)}_{j});Z) \right] \right] \right\|_{2} + \left\| \mathbb{E}_{z_{i}'} \left[ \mathbb{E}_{Z} \left[ \nabla f(A(S^{(i)});Z) \right] - \nabla f(A(S^{(i)}_{j});z_{i}) \right] \right\|_{2} \\ & \leq 2\beta, \end{aligned}$$

where  $S^{(i)} = \{z_i, \ldots, z_{i-1}, z'_i, z_{i+1}, \ldots, z_n\}$ . Thus, we have verified that three conditions in Theorem 1 are satisfied for  $\mathbf{g}_i(S)$ . We have the following result for any p > 2

$$\left\| \left\| \sum_{i=1}^{n} \mathbf{g}_{i}(S) \right\| \right\|_{p} \leq 4(\sqrt{2p}+1)\sqrt{n}M + 8 \times 2^{\frac{1}{4}} \left(\sqrt{\frac{p}{e}}\right) (\sqrt{2p}+1)n\beta \left\lceil \log_{2} n \right\rceil.$$

576 We can combine the above inequality and (16) to derive the following inequality

$$n \left\| \left\| \nabla F(A(S)) - n \nabla F_S(A(S)) \right\| \right\|_p$$
  
$$\leq 2n\beta + 4(\sqrt{2p} + 1)\sqrt{n}M + 8 \times 2^{\frac{1}{4}} \left( \sqrt{\frac{p}{e}} \right) (\sqrt{2p} + 1)n\beta \left[ \log_2 n \right].$$

According to Lemma 6 for any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ , we have

$$\begin{split} n\|\nabla F(A(S)) - \nabla F_S(A(S))\|_2 \\ \leq & 2n\beta + 4\sqrt{n}M + 8 \times 2^{\frac{3}{4}}\sqrt{e}n\beta \left\lceil \log_2 n \right\rceil \log\left(e/\delta\right) + (4e\sqrt{2n}M + 8 \times 2^{\frac{1}{4}}\sqrt{e}n\beta \left\lceil \log_2 n \right\rceil)\sqrt{\log e/\delta}. \end{split}$$

578 This implies that

$$\begin{aligned} \|\nabla F(A(S)) - \nabla F_S(A(S))\|_2 \\ \leq 2\beta + \frac{4M\left(1 + e\sqrt{2\log\left(e/\delta\right)}\right)}{\sqrt{n}} + 8 \times 2^{\frac{1}{4}}(\sqrt{2}+1)\sqrt{e\beta}\left\lceil\log_2 n\right\rceil\log\left(e/\delta\right). \end{aligned}$$

579 The proof is completed.

580 Proof of Theorem 3. We can firstly write the following decomposition

$$\begin{split} n\nabla F(A(S)) &- n\nabla F_S(A(S)) \\ = \sum_{i=1}^n \mathbb{E}_Z \left[ \nabla f(A(S); Z) \right] - \mathbb{E}_{z'_i} \left[ \nabla f(A(S^{(i)}), Z) \right] \right] \\ &+ \sum_{i=1}^n \mathbb{E}_{z'_i} \left[ \mathbb{E}_Z \left[ \nabla f(A(S^{(i)}), Z) \right] - \nabla f(A(S^{(i)}), z_i) \right] \\ &+ \sum_{i=1}^n \mathbb{E}_{z'_i} \left[ \nabla f(A(S^{(i)}), z_i) \right] - \sum_{i=1}^n \nabla f(A(S), z_i). \end{split}$$

We denote that  $\mathbf{h}_i(S) = \mathbb{E}_{z'_i} \left[ \mathbb{E}_Z \left[ \nabla f(A(S^{(i)}), Z) \right] - \nabla f(A(S^{(i)}), z_i) \right]$ , we have

$$n\nabla F(A(S)) - n\nabla F_S(A(S)) - \sum_{i=1}^n \mathbf{h}_i(S)$$
  
=  $\sum_{i=1}^n \mathbb{E}_Z \left[ \nabla f(A(S); Z) \right] - \mathbb{E}_{z'_i} \left[ \nabla f(A(S^{(i)}), Z) \right] \right]$   
+  $\sum_{i=1}^n \mathbb{E}_{z'_i} \left[ \nabla f(A(S^{(i)}), z_i) \right] - \sum_{i=1}^n \nabla f(A(S), z_i)$ 

582 which implies that

$$\left\| n\nabla F(A(S)) - n\nabla F_S(A(S)) - \sum_{i=1}^n \mathbf{h}_i(S) \right\|_2$$
  
= 
$$\left\| \sum_{i=1}^n \mathbb{E}_Z \left[ \nabla f(A(S); Z) \right] - \mathbb{E}_{z'_i} \left[ \nabla f(A(S^{(i)}), Z) \right] \right]$$
  
+ 
$$\sum_{i=1}^n \mathbb{E}_{z'_i} \left[ \nabla f(A(S^{(i)}), z_i) \right] - \sum_{i=1}^n \nabla f(A(S), z_i) \right\|_2$$
  
$$\leq 2n\beta,$$
  
(17)

<sup>583</sup> where the inequality holds from the definition of uniform stability in gradients.

Then, for any i = 1, ..., n, we define  $\mathbf{q}_i(S) = \mathbf{h}_i(S) - \mathbb{E}_{S\{z_i\}}[\mathbf{h}_i(S)]$ . It is easy to verify that  $\mathbb{E}_{S \setminus \{z_i\}}[\mathbf{q}_i(S)] = \mathbf{0}$  and  $\mathbb{E}_{z_i}[\mathbf{h}_i(S)] = \mathbb{E}_{z_i}[\mathbf{q}_i(S)] - \mathbb{E}_{z_i}\mathbb{E}_{S \setminus \{z_i\}}[\mathbf{q}_i(S)] = \mathbf{0} - \mathbf{0} = \mathbf{0}$ . Also, for any  $j \in [n]$  with  $j \neq i$  and  $z''_j \in \mathbb{Z}$ , we have the following inequality

$$\begin{aligned} \|\mathbf{q}_{i}(S) - \mathbf{q}_{i}(z_{1}, \dots, z_{j-1}, z_{j}'', z_{j+1}, \dots, z_{n})\|_{2} \\ \leq \|\mathbf{h}_{i}(S) - \mathbf{h}_{i}(z_{1}, \dots, z_{j-1}, z_{j}'', z_{j+1}, \dots, z_{n})\|_{2} \\ + \|\mathbb{E}_{S \setminus \{z_{i}\}}[\mathbf{h}_{i}(S)] - \mathbb{E}_{S \setminus \{z_{i}\}}[\mathbf{h}_{i}(1, \dots, z_{j-1}, z_{j}'', z_{j+1}, \dots, z_{n})]\|_{2}. \end{aligned}$$

For the first term  $\|\mathbf{h}_i(S) - \mathbf{h}_i(z_1, \dots, z_{j-1}, z''_j, z_{j+1}, \dots, z_n)\|_2$ , it can be bounded by  $2\beta$  according to the definition of uniform stability. Similar result holds for the second term  $\|\mathbb{E}_{S \setminus \{z_i\}}[\mathbf{h}_i(S)] - \mathbb{E}_{S \setminus \{z_i\}}[\mathbf{h}_i(1, \dots, z_{j-1}, z''_j, z_{j+1}, \dots, z_n)]\|_2$  according to the uniform stability. By a combination of the above analysis, we get  $\|\mathbf{q}_i(S) - \mathbf{q}_i(1, \dots, z_{j-1}, z''_j, z_{j+1}, \dots, z_n)\|_2 \leq \|\mathbf{h}_i(S) - \mathbf{h}_i(1, \dots, z_{j-1}, z''_j, z_{j+1}, \dots, z_n)\|_2 \leq 4\beta$ .

Thus, we have verified that three conditions in Theorem 1 are satisfied for  $\mathbf{q}_i(S)$ . We have the following result for any  $p \ge 2$ 

$$\left\| \left\| \sum_{i=1}^{n} \mathbf{q}_{i}(S) \right\| \right\|_{p} \leq 2^{4+\frac{1}{4}} \left( \sqrt{\frac{p}{e}} \right) (\sqrt{2p} + 1)n\beta \left\lceil \log_{2} n \right\rceil.$$
(18)

594 Furthermore, we can derive that

$$\begin{split} n\nabla F(A(S)) &- n\nabla F_S(A(S)) - \sum_{i=1}^n \mathbf{h}_i(S) + \sum_{i=1}^n \mathbf{q}_i(S) \\ &= n\nabla F(A(S)) - n\nabla F_S(A(S)) - \sum_{i=1}^n \mathbb{E}_{S \setminus \{z_i\}}[\mathbf{h}_i(S)] \\ &= n\nabla F(A(S)) - n\nabla F_S(A(S)) - n\mathbb{E}_{S'}[\nabla F(A(S'))] + n\mathbb{E}_S[\nabla F(A(S))]. \end{split}$$

<sup>595</sup> Due to the i.i.d. property between S and S', we know that  $\mathbb{E}_{S'}[\nabla F(A(S'))] = \mathbb{E}_{S}[\nabla F(A(S))]$ . <sup>596</sup> Thus, combined above equality, (17) and (18), we have

$$\begin{split} & \|\|n\nabla F(A(S)) - n\nabla F_S(A(S)) - n\mathbb{E}_S[\nabla F(A(S))] + n\mathbb{E}_{S'}[\nabla F_S(A(S'))]\|\|_p \\ & \leq \left\| \left\| n\nabla F(A(S)) - n\nabla F_S(A(S)) - \sum_{i=1}^n \mathbf{h}_i(S) \right\| \right\|_p \\ & + \left\| \left\| \sum_{i=1}^n \mathbf{h}_i(S) - n\mathbb{E}_S[\nabla F(A(S))] + n\mathbb{E}_{S'}F_S[A(S')] \right\| \right\|_p \\ & = \left\| \left\| n\nabla F(A(S)) - n\nabla F_S(A(S)) - \sum_{i=1}^n \mathbf{h}_i(S) \right\| \right\|_p + \left\| \left\| \sum_{i=1}^n \mathbf{q}_i(S) \right\| \right\|_p \\ & \leq 2n\beta + 2^{4+\frac{1}{4}} \left( \sqrt{\frac{p}{e}} \right) (\sqrt{2p} + 1)n\beta \left\lceil \log_2 n \right\rceil \\ & \leq 16 \times 2^{\frac{3}{4}} \left( \sqrt{\frac{1}{e}} \right) pn\beta \left\lceil \log_2 n \right\rceil + 32 \left( \sqrt{\frac{1}{e}} \right) \sqrt{p}n\beta \left\lceil \log_2 n \right\rceil . \end{split}$$

597 According to Lemma 6 for any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta/3$ , we have

$$\|\nabla F(A(S)) - \nabla F_S(A(S))\|_2$$
  

$$\leq \|\mathbb{E}_{S'}[\nabla F_S(A(S'))] - \mathbb{E}_S[\nabla F(A(S))]\|_2$$
  

$$+ 16 \times 2^{\frac{3}{4}} \sqrt{e\beta} \left\lceil \log_2 n \right\rceil \log (3e/\delta) + 32\sqrt{e\beta} \left\lceil \log_2 n \right\rceil \sqrt{\log 3e/\delta}.$$
(19)

Next, we need to bound the term  $\|\mathbb{E}_{S'}[\nabla F_S(A(S'))] - \mathbb{E}_S[\nabla F(A(S))]\|_2$ . There holds that  $\|\mathbb{E}_S\mathbb{E}_{S'}[\nabla F_S(A(S'))]\|_2 = \|\mathbb{E}_S[\nabla F(A(S))]\|_2$ . Then, by the Bernstein inequality in Lemma 7, we obtain the following inequality with probability at least  $1 - \delta/3$ ,

$$\left\|\mathbb{E}_{S'}[\nabla F_S(A(S'))] - \mathbb{E}_S[\nabla F(A(S))]\right\|_2 \le \sqrt{\frac{2\mathbb{E}_{z_i}[\|\mathbb{E}_{S'}\nabla f(A(S'); z_i)\|_2^2]\log\frac{6}{\delta}}{n}} + \frac{M\log\frac{6}{\delta}}{n}.$$
 (20)

<sup>601</sup> Then using Jensen's inequality, we have

$$\mathbb{E}_{z_{i}}[\|\mathbb{E}_{S'}\nabla f(A(S'); z_{i})\|_{2}^{2}] \leq \mathbb{E}_{z_{i}}\mathbb{E}_{S'}\|\nabla f(A(S'); z_{i})\|_{2}^{2} \\
= \mathbb{E}_{Z}\mathbb{E}_{S'}\|\nabla f(A(S'); Z)\|_{2}^{2} \\
= \mathbb{E}_{Z}\mathbb{E}_{S}\|\nabla f(A(S); Z)\|_{2}^{2}.$$
(21)

<sup>602</sup> Combing (19), (20) with (21), we finally obtain that with probability at least  $1 - 2\delta/3$ ,

$$\|\nabla F(A(S)) - \nabla F_{S}(A(S))\|_{2}$$

$$\leq \sqrt{\frac{2\mathbb{E}_{Z}\mathbb{E}_{S}\|\nabla f(A(S);Z)\|_{2}^{2}\log\frac{6}{\delta}}{n}} + \frac{M\log\frac{6}{\delta}}{n} + \frac{M\log\frac{6}{\delta}}{n}$$

$$+ 16 \times 2^{\frac{3}{4}}\sqrt{e\beta} \left\lceil \log_{2}n \right\rceil \log(3e/\delta) + 32\sqrt{e\beta} \left\lceil \log_{2}n \right\rceil \sqrt{\log 3e/\delta}.$$

$$(22)$$

603 Next, since  $S = \{z_i, ..., z_n\}$ , we define  $p = p(z_1, ..., z_n) = \mathbb{E}_Z[\|\nabla f(A(S); Z)\|_2^2]$  and  $p_i = p_i(z_1, ..., z_n) = \sup_{z_i \in \mathcal{Z}} p(z_i, ..., z_n)$ . So there holds  $p_i \ge p$  for any i = 1, ..., n and any

605  $\{z_1,\ldots,z_n\}\in \mathcal{Z}^n$ . Also, there holds that

$$\begin{split} &\sum_{i=1}^{n} (p_{i} - p)^{2} \\ &= \sum_{i=1}^{n} \left( \sup_{z_{i} \in \mathcal{Z}} \mathbb{E}_{Z} [\|\nabla f(A(S'); Z)\|_{2}^{2}] - \mathbb{E}_{Z} [\|\nabla f(A(S); Z)\|_{2}^{2}] \right)^{2} \\ &\leq \sum_{i=1}^{n} \left( \mathbb{E}_{Z} \left[ \sup_{z_{i} \in \mathcal{Z}} \|\nabla f(A(S'); Z)\|_{2}^{2} - \|\nabla f(A(S); Z)\|_{2}^{2} \right] \right)^{2} \\ &= \sum_{i=1}^{n} \left( \mathbb{E}_{Z} \left[ \left( \sup_{z_{i} \in \mathcal{Z}} \|\nabla f(A(S'); Z)\|_{2} - \|\nabla f(A(S); Z)\|_{2} \right) \left( \sup_{z_{i} \in \mathcal{Z}} \|\nabla f(A(S'); Z)\|_{2} + \|\nabla f(A(S); Z)\|_{2} \right) \right] \right)^{2} \\ &\leq \sum_{i=1}^{n} \beta^{2} \left( \mathbb{E}_{Z} \left[ \|\nabla f(A(S); Z)\|_{2} + \sup_{z_{i} \in \mathcal{Z}} \|\nabla f(A(S); Z)\|_{2} \right] \right)^{2} \\ &\leq n\beta^{2} \left( 2\mathbb{E}_{Z} [\|\nabla f(A(S); Z)\|_{2} + \beta] \right)^{2} \\ &\leq 8n\beta^{2}p + 2n\beta^{4}, \end{split}$$
(23)

where the first inequality follows from the Jensen's inequality. The second and third inequalities follow from the definition of uniform stability in gradients. The last inequality holds from that  $(a + b)^2 \le 2a^2 + 2b^2$ .

From (23), we know that p is  $(8n\beta^2, 2n\beta^4)$  weakly self-bounded. Thus, by Lemma 8, we obtain that with probability at least  $1 - \delta/3$ ,

$$\mathbb{E}_{Z}\mathbb{E}_{S}[\|\nabla f(A(S); Z)\|_{2}^{2}] - \mathbb{E}_{Z}[\|\nabla f(A(S); Z)\|_{2}^{2}]$$

$$\leq \sqrt{(16n\beta^{2}\mathbb{E}_{S}\mathbb{E}_{Z}[\|\nabla f(A(S); Z)\|_{2}^{2}] + 4n\beta^{4})\log(3/\delta)}$$

$$= \sqrt{(\mathbb{E}_{S}\mathbb{E}_{Z}[\|\nabla f(A(S); Z)\|_{2}^{2}] + \frac{1}{4}\beta^{2})16n\beta^{2}\log(3/\delta)}$$

$$\leq \frac{1}{2}(\mathbb{E}_{S}\mathbb{E}_{Z}[\|\nabla f(A(S); Z)\|_{2}^{2}] + \frac{1}{4}\beta^{2}) + 8n\beta^{2}\log(3/\delta),$$

611 where the last inequality follows from that  $\sqrt{ab} \leq \frac{a+b}{2}$  for all a, b > 0. Thus, we have

$$\mathbb{E}_{Z}\mathbb{E}_{S}[\|\nabla f(A(S);Z)\|_{2}^{2}] \leq 2\mathbb{E}_{Z}[\|\nabla f(A(S);Z)\|_{2}^{2}] + \frac{1}{4}\beta^{2} + 16n\beta^{2}\log(3/\delta).$$
(24)

Substituting (24) into (22), we finally obtain that with probability at least  $1 - \delta$ 

$$\begin{aligned} \|\nabla F(A(S)) - \nabla F_S(A(S))\|_2 \\ \leq & \sqrt{\frac{2\left(2\mathbb{E}_Z[\|\nabla f(A(S);Z)\|_2^2] + \frac{1}{4}\beta^2 + 16n\beta^2\log(3/\delta)\right)\log\frac{6}{\delta}}{n}} + \frac{M\log\frac{6}{\delta}}{n} \qquad (25) \\ & + 16 \times 2^{\frac{3}{4}}\sqrt{e}\beta\left\lceil\log_2 n\right\rceil\log(3e/\delta) + 32\sqrt{e}\beta\left\lceil\log_2 n\right\rceil\sqrt{\log 3e/\delta}. \end{aligned}$$

According to inequality  $\sqrt{a+b} = \sqrt{a} + \sqrt{b}$  for any a, b > 0, with probability at least  $1 - \delta$ , we have

$$\begin{split} \|\nabla F(A(S)) - \nabla F_S(A(S))\|_2 \\ \leq & \sqrt{\frac{4\mathbb{E}_Z[\|\nabla f(A(S);Z)\|_2^2]\log\frac{6}{\delta}}{n}} + \sqrt{\frac{\left(\frac{1}{2}\beta^2 + 32n\beta^2\log(3/\delta)\right)\log\frac{6}{\delta}}{n}} + \frac{M\log\frac{6}{\delta}}{n} \\ & + 16 \times 2^{\frac{3}{4}}\sqrt{e}\beta\left\lceil\log_2 n\right\rceil\log\left(3e/\delta\right) + 32\sqrt{e}\beta\left\lceil\log_2 n\right\rceil\sqrt{\log 3e/\delta}. \end{split}$$

614 The proof is complete.

615

<sup>616</sup> *Proof of Remark 4.* According to the proof in Theorem 3, we have the following inequality that with <sup>617</sup> probability at least  $1 - \delta$ 

$$\begin{aligned} \|\nabla F(A(S)) - \nabla F_S(A(S))\|_2 \\ \leq & \sqrt{\frac{4\mathbb{E}_Z[\|\nabla f(A(S);Z)\|_2^2]\log\frac{6}{\delta}}{n}} + \sqrt{\frac{\left(\frac{1}{2}\beta^2 + 32n\beta^2\log(3/\delta)\right)\log\frac{6}{\delta}}{n}} + \frac{M\log\frac{6}{\delta}}{n} \quad (26) \\ & + 16 \times 2^{\frac{3}{4}}\sqrt{e}\beta\left\lceil\log_2 n\right\rceil\log\left(3e/\delta\right) + 32\sqrt{e}\beta\left\lceil\log_2 n\right\rceil\sqrt{\log 3e/\delta}. \end{aligned}$$

618 Since  $f(\mathbf{w})$  is  $\gamma$ -smooth, we have

$$\mathbb{E}_{Z}[\|\nabla f(A(S); Z)\|_{2}^{2}] \leq \mathbb{E}_{Z}[\|\nabla f(A(S); Z) - \nabla f(\mathbf{w}^{*}; Z)\|_{2}^{2} + \|\nabla f(\mathbf{w}^{*}; Z)\|_{2}^{2}] \leq \gamma^{2}\|A(S) - \mathbf{w}^{*}\|_{2}^{2} + \mathbb{E}_{Z}[\|\nabla f(\mathbf{w}^{*}; Z)\|_{2}^{2}]$$
(27)

619 Plugging (27) into (26), we have

$$\begin{aligned} \|\nabla F(A(S)) - \nabla F_{S}(A(S))\|_{2} \\ \leq & \sqrt{\frac{4(\gamma^{2}\|A(S) - \mathbf{w}^{*}\|_{2}^{2} + \mathbb{E}_{Z}[\|\nabla f(\mathbf{w}^{*}; Z)\|_{2}^{2}])\log\frac{6}{\delta}}{n}} + \sqrt{\frac{\left(\frac{1}{2}\beta^{2} + 32n\beta^{2}\log(3/\delta)\right)\log\frac{6}{\delta}}{n}}{n}} \\ & + \frac{M\log\frac{6}{\delta}}{n} + 16 \times 2^{\frac{3}{4}}\sqrt{e}\beta\left\lceil\log_{2}n\right\rceil\log\left(3e/\delta\right) + 32\sqrt{e}\beta\left\lceil\log_{2}n\right\rceil\sqrt{\log 3e/\delta}} \\ \leq & 2\gamma\|A(S) - \mathbf{w}^{*}\|_{2}\sqrt{\frac{\log\frac{6}{\delta}}{n}} + \sqrt{\frac{4\mathbb{E}_{Z}[\|\nabla f(\mathbf{w}^{*}; Z)\|_{2}^{2}]\log\frac{6}{\delta}}{n}} \\ & + \sqrt{\frac{\left(\frac{1}{2}\beta^{2} + 32n\beta^{2}\log(3/\delta)\right)\log\frac{6}{\delta}}{n}} + \frac{M\log\frac{6}{\delta}}{n} \\ & + 16 \times 2^{\frac{3}{4}}\sqrt{e}\beta\left\lceil\log_{2}n\right\rceil\log\left(3e/\delta\right) + 32\sqrt{e}\beta\left\lceil\log_{2}n\right\rceil\sqrt{\log 3e/\delta}, \end{aligned}$$
(28)

where the second inequality holds because  $\sqrt{a+b} + \sqrt{a} + \sqrt{b}$  for any a, b > 0, which means that  $\|\nabla F(A(S)) - \nabla F_S(A(S))\|_2$ 

$$\lesssim \beta \log n \log(1/\delta) + \frac{\log(1/\delta)}{n} + \sqrt{\frac{\mathbb{E}_Z\left[\nabla \|f(\mathbf{w}^*;Z)\|_2^2\right]\log(1/\delta)}{n}} + \|A(S) - \mathbf{w}^*\|\sqrt{\frac{\log(1/\delta)}{n}}.$$

621 The proof is complete.

622

623 Proof of Lemma 1. Inequality (28) implies that

$$\begin{split} \|\nabla F(A(S))\|_{2} &- \|\nabla F_{S}(A(S))\|_{2} \\ \leq & \sqrt{\frac{4(\gamma^{2}\|A(S) - \mathbf{w}^{*}\|_{2}^{2} + \mathbb{E}_{Z}[\|\nabla f(\mathbf{w}^{*};Z)\|_{2}^{2}])\log\frac{6}{\delta}}{n}} + \sqrt{\frac{\left(\frac{1}{2}\beta^{2} + 32n\beta^{2}\log(3/\delta)\right)\log\frac{6}{\delta}}{n}} \\ &+ \frac{M\log\frac{6}{\delta}}{n} + 16 \times 2^{\frac{3}{4}}\sqrt{e}\beta\left\lceil\log_{2}n\right\rceil\log\left(3e/\delta\right) + 32\sqrt{e}\beta\left\lceil\log_{2}n\right\rceil\sqrt{\log 3e/\delta} \\ \leq & 2\gamma\|A(S) - \mathbf{w}^{*}\|_{2}\sqrt{\frac{\log\frac{6}{\delta}}{n}} + \sqrt{\frac{4\mathbb{E}_{Z}[\|\nabla f(\mathbf{w}^{*};Z)\|_{2}^{2}]\log\frac{6}{\delta}}{n}} + \sqrt{\frac{\left(\frac{1}{2}\beta^{2} + 32n\beta^{2}\log(3/\delta)\right)\log\frac{6}{\delta}}{n}} \\ &+ \frac{M\log\frac{6}{\delta}}{n} + 16 \times 2^{\frac{3}{4}}\sqrt{e}\beta\left\lceil\log_{2}n\right\rceil\log\left(3e/\delta\right) + 32\sqrt{e}\beta\left\lceil\log_{2}n\right\rceil\sqrt{\log 3e/\delta}, \end{split}$$

<sup>624</sup> When  $F(\mathbf{w})$  satisfies the PL condition, there holds the following error bound property (refer to <sup>625</sup> Theorem 2 in Karimi et al. [2016])

$$\|\nabla F(A(S))\|_2 \ge \mu \|A(S) - \mathbf{w}^*\|_2.$$

#### 626 Thus, we have

$$\begin{split} \mu \|A(S) - \mathbf{w}^*\|_2 &\leq \|\nabla F(A(S))\|_2 \\ \leq \|\nabla F_S(A(S))\|_2 + 2\gamma \|A(S) - \mathbf{w}^*\|_2 \sqrt{\frac{\log \frac{6}{\delta}}{n}} + \sqrt{\frac{4\mathbb{E}_Z[\|\nabla f(\mathbf{w}^*; Z)\|_2^2]\log \frac{6}{\delta}}{n}} \\ &+ \sqrt{\frac{\left(\frac{1}{2}\beta^2 + 32n\beta^2\log(3/\delta)\right)\log \frac{6}{\delta}}{n}} + \frac{M\log \frac{6}{\delta}}{n} \\ &+ 16 \times 2^{\frac{3}{4}}\sqrt{e}\beta \left\lceil \log_2 n \right\rceil \log (3e/\delta) + 32\sqrt{e}\beta \left\lceil \log_2 n \right\rceil \sqrt{\log 3e/\delta}. \end{split}$$

627 When  $n \ge \frac{16\gamma^2 \log \frac{6}{\delta}}{\mu^2}$ , we have  $2\gamma \sqrt{\frac{\log \frac{6}{\delta}}{n}} \le \frac{\mu}{2}$ , then we can derive that

$$\begin{split} & \mu \|A(S) - \mathbf{w}^*\|_2 \le \|\nabla F(A(S))\|_2 \\ \le \|\nabla F_S(A(S))\|_2 + \frac{\mu}{2} \|A(S) - \mathbf{w}^*\|_2 + \sqrt{\frac{4\mathbb{E}_Z[\|\nabla f(\mathbf{w}^*; Z)\|_2^2]\log\frac{6}{\delta}}{n}} \\ & + \sqrt{\frac{\left(\frac{1}{2}\beta^2 + 32n\beta^2\log(3/\delta)\right)\log\frac{6}{\delta}}{n}} + \frac{M\log\frac{6}{\delta}}{n} \\ & + 16 \times 2^{\frac{3}{4}}\sqrt{e}\beta \left\lceil \log_2 n \right\rceil \log(3e/\delta) + 32\sqrt{e}\beta \left\lceil \log_2 n \right\rceil \sqrt{\log 3e/\delta}. \end{split}$$

628 This implies that

$$||A(S) - \mathbf{w}^*||_2$$

$$\leq \frac{2}{\mu} \Big( ||\nabla F_S(A(S))||_2 + \sqrt{\frac{4\mathbb{E}_Z[||\nabla f(\mathbf{w}^*; Z)||_2^2]\log\frac{6}{\delta}}{n}} + \sqrt{\frac{(\frac{1}{2}\beta^2 + 32n\beta^2\log(3/\delta))\log\frac{6}{\delta}}{n}} + \frac{M\log\frac{6}{\delta}}{n} + 16 \times 2^{\frac{3}{4}}\sqrt{e}\beta \lceil \log_2 n \rceil \log(3e/\delta) + 32\sqrt{e}\beta \lceil \log_2 n \rceil \sqrt{\log 3e/\delta} \Big).$$

$$(29)$$

Then, substituting (29) into (28), when  $n \ge \frac{16\gamma^2 \log \frac{6}{\delta}}{\mu^2}$ , with probability at least  $1 - \delta$ 

$$\begin{aligned} \|\nabla F(A(S)) - \nabla F_S(A(S))\| \\ \leq \|\nabla F_S(A(S))\| + 4\sqrt{\frac{\mathbb{E}_Z[\|\nabla f(\mathbf{w}^*; Z)\|^2]\log\frac{6}{\delta}}{n}} + 2\sqrt{\frac{\left(\frac{1}{2}\beta^2 + 32n\beta^2\log(3/\delta)\right)\log\frac{6}{\delta}}{n}} \\ + \frac{2M\log\frac{6}{\delta}}{n} + 32 \times 2^{\frac{3}{4}}\sqrt{e\beta} \left\lceil \log_2 n \right\rceil \log(3e/\delta) + 64\sqrt{e\beta} \left\lceil \log_2 n \right\rceil \sqrt{\log 3e/\delta}. \end{aligned}$$

630 The proof is complete.

Proof of Remark 5. Here we briefly prove the results given in Remark 5. Since F satisfies the PL condition with  $\mu$ , we have

$$F(A(S)) - F(\mathbf{w}^*) \le \frac{\|\nabla F(A(S))\|^2}{2\mu}, \quad \forall \mathbf{w} \in \mathcal{W}.$$
(30)

So to bound F(A(S)) - F(A(S)), we need to bound the term  $\|\nabla F(A(S))\|^2$ . And there holds

$$\|\nabla F(A(S))\|_{2}^{2} = 2 \|\nabla F(A(S)) - \nabla F_{S}(A(S))\|^{2} + 2\|\nabla F_{S}(A(S))\|_{2}^{2}.$$
 (31)

From Lemma 1, if f is M-Lipschitz and  $\gamma$ -smooth and F satisfies PL condition with  $\mu$ , for any  $\delta > 0$ , when  $n \ge \frac{16\gamma^2 \log \frac{6}{\delta}}{\mu^2}$ , with probability at least  $1 - \delta$ , there holds

$$\begin{aligned} \|\nabla F(A(S)) - \nabla F_S(A(S))\|_2 \\ &\leq \|\nabla F_S(A(S))\|_2 + C\left(\sqrt{\frac{2\mathbb{E}_Z[\|\nabla f(\mathbf{w}^*;Z)\|_2^2]\log\frac{6}{\delta}}{n}} + \frac{M\log\frac{6}{\delta}}{n} + e\beta \left\lceil \log_2 n \right\rceil \log\left(3e/\delta\right)\right) \\ &\leq \|\nabla F_S(A(S))\|_2 + C\left(\sqrt{\frac{8\gamma F(\mathbf{w}^*)\log\frac{6}{\delta}}{n}} + \frac{M\log\frac{6}{\delta}}{n} + e\beta \left\lceil \log_2 n \right\rceil \log\left(3e/\delta\right)\right), \end{aligned}$$

where C is a positive constant and the last inequality follows from Lemma 4.1 of Srebro et al. [2010] when f is nonegative and  $\gamma$ -smooth (see (44)).

638 Combing above inequality with (30), (31), we can derive that

$$F(A(S)) - F(\mathbf{w}^*) \lesssim \|\nabla F_S(A(S))\|_2 + \frac{F(\mathbf{w}^*)\log(1/\delta)}{n} + \frac{M\log^2(1/\delta)}{n^2} + \beta^2\log^2 n\log^2(1/\delta).$$

639 The proof is complete.

640

<sup>641</sup> *Proof of Lemma 2.* According to the proof in Theorem 3, we have the following inequality with <sup>642</sup> probability at least  $1 - \delta$ 

$$\begin{aligned} \|\nabla F(A(S)) - \nabla F_S(A(S))\|_2 \\ \leq & \sqrt{\frac{2\left(2\mathbb{E}_Z[\|\nabla f(A(S);Z)\|_2^2] + \frac{1}{4}\beta^2 + 16n\beta^2\log(3/\delta)\right)\log\frac{6}{\delta}}{n}} \\ & + \frac{M\log\frac{6}{\delta}}{n} + 16 \times 2^{\frac{3}{4}}\sqrt{e}\beta\left\lceil\log_2 n\right\rceil\log\left(3e/\delta\right) + 32\sqrt{e}\beta\left\lceil\log_2 n\right\rceil\sqrt{\log 3e/\delta}. \end{aligned}$$
(32)

Since SGC implies that  $\mathbb{E}_Z[\|\nabla f(\mathbf{w}; Z)\|_2^2] \le \rho \|\nabla F(\mathbf{w})\|_2^2$ , according to inequalities  $\sqrt{ab} \le \eta a + \frac{1}{\eta}b$ and  $\sqrt{a+b} \le \sqrt{a} + \sqrt{b}$  for any  $a, b, \eta > 0$ , we have the following inequality with probability at least  $1 - \delta$ 

$$\begin{split} \|\nabla F(A(S)) - \nabla F_S(A(S))\|_2 \\ \leq & \sqrt{\frac{2\left(2\rho\|\nabla F(A(S))\|_2^2 + \frac{1}{4}\beta^2 + 16n\beta^2\log(3/\delta)\right)\log\frac{6}{\delta}}{n}} \\ & + \frac{M\log\frac{6}{\delta}}{n} + 16 \times 2^{\frac{3}{4}}\sqrt{e}\beta\left\lceil\log_2 n\right\rceil\log\left(3e/\delta\right) + 32\sqrt{e}\beta\left\lceil\log_2 n\right\rceil\sqrt{\log 3e/\delta} \\ \leq & \sqrt{\frac{\left(\frac{1}{2}\beta^2 + 32n\beta^2\log(3/\delta)\right)\log\frac{6}{\delta}}{n}} + \frac{\eta}{1+\eta}\|\nabla F(A(S))\| + \frac{1+\eta}{\eta}\frac{4\rho M\log\frac{6}{\delta}}{n} \\ & + \frac{M\log\frac{6}{\delta}}{n} + 16 \times 2^{\frac{3}{4}}\sqrt{e}\beta\left\lceil\log_2 n\right\rceil\log\left(3e/\delta\right) + 32\sqrt{e}\beta\left\lceil\log_2 n\right\rceil\sqrt{\log 3e/\delta}. \end{split}$$

646 which implies that

$$\|\nabla F(A(S))\|_{2} \leq (1+\eta)\|\nabla F_{S}(A(S))\|_{2} + C\frac{1+\eta}{\eta}\left(\frac{M}{n}\log\frac{6}{\delta} + \beta\log n\log\frac{1}{\delta}\right).$$

647 The proof is complete.

## 648 **D** Proofs of ERM

649 Proof of Lemma 3. Since  $F_{S^{(i)}}(\mathbf{w}) = \frac{1}{n} \left( f(\mathbf{w}; z'_i) + \sum_{j \neq i} f(\mathbf{w}, z_j) \right)$ , we have

$$\begin{aligned} F_{S}(\hat{\mathbf{w}}^{*}(S^{(i)})) &- F_{S}(\hat{\mathbf{w}}^{*}(S)) \\ &= \frac{f(\hat{\mathbf{w}}^{*}(S^{(i)}); z_{i}) - f(\hat{\mathbf{w}}^{*}(S); z_{i})}{n} + \frac{\sum_{j \neq i} (f(\hat{\mathbf{w}}^{*}(S^{(i)}); z_{j}) - f(\hat{\mathbf{w}}^{*}(S); z_{j}))}{n} \\ &= \frac{f(\hat{\mathbf{w}}^{*}(S^{(i)}); z_{i}) - f(\hat{\mathbf{w}}^{*}(S); z_{i})}{n} + \frac{f(\hat{\mathbf{w}}^{*}(S); z_{i}') - f(\hat{\mathbf{w}}^{*}(S^{(i)}); z_{i}')}{n} \\ &+ \left(F_{S^{(i)}}(\hat{\mathbf{w}}^{*}(S^{(i)})) - F_{S^{(i)}}(\hat{\mathbf{w}}^{*}(S))\right) \\ &\leq \frac{f(\hat{\mathbf{w}}^{*}(S^{(i)}); z_{i}) - f(\hat{\mathbf{w}}^{*}(S); z_{i})}{n} + \frac{f(\hat{\mathbf{w}}^{*}(S); z_{i}') - f(\hat{\mathbf{w}}^{*}(S^{(i)}); z_{i}')}{n} \\ &\leq \frac{2M}{n} \|\hat{\mathbf{w}}^{*}(S^{(i)}) - \hat{\mathbf{w}}^{*}(S)\|_{2}, \end{aligned}$$

where the first inequality follows from the fact that  $\hat{\mathbf{w}}^*(S^{(i)})$  is the ERM of  $F_{S^{(i)}}$  and the second inequality follows from the Lipschitz property. Furthermore, for  $\hat{\mathbf{w}}^*(S^{(i)})$ , the convexity of f and the strongly-convex property of  $F_S$  imply that its closest optima point of  $F_S$  is  $\hat{\mathbf{w}}^*(S)$  (the global minimizer of  $F_S$  is unique). Then, there holds that

$$F_{S}(\hat{\mathbf{w}}^{*}(S^{(i)})) - F_{S}(\hat{\mathbf{w}}^{*}(S)) \geq \frac{\mu}{2} \|\hat{\mathbf{w}}^{*}(S^{(i)}) - \hat{\mathbf{w}}^{*}(S)\|_{2}^{2}.$$

654 Then we get

$$\frac{\mu}{2} \|\hat{\mathbf{w}}^*(S^{(i)}) - \hat{\mathbf{w}}^*(S)\|_2^2 \le F_S(\hat{\mathbf{w}}^*(S^{(i)})) - F_S(\hat{\mathbf{w}}^*(S)) \le \frac{2M}{n} \|\hat{\mathbf{w}}^*(S^{(i)}) - \hat{\mathbf{w}}^*(S)\|_2,$$

which implies that  $\|\hat{\mathbf{w}}^*(S^{(i)}) - \hat{\mathbf{w}}^*(S)\|_2 \le \frac{4M}{n\mu}$ . Combined with the smoothness property of f we obtain that for any  $S^{(i)}$  and S

$$\forall z \in \mathcal{Z}, \quad \left\| \nabla f(\hat{\mathbf{w}}^*(S^{(i)}); z) - \nabla f(\hat{\mathbf{w}}^*(S); z) \right\|_2 \le \frac{4M\gamma}{n\mu}.$$
ete.  $\Box$ 

657 The proof is complete.

658 Proof of Theorem 4. Since F is  $\mu$ -strongly convex, we have

$$F(\mathbf{w}) - F(\mathbf{w}^*) \le \frac{\|\nabla F(\mathbf{w})\|_2^2}{2\mu}, \quad \forall \mathbf{w} \in \mathcal{W}.$$
(33)

So to bound  $F(\hat{\mathbf{w}}^*) - F(\mathbf{w}^*)$ , we need to bound the term  $\|\nabla F(\hat{\mathbf{w}}^*)\|_2^2$ . And there holds

$$\|\nabla F(\hat{\mathbf{w}}^*)\|_2^2 = 2 \|\nabla F(\hat{\mathbf{w}}^*) - \nabla F_S(\hat{\mathbf{w}}^*)\|_2^2 + 2 \|\nabla F_S(\hat{\mathbf{w}}^*)\|_2^2.$$
(34)

From Lemma 1, if f is M-Lipschitz and  $\gamma$ -smooth and  $F_S$  is  $\mu$ -strongly convex, for any  $\delta > 0$ , when  $n \ge \frac{16\gamma^2 \log \frac{6}{\delta}}{\mu^2}$ , with probability at least  $1 - \delta$ , there holds

$$\begin{aligned} \|\nabla F(\hat{\mathbf{w}}^{*}) - \nabla F_{S}(\hat{\mathbf{w}}^{*})\|_{2} \\ &\leq \|\nabla F_{S}(\hat{\mathbf{w}}^{*})\|_{2} + C\left(\sqrt{\frac{2\mathbb{E}_{Z}[\|\nabla f(\mathbf{w}^{*};Z)\|_{2}^{2}]\log\frac{6}{\delta}}{n}} + \frac{M\log\frac{6}{\delta}}{n} + e\hat{\beta}\left\lceil\log_{2}n\right\rceil\log\left(3e/\delta\right)\right) \\ &\leq \|\nabla F_{S}(\hat{\mathbf{w}}^{*})\|_{2} + C\left(\sqrt{\frac{8\gamma F(\mathbf{w}^{*})\log\frac{6}{\delta}}{n}} + \frac{M\log\frac{6}{\delta}}{n} + e\hat{\beta}\left\lceil\log_{2}n\right\rceil\log\left(3e/\delta\right)\right), \end{aligned}$$

$$(35)$$

where the last inequality follows from Lemma 4.1 of Srebro et al. [2010] when f is nonegative (see (44)) and  $\gamma$ -smooth and  $\hat{\beta} = \|\nabla f(\hat{\mathbf{w}}^*(S); z) - \nabla f(\hat{\mathbf{w}}^*(S'); z)\|_2$ . C is a positive constant.

From Lemma 3, we have  $\|\nabla f(\hat{\mathbf{w}}^*(S); z) - \nabla f(\hat{\mathbf{w}}^*(S'); z)\|_2 \le \frac{4M\gamma}{n\mu}$ . Since  $\nabla F_S(\hat{\mathbf{w}}^*) = 0$ , we have  $\|\nabla F_S(\hat{\mathbf{w}}^*)\|_2 = 0$ , then we can derive that

$$F(\mathbf{w}) - F(\mathbf{w}^*) \lesssim \frac{F(\mathbf{w}^*) \log (1/\delta)}{n} + \frac{\log^2 n \log^2(1/\delta)}{n^2}.$$

666

# 667 E Proofs of PGD

Proof of Theorem 5. According to smoothness assumption and  $\eta = 1/\gamma$ , we can derive that

$$F_{S}(\mathbf{w}_{t+1}) - F_{S}(\mathbf{w}_{t})$$

$$\leq \langle \mathbf{w}_{t+1} - \mathbf{w}_{t}, \nabla F_{S}(\mathbf{w}_{t}) \rangle + \frac{\gamma}{2} \|\mathbf{w}_{t+1} - \mathbf{w}_{t}\|_{2}^{2}$$

$$= -\eta_{t} \|\nabla F_{S}(\mathbf{w}_{t})\|_{2}^{2} + \frac{\gamma}{2} \eta_{t}^{2} \|\nabla F_{S}(\mathbf{w}_{t})\|_{2}^{2}$$

$$= \left(\frac{\gamma}{2} \eta_{t}^{2} - \eta_{t}\right) \|\nabla F_{S}(\mathbf{w}_{t})\|_{2}^{2}$$

$$\leq -\frac{1}{2} \eta_{t} \|\nabla F_{S}(\mathbf{w}_{t})\|_{2}^{2}.$$

According to above inequality and the assumptions that  $F_S$  is  $\mu$ -strongly convex, we can prove that

$$F_{S}(\mathbf{w}_{t+1}) - F_{S}(\mathbf{w}_{t}) \leq -\frac{1}{2}\eta_{t} \|\nabla F_{S}(\mathbf{w}_{t})\|_{2}^{2} \leq -\mu\eta_{t}(F_{S}(\mathbf{w}_{t}) - F_{S}(\hat{\mathbf{w}}^{*})),$$

670 which implies that

$$F_S(\mathbf{w}_{t+1}) - F_S(\hat{\mathbf{w}}^*) \le (1 - \mu \eta_t)(F_S(\mathbf{w}_t) - F_S(\hat{\mathbf{w}}^*)).$$

According to the property for  $\gamma$ -smooth for  $F_S$  and the property for  $\mu$ -strongly convex for  $F_S$ , we have

$$\frac{1}{2\gamma} \|\nabla F_S(\mathbf{w})\|_2^2 \le F_S(\mathbf{w}) - F_S(\hat{\mathbf{w}}^*) \le \frac{1}{2\mu} \|\nabla F_S(\mathbf{w})\|_2^2$$

- 673 which means that  $\frac{\mu}{\gamma} \leq 1$ .
- <sup>674</sup> Then If  $\eta_t = 1/\gamma, 0 \le 1 \mu \eta_t < 1$ , taking over T iterations, we get

$$F_{S}(\mathbf{w}_{t+1}) - F_{S}(\hat{\mathbf{w}}^{*}) \le (1 - \mu \eta_{t})^{T} (F_{S}(\mathbf{w}_{t}) - F_{S}(\hat{\mathbf{w}}^{*})).$$
(36)

<sup>675</sup> Combined (36), the smoothness of  $F_S$  and the nonnegative property of f, it can be derive that

$$\|\nabla F_S(\mathbf{w}_{T+1}))\|_2^2 = O\left((1-\frac{\mu}{\gamma})^T\right).$$

Furthermore, since F is  $\mu$ -strongly convex, we have

$$F(\mathbf{w}) - F(\mathbf{w}^*) \le \frac{\|\nabla F(\mathbf{w})\|_2^2}{2\mu}, \quad \forall \mathbf{w} \in \mathcal{W}.$$
(37)

So to bound  $F(\mathbf{w}_{T+1}) - F(\mathbf{w}^*)$ , we need to bound the term  $\|\nabla F(\mathbf{w}_{T+1})\|_2^2$ . And there holds

$$\|\nabla F(\mathbf{w}_{T+1})\|_2^2 = 2 \|\nabla F(\mathbf{w}_{T+1}) - \nabla F_S(\mathbf{w}_{T+1})\|_2^2 + 2\|\nabla F_S(\mathbf{w}_{T+1})\|_2^2.$$
(38)

From Lemma 1, if f is M-Lipschitz and  $\gamma$ -smooth and  $F_S$  is  $\mu$ -strongly convex, for any  $\delta > 0$ , when  $n \ge \frac{16\gamma^2 \log \frac{6}{\delta}}{\mu^2}$ , with probability at least  $1 - \delta$ , there holds

$$\|\nabla F(\mathbf{w}_{T+1}) - \nabla F_S(\mathbf{w}_{T+1})\|_2$$

$$\leq \|\nabla F_S(\mathbf{w}_{T+1})\|_2 + C\left(\sqrt{\frac{2\mathbb{E}_Z[\|\nabla f(\mathbf{w}^*;Z)\|_2^2]\log\frac{6}{\delta}}{n}} + \frac{M\log\frac{6}{\delta}}{n} + e\beta\left\lceil\log_2 n\right\rceil\log\left(3e/\delta\right)\right)$$

$$\leq \|\nabla F_S(\mathbf{w}_{T+1})\|_2 + C\left(\sqrt{\frac{8\gamma F(\mathbf{w}^*)\log\frac{6}{\delta}}{n}} + \frac{M\log\frac{6}{\delta}}{n} + e\beta\left\lceil\log_2 n\right\rceil\log\left(3e/\delta\right)\right),$$
(30)

where the last inequality follows from Lemma 4.1 of Srebro et al. [2010] when f is nonegative and  $\gamma$ -smooth (see (44)) and  $\beta = \|\nabla f(\mathbf{w}_{T+1}(S); z) - \nabla f(\mathbf{w}_{T+1}(S'); z)\|_2$ . C is a positive constant.

From Lemma 4, we have  $\beta = \|\nabla f(\mathbf{w}_{T+1}(S); z) - \nabla f(\mathbf{w}_{T+1}(S'); z)\|_2 \leq \frac{2M\gamma}{n\mu}$ . Since  $\|\nabla F_S(\mathbf{w}_{T+1})\|_2 = O\left((1-\frac{\mu}{\gamma})^T\right)$ , then we can derive that

$$F(\mathbf{w}) - F(\mathbf{w}^*) \lesssim \left(1 - \frac{\mu}{\gamma}\right)^{2T} + \frac{F(\mathbf{w}^*)\log\left(1/\delta\right)}{n} + \frac{\log^2 n \log^2(1/\delta)}{n^2}$$

684 Let  $T \asymp \log n$ , we have

$$F(\mathbf{w}) - F(\mathbf{w}^*) \lesssim \frac{F(\mathbf{w}^*) \log(1/\delta)}{n} + \frac{\log^2 n \log^2(1/\delta)}{n^2}$$

685 The proof is complete.

## 686 F Proofs of SGD

- <sup>687</sup> We first introduce some necessary lemmata on the empirical risk.
- **Lemma 14** ([Lei and Tang, 2021]). Let  $\{\mathbf{w}_t\}_t$  be the sequence produced by SGD with  $\eta_t \leq \frac{1}{2\gamma}$  for all  $t \in \mathbb{N}$ . Suppose Assumption 1 hold. Assume for all z, the function  $\mathbf{w} \mapsto f(\mathbf{w}; z)$  is *M*-Lipschitz and  $\gamma$ -smooth. Then, for any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ , there holds that

$$\sum_{k=1}^t \eta_k \|\nabla F_S(\mathbf{w}_k)\|_2^2 = O\left(\log\frac{1}{\delta} + \sum_{k=1}^t \eta_k^2\right).$$

**Lemma 15** ([Lei and Tang, 2021]). Let  $\{\mathbf{w}_t\}_t$  be the sequence produced by SGD with  $\eta_t = \frac{2}{\mu(t+t_0)}$ such that  $t_0 \ge \max\{\frac{4\gamma}{\mu}, 1\}$  for all  $t \in \mathbb{N}$ . Suppose Assumption 1 hold. Assume for all z, the function  $\mathbf{w} \mapsto f(\mathbf{w}; z)$  is *M*-Lipschitz and  $\gamma$ -smooth and assume  $F_S$  satisfies *PL* condition with parameter  $\mu$ . Then, for any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ , there holds that

$$F_S(\mathbf{w}_{T+1}) - F_S(\hat{\mathbf{w}}^*) = O\left(\frac{\log(T)\log^3(1/\delta)}{T}\right)$$

**Lemma 16** ([Lei and Tang, 2021]). Let *e* be the base of the natural logarithm. There holds the following elementary inequalities.

697 • If 
$$\theta \in (0, 1)$$
, then  $\sum_{k=1}^{t} k^{-\theta} \le t^{1-\theta}/(1-\theta)$ ;

698

• If 
$$\theta = 1$$
, then  $\sum_{k=1}^{t} k^{-\theta} \le \log(et)$ ;

699 • If  $\theta > 1$ , then  $\sum_{k=1}^{t} k^{-\theta} \leq \frac{\theta}{\theta - 1}$ .

Proof of Lemma 5. We have known that  $F_{S^{(i)}}(\mathbf{w}) = \frac{1}{n} \left( f(\mathbf{w}; z'_i) + \sum_{j \neq i} f(\mathbf{w}; z_j) \right)$ . We denote  $\hat{\mathbf{w}}^*(S^{(i)})$  be the ERM of  $F_{S^{(i)}}(\mathbf{w})$  and  $\hat{\mathbf{w}}^*_S$  be the ERM of  $F_S(\mathbf{w})$ . From Lemma 3, we know that

$$\forall z \in \mathcal{Z}, \quad \left\| \nabla f(\hat{\mathbf{w}}^*(S^{(i)}); z) - f(\hat{\mathbf{w}}^*(S); z) \right\|_2 \le \frac{4M\gamma}{n\mu}.$$

Also, for  $\mathbf{w}_t$ , the convexity of f and the strongly-convex property implies that its closest optima point of  $F_S$  is  $\hat{\mathbf{w}}^*(S)$  (the global minimizer of  $F_S$  is unique). Then, there holds that

$$\frac{\mu}{2} \|\mathbf{w}_t - \hat{\mathbf{w}}^*(S)\|_2^2 \le F_S(\mathbf{w}_t) - F_S(\hat{\mathbf{w}}^*(S)) = \epsilon_{opt}(\mathbf{w}_t).$$

Thus we have  $\|\mathbf{w}_t - \hat{\mathbf{w}}^*(S)\|_2 \le \sqrt{\frac{2\epsilon_{opt}(\mathbf{w}_t)}{\mu}}$ . A similar relation holds between  $\hat{\mathbf{w}}^*(S^{(i)})$  and  $\mathbf{w}_t^i$ . Combined with the Lipschitz property of f we obtain that for  $\forall z \in \mathcal{Z}$ , there holds that

$$\begin{aligned} \left\|\nabla f(\mathbf{w}_{t};z) - \nabla f(\mathbf{w}_{t}^{i};z)\right\|_{2} \\ &\leq \left\|\nabla f(\mathbf{w}_{t};z) - \nabla f(\hat{\mathbf{w}}^{*}(S);z)\right\|_{2} + \left\|\nabla f(\hat{\mathbf{w}}^{*}(S);z) - \nabla f(\hat{\mathbf{w}}^{*}(S^{(i)});z)\right\|_{2} \\ &+ \left\|\nabla f(\hat{\mathbf{w}}^{*}(S^{(i)});z) - \nabla f(\mathbf{w}_{t}^{i};z)\right\|_{2} \\ &\leq \gamma \|\mathbf{w}_{t} - \hat{\mathbf{w}}^{*}(S)\|_{2} + \frac{4M\gamma}{n\mu} + \gamma \|\hat{\mathbf{w}}^{*}(S^{(i)}) - \mathbf{w}_{t}^{i}\|_{2} \\ &\leq \gamma \sqrt{\frac{2\epsilon_{opt}(\mathbf{w}_{t})}{\mu}} + \frac{4M\gamma}{n\mu} + \gamma \sqrt{\frac{2\epsilon_{opt}(\mathbf{w}_{t}^{i})}{\mu}}. \end{aligned}$$

According to Lemma 15, for any dataset S, the optimization error  $\epsilon_{opt}(\mathbf{w}_t)$  is uniformly bounded by the same upper bound. Therefore, we write  $\|\nabla f(\mathbf{w}_t; z) - \nabla f(\mathbf{w}_t^i; z)\|_2 \le 2\gamma \sqrt{\frac{2\epsilon_{opt}(\mathbf{w}_t)}{\mu}} + \frac{4M\gamma}{n\mu}$ here.

709 The proof is complete.

- 710 Now We begin to prove Lemma 6.
- *Proof of Lemma 6.* If f is L-Lipschitz and  $\gamma$ -smooth and  $F_S$  is  $\mu$ -strongly convex. According to Lemma 1, we know that for all  $\mathbf{w} \in \mathcal{W}$  and any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta/2$ , when  $n > \frac{16\gamma^2 \log \frac{6}{\delta}}{2}$  we have

$$\begin{aligned} & \left(\sum_{t=1}^{T} \eta_{t}\right)^{-1} \sum_{t=1}^{T} \eta_{t} \|\nabla F(\mathbf{w}_{t})\|_{2}^{2} \\ & \leq & 16 \left(\sum_{t=1}^{T} \eta_{t}\right)^{-1} \sum_{t=1}^{T} \eta_{t} \|\nabla F_{S}(\mathbf{w}_{t})\|_{2}^{2} + \frac{4C^{2}L^{2}\log^{2}\frac{6}{\delta}}{n^{2}} + \frac{8C^{2}\mathbb{E}_{Z}[\|\nabla f(\mathbf{w}^{*};Z)\|_{2}^{2}]\log^{2}\frac{6}{\delta}}{n} \quad (40) \\ & + \left(\sum_{t=1}^{T} \eta_{t}\right)^{-1} \sum_{t=1}^{T} \eta_{t}C^{2}e^{2}\beta_{t}^{2} \left\lceil \log_{2}n \right\rceil^{2}\log^{2}(3e/\delta), \end{aligned}$$

where  $\beta_t = \left\| \nabla f(\mathbf{w}_t; z) - \nabla f(\mathbf{w}_t^i; z) \right\|_2$  and C is a positive constant.

ŀ

From Lemma 5, we have 
$$\left\|\nabla f(\mathbf{w}_t; z) - \nabla f(\mathbf{w}_t^i; z)\right\|_2 \le 2\gamma \sqrt{\frac{2\epsilon_{opt}(\mathbf{w}_t)}{\mu} + \frac{4M\gamma}{n\mu}}$$
, thus

$$\begin{aligned} \beta_t^2 &= \left\| \nabla f(\mathbf{w}_t; z) - \nabla f(\mathbf{w}_t^i; z) \right\|_2^2 \\ &\leq \left( 2\gamma \sqrt{\frac{2\epsilon_{opt}(\mathbf{w}_t)}{\mu}} + \frac{4M\gamma}{n\mu} \right)^2 \\ &\leq \frac{16\gamma^2 (F_S(\mathbf{w}_t) - F_S(\hat{\mathbf{w}}^*(S)))}{\mu} + \frac{32M^2\gamma^2}{n^2\mu^2} \\ &\leq \frac{8\gamma^2 \|\nabla F_S(\mathbf{w}_t)\|_2^2}{\mu^2} + \frac{32M^2\gamma^2}{n^2\mu^2}, \end{aligned}$$
(41)

where the second inequality holds from Cauchy-Bunyakovsky-Schwarz inequality and the second inequality satisfies because  $F_S$  is  $\mu$ -strongly convex. Plugging (41) into (40), with probability at least  $1 - \delta/2$ , when  $n > \frac{16\gamma^2 \log \frac{6}{\delta}}{\mu^2}$ , we have

$$\left(\sum_{t=1}^{T} \eta_{t}\right)^{-1} \sum_{t=1}^{T} \eta_{t} \|\nabla F(\mathbf{w}_{t})\|_{2}^{2} \\
\leq \left(16 + \frac{8\gamma^{2}C^{2}e^{2}\left\lceil\log_{2}n\right]^{2}\log^{2}\left(6e/\delta\right)}{\mu^{2}}\right) \left(\sum_{t=1}^{T} \eta_{t}\right)^{-1} \sum_{t=1}^{T} \eta_{t} \|\nabla F_{S}(\mathbf{w}_{t})\|_{2}^{2} \\
+ \frac{4C^{2}L^{2}\log^{2}\frac{12}{\delta}}{n^{2}} + \frac{8C^{2}\mathbb{E}_{Z}[\|\nabla f(\mathbf{w}^{*};Z)\|_{2}^{2}]\log^{2}\frac{12}{\delta}}{n} + \frac{32L^{2}\gamma^{2}C^{2}e^{2}\left\lceil\log_{2}n\right]^{2}\log^{2}\left(6e/\delta\right)}{n^{2}\mu^{2}},$$
(42)

When  $\eta_t = \eta_1 t^{-\theta}, \theta \in (0, 1)$ , with  $\eta_1 \leq \frac{1}{2\beta}$  and Assumption 1, according to Lemma 14 and Lemma 16, we obtain the following inequality with probability at least  $1 - \delta/2$ ,

$$\left(\sum_{t=1}^{T} \eta_t\right)^{-1} \sum_{t=1}^{T} \eta_t \|\nabla F_S(\mathbf{w}_t)\|^2 = \begin{cases} O\left(\frac{\log(1/\delta)}{T-\theta}\right), & \text{if } \theta < 1/2\\ O\left(\frac{\log(1/\delta)}{T^{-\frac{1}{2}}}\right), & \text{if } \theta = 1/2\\ O\left(\frac{\log(1/\delta)}{T^{\theta-1}}\right), & \text{if } \theta > 1/2. \end{cases}$$
(43)

On the other hand, when f is nonegative and  $\gamma$ -smooth, from Lemma 4.1 of Srebro et al. [2010], we have

$$\|\nabla f(\mathbf{w}^*; z)\|_2^2 \le 4\gamma f(\mathbf{w}^*; z)$$

723 which implies that

$$\mathbb{E}_{Z}[\|\nabla f(\mathbf{w}^{*}; Z)\|_{2}^{2}] \leq 4\gamma \mathbb{E}_{Z} f(\mathbf{w}^{*}; Z) = 4\gamma F(\mathbf{w}^{*}).$$
(44)

Plugging (44), (43) into (42), with probability at least  $1 - \delta$ , we derive that

$$\begin{pmatrix} \sum_{t=1}^{T} \eta_t \end{pmatrix}^{-1} \sum_{t=1}^{T} \eta_t \|\nabla F(\mathbf{w}_t)\|_2^2 \\ = \begin{cases} O\left(\frac{\log^2 n \log^3(1/\delta)}{T^{-\theta}}\right) + O\left(\frac{\log^2 n \log^2(1/\delta)}{n^2} + \frac{F(\mathbf{w}^*) \log^2(1/\delta)}{n}\right), & \text{if } \theta < 1/2 \\ O\left(\frac{\log^2 n \log^3(1/\delta)}{T^{-\frac{1}{2}}}\right) + O\left(\frac{\log^2 n \log^2(1/\delta)}{n^2} + \frac{F(\mathbf{w}^*) \log^2(1/\delta)}{n}\right), & \text{if } \theta = 1/2 \\ O\left(\frac{\log^2 n \log^3(1/\delta)}{T^{\theta-1}}\right) + O\left(\frac{\log^2 n \log^2(1/\delta)}{n^2} + \frac{F(\mathbf{w}^*) \log^2(1/\delta)}{n}\right), & \text{if } \theta > 1/2. \end{cases}$$

When  $\theta < 1/2$ , we set  $T \simeq n^{\frac{2}{\theta}}$  and assume  $F(\mathbf{w}^*) = O(\frac{1}{n})$ , then we obtain the following result with probability at least  $1 - \delta$ 

$$\left(\sum_{t=1}^{T} \eta_t\right)^{-1} \sum_{t=1}^{T} \eta_t \|\nabla F(\mathbf{w}_t)\|_2^2 = O\left(\frac{\log^2 n \log^3(1/\delta)}{n^2}\right)$$

When  $\theta = 1/2$ , we set  $T \simeq n^4$  and assume  $F(\mathbf{w}^*) = O(\frac{1}{n})$ , then we obtain the following result with probability at least  $1 - \delta$ 

$$\left(\sum_{t=1}^{T} \eta_t\right)^{-1} \sum_{t=1}^{T} \eta_t \|\nabla F(\mathbf{w}_t)\|_2^2 = O\left(\frac{\log^2 n \log^3(1/\delta)}{n^2}\right)$$

When  $\theta > 1/2$ , we set  $T \approx n^{\frac{2}{1-\theta}}$  and assume  $F(\mathbf{w}^*) = O(\frac{1}{n})$ , then we obtain the following result with probability at least  $1 - \delta$ 

$$\left(\sum_{t=1}^{T} \eta_t\right)^{-1} \sum_{t=1}^{T} \eta_t \|\nabla F(\mathbf{w}_t)\|_2^2 = O\left(\frac{\log^2 n \log^3(1/\delta)}{n^2}\right)$$

731 The proof is complete.

732

Proof of Theorem 7. Since F is  $\mu$ -strongly convex, we have

$$F(\mathbf{w}) - F(\mathbf{w}^*) \le \frac{\|\nabla F(\mathbf{w})\|_2^2}{2\mu}, \quad \forall \mathbf{w} \in \mathcal{W}.$$
 (45)

So to bound  $F(\mathbf{w}_{T+1}) - F(\mathbf{w}^*)$ , we need to bound the term  $\|\nabla F(\mathbf{w}_{T+1})\|_2^2$ . And there holds

$$\|\nabla F(\mathbf{w}_{T+1})\|_{2}^{2} = 2 \|\nabla F(\mathbf{w}_{T+1}) - \nabla F_{S}(\mathbf{w}_{T+1})\|^{2} + 2 \|\nabla F_{S}(\mathbf{w}_{T+1})\|_{2}^{2}.$$
 (46)

From Lemma 1, if f is L-Lipschitz and  $\gamma$ -smooth and  $F_S$  is  $\mu$ -strongly convex, for all  $\mathbf{w} \in \mathcal{W}$  and any  $\delta > 0$ , when  $n \ge \frac{16\gamma^2 \log \frac{6}{\delta}}{\mu^2}$ , with probability at least  $1 - \delta/2$ , there holds

$$\begin{aligned} \|\nabla F(\mathbf{w}_{T+1}) - \nabla F_S(\mathbf{w}_{T+1})\|_2 \\ &\leq \|\nabla F_S(\mathbf{w}_{T+1})\|_2 + C\left(\sqrt{\frac{2\mathbb{E}_Z[\|\nabla f(\mathbf{w}^*;Z)\|_2^2]\log\frac{12}{\delta}}{n}} + \frac{M\log\frac{12}{\delta}}{n} + e\beta\left\lceil\log_2 n\right\rceil\log\left(6e/\delta\right)\right) \\ &\leq \|\nabla F_S(\mathbf{w}_{T+1})\|_2 + C\left(\sqrt{\frac{8\gamma F(\mathbf{w}^*)\log\frac{12}{\delta}}{n}} + \frac{M\log\frac{12}{\delta}}{n} + e\beta\left\lceil\log_2 n\right\rceil\log\left(6e/\delta\right)\right), \end{aligned}$$

$$(47)$$

where the last inequality follows from Lemma 4.1 of Srebro et al. [2010] when f is nonegative and  $\gamma$ -smooth (see (44)) and C is a positive constant. Then we can derive that

$$\begin{aligned} \|\nabla F(\mathbf{w}_{T+1}) - \nabla F_S(\mathbf{w}_{T+1})\|_2^2 \\ \leq 4\|\nabla F_S(\mathbf{w}_{T+1})\|_2^2 + \frac{32C^2\gamma F(\mathbf{w}^*)\log\frac{12}{\delta}}{n} + \frac{4M^2C^2\log^2\frac{12}{\delta}}{n^2} + 4e^2\beta_{T+1}^2\left\lceil\log_2 n\right\rceil^2\log^2(6e/\delta). \end{aligned}$$
(48)

From Lemma 5, we have  $\left\|\nabla f(\mathbf{w}_t; z) - \nabla f(\mathbf{w}_t^i; z)\right\|_2 \le 2\gamma \sqrt{\frac{2\epsilon_{opt}(\mathbf{w}_t)}{\mu}} + \frac{4M\gamma}{n\mu}$ , thus

$$\beta_t^2 = \left\| \nabla f(\mathbf{w}_t; z) - \nabla f(\mathbf{w}_t^i; z) \right\|_2^2$$

$$\leq \left( 2\gamma \sqrt{\frac{2\epsilon_{opt}(\mathbf{w}_t)}{\mu}} + \frac{4M\gamma}{n\mu} \right)^2$$

$$\leq \frac{16\gamma^2 (F_S(\mathbf{w}_t) - F_S(\hat{\mathbf{w}}^*(S)))}{\mu} + \frac{32M^2\gamma^2}{n^2\mu^2}$$

$$\leq \frac{8\gamma^2 \|\nabla F_S(\mathbf{w}_t)\|_2^2}{\mu^2} + \frac{32M^2\gamma^2}{n^2\mu^2},$$
(49)

- where the second inequality holds from Cauchy-Bunyakovsky-Schwarz inequality and the second inequality satisfies because  $F_S$  is  $\mu$ -strongly convex.
- Plugging (49) into (48), with probability at least  $1 \delta/2$ , when , we have

$$\begin{aligned} \|\nabla F(\mathbf{w}_{T+1}) - \nabla F_S(\mathbf{w}_{T+1})\|_2^2 \\ &\leq \left(4 + 32e^2 \left\lceil \log_2 n \right\rceil^2 \log^2 (6e/\delta) \right) \|\nabla F_S(\mathbf{w}_{T+1})\|_2^2 + \frac{32C^2 \gamma F(\mathbf{w}^*) \log \frac{6}{\delta}}{n} \\ &+ \frac{4L^2 C^2 \log^2 \frac{12}{\delta}}{n^2} + \frac{128M^2 \gamma^2 e^2 \left\lceil \log_2 n \right\rceil^2 \log^2 (6e/\delta)}{n^2 \mu^2}. \end{aligned}$$
(50)

According to the smoothness property of  $F_S$  and Lemma 15, it can be derived that with propability at least  $1 - \delta/2$ 

$$\|\nabla F_S(\mathbf{w}_{T+1})\|_2^2 = O\left(\frac{\log T \log^3(1/\delta)}{T}\right).$$
 (51)

<sup>745</sup> Substituting (51), (50) into (46), we derive that

$$\|\nabla F(\mathbf{w}_{T+1})\|_{2}^{2} = O\left(\frac{\left\lceil \log_{2} n \right\rceil^{2} \log T \log^{5}(1/\delta)}{T}\right) + O\left(\frac{\left\lceil \log_{2} n \right\rceil^{2} \log^{2}(1/\delta)}{n^{2}} + \frac{F(\mathbf{w}^{*}) \log(1/\delta)}{n}\right).$$
(52)

Further substituting (52) into (45) and choosing  $T \asymp n^2$ , we finally obtain that when n, with probability at least  $1 - \delta$ 

$$F(\mathbf{w}_{T+1}) - F(\mathbf{w}^*) = O\left(\frac{\log^4 n \log^5(1/\delta)}{n^2} + \frac{F(\mathbf{w}^*) \log(1/\delta)}{n}\right).$$

748

# 749 NeurIPS Paper Checklist

750	1.	Claims
751		Question: Do the main claims made in the abstract and introduction accurately reflect the
752		paper's contributions and scope?
753		Answer: [Yes]
754		Institution: We have checked that the abstract and introduction accurately reflect our
755		contributions and scope.
756		Guidelines:
757		• The answer NA means that the abstract and introduction do not include the claims
758		made in the paper.
759		• The abstract and/or introduction should clearly state the claims made, including the
760		contributions made in the paper and important assumptions and limitations. A No or
761		NA answer to this question will not be perceived well by the reviewers.
762		• The claims made should match theoretical and experimental results, and reflect how
763		much the results can be expected to generalize to other settings.
764		• It is fine to include aspirational goals as motivation as long as it is clear that these goals
765		are not attained by the paper.
766	2.	Limitations
767		Question: Does the paper discuss the limitations of the work performed by the authors?
768		Answer: [Yes]
769		Justification: We have clearly stated the required assumptions for each theorem and lemma,
770		and the conditions for the assumptions to hold are also stated in the main text.
771		Guidelines:
772		• The answer NA means that the paper has no limitation while the answer No means that
773		the paper has limitations, but those are not discussed in the paper.
774		• The authors are encouraged to create a separate "Limitations" section in their paper.
775		• The paper should point out any strong assumptions and how robust the results are to
776		violations of these assumptions (e.g., independence assumptions, noiseless settings,
777		model well-specification, asymptotic approximations only holding locally). The authors
778		should reflect on how these assumptions might be violated in practice and what the
779		implications would be.
780		• The authors should reflect on the scope of the claims made, e.g., if the approach was
781 782		depend on implicit assumptions, which should be articulated.
783		• The authors should reflect on the factors that influence the performance of the approach.
784		For example, a facial recognition algorithm may perform poorly when image resolution
785		is low or images are taken in low lighting. Or a speech-to-text system might not be
786		used reliably to provide closed captions for online lectures because it fails to handle
787		technical jargon.
788		• The authors should discuss the computational efficiency of the proposed algorithms
789		and how they scale with dataset size.
790		• If applicable, the authors should discuss possible limitations of their approach to
791		address problems of privacy and fairness.
792		• While the authors might fear that complete honesty about limitations might be used by
793		reviewers as grounds for rejection, a worse outcome might be that reviewers discover
794		limitations that aren't acknowledged in the paper. The authors should use their best
795		judgment and recognize that individual actions in favor of transparency play an impor-
796 797		will be specifically instructed to not penalize honesty concerning limitations.
709	3	Theory Assumptions and Proofs
1 90	5.	Ouestion: For each theoretical result does the nerver provide the full set of accurations and
/99		Question. For each medical result, does the paper provide the full set of assumptions and

a complete (and correct) proof? 800

801	Answer: [Yes]
802 803	Justification: We have clearly stated the required assumptions for each theorem and lemma, and all proofs are provided in the appendix.
804	Guidelines:
805	• The answer NA means that the paper does not include theoretical results
806	• All the theorems formulas and proofs in the paper should be numbered and cross-
807	referenced.
808	• All assumptions should be clearly stated or referenced in the statement of any theorems.
809	• The proofs can either appear in the main paper or the supplemental material, but if
810	they appear in the supplemental material, the authors are encouraged to provide a short
811	proof sketch to provide intuition.
812	• Inversely, any informal proof provided in the core of the paper should be complemented
813	by formal proofs provided in appendix or supplemental material.
814	<ul> <li>Theorems and Lemmas that the proof relies upon should be properly referenced.</li> </ul>
815 4	. Experimental Result Reproducibility
816	Ouestion: Does the paper fully disclose all the information needed to reproduce the main ex-
817	perimental results of the paper to the extent that it affects the main claims and/or conclusions
818	of the paper (regardless of whether the code and data are provided or not)?
819	Answer: [NA]
820	Justification: This paper focuses on learning theory.
821	Guidelines:
822	• The answer NA means that the paper does not include experiments.
823	• If the paper includes experiments, a No answer to this question will not be perceived
824	well by the reviewers: Making the paper reproducible is important, regardless of
825	whether the code and data are provided or not.
826	• If the contribution is a dataset and/or model, the authors should describe the steps taken
827	to make their results reproducible or verifiable.
828	• Depending on the contribution, reproducibility can be accomplished in various ways.
829	For example, if the contribution is a novel architecture, describing the architecture fully
830	might suffice, or if the contribution is a specific model and empirical evaluation, it may
831	dataset or provide access to the model. In general releasing code and data is often
833	one good way to accomplish this, but reproducibility can also be provided via detailed
834	instructions for how to replicate the results, access to a hosted model (e.g., in the case
835	of a large language model), releasing of a model checkpoint, or other means that are
836	appropriate to the research performed.
837	• While NeurIPS does not require releasing code, the conference does require all submis-
838	sions to provide some reasonable avenue for reproducibility, which may depend on the
839	nature of the contribution. For example
840	(a) If the contribution is primarily a new algorithm, the paper should make it clear how
841	to reproduce that algorithm.
842	(b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully
843	(a) If the contribution is a new model (a $g_{-}$ a large language model), then there should
044 845	either be a way to access this model for reproducing the results or a way to reproduce
846	the model (e.g., with an open-source dataset or instructions for how to construct
847	the dataset).
848	(d) We recognize that reproducibility may be tricky in some cases, in which case
849	authors are welcome to describe the particular way they provide for reproducibility.
850	In the case of closed-source models, it may be that access to the model is limited in
851	some way (e.g., to registered users), but it should be possible for other researchers
852	to have some path to reproducing or verifying the results.
853	5. Open access to data and code

854 855 856	Question: Does the paper provide open access to the data and code, with sufficient instruc- tions to faithfully reproduce the main experimental results, as described in supplemental material?
857	Answer: [NA]
007	Justification: This paper focuses on learning theory and does not include experiments
858	Justification. This paper focuses on learning meory and does not include experiments.
859	Guidelines:
860	• The answer NA means that paper does not include experiments requiring code.
861 862	• Please see the NeurIPS code and data submission guidelines (https://nips.cc/ public/guides/CodeSubmissionPolicy) for more details.
863	• While we encourage the release of code and data, we understand that this might not be
864	possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not
865 866	including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
867	• The instructions should contain the exact command and environment needed to run to
868 869	reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details
870	• The authors should provide instructions on data access and preparation including how
871	to access the raw data, preprocessed data, intermediate data, and generated data, etc.
872	• The authors should provide scripts to reproduce all experimental results for the new
873	proposed method and baselines. If only a subset of experiments are reproducible, they
874	should state which ones are omitted from the script and why.
875	• At submission time, to preserve anonymity, the authors should release anonymized
876	versions (if applicable).
877 878	• Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.
879	6. Experimental Setting/Details
880	Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
881	parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
882	results?
883	Answer: [NA]
884	Justification: This paper focuses on learning theory and does not include experiments.
885	Guidelines:
886	• The answer NA means that the paper does not include experiments.
887	• The experimental setting should be presented in the core of the paper to a level of detail
888	that is necessary to appreciate the results and make sense of them.
889 890	• The full details can be provided either with the code, in appendix, or as supplemental material.
891	7. Experiment Statistical Significance
892	Question: Does the paper report error bars suitably and correctly defined or other appropriate
893	information about the statistical significance of the experiments?
894	Answer: [NA]
895	Justification: This paper focuses on learning theory and does not include experiments.
896	Guidelines:
897	• The answer NA means that the paper does not include experiments.
898	• The authors should answer "Yes" if the results are accompanied by error bars, confi-
899	dence intervals, or statistical significance tests, at least for the experiments that support
900	the main claims of the paper.
901	• The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or every
902 903	run with given experimental conditions).
904	• The method for calculating the error bars should be explained (closed form formula.
905	call to a library function, bootstrap, etc.)

906		• The assumptions made should be given (e.g., Normally distributed errors).
907		• It should be clear whether the error bar is the standard deviation or the standard error
908		of the mean.
909		• It is OK to report 1-sigma error bars, but one should state it. The authors should
910		preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
911		of Normality of errors is not verified.
912		• For asymmetric distributions, the authors should be careful not to show in tables or
913		figures symmetric error bars that would yield results that are out of range (e.g. negative
914		error rates).
915 916		• If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.
917	8.	Experiments Compute Resources
918		Question: For each experiment, does the paper provide sufficient information on the com-
919		puter resources (type of compute workers, memory, time of execution) needed to reproduce
920		the experiments?
921		Answer: [NA]
922		Justification: This paper focuses on learning theory and does not include experiments.
923		Guidelines:
924		• The answer NA means that the paper does not include experiments.
925		• The paper should indicate the type of compute workers CPU or GPU, internal cluster,
926		or cloud provider, including relevant memory and storage.
927		• The paper should provide the amount of compute required for each of the individual
928		experimental runs as well as estimate the total compute.
929		• The paper should disclose whether the full research project required more compute
930		than the experiments reported in the paper (e.g., preliminary or failed experiments that
021		
551		didn't make it into the paper).
932	9.	didn't make it into the paper). Code Of Ethics
932 933	9.	didn't make it into the paper). <b>Code Of Ethics</b> Question: Does the research conducted in the paper conform, in every respect, with the
932 933 934	9.	didn't make it into the paper). <b>Code Of Ethics</b> Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?
932 933 934 935	9.	didn't make it into the paper). <b>Code Of Ethics</b> Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines? Answer: [Yes]
932 933 934 935 936	9.	didn't make it into the paper). <b>Code Of Ethics</b> Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines? Answer: [Yes] Justification: We have reviewed the NeurIPS Code of Ethics.
932 933 934 935 936 937	9.	didn't make it into the paper).         Code Of Ethics         Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?         Answer: [Yes]         Justification: We have reviewed the NeurIPS Code of Ethics.         Guidelines:
932 933 934 935 936 937 938	9.	didn't make it into the paper).         Code Of Ethics         Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?         Answer: [Yes]         Justification: We have reviewed the NeurIPS Code of Ethics.         Guidelines:         • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
932 933 934 935 936 937 938 939	9.	<ul> <li>didn't make it into the paper).</li> <li>Code Of Ethics</li> <li>Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?</li> <li>Answer: [Yes]</li> <li>Justification: We have reviewed the NeurIPS Code of Ethics.</li> <li>Guidelines: <ul> <li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.</li> <li>If the authors answer No, they should explain the special circumstances that require a</li> </ul> </li> </ul>
932 933 934 935 936 937 938 939 940	9.	<ul> <li>didn't make it into the paper).</li> <li>Code Of Ethics</li> <li>Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?</li> <li>Answer: [Yes]</li> <li>Justification: We have reviewed the NeurIPS Code of Ethics.</li> <li>Guidelines: <ul> <li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.</li> <li>If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.</li> </ul> </li> </ul>
932 933 934 935 936 937 938 939 940 941	9.	<ul> <li>didn't make it into the paper).</li> <li>Code Of Ethics</li> <li>Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?</li> <li>Answer: [Yes]</li> <li>Justification: We have reviewed the NeurIPS Code of Ethics.</li> <li>Guidelines: <ul> <li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.</li> <li>If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.</li> <li>The authors should make sure to preserve anonymity (e.g., if there is a special consid-</li> </ul> </li> </ul>
932 933 934 935 936 937 938 939 940 941 942	9.	<ul> <li>didn't make it into the paper).</li> <li>Code Of Ethics</li> <li>Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?</li> <li>Answer: [Yes]</li> <li>Justification: We have reviewed the NeurIPS Code of Ethics.</li> <li>Guidelines: <ul> <li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.</li> <li>If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.</li> <li>The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).</li> </ul> </li> </ul>
932 933 934 935 936 937 938 939 940 941 942 943	9.	<ul> <li>didn't make it into the paper).</li> <li>Code Of Ethics</li> <li>Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?</li> <li>Answer: [Yes]</li> <li>Justification: We have reviewed the NeurIPS Code of Ethics.</li> <li>Guidelines: <ul> <li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.</li> <li>If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.</li> <li>The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).</li> </ul> </li> <li>Broader Impacts</li> </ul>
932 932 933 934 935 936 937 938 939 940 941 942 943 944	9.	<ul> <li>didn't make it into the paper).</li> <li>Code Of Ethics</li> <li>Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?</li> <li>Answer: [Yes]</li> <li>Justification: We have reviewed the NeurIPS Code of Ethics.</li> <li>Guidelines: <ul> <li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.</li> <li>If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.</li> <li>The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).</li> </ul> </li> <li>Broader Impacts <ul> <li>Question: Does the paper discuss both potential positive societal impacts and negative</li> </ul> </li> </ul>
932 933 934 935 936 937 938 939 940 941 942 943 944 945	9.	<ul> <li>didn't make it into the paper).</li> <li>Code Of Ethics</li> <li>Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?</li> <li>Answer: [Yes]</li> <li>Justification: We have reviewed the NeurIPS Code of Ethics.</li> <li>Guidelines: <ul> <li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.</li> <li>If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.</li> <li>The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).</li> </ul> </li> <li>Broader Impacts <ul> <li>Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?</li> </ul> </li> </ul>
932 933 934 935 936 937 938 939 940 941 942 943 944 945 946	9.	<ul> <li>didn't make it into the paper).</li> <li>Code Of Ethics</li> <li>Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?</li> <li>Answer: [Yes]</li> <li>Justification: We have reviewed the NeurIPS Code of Ethics.</li> <li>Guidelines: <ul> <li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.</li> <li>If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.</li> <li>The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).</li> </ul> </li> <li>Broader Impacts <ul> <li>Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?</li> </ul> </li> </ul>
932 933 934 935 936 937 938 939 940 941 942 943 944 945 946 947	9.	<ul> <li>didn't make it into the paper).</li> <li>Code Of Ethics</li> <li>Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?</li> <li>Answer: [Yes]</li> <li>Justification: We have reviewed the NeurIPS Code of Ethics.</li> <li>Guidelines: <ul> <li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.</li> <li>If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.</li> <li>The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).</li> </ul> </li> <li>Broader Impacts <ul> <li>Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?</li> </ul> </li> <li>Answer: [NA]</li> <li>Justification: This paper focuses on learning theory and there is no societal impact of the work performed.</li> </ul>
932 932 933 934 935 936 937 938 939 940 941 942 943 944 945 946 947 948 949	9.	<ul> <li>didn't make it into the paper).</li> <li>Code Of Ethics</li> <li>Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?</li> <li>Answer: [Yes]</li> <li>Justification: We have reviewed the NeurIPS Code of Ethics.</li> <li>Guidelines: <ul> <li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.</li> <li>If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.</li> <li>The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).</li> </ul> </li> <li>Broader Impacts <ul> <li>Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?</li> <li>Answer: [NA]</li> <li>Justification: This paper focuses on learning theory and there is no societal impact of the work performed.</li> <li>Guidelines:</li> </ul> </li> </ul>
932 933 934 935 936 937 938 939 940 941 942 943 944 945 946 947 948 949	9.	<ul> <li>didn't make it into the paper).</li> <li>Code Of Ethics</li> <li>Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?</li> <li>Answer: [Yes]</li> <li>Justification: We have reviewed the NeurIPS Code of Ethics.</li> <li>Guidelines: <ul> <li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.</li> <li>If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.</li> <li>The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).</li> </ul> </li> <li>Broader Impacts <ul> <li>Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?</li> <li>Answer: [NA]</li> <li>Justification: This paper focuses on learning theory and there is no societal impact of the work performed.</li> <li>Guidelines: <ul> <li>The answer NA means that there is no societal impact of the work performed.</li> </ul> </li> </ul></li></ul>
932 933 934 935 936 937 938 939 940 941 942 943 944 945 944 945 946 947 948 949 950	9.	<ul> <li>didn't make it into the paper).</li> <li>Code Of Ethics</li> <li>Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?</li> <li>Answer: [Yes]</li> <li>Justification: We have reviewed the NeurIPS Code of Ethics.</li> <li>Guidelines: <ul> <li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.</li> <li>If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.</li> <li>The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).</li> </ul> </li> <li>Broader Impacts <ul> <li>Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?</li> <li>Answer: [NA]</li> <li>Justification: This paper focuses on learning theory and there is no societal impact of the work performed.</li> <li>Guidelines: <ul> <li>The answer NA means that there is no societal impact of the work performed.</li> </ul> </li> </ul></li></ul>
932 932 933 934 935 936 937 938 939 940 941 942 943 944 945 944 945 946 947 948 949 949	9.	<ul> <li>didn't make it into the paper).</li> <li>Code Of Ethics</li> <li>Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?</li> <li>Answer: [Yes]</li> <li>Justification: We have reviewed the NeurIPS Code of Ethics.</li> <li>Guidelines: <ul> <li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.</li> <li>If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.</li> <li>The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).</li> </ul> </li> <li>Broader Impacts <ul> <li>Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?</li> <li>Answer: [NA]</li> <li>Justification: This paper focuses on learning theory and there is no societal impact of the work performed.</li> <li>If the authors answer NA or No, they should explain why their work has no societal impact or why the maner does not address societal impact</li> </ul> </li> </ul>
932 932 933 934 935 936 937 938 939 940 941 942 943 944 945 944 945 946 947 948 949 949 950 951 952	9.	<ul> <li>didn't make it into the paper).</li> <li>Code Of Ethics</li> <li>Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?</li> <li>Answer: [Yes]</li> <li>Justification: We have reviewed the NeurIPS Code of Ethics.</li> <li>Guidelines: <ul> <li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.</li> <li>If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.</li> <li>The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).</li> </ul> </li> <li>Broader Impacts <ul> <li>Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?</li> <li>Answer: [NA]</li> <li>Justification: This paper focuses on learning theory and there is no societal impact of the work performed.</li> <li>Guidelines: <ul> <li>The answer NA means that there is no societal impact of the work performed.</li> <li>If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.</li> </ul> </li> </ul></li></ul>
932 932 933 934 935 936 937 938 939 940 941 942 943 944 945 946 947 946 947 948 949 949 950 951 952 953	9.	<ul> <li>didn't make it into the paper).</li> <li>Code Of Ethics</li> <li>Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?</li> <li>Answer: [Yes]</li> <li>Justification: We have reviewed the NeurIPS Code of Ethics.</li> <li>Guidelines: <ul> <li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.</li> <li>If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.</li> <li>The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).</li> </ul> </li> <li>Broader Impacts <ul> <li>Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?</li> <li>Answer: [NA]</li> <li>Justification: This paper focuses on learning theory and there is no societal impact of the work performed.</li> <li>Guidelines: <ul> <li>The answer NA means that there is no societal impact of the work performed.</li> <li>If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.</li> <li>Examples of negative societal impacts include potential malicious or unintended uses (e.g., digingraphic profiles, surveillance) foirness considerations</li> </ul> </li> </ul></li></ul>
932 932 933 934 935 936 937 938 939 940 941 942 943 944 945 946 944 945 946 947 948 949 950 951 952 953 954 955	9.	<ul> <li>didn't make it into the paper).</li> <li>Code Of Ethics</li> <li>Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?</li> <li>Answer: [Yes]</li> <li>Justification: We have reviewed the NeurIPS Code of Ethics.</li> <li>Guidelines: <ul> <li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.</li> <li>If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.</li> <li>The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).</li> </ul> </li> <li>Broader Impacts <ul> <li>Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?</li> </ul> </li> <li>Answer: [NA]</li> <li>Justification: This paper focuses on learning theory and there is no societal impact of the work performed.</li> <li>Guidelines: <ul> <li>The answer NA means that there is no societal impact of the work performed.</li> <li>If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.</li> <li>Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., denlowment of technologies that could make decisions that unfairly impact specific</li> </ul> </li> </ul>

957 958 959 960 961 962 963 964 965 966 966 966 967 968 969 970 971	<ul> <li>The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.</li> <li>The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.</li> <li>If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).</li> </ul>
972	11. Safeguards
973	Question: Does the paper describe safeguards that have been put in place for responsible
974	release of data or models that have a high risk for misuse (e.g., pretrained language models,
975	image generators, or scraped datasets)?
976	Answer: [NA]
977	Justification: This paper focuses on learning theory and poses no such risks.
978	Guidelines:
979	• The answer NA means that the paper poses no such risks.
980	• Released models that have a high risk for misuse or dual-use should be released with
981	necessary safeguards to allow for controlled use of the model, for example by requiring
982 983	that users adhere to usage guidelines or restrictions to access the model or implementing safety filters
984	• Datasets that have been scraped from the Internet could nose safety risks. The authors
985	should describe how they avoided releasing unsafe images.
986	• We recognize that providing effective safeguards is challenging, and many papers do
987 988	not require this, but we encourage authors to take this into account and make a best faith effort.
989	12. Licenses for existing assets
990	Question: Are the creators or original owners of assets (e.g., code, data, models), used in
991	the paper, properly credited and are the license and terms of use explicitly mentioned and
992	properly respected?
993	Answer: [NA]
994	Justification: This paper focuses on learning theory and does not use existing assets.
995	Guidelines:
996	• The answer NA means that the paper does not use existing assets.
997	• The authors should cite the original paper that produced the code package or dataset.
998	• The authors should state which version of the asset is used and, if possible, include a
999	
1000	• The name of the license (e.g., CC-BY 4.0) should be included for each asset.
1001 1002	• For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
1003	• If assets are released, the license, copyright information, and terms of use in the
1004	package should be provided. For popular datasets, paperswithcode.com/datasets
1005	has curated licenses for some datasets. Their licensing guide can help determine the
1006	license of a dataset.
1007	• For existing datasets that are re-packaged, both the original license and the license of the derived exect (if it has always of blowled be area it.)
1008	the derived asset (if it has changed) should be provided.

1009 1010		• If this information is not available online, the authors are encouraged to reach out to the asset's creators.
1011	13.	New Assets
1012 1013		Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?
1014		Answer: [NA]
1015		Justification: The paper does not release new assets.
1016		Guidelines:
1017		• The answer NA means that the paper does not release new assets
1018		<ul> <li>Researchers should communicate the details of the dataset/code/model as part of their</li> </ul>
1019		submissions via structured templates. This includes details about training, license,
1020		limitations, etc.
1021 1022		• The paper should discuss whether and how consent was obtained from people whose asset is used.
1023 1024		• At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.
1025	14.	Crowdsourcing and Research with Human Subjects
1026		Ouestion: For crowdsourcing experiments and research with human subjects, does the paper
1027		include the full text of instructions given to participants and screenshots, if applicable, as
1028		well as details about compensation (if any)?
1029		Answer: [NA]
1030		Justification: This paper does not involve crowdsourcing nor research with human subjects.
1031		Guidelines:
1032		• The answer NA means that the paper does not involve crowdsourcing nor research with
1033		human subjects.
1034		• Including this information in the supplemental material is fine, but if the main contribu-
1035		tion of the paper involves human subjects, then as much detail as possible should be
1036		Included in the main paper.
1037		or other labor should be paid at least the minimum wage in the country of the data
1039		collector.
1040	15.	Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
1041		Subjects
1042		Question: Does the paper describe potential risks incurred by study participants, whether
1043		approvals (or an equivalent approval/review based on the requirements of your country or
1045		institution) were obtained?
1046		Answer: [NA]
1047		Justification: This paper does not involve crowdsourcing nor research with human subjects.
1048		Guidelines:
1049		• The answer NA means that the paper does not involve crowdsourcing nor research with
1050		human subjects.
1051		• Depending on the country in which research is conducted, IRB approval (or equivalent)
1052		may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper
1053		• We recognize that the procedures for this may vary significantly between institutions
1054		and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
1056		guidelines for their institution.
1057		• For initial submissions, do not include any information that would break anonymity (if
1058		applicable), such as the institution conducting the review.