

High-Probability Generalization Bounds for Pointwise Uniformly Stable Algorithms

Jun Fan^a, Yunwen Lei^{b,*}

^a*Department of Mathematics, Hong Kong Baptist University, Kowloon, Hong Kong, China*

^b*Department of Mathematics, The University of Hong Kong, Pokfulam, Hong Kong, China*

Abstract

Algorithmic stability is a fundamental concept in statistical learning theory to understand the generalization behavior of optimization algorithms. Existing high-probability bounds are developed for the generalization gap as measured by function values and require the algorithm to be uniformly stable. In this paper, we introduce a novel stability measure called pointwise uniform stability by considering the sensitivity of the algorithm with respect to the perturbation of each training example. We show this weaker pointwise uniform stability guarantees almost optimal bounds, and gives the first high-probability bound for the generalization gap as measured by gradients. Sharper bounds are given for strongly convex and smooth problems. We further apply our general result to derive improved generalization bounds for stochastic gradient descent. As a byproduct, we develop concentration inequalities for a summation of weakly-dependent vector-valued random variables.

Keywords:

Learning Theory, Algorithmic Stability, Stochastic Gradient Descent, Generalization Analysis

1. Introduction

How to understand the generalization behavior of a learning algorithm is a central problem in statistical learning theory. A popular approach to developing generalization bounds is based on the uniform convergence, which controls the uniform deviation between population risks and empirical risks over a function space [39, 2, 36, 9]. This approach ignores how an algorithm explores over the function space, and leads to generalization bounds depending on the complexity of function spaces such as VC dimension [39], covering numbers [45, 36] and Rademacher complexities [2].

An alternative approach for generalization analysis is based on a fundamental concept of algorithmic stability. Roughly speaking, we say a learning algorithm is algorithmically stable if a change of a single example in the training dataset brings only a small change in the output model, i.e., the algorithm is insensitive with respect to (w.r.t.) the perturbation of training datasets [32, 5]. Algorithmic stability was introduced in 1970s to derive leave-one-out bounds for certain nonparametric local learning algorithms (such as nearest-neighbor rules) [11, 32]. The modern framework of stability analysis

*Corresponding author

Email addresses: junfan@hkbu.edu.hk (Jun Fan), leiyw@hku.hk (Yunwen Lei)

14 was established in [5], where a celebrated concept called the uniform stability has been introduced to
 15 study regularization methods.

16 We need to answer two questions in applying algorithmic stability to get generalization bounds for
 17 an algorithm. The first question is how to guarantee the generalization by stability, i.e., whether a
 18 stable algorithm can always produce models with good generalization behavior. The second question
 19 is how to develop stability bounds for an algorithm in terms of algorithm parameters such as the
 20 regularization parameter, the step size and the number of iterations.

21 The second question is algorithm-dependent, which allows us to exploit the special property of
 22 algorithms to get bounds better than algorithm-independent bounds based on complexity measures [2].
 23 The stability of various optimization algorithms has been developed in the literature. For example,
 24 the uniform stability has been developed for stochastic gradient descent (SGD) [17], which is one of
 25 the most widely used optimization methods to solve large-scale problems in machine learning.

26 For the first question, quantitative connection either in expectation or with high probability has
 27 been established. In particular, with probability at least $1 - \delta$ the following generalization bounds were
 28 developed for β -uniformly stable algorithms¹ [6, 14]

$$|F(A(S)) - F_S(A(S))| \lesssim \beta \log n \log(1/\delta) + \log^{\frac{1}{2}}(1/\delta)/n^{\frac{1}{2}}, \quad (1.1)$$

29 where $A(S)$ denotes the output model by applying an algorithm A to the dataset S , $F(\mathbf{w})$ denotes
 30 the population risk of a model \mathbf{w} , $F_S(\mathbf{w})$ denotes the empirical risk of \mathbf{w} (definitions are given in
 31 Section 3.1) and n is the sample size. Eq. (1.1) is a breakthrough result on the high-probability
 32 generalization analysis for uniformly stable algorithms initialized in 2002 [5]. However, some questions
 33 on Eq. (1.1) still remain.

- 34 • Eq. (1.1) provides generalization bounds in terms of function values. For nonconvex problems,
 35 optimization algorithms can only find a local minimizer and therefore we can only get optimiza-
 36 tion error bounds for $\|\nabla F(A(S))\|$ [15], where ∇ denotes the gradient operator. Therefore, it is
 37 interesting to study the generalization behavior of $A(S)$ as measured by $\nabla F(A(S))$, which moti-
 38 vates the question of developing high-probability bounds on the generalization gap as measured
 39 by gradients, i.e., $\|\nabla F(A(S)) - \nabla F_S(A(S))\|$.
- 40 • Eq. (1.1) requires the algorithm to be uniformly stable, which is arguably the strongest concept
 41 of algorithmic stability. Is it possible to relax this uniform stability to a weaker version of uniform
 42 stability, and can we develop better bounds on this weaker stability for popular algorithms such
 43 as SGD?
- 44 • The recent sharper generalization bounds in [18] require the loss function to be simultaneously
 45 Lipschitz continuous and λ -strongly convex, which cannot be satisfied globally due to the conflict

¹We use the notation \lesssim to ignore constant factors in an inequality.

46 between the Lipschitz continuity and strong convexity. Furthermore, their generalization bounds
 47 involve $\hat{\Delta}_\lambda^{\frac{1}{2}}$, where $\hat{\Delta}_\lambda = \lambda^{-1}(F_S(A(S)) - \min_{\mathbf{w}} F_S(\mathbf{w}))$ denotes a weighted suboptimality of
 48 the output model in terms of the empirical risk. This square-root dependency on $\hat{\Delta}_\lambda$ is slow in
 49 practice. Can we address the above conflict and improve the dependency on $\hat{\Delta}_\lambda$?

50 In this paper, we aim to provide affirmative answers to the above questions. Our main contributions
 51 are as follows.

- 52 • We develop a concentration inequality for a summation of weakly-dependent vector-valued ran-
 53 dom variables, which generalizes a similar result in Bousquet et al. [6] from real-valued random
 54 variables to random variables taking values in a Hilbert space.
- 55 • We introduce a new stability measure termed as the pointwise uniform stability. While this
 56 stability is weaker than the uniform stability, we show it guarantees high-probability gener-
 57 alization bounds on $F(A(S)) - F_S(A(S))$. We also give the first high-probability bound for
 58 $\|\nabla F(A(S)) - \nabla F_S(A(S))\|$ based on stability analysis.
- 59 • We improve the high-probability bound in [18] by considering a loss function of a structure,
 60 which reconciles the conflict between Lipschitz continuity and strong convexity. Furthermore,
 61 we derive a sharper bound involving $\hat{\Delta}_\lambda^{\frac{1+\alpha}{2}}$ to exploit the α -Hölder continuity of gradients. In
 62 particular, if $\alpha = 1$, the term $\hat{\Delta}_\lambda^{\frac{1+\alpha}{2}}$ decays quadratically faster than $\hat{\Delta}_\lambda^{\frac{1}{2}}$ in [18].
- 63 • We study the pointwise uniform stability of SGD for convex and strongly convex problems,
 64 covering smooth and nonsmooth problems. We then apply our connection between stability and
 65 generalization to give high-probability generalization bounds.

66 The paper is organized as follows. We review the related work in Section 2. We present our main
 67 results in Section 3, and give applications to SGD in Section 4. We present the proofs on connecting
 68 stability and generalization in Section 5, and the proofs on SGD in Section 6. The conclusion is given
 69 in Section 7. Some lemmas and proofs are given in the Appendix.

70 2. Related Work

71 2.1. Connection on Stability and Generalization

72 Algorithmic stability can imply generalization bounds in expectation and with high probability. We
 73 first consider generalization bounds in expectation. On-average stability can imply generalization under
 74 a Lipschitz condition of loss functions [34]. For non-Lipschitz problems, an on-average model stability
 75 was proposed to give generalization bounds by exploiting the smoothness of loss functions [22], which
 76 can further imply fast rates under a low-noise condition. On-average stability can imply generalization
 77 bounds for any learning algorithms to solve gradient-dominated problems [23, 7]. For nonconvex
 78 and smooth problems, generalization as measured by gradients can be guaranteed by stability in
 79 gradients [21].

80 We now consider generalization bounds with high probability. In a seminal paper [5], β -uniform
81 stability was introduced to give bounds of order $O((\beta + 1/n)\sqrt{n\log(1/\delta)})$, which was extended to ran-
82 domized learning algorithms [12]. These results were significantly improved to $O(\sqrt{(\beta + 1/n)\log(1/\delta)})$
83 in [13] by techniques in adaptive data analysis. Almost optimal generalization bounds in Eq. (1.1)
84 were further derived by developing concentration inequalities for a summation of weakly-dependent
85 random variables [6, 14]. The above-mentioned high-probability analysis can imply bounds of the order
86 at most $O(1/\sqrt{n})$. Under a Bernstein condition on variances, it was shown that β -uniformly stable
87 algorithms can enjoy high-probability bounds of the order $O((\beta \log n + 1/n)\log(1/\delta))$ [18].

88 2.2. Stability of Learning Algorithms

89 Algorithmic stability has been studied for various learning algorithms. Uniform stability bounds of
90 order $O(1/(n\lambda))$ were developed for empirical risk minimization to solve λ -strongly convex problems [5].
91 In a seminal paper, uniform stability bounds of order $O(G^2 \sum_{t=1}^T \eta_t/n)$ were developed for SGD with
92 T iterations and step size sequences $\{\eta_t\}$ for convex, smooth and G -Lipschitz problems [17]. Data-
93 dependent stability bounds reflecting the effect of initialization point were established for SGD [20].
94 For nonsmooth and convex problems, stability bounds of order $O(\eta\sqrt{T} + \eta T/n)$ were developed for
95 SGD with $\eta_t = \eta$ either in expectation [22] or with high probability [3]. The Lipschitz constant G in the
96 existing stability bounds [17] was replaced by the training error based on on-average model stability,
97 which can imply fast excess risk bounds under a low-noise condition [22]. On-average model stability
98 was also used to understand the benefit of overparameterization for shallow neural networks [30, 37, 24],
99 and the implicit bias of gradient methods for separable data and self-bounding loss functions [33].
100 Other than the standard SGD/GD, the stability of differentially private SGD [42, 3, 21], gradient-free
101 optimization methods [28], accelerated methods [41] and noisy SGD [46, 25, 38, 27] was studied in the
102 literature. Lower bounds on the stability of gradient methods were also developed [3, 19, 1].

103 3. Main Results

104 3.1. Problem Setup

105 Let ρ be a probability measure defined on a sample space $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$, where \mathcal{X} is an input space and
106 \mathcal{Y} is an output space. Let $S = (z_1, \dots, z_n)$ be a training dataset drawn independently from ρ , based on
107 which we aim to find a model $h : \mathcal{X} \mapsto \mathcal{Y}$ for further prediction. We consider a parametric model, i.e., a
108 model can be indexed by a parameter $\mathbf{w} \in \mathcal{W}$, where $\mathcal{W} \subset \mathbb{R}^d$ is the parameter space. The performance
109 of a model \mathbf{w} on an example z can be measured by $f(\mathbf{w}; z)$, where $f : \mathcal{W} \times \mathcal{Z} \mapsto \mathbb{R}^+$ is the loss function.
110 The empirical behavior of a model \mathbf{w} can be quantified by the empirical risk $F_S(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n f(\mathbf{w}; z_i)$,
111 while the prediction behavior can be quantified by the population risk $F(\mathbf{w}) = \mathbb{E}_z[f(\mathbf{w}; z)]$, where $\mathbb{E}_z[\cdot]$
112 denotes the expectation w.r.t. z . We often apply an algorithm A onto S to get a model $A(S) \in \mathcal{W}$
113 with a small empirical risk. However, this does not necessarily imply a small population risk referred
114 to as the overfitting phenomenon. To this aim, we need to handle an important concept called the

115 generalization gap $F(A(S)) - F_S(A(S))$, i.e., the difference between population risk and empirical risk
 116 at the output model $A(S)$. In this paper, we will leverage the celebrated concept called algorithmic
 117 stability to develop high-probability bounds on the generalization gap.

118 3.2. Concentration Inequality

119 We give a p -norm bound for a summation of weakly-dependent random variables taking values
 120 in a Hilbert space, whose proof is given in Section 5.1. It will play a fundamental role in deriving
 121 the connection between stability and generalization. The L_p -norm of a real-valued random variable
 122 Z is denoted by $\|Z\|_p := (\mathbb{E}[|Z|^p])^{\frac{1}{p}}$, $p \geq 1$. Let $\|\cdot\|$ denote the norm in a Hilbert space \mathcal{H} . Then
 123 $\|\nabla f(\mathbf{w}; Z)\|$ is a real-valued random variable (as a function of Z). According to our notation, we have

$$\|\|\nabla f(\mathbf{w}; Z)\|\|_p = \left(\mathbb{E}_Z[\|\nabla f(\mathbf{w}; Z)\|^p]\right)^{\frac{1}{p}}, \quad \forall p \geq 1.$$

124 **Theorem 1.** Let $\mathbf{Z} = (Z_1, \dots, Z_n)$ be a sequence of independent random variables taking values in a
 125 Hilbert space \mathcal{H} . Let g_1, \dots, g_n be functions $g_i : \mathcal{Z}^n \mapsto \mathcal{H}$ such that the following holds.

- 126 1. For any $i \in [n]$, almost surely we have $\sup_{z_i} \|\mathbb{E}[g_i(\mathbf{Z})|Z_i = z_i]\| \leq M$.
- 127 2. For any $i \in [n]$, almost surely we have $\mathbb{E}[g_i(\mathbf{Z})|\mathbf{Z}_{[n]\setminus\{i\}} = (z_j)_{j \neq i}] = 0, \forall z_j \in \mathcal{Z}, j \neq i$.
- 128 3. For any $i \in [n]$, the following inequality holds

$$\sup_{z_1, \dots, z_n, z'_j: j \neq i} \|g_i(z_1, \dots, z_{j-1}, z_j, z_{j+1}, \dots, z_n) - g_i(z_1, \dots, z_{j-1}, z'_j, z_{j+1}, \dots, z_n)\| \leq \beta_j. \quad (3.1)$$

129 Then, for any $p \geq 2$ we have

$$\left\| \left\| \sum_{i=1}^n g_i \right\| \right\|_p \leq 2(\sqrt{2} + 1)M\sqrt{np} + 2(\sqrt{2} + 1)p \lceil \log_2 n \rceil \left(n \sum_{i=1}^n \beta_i^2 \right)^{\frac{1}{2}}.$$

130 **Remark 1.** If $\mathcal{H} = \mathbb{R}$, a similar bound was established in [6]. That is, let $\tilde{g}_1, \dots, \tilde{g}_n$ be real-valued
 131 functions such that $\|\mathbb{E}[\tilde{g}_i(Z)|Z_i]\| \leq M$, $\mathbb{E}[\tilde{g}_i(Z)|\mathbf{Z}_{[n]\setminus\{i\}}] = 0$ and

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

132 Then, the following inequality was established for any $p \geq 2$ [6]

$$\left\| \sum_{i=1}^n \tilde{g}_i(Z) \right\|_p \leq 4M\sqrt{np} + 12\sqrt{2}pn\beta \lceil \log_2 n \rceil. \quad (3.2)$$

133 There are two differences between our result and Eq. (3.2). First, we extend the discussion in [6]
 134 from real-valued random variables to random variables taking values in a general Hilbert space, and
 135 slightly improve the constant factor. Second, the discussions [6] assume the change of j -th example
 136 in $\mathbf{z} = (z_1, \dots, z_n)$ would lead to a change of value uniformly bounded by β . As a comparison, we
 137 allow different β_j for different $j \in [n]$. As we will show, this is useful for us to get a new generalization
 138 bound based on our pointwise uniform stability.

139 *3.3. Stability and Generalization*

140 Stability measures the sensitivity of an algorithm up to the perturbation of the training dataset
 141 by a single example. A very popular stability measure is the uniform stability, which considers the
 142 change of any single example of any training dataset by any $z \in \mathcal{Z}$.

143 **Definition 1** (Uniform Stability). Let A be an algorithm and $\beta > 0$.

- 144 1. We say A is β -uniformly-stable in *function values* if for all datasets S, S' such that S and S'
 145 differ by a single example, we have

$$\sup_z |f(A(S); z) - f(A(S'); z)| \leq \beta. \quad (3.3)$$

- 146 2. We say A is β -uniformly-stable in *gradients* if for all datasets S, S' such that S and S' differ by
 147 a single example, we have

$$\sup_z \|\nabla f(A(S); z) - \nabla f(A(S'); z)\| \leq \beta. \quad (3.4)$$

148 In this paper, we introduce a new stability measure which we call the pointwise uniform stability.
 149 The basic idea is to give a single stability parameter β_i for perturbing the i -th example of the dataset. It
 150 is clear that if A is β -uniformly stable, then it is also (β, \dots, β) -pointwise uniformly stable. Therefore,
 151 pointwise uniform stability is weaker than the uniform stability. In this paper, we will show that this
 152 weaker stability can also imply high-probability generalization bounds. We say two datasets S and
 153 $S^{(i)}$ differ only by the i -th example if $S = (z_1, \dots, z_n)$ and $S^{(i)} = (z_1, \dots, z_{i-1}, z'_i, z_{i+1}, \dots, z_n)$ for
 154 some $z'_i \in \mathcal{Z}$.

155 **Definition 2** (Pointwise Uniform Stability). Let A be an algorithm and $\beta = (\beta_1, \dots, \beta_n), \beta_i > 0$.

- 156 1. We say A is β -pointwise uniformly-stable in *function values* if for all $S, S^{(i)}$ such that S and $S^{(i)}$
 157 differ by the i -th example, we have

$$\sup_z |f(A(S); z) - f(A(S^{(i)}); z)| \leq \beta_i. \quad (3.5)$$

- 158 2. We say A is β -pointwise uniformly-stable in *gradients* if for all $S, S^{(i)}$ such that S and $S^{(i)}$ differ
 159 by the i -th example, we have

$$\sup_z \|\nabla f(A(S); z) - \nabla f(A(S^{(i)}); z)\| \leq \beta_i. \quad (3.6)$$

160 Theorem 2 gives a high-probability bound on the generalization gap $F(A(S)) - F_S(A(S))$ for
 161 pointwise uniformly stable algorithms. We omit the proof due to its similarity with Theorem 3.

162 **Theorem 2** (Generalization via Function Values). Let $\beta = (\beta_1, \dots, \beta_n)$. Consider an algorithm A
 163 and $\delta \in (0, 1)$. Assume for any S and any z , $|f(A(S); z)| \leq M$. If A is β -pointwise uniformly-stable
 164 in function values, then the following inequality holds with probability at least $1 - \delta$

$$|F(A(S)) - F_S(A(S))| \lesssim \frac{M \log^{\frac{1}{2}}(1/\delta)}{\sqrt{n}} + \left(\frac{1}{n} \sum_{i=1}^n \beta_i^2\right)^{\frac{1}{2}} \log n \log(1/\delta).$$

165 **Remark 2.** If A is β -uniformly stable in function values, the generalization bound in Eq. (1.1)
 166 was developed [5, 14]. As a comparison, our bound involves an average of stability parameters over all
 167 indices, i.e., the term $(\frac{1}{n} \sum_{i=1}^n \beta_i^2)^{\frac{1}{2}}$, which is smaller than the uniform stability parameter $\beta = \max_i \beta_i$
 168 considered in [5, 14]. As we will show, for SGD we can establish a bound for $(\frac{1}{n} \sum_{i=1}^n \beta_i^2)^{\frac{1}{2}}$ which is
 169 smaller than that for $\max_i \beta_i$.

170 Our next result is a high-probability bound on the generalization gap in terms of gradients, which
 171 extends the high-probability generalization bound in function values in [5, 14]. We show that the
 172 deviation between population gradients and empirical gradients at the output model can be bounded
 173 by the stability parameter in gradients. We require f to be differentiable, and do not require a convexity
 174 or smoothness assumption in Theorem 3. The proof is given in Section 5.2.

175 **Theorem 3** (Generalization via Gradients). *Let $\boldsymbol{\beta} = (\beta_1, \dots, \beta_n)$. Consider an algorithm A and*
 176 *$\delta \in (0, 1)$. Assume for any S and any z , $\|\nabla f(A(S); z)\| \leq M$. If A is $\boldsymbol{\beta}$ -pointwise uniformly-stable in*
 177 *gradients, then the following inequality holds with probability at least $1 - \delta$*

$$\|\nabla F(A(S)) - \nabla F_S(A(S))\| \lesssim \frac{M \log^{\frac{1}{2}}(1/\delta)}{\sqrt{n}} + \left(\frac{1}{n} \sum_{i=1}^n \beta_i^2\right)^{\frac{1}{2}} \log n \log(1/\delta).$$

178 **Remark 3.** If A is β -uniformly stable in gradients with, then it was shown [21]

$$\mathbb{E}[\|\nabla F(A(S)) - \nabla F_S(A(S))\|] \lesssim \beta + \sqrt{\frac{1}{n} \mathbb{E}[\mathbb{V}_Z[f(A(S); Z)]]},$$

179 where $\mathbb{V}_Z[f(A(S); Z)]$ is the variance of $\nabla f(A(S); Z)$ as a function of Z . This bound was established in
 180 expectation. As a comparison, we develop high-probability bounds on the generalization gap between
 181 population and empirical gradients. High-probability bounds of order $\sqrt{d \log(1/\delta)/n}$ were also estab-
 182 lished for $\sup_{\mathbf{w}} \|\nabla F(\mathbf{w}) - \nabla F_S(\mathbf{w})\|$ based on complexity measures of function spaces, which, however,
 183 depend on the dimensionality d of the problem and are not appealing for high-dimensional learning
 184 problems. As a comparison, our stability analysis implies dimension-free generalization bounds.

185 3.4. Sharper Generalization Bounds

186 Theorem 2 implies generalization bounds of the order $O(1/\sqrt{n})$. In this section, we improve this
 187 dependency to $O(1/n)$ for pointwise uniformly stable algorithms. The following theorem is an extension
 188 of the stability analysis in [18]. We consider functions with a composite structure.

189 **Definition 3** (Lipschitzness, Smoothness and Convexity). Let $G, L_\alpha, L > 0, \lambda \geq 0$ and $g : \mathcal{W} \mapsto \mathbb{R}$.

- 190 • We say g is G -Lipschitz continuous if $|g(\mathbf{w}) - g(\mathbf{w}')| \leq G\|\mathbf{w} - \mathbf{w}'\|, \forall \mathbf{w}, \mathbf{w}' \in \mathcal{W}$.
- 191 • We say g has (α, L_α) -Hölder continuous gradients ($\alpha \in [0, 1]$) if

$$\|\nabla g(\mathbf{w}) - \nabla g(\mathbf{w}')\| \leq L_\alpha \|\mathbf{w} - \mathbf{w}'\|^\alpha, \quad \forall \mathbf{w}, \mathbf{w}' \in \mathcal{W}.$$

192 We say g is L -smooth if g has $(1, L)$ -Hölder continuous gradients.

193 • We say g is λ -strongly convex if

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

194 We say g is convex if the above inequality holds with $\lambda = 0$.

195 **Assumption 1.** Let $\lambda > 0, \ell : \mathcal{W} \times \mathcal{Z} \mapsto \mathbb{R}_+$ and $r : \mathcal{W} \mapsto \mathbb{R}_+$. Assume $f : \mathcal{W} \times \mathcal{Z} \mapsto \mathbb{R}_+$ has the
196 following structure

$$f(\mathbf{w}; z) = \ell(\mathbf{w}; z) + r(\mathbf{w}). \quad (3.7)$$

197 Assume for any z , the function $\mathbf{w} \mapsto \ell(\mathbf{w}; z)$ is nonnegative and has (α, L_α) -Hölder continuous gradi-
198 ents. Assume r is L_r -smooth, and for any z , the function $\mathbf{w} \mapsto f(\mathbf{w}; z)$ is λ -strongly convex.

199 For non-composite problems, our analysis can still imply faster rates if f is strongly convex, smooth
200 and $\|\nabla f(A(S); z)\| \leq G, \|\nabla f(A_e(S); z)\| \leq G$, where we denote by A_e the empirical risk minimization
201 (ERM) algorithm, i.e.,

$$A_e(S) = \arg \min_{\mathbf{w} \in \mathcal{W}} F_S(\mathbf{w}).$$

202 Let $L_S(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}; z_i)$ and $L(\mathbf{w}) = \mathbb{E}_z[\ell(\mathbf{w}; z)]$. Let $\mathbf{w}^* = \arg \min_{\mathbf{w} \in \mathcal{W}} F(\mathbf{w})$ be a minimizer of
203 the population risk. The proof is given in Section C.

204 **Theorem 4.** Let $\beta = (\beta_1, \dots, \beta_n)$ and f take a structure in Eq. (3.7). Assume A is β -pointwise
205 uniformly-stable in function values (measured by ℓ), i.e., Eq. (3.5) holds with f replaced by ℓ . Let
206 $M = \sup_z |\mathbb{E}_S[\ell(A(S)); z] - \ell(\mathbf{w}^*; z)|$. Then for any $\delta \in (0, 1)$, the following inequality holds with
207 probability at least $1 - \delta$

$$F(A(S)) - F_S(A(S)) - F(\mathbf{w}^*) + F_S(\mathbf{w}^*) \lesssim \left(\frac{1}{n} \sum_{i=1}^n \beta_i^2 \right)^{\frac{1}{2}} \log n \log(1/\delta) + \frac{M \log \frac{1}{\delta}}{n} + \left(\frac{\sigma_A^2 \log(1/\delta)}{n} \right)^{\frac{1}{2}},$$

208 where

$$\sigma_A^2 = \mathbb{E}_Z \left[\left(\mathbb{E}_S[\ell(A(S); Z)] - \ell(\mathbf{w}^*; Z) \right)^2 \right] - \left(\mathbb{E}_S[L(A(S))] - L(\mathbf{w}^*) \right)^2.$$

209 **Remark 4.** A key difference between Theorem 4 and Theorem 2 is that the term $n^{-\frac{1}{2}} M \log^{\frac{1}{2}}(1/\delta)$ in
210 Theorem 2 is replaced by $n^{-1} M \log(1/\delta)$ in Theorem 4, at the cost of introducing $\sigma_A n^{-\frac{1}{2}} \log^{\frac{1}{2}}(1/\delta)$.
211 Then, Theorem 4 can imply fast excess risk bounds if the variance σ_A^2 is small. Similar bounds were
212 derived in [18] under the following Bernstein assumption

$$\mathbb{E}_Z [(f(\mathbf{w}; Z) - f(\mathbf{w}^*; Z))^2] \leq B(F(\mathbf{w}) - F(\mathbf{w}^*)), \quad \forall \mathbf{w} \in \mathcal{W}. \quad (3.8)$$

213 The bound in [18] involves the uniform stability. As a comparison, our analysis uses the pointwise
214 uniform stability. Furthermore, we do not impose a Bernstein assumption, and instead include the
215 variance term σ_A^2 in the upper bound. Finally, we consider a problem with a composite structure and
216 our stability assumption is imposed to ℓ instead of f . The underlying reason is that it is possible ℓ
217 is Lipschitz continuous but f not. In this case, if we can derive a bound on $\|A(S) - A(S^{(i)})\|$, we
218 can use the Lipschitz continuity of ℓ to get a bound on $\ell(A(S); z) - \ell(A(S^{(i)}); z)$ but not a bound
219 on $f(A(S); z) - f(A(S^{(i)}); z)$. As a comparison, the analysis in [18] does not consider this composite
220 structure.

221 To apply Theorem 4, we need to estimate the variance term σ_A^2 , which can be related to the excess
 222 risk $F(A(S)) - F(\mathbf{w}^*)$. In the following theorem to be proved in Section 5.3, we show that the Bernstein
 223 condition holds if Assumption 1 holds. Furthermore, we also give generalization bounds in expectation,
 224 which involves optimization error $F_S(A(S)) - F_S(A_e(S))$ and the strong convexity parameter λ . Define

$$c_\alpha = \begin{cases} (1 + 1/\alpha)^{\frac{\alpha}{1+\alpha}} L_\alpha^{\frac{1}{1+\alpha}}, & \text{if } \alpha \in (0, 1], \\ \sup_z \|\nabla \ell(0; z)\| + L_\alpha, & \text{if } \alpha = 0. \end{cases} \quad (3.9)$$

225 The proof is given in Section D.

226 **Lemma 5.** *Let Assumption 1 hold. Then*

$$\sigma_A^2 \leq C \lambda^{-1} \mathbb{E}_S[F(A(S)) - F(\mathbf{w}^*)], \quad (3.10)$$

227 where

$$C = 2c_\alpha^2 \mathbb{E}_{S,Z}[\max\{\ell^{\frac{2\alpha}{1+\alpha}}(A(S); Z), \ell^{\frac{2\alpha}{1+\alpha}}(\mathbf{w}^*; Z)\}]. \quad (3.11)$$

228 Furthermore, if $(F(A_e(S)) - F(\mathbf{w}^*))^{\frac{1-\alpha}{1+\alpha}} \leq \tilde{C}$ for some $\tilde{C} > 0$ independent of n or λ , then any
 229 algorithm A satisfies

$$\mathbb{E}[F(A(S))] - F(\mathbf{w}^*) \leq \mathfrak{C} \left(\Delta_\lambda^{\frac{1+\alpha}{2}} + \Delta_\lambda + \nabla_\lambda \right), \quad (3.12)$$

230 where \mathfrak{C} is a constant independent of λ or n (may depend on α, L_α, L_r and is explicitly given in Eq.
 231 (D.5)) and

$$\Delta_\lambda = \lambda^{-1} \mathbb{E}[F_S(A(S)) - F_S(A_e(S))], \quad \nabla_\lambda = \frac{1}{n\lambda} \mathbb{E}\left[L_S^{\frac{2\alpha}{1+\alpha}}(A_e(S)) + L^{\frac{2\alpha}{1+\alpha}}(A_e(S))\right].$$

232 The assumption $(F(A_e(S)) - F(\mathbf{w}^*))^{\frac{1-\alpha}{1+\alpha}} \leq \tilde{C}$ is introduced just for simplifying the analysis, and
 233 can be removed with more complicated computation. This assumption holds automatically if $\alpha = 1$.
 234 We can combine Eq. (3.10) and Eq. (3.12) to derive

$$\begin{aligned} \sigma_A^2 \leq 2c_\alpha^2 \mathfrak{C} (\mathbb{E}[F(A(S))])^{\frac{2\alpha}{1+\alpha}} & \left(\frac{1}{\lambda^{1+\frac{1+\alpha}{2}}} \left(\mathbb{E}[F_S(A(S)) - F_S(A_e(S))] \right)^{\frac{1+\alpha}{2}} \right. \\ & \left. + \frac{\mathbb{E}[F_S(A(S)) - F_S(A_e(S))]}{\lambda^2} + \frac{2(\mathbb{E}[L(A_e(S))])^{\frac{2\alpha}{1+\alpha}}}{n\lambda^2} \right), \end{aligned}$$

235 where we have used the Jensen's inequality. We can plug the above bound back into Theorem 4, and
 236 get the following high-probability bound. We omit the proof for simplicity. For simplicity, we assume
 237 $\Delta_\lambda = O(1)$ and absorb all constant factors independent of β_i, n, λ (e.g., α, L_α, L_r) into the \lesssim notation.

238 **Corollary 6.** *Let Assumptions in Lemma 5 hold. Let $\boldsymbol{\beta} = (\beta_1, \dots, \beta_n)$ and assume A is $\boldsymbol{\beta}$ -pointwise
 239 uniformly-stable in function values (measured by ℓ). Let $M = \sup_z |\mathbb{E}_S[\ell(A(S)); z] - \ell(\mathbf{w}^*; z)|$. Then
 240 for any $\delta \in (0, 1)$, the following inequality holds with probability at least $1 - \delta$*

$$\begin{aligned} F(A(S)) - F_S(A(S)) - F(\mathbf{w}^*) + F_S(\mathbf{w}^*) & \lesssim \left(\frac{1}{n} \sum_{i=1}^n \beta_i^2 \right)^{\frac{1}{2}} \log n \log(1/\delta) + \frac{M \log \frac{1}{\delta}}{n} \\ & + \frac{\log^{1/2}(1/\delta) (\mathbb{E}[F(A(S))])^{\frac{\alpha}{1+\alpha}}}{\sqrt{n}} \left(\frac{1}{\lambda^{\frac{3+\alpha}{4}}} \left(\mathbb{E}[F_S(A(S)) - F_S(A_e(S))] \right)^{\frac{1+\alpha}{4}} + \frac{(\mathbb{E}[L(A_e(S))])^{\frac{\alpha}{1+\alpha}}}{\sqrt{n\lambda}} \right). \end{aligned}$$

241 If we further assume ℓ is Lipschitz continuous, we can have the following high-probability bounds
 242 for any algorithm to solve strongly convex problems. The proof is given in Section 5.3.

243 **Theorem 7.** *Let Assumptions in Lemma 5 hold and ℓ be G -Lipschitz continuous. If $\sup_z |\mathbb{E}_S[\ell(A(S)); z] -$
 244 $\ell(\mathbf{w}^*; z)| < \infty$, then for any $\delta \in (0, 1)$ with probability at least $1 - \delta$ we have*

$$F(A(S)) - F(\mathbf{w}^*) \lesssim (n\lambda)^{-1} \log n \log(1/\delta) + \hat{\Delta}_\lambda^{\frac{1+\alpha}{2}},$$

245 where $\hat{\Delta}_\lambda = \lambda^{-1}(F_S(A(S)) - F_S(A_e(S)))$.

246 **Remark 5.** The term $F_S(A(S)) - F_S(A_e(S))$ is the optimization error, which measures the subopti-
 247 mality of $A(S)$ to the minimal empirical risk. The recent work [18] gives the following high probability
 248 bound if F_S is λ -strongly convex and f is Lipschitz continuous

$$F(A(S)) - F(\mathbf{w}^*) \lesssim \left(\frac{1}{n\lambda} + \bar{\Delta}_\lambda^{\frac{1}{2}} \right) \log n \log(1/\delta), \quad (3.13)$$

249 where $\bar{\Delta}_\lambda$ is a deterministic number and an upper bound of $\hat{\Delta}_\lambda$. However, a strongly convex function
 250 cannot be Lipschitz continuous in the whole region. Therefore, the strong convexity assumption is
 251 contradictory to the Lipschitz condition. As a comparison, we consider an objective with a composite
 252 structure where ℓ has α -Hölder continuous gradients and is Lipschitz continuous. Our assumption is
 253 satisfied by various machine learning problems. For example, for logistic regression we have

$$f(\mathbf{w}; z) = \log(1 + \exp(-y\mathbf{w}^\top x)) + 2^{-1}\lambda\|\mathbf{w}\|^2,$$

254 which satisfies Assumption 1 with $\alpha = 1$. Moreover, the function $z \mapsto \log(1 + \exp(-y\mathbf{w}^\top x))$ is Lipschitz
 255 continuous.

256 Furthermore, we show that the term $\bar{\Delta}_\lambda^{\frac{1}{2}}$ in Eq. (3.13) can be replaced by a faster-decaying term
 257 $\hat{\Delta}_\lambda^{\frac{1+\alpha}{2}}$. In particular, if ℓ is smooth, we have $\hat{\Delta}_\lambda^{\frac{1+\alpha}{2}} = \hat{\Delta}_\lambda$, which decays quadratically faster than $\bar{\Delta}_\lambda^{\frac{1}{2}}$
 258 in Eq. (3.13). This shows that we can stop the algorithm earlier if we impose a stronger assumption
 259 on the smoothness, and shows the benefit of smoothness in improving the generalization. Indeed, the
 260 analysis in [18] first shows that the algorithm A is β -uniformly stable with $\beta = 4G^2/(\lambda n) + \sqrt{8G^2\bar{\Delta}_\lambda}$.
 261 Then, they apply the high-probability bound on uniform stability to A and give the bound in Eq.
 262 (3.13). Since a smoothness assumption would not affect the uniform stability, the uniform stability
 263 parameter there involves $\bar{\Delta}_\lambda^{\frac{1}{2}}$, and the strategy fails to use the smoothness assumption to improve the
 264 bound. We take a different strategy. We apply Theorem 4 to the algorithm A_e to first give a bound
 265 on $F(A_e(S)) - F_S(A_e(S)) - F(\mathbf{w}^*) + F_S(\mathbf{w}^*)$, which does not involve $\bar{\Delta}_\lambda$ since A_e outputs the ERM
 266 model. Then we control $F(A(S)) - F_S(A(S))$ in terms of $\hat{\Delta}_\lambda$, and use the smoothness assumption to
 267 show this bound improves as f is becoming more and more smooth. Finally, Eq. (3.13) requires $\hat{\Delta}_\lambda$ to
 268 be upper bounded by a deterministic number $\bar{\Delta}_\lambda$. As a comparison, our result directly involves $\hat{\Delta}_\lambda$.

269 4. Applications to Stochastic Gradient Descent

270 In this section, we apply our connection between stability and generalization to derive generalization
 271 bounds for SGD.

272 **Definition 4** (SGD). Let $\mathbf{w}_1 \in \mathcal{W}$ and $\{\eta_t\}$ be a sequence of positive step sizes. At each iteration, we
 273 first randomly select an index j_t according to the uniform distribution over $[n]$ and update the model
 274 as follows

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \nabla f(\mathbf{w}_t; z_{j_t}).$$

275 4.1. Stability Bounds

276 We first develop the pointwise uniform stability bounds for SGD. We consider three classes of
 277 problems: convex and smooth problems, convex and nonsmooth problems, and strongly convex and
 278 smooth problems. Let $\mathbb{I}_{[E]}$ be the indicator function, i.e., $\mathbb{I}_{[E]} = 1$ if the event E happens, and 0
 279 otherwise. The proofs are given in Section 6.

280 **Theorem 8** (Stability of SGD: Smooth Case). *Assume $f : \mathcal{W} \times \mathcal{Z} \mapsto \mathbb{R}^+$ is convex, L -smooth and
 281 G -Lipschitz. Let $\{\mathbf{w}_t\}_{t \in \mathbb{N}}$ be produced by SGD with $\eta_t = \eta \leq 2/L$. Then SGD with T iterations is
 282 β -pointwise uniformly stable in function values, where $\frac{1}{n} \sum_{i=1}^n \beta_i = \frac{2G^2 T \eta}{n}$ and*

$$\frac{1}{n} \sum_{i=1}^n \beta_i^2 = \frac{4G^4 \eta^2}{n} \sum_{i=1}^n \left(\sum_{t=1}^T \mathbb{I}_{[j_t=i]} \right)^2. \quad (4.1)$$

283 **Remark 6.** Under the same condition, one can show that SGD is β_{unif} -uniformly stable in function
 284 values with (implicitly shown in the proof of Theorem 8)

$$\beta_{\text{unif}} = 2G^2 \eta \max_{i \in [n]} \sum_{t=1}^T \mathbb{I}_{[j_t=i]}. \quad (4.2)$$

285 To see the comparison between the uniform stability bound in Eq. (4.2) and the pointwise stability
 286 bound in Eq. (4.1), we introduce

$$\tilde{\beta}_{\text{unif}} = \max_{i \in [n]} \sum_{t=1}^T \mathbb{I}_{[j_t=i]}, \quad \tilde{\beta}_{\text{point}} = \left(\frac{1}{n} \sum_{i=1}^n \left(\sum_{t=1}^T \mathbb{I}_{[j_t=i]} \right)^2 \right)^{\frac{1}{2}}. \quad (4.3)$$

287 It is clear that $\tilde{\beta}_{\text{unif}}, \tilde{\beta}_{\text{point}}$ differ from the above uniform/pointwise stability bounds by a factor of
 288 $2G^2 \eta$. For simplicity, we set $T = n$ as this implies the optimal excess risk bounds [17]. Then, we have

$$\tilde{\beta}_{\text{point}} \leq \left(\frac{1}{n} \sum_{i=1}^n \left(\sum_{t=1}^n \mathbb{I}_{[j_t=i]} \right) \max_{i \in [n]} \left(\sum_{t=1}^n \mathbb{I}_{[j_t=i]} \right) \right)^{\frac{1}{2}} = \left(\frac{1}{n} \sum_{t=1}^n \sum_{i=1}^n \mathbb{I}_{[j_t=i]} \right)^{\frac{1}{2}} \tilde{\beta}_{\text{unif}}^{\frac{1}{2}} = \tilde{\beta}_{\text{unif}}^{\frac{1}{2}}, \quad (4.4)$$

289 where we have used the identity $\sum_{i=1}^n \mathbb{I}_{[j_t=i]} = 1$ for any t . The term $\tilde{\beta}_{\text{unif}}$ is related to the balls and
 290 bins problem [29]. It was shown that with probability at least $1 - 1/n$, $\tilde{\beta}_{\text{unif}} = \Theta\left(\frac{\log n}{\log \log n}\right)$ [29]. Then,
 291 by Eq. (4.4), with probability at least $1 - 1/n$ we have $\tilde{\beta}_{\text{point}} = O\left(\frac{\log^{\frac{1}{2}} n}{(\log \log n)^{\frac{1}{2}}}\right)$. Note Eq. (4.4) is not
 292 tight, and we expect that $\tilde{\beta}_{\text{point}}$ has a tighter upper bound. For example, we can show that the second
 293 moment of $\tilde{\beta}_{\text{point}}$ is bounded by a constant independent of n :

$$\begin{aligned} \mathbb{E}[\tilde{\beta}_{\text{point}}^2] &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}\left[\left(\sum_{t=1}^n \mathbb{I}_{[j_t=i]}\right)^2\right] = \frac{1}{n} \sum_{i=1}^n \sum_{t=1}^n \mathbb{E}[\mathbb{I}_{[j_t=i]}^2] + \frac{1}{n} \sum_{i=1}^n \sum_{t \neq t' \in [n]} \mathbb{E}[\mathbb{I}_{[j_t=i]} \mathbb{I}_{[j_{t'}=i]}] \\ &= \frac{1}{n} \sum_{t=1}^n \sum_{i=1}^n \mathbb{E}[\mathbb{I}_{[j_t=i]}] + \frac{1}{n} \sum_{i=1}^n \sum_{t \neq t' \in [n]} \mathbb{E}[\mathbb{I}_{[j_t=i]}] \mathbb{E}[\mathbb{I}_{[j_{t'}=i]}] = 1 + \frac{n^2 - n}{n^2} \leq 2. \end{aligned}$$

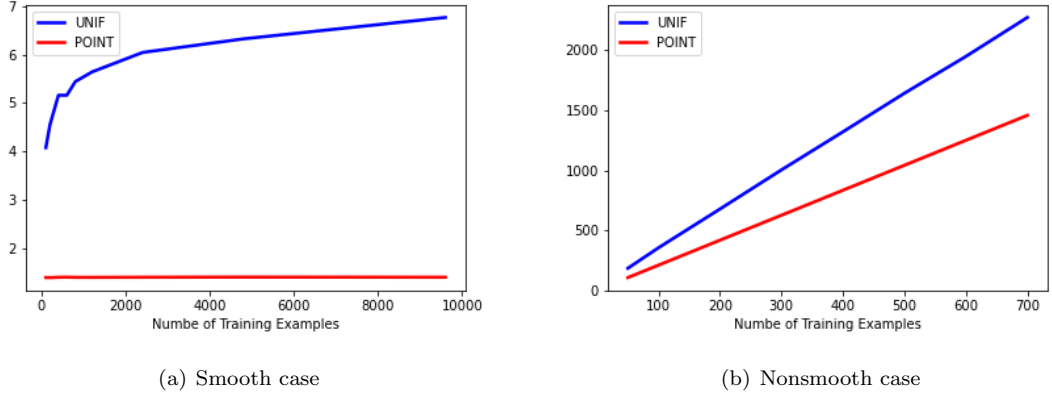


Figure 1: $\tilde{\beta}_{\text{unif}}$ (blue curve) and $\tilde{\beta}_{\text{point}}$ (red curve) as a function of n for SGD applied to convex and Lipschitz problems. Left panel considers the smooth case with $T = n$, where $\tilde{\beta}_{\text{unif}}$ and $\tilde{\beta}_{\text{point}}$ are defined in Eq. (4.3). Right panel considers the nonsmooth case with $T = n^2$, where $\tilde{\beta}_{\text{unif}}$ and $\tilde{\beta}_{\text{point}}$ are defined in Eq. (4.6).

294 As a comparison, $\mathbb{E}[\tilde{\beta}_{\text{unif}}] = \Theta\left(\frac{\log n}{\log \log n}\right)$ [29], which grows as n increases. We perform a simulation to
 295 compare $\tilde{\beta}_{\text{unif}}$ and $\tilde{\beta}_{\text{point}}$. We set $T = n$, and get a sequence of indices $\{j_t\}_{t \in [T]}$ by drawing j_t from the
 296 uniform distribution over $[n]$. Then, we compute $\tilde{\beta}_{\text{unif}}$ and $\tilde{\beta}_{\text{point}}$ according to Eq. (4.3). We repeat
 297 the experiments 25 times, and report the average of the experimental results. In Figure 1 (left panel),
 298 we plot $\tilde{\beta}_{\text{unif}}$ and $\tilde{\beta}_{\text{point}}$ as functions of n . The plot shows that $\tilde{\beta}_{\text{unif}}$ is substantially larger than $\tilde{\beta}_{\text{point}}$,
 299 and the difference grows as n increases. This shows the benefit of using pointwise uniform stability to
 300 study generalization.

301 **Remark 7.** Recently, fast excess risk bounds were derived for SGD based on the on-average model
 302 stability in the realizable (low-noise) setting [22]. Their bounds are stated in expectation, and their
 303 key idea is to incorporate the empirical risk in the stability bounds by using the expectation over S .
 304 For example, for SGD in a convex and smooth case, we can build the following inequality for two
 305 datasets $S, S^{(i)}$ differing by the i -th example

$$\|\mathbf{w}_{t+1} - \mathbf{w}_{t+1}^{(i)}\| \leq \|\mathbf{w}_t - \mathbf{w}_t^{(i)}\| + \eta_t \mathbb{1}_{[j_t=i]} \left(\|\nabla f(\mathbf{w}_t; z_i)\| + \|\nabla f(\mathbf{w}_t^{(i)}; z'_i)\| \right), \quad (4.5)$$

306 where z_i and z'_i are respectively the i -th example in S and $S^{(i)}$, and $\{\mathbf{w}_t\}, \{\mathbf{w}_t^{(i)}\}$ are SGD iterates
 307 on S and $S^{(i)}$, respectively. Then, the self-bounding property of smooth functions and the symmetry
 308 between z_i and z'_i imply

$$\mathbb{E}[\|\mathbf{w}_{t+1} - \mathbf{w}_{t+1}^{(i)}\|] \leq \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}_t^{(i)}\|] + \frac{\sqrt{2L}\eta_t}{n} \mathbb{E}[f^{\frac{1}{2}}(\mathbf{w}_t; z_i) + f^{\frac{1}{2}}(\mathbf{w}_t^{(i)}; z'_i)] = \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}_t^{(i)}\|] + \frac{2\sqrt{2L}\eta_t}{n} \mathbb{E}[f^{\frac{1}{2}}(\mathbf{w}_t; z_i)].$$

309 An average over $i \in [n]$ further includes the empirical risks in the stability bounds

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E}[\|\mathbf{w}_{t+1} - \mathbf{w}_{t+1}^{(i)}\|] \leq \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}_t^{(i)}\|] + \frac{2\sqrt{2L}\eta_t}{n} \mathbb{E}[F_S^{\frac{1}{2}}(\mathbf{w}_t)],$$

310 which implies fast rates if $F_S(\mathbf{w}_t)$ are small.

311 As a comparison, the pointwise uniform stability takes a supremum over all neighboring datasets,
 312 and this supremum comes from the bounded increment condition in Eq. (A.1), which takes the
 313 supremum over all z_j . Then, we need to take supremum over S on both sides of Eq. (4.5) to get

$$\sup_{S, z'_i} \|\mathbf{w}_{t+1} - \mathbf{w}_{t+1}^{(i)}\| \leq \sup_{S, z'_i} \|\mathbf{w}_t - \mathbf{w}_t^{(i)}\| + \eta_t \mathbb{I}_{[j_t=i]} \sup_{S, z'_i} \left(\|\nabla f(\mathbf{w}_t; z_i)\| + \|\nabla f(\mathbf{w}_t^{(i)}; z'_i)\| \right),$$

314 from which we fail to incorporate empirical risks in the stability bounds for a fast rate.

315 We now consider the convex and nonsmooth case. The following theorem shows that the stability
 316 of SGD in the nonsmooth case is worse than that in the smooth case.

317 **Theorem 9** (Stability of SGD: Nonsmooth Case). *Assume $f : \mathcal{W} \times \mathcal{Z} \mapsto \mathbb{R}^+$ is convex and G -*
 318 *Lipschitz. Let $\{\mathbf{w}_t\}_{t \in \mathbb{N}}$ be produced by SGD with $\eta_t = \eta$. Then SGD with T iterations is β -pointwise*
 319 *uniformly-stable in function values, where*

$$\frac{1}{n} \sum_{i=1}^n \beta_i^2 \leq \frac{4G^4\eta^2}{n} \left(Tn + 4(T+1)^{\frac{3}{2}}/3 + 4 \sum_{i=1}^n \sum_{t=1}^T \sum_{k=1}^t \mathbb{I}_{[j_t=i]} \mathbb{I}_{[j_k=i]} \right).$$

320 **Remark 8.** As we will show in the proof of Theorem 9, we can show that SGD is β_{unif} -uniformly
 321 stable with

$$\beta_{\text{unif}} \leq 2G^2\eta\sqrt{T} + 4G^2\eta \max_{i \in [n]} \sum_{t=1}^T \mathbb{I}_{[j_t=i]}.$$

322 Analogous to Remark 6, we introduce

$$\tilde{\beta}_{\text{unif}} = \sqrt{T} + 2 \max_{i \in [n]} \sum_{t=1}^T \mathbb{I}_{[j_t=i]}, \quad \tilde{\beta}_{\text{point}} = \left(T + \frac{4(T+1)^{\frac{3}{2}}}{3n} + \frac{4}{n} \sum_{i=1}^n \sum_{t=1}^T \sum_{k=1}^t \mathbb{I}_{[j_t=i]} \mathbb{I}_{[j_k=i]} \right)^{\frac{1}{2}}. \quad (4.6)$$

323 It is clear that $\tilde{\beta}_{\text{unif}}, \tilde{\beta}_{\text{point}}$ differ from the above uniform/pointwise stability bounds by a factor of
 324 $2G^2\eta$. Since $\sum_{t=1}^T \sum_{k=1}^{t-1} \mathbb{I}_{[j_t=i]} \mathbb{I}_{[j_k=i]} = \sum_{k=1}^T \sum_{t=k+1}^T \mathbb{I}_{[j_t=i]} \mathbb{I}_{[j_k=i]} = \sum_{t=1}^T \sum_{k=t+1}^T \mathbb{I}_{[j_t=i]} \mathbb{I}_{[j_k=i]}$, we
 325 know

$$\begin{aligned} 2 \sum_{i=1}^n \sum_{t=1}^T \sum_{k=1}^t \mathbb{I}_{[j_t=i]} \mathbb{I}_{[j_k=i]} &= \sum_{i=1}^n \sum_{t=1}^T \mathbb{I}_{[j_t=i]}^2 + \sum_{i=1}^n \sum_{t=1}^T \sum_{k=1}^t \mathbb{I}_{[j_t=i]} \mathbb{I}_{[j_k=i]} = T + \sum_{i=1}^n \left(\sum_{t=1}^T \mathbb{I}_{[j_t=i]} \right)^2 \\ &\leq T + \sum_{i=1}^n \left(\sum_{t=1}^T \mathbb{I}_{[j_t=i]} \right) \max_{i \in [n]} \sum_{t=1}^T \mathbb{I}_{[j_t=i]} = T + T \max_{i \in [n]} \sum_{t=1}^T \mathbb{I}_{[j_t=i]}. \end{aligned}$$

326 It then follows that

$$\tilde{\beta}_{\text{point}} \leq \left(T + \frac{4(T+1)^{\frac{3}{2}}}{3n} + \frac{2T}{n} + \frac{2T}{n} \max_{i \in [n]} \sum_{t=1}^T \mathbb{I}_{[j_t=i]} \right)^{\frac{1}{2}} \leq \tilde{\beta}_{\text{unif}}. \quad (4.7)$$

327 For the nonsmooth case, $\tilde{\beta}_{\text{point}}$ and $\tilde{\beta}_{\text{unif}}$ are of similar order. Indeed, if $T = O(n^2)$, the dominating
 328 term in both $\tilde{\beta}_{\text{point}}$ and $\tilde{\beta}_{\text{unif}}$ is \sqrt{T} . Furthermore, if $T = \Omega(n^2)$, then $\max_{i \in [n]} \sum_{t=1}^T \mathbb{I}_{[j_t=i]} = \Theta(T/n)$
 329 with high probability [29], which implies that $\tilde{\beta}_{\text{point}} = \Theta(T/n)$ and $\tilde{\beta}_{\text{unif}} = \Theta(T/n)$ in this case. In
 330 Figure 1 (right panel), we also plot $\tilde{\beta}_{\text{point}}$ and $\tilde{\beta}_{\text{unif}}$ as a function of n . We set $T = n^2$, and get a
 331 sequence of indices $\{j_t\}_{t \in [T]}$ by drawing j_t from the uniform distribution over $[n]$. Then, we compute

332 $\tilde{\beta}_{\text{unif}}$ and $\tilde{\beta}_{\text{point}}$ according to Eq. (4.6). We repeat the experiments 25 times, and report the average
 333 of the experimental results. The experimental results show that $\tilde{\beta}_{\text{point}}$ and $\tilde{\beta}_{\text{unif}}$ behave as linear
 334 functions of n in the nonsmooth case (note $T/n = n$ in our experiments), which is consistent with our
 335 theoretical analysis.

336 Finally, we consider SGD for strongly convex and smooth problems.

337 **Theorem 10** (Stability of SGD: Strongly Convex Case). *Let Assumption 1 hold with $\alpha = 1$ and ℓ be*
 338 *G -Lipschitz. Let $\{\mathbf{w}_t\}_{t \in \mathbb{N}}$ be produced by SGD with $\eta_t \leq 1/L$, where $L := L_\alpha + L_r$. Then SGD with T*
 339 *iterations is β -pointwise uniformly-stable in function values (measured by ℓ), where $\frac{1}{n} \sum_{i=1}^n \beta_i \leq \frac{4G^2}{n\lambda}$*
 340 *and*

$$\frac{1}{n} \sum_{i=1}^n \beta_i^2 = \frac{4G^4}{n} \sum_{i=1}^n \left(\sum_{t=1}^T \eta_t \mathbb{I}_{[j_t=i]} \prod_{t'=t+1}^T (1 - \eta_{t'} \lambda/2) \right)^2.$$

341 **Remark 9.** Under the same condition, one can show that SGD is β_{unif} -uniformly stable in function val-
 342 ues with (implicitly shown in the proof of Theorem 10) $\beta_{\text{unif}} = 2G^2 \max_{i \in [n]} \sum_{t=1}^T \eta_t \mathbb{I}_{[j_t=i]} \prod_{t'=t+1}^T (1 -$
 343 $\eta_{t'} \lambda/2)$. It is clear that $\beta_{\text{unif}}^2 \geq \frac{1}{n} \sum_{i=1}^n \beta_i^2$ for β_i in Theorem 10.

344 **Remark 10** (Lower bounds). Recently, lower bounds on the uniform stability were also developed for
 345 Lipschitz problems [3, 44, 19, 1]. A lower bound of order $\Omega(\min\{1, t/n\} \eta \sqrt{t} + \eta t/n)$ was established
 346 for SGD with convex and nonsmooth problems [3], a lower bound of order $\Omega(\eta t/n)$ was established for
 347 convex and smooth problems [44], and a lower bound of order $\Omega(\eta^2 n)$ was established for nonconvex
 348 problems [19]. These bounds are developed for uniform stability and are stated in expectation. As a
 349 comparison, this paper considers pointwise uniform stability. It is interesting to develop lower bounds
 350 on pointwise uniform stability with high probability.

351 4.2. Generalization Bounds

352 We now apply the above stability bounds to get high-probability generalization bounds of SGD.
 353 To our knowledge, Corollary 11 gives the first high-probability bounds on $\|\nabla F(A(S)) - \nabla F_S(A(S))\|$
 354 based on algorithmic stability. The bounds can be directly derived by plugging the stability bounds
 355 in Section 4.1 to Theorem 2 (Theorem 3). We omit the proofs for simplicity.

356 **Corollary 11** (Generalization of SGD: Smooth Case). *Assume $f : \mathcal{W} \times \mathcal{Z} \mapsto \mathbb{R}^+$ is convex, L -smooth*
 357 *and G -Lipschitz. Let $\{\mathbf{w}_t\}_{t \in \mathbb{N}}$ be produced by SGD with $\eta_t = \eta$. Let $\delta \in (0, 1)$. Then with probability*
 358 *at least $1 - \delta$ we have*

$$|F(A(S)) - F_S(A(S))| \lesssim \mathfrak{F}_1 \quad \text{and} \quad \|\nabla F(A(S)) - \nabla F_S(A(S))\| \lesssim \mathfrak{F}_1,$$

359 where

$$\mathfrak{F}_1 = \frac{\log^{\frac{1}{2}}(1/\delta)}{\sqrt{n}} + \frac{\eta \log n \log(1/\delta)}{\sqrt{n}} \left(\sum_{i=1}^n \left(\sum_{k=1}^T \mathbb{I}_{[j_k=i]} \right)^2 \right)^{\frac{1}{2}}.$$

360 We now turn to high-probability bounds for SGD applied to nonsmooth problems.

361 **Corollary 12** (Generalization of SGD: Nonsmooth Case). *Assume $f : \mathcal{W} \times \mathcal{Z} \mapsto \mathbb{R}^+$ is convex and*
362 *G -Lipschitz. Let $\{\mathbf{w}_t\}_{t \in \mathbb{N}}$ be produced by SGD with $\eta_t = \eta$ and $\delta \in (0, 1)$. With probability at least*
363 *$1 - \delta$, we have*

$$|F(A(S)) - F_S(A(S))| \lesssim \frac{\log^{\frac{1}{2}}(1/\delta)}{\sqrt{n}} + \frac{\eta \log n \log(1/\delta)}{\sqrt{n}} \left(Tn + T^{\frac{3}{2}} + \sum_{i=1}^n \sum_{t=1}^T \sum_{k=1}^t \mathbb{I}_{[j_t=i]} \mathbb{I}_{[j_k=i]} \right)^{\frac{1}{2}}.$$

364 Finally, we can directly plug Theorem 10 to Corollary 6 to get high-probability bounds for SGD
365 applied to strongly convex problems.

366 **Corollary 13** (Generalization of SGD: Strongly Convex Case). *Let Assumptions in Lemma 5 hold*
367 *with $\alpha = 1$ and ℓ be G -Lipschitz continuous. Let A be SGD with T iterations and $\eta_t \leq 1/L$. If*
368 *$\sup_z |\mathbb{E}_S[\ell(A(S)); z] - \ell(\mathbf{w}^*; z)| < \infty$, then for any $\delta \in (0, 1)$ with probability at least $1 - \delta$ we have*

$$F(A(S)) - F_S(A(S)) - F(\mathbf{w}^*) + F_S(\mathbf{w}^*) \lesssim \mathfrak{T}_2 \quad \text{and} \quad \|\nabla F(A(S)) - \nabla F_S(A(S)) + \nabla F_S(\mathbf{w}^*)\| \lesssim \mathfrak{T}_2,$$

369 where

$$\begin{aligned} \mathfrak{T}_2 = & \left(\frac{1}{n} \sum_{i=1}^n \left(\sum_{t=1}^T \eta_t \mathbb{I}_{[j_t=i]} \prod_{t'=t+1}^T (1 - \eta_{t'} \lambda/2) \right)^2 \right)^{\frac{1}{2}} \log n \log(1/\delta) \\ & + \frac{\log^{1/2}(1/\delta) (\mathbb{E}[F(A(S))])^{\frac{1}{2}}}{\lambda \sqrt{n}} \left((\mathbb{E}[F_S(A(S)) - F_S(A_e(S))])^{\frac{1}{2}} + \frac{\mathbb{E}[L(A_e(S))]}{\sqrt{n}} \right). \end{aligned}$$

370 5. Proofs on Connecting Stability and Generalization

371 5.1. Proof of Theorem 1

372 To prove Theorem 1, we need the following Marcinkiewicz-Zygmund's inequality for random vari-
373 ables taking values in a Hilbert space. It shows that the p -norm of a summation of independent random
374 variables can be bounded by the summation of the p -norm of random variables.

375 **Lemma 14.** *Let X_1, \dots, X_n be independent random variables taking values in a Hilbert space with*
376 *$\mathbb{E}[X_i] = 0$ for all $i \in [n]$. Then for any $p \geq 2$ we have*

$$\left\| \left\| \sum_{i=1}^n X_i \right\| \right\|_p \leq 2\sqrt{np} \left(\frac{1}{n} \sum_{i=1}^n \left\| \|X_i\| \right\|_p^p \right)^{\frac{1}{p}}.$$

377 The Marcinkiewicz-Zygmund's inequality can be proved by using its connection to Khintchine-
378 Kahane's inequality [4, page 441], where the Marcinkiewicz-Zygmund's inequality was established for
379 real-valued random variables. To get Marcinkiewicz-Zygmund's inequality for vector-valued random
380 variables, we need to use the following Khintchine-Kahane's inequality [10, Theorem 1.3.1]

$$\mathbb{E}_\epsilon \left\| \sum_{i=1}^n \epsilon_i X_i \right\|^p \leq \max((p-1)^{\frac{p}{2}}, 1) \left(\sum_{i=1}^n \|X_i\|^2 \right)^{\frac{p}{2}} \quad p \geq 2,$$

381 where X_1, \dots, X_n are elements in a Hilbert space, and $\epsilon_1, \dots, \epsilon_n$ are independent Rademacher variables
382 (i.e., taking values in $\{1, -1\}$ with the same probability). For brevity, we omit the proof of Lemma 14.

383 We now give the proof of Theorem 1, which is motivated by the analysis in [6]. For $f(Z_1, \dots, Z_n)$
384 and $A \subset [n]$, we write $\|f\|_p(Z_A) = (\mathbb{E}[|f|^p | Z_A])^{\frac{1}{p}}$.

385 *Proof of Theorem 1.* For simplicity, we assume $n = 2^k$. Define a sequence of partitions $\mathcal{B}_0, \dots, \mathcal{B}_k$ with
386 $\mathcal{B}_k = \{1, 2, \dots, 2^k\}$, where \mathcal{B}_l is derived from \mathcal{B}_{l+1} by splitting each subset in \mathcal{B}_{l+1} into two equal parts.
387 Then, there holds

$$\mathcal{B}_0 = \{\{1\}, \{2\}, \dots, \{2^k\}\}, \mathcal{B}_1 = \{\{1, 2\}, \{3, 4\}, \dots, \{2^k - 1, 2^k\}\}, \dots, \mathcal{B}_k = \{[n]\}.$$

388 For each $i \in [n]$ and $l = 0, 1, \dots, k$, denote by $B^l(i) \in \mathcal{B}_l$ the only set from \mathcal{B}_l containing i . According
389 to this definition, we know $B^0(i) = \{i\}$ and $B^k(i) = [n]$.

390 For each $i \in [n]$ and each $l = 0, 1, \dots, k$, we introduce random vectors as follows

$$g_i^l := g_i^l(Z_i, Z_{[n] \setminus B^l(i)}) = \mathbb{E}[g_i | Z_i, Z_{[n] \setminus B^l(i)}].$$

391 That is, we condition on Z_i and all the variables that are not in the same set as Z_i in \mathcal{B}_l . This definition
392 shows that $g_i^0 = g_i$ and $g_i^k = \mathbb{E}[g_i | Z_i]$. For each $i \in [n]$, we can decompose g_i as follows

$$g_i = \mathbb{E}[g_i | Z_i] + \sum_{l=0}^{k-1} (g_i^l - g_i^{l+1}).$$

393 It then follows from the triangle inequality that

$$\begin{aligned} \left\| \sum_{i=1}^n g_i \right\|_p &= \left\| \sum_{i=1}^n \left(\mathbb{E}[g_i | Z_i] + \sum_{l=0}^{k-1} (g_i^l - g_i^{l+1}) \right) \right\|_p \\ &\leq \left\| \sum_{i=1}^n \mathbb{E}[g_i | Z_i] \right\|_p + \sum_{l=0}^{k-1} \left\| \sum_{i=1}^n (g_i^l - g_i^{l+1}) \right\|_p. \end{aligned} \quad (5.1)$$

394 Since $\|\mathbb{E}[g_i | Z_i]\| \leq M$, one can check that $f(Z_1, \dots, Z_n) = \sum_{i=1}^n \mathbb{E}[g_i | Z_i]$ satisfies Eq. (A.1) with
395 $\beta_i = 2M$. Furthermore, we have $\mathbb{E}[\mathbb{E}[g_i | Z_i]] = 0$. Now we can apply Lemma A.3 with $\beta_i = 2M$ to
396 derive the following inequality

$$\left\| \sum_{i=1}^n \mathbb{E}[g_i | Z_i] \right\|_p \leq 2(\sqrt{2} + 1)\sqrt{np}M. \quad (5.2)$$

397 The definition of g_i^l implies that

$$\mathbb{E}_{Z_{B^{l+1}(i) \setminus B^l(i)}}[g_i^l] = g_i^{l+1}.$$

398 We view g_i^l as a function of $Z_j, j \in B^{l+1}(i) \setminus B^l(i)$. Changing any Z_j would change g_i^l by β_j . Therefore,
399 one can apply Lemma A.3 with $f = g_i^l$ to derive the following inequality with (there are 2^l random
400 variables)

$$\left\| g_i^l - g_i^{l+1} \right\|_p(Z_i, Z_{[n] \setminus B^{l+1}(i)}) \leq (\sqrt{2} + 1) \left(p \sum_{j \in B^{l+1}(i) \setminus B^l(i)} \beta_j^2 \right)^{\frac{1}{2}}. \quad (5.3)$$

401 We now turn to the sum $\sum_{i \in B} (g_i^l - g_i^{l+1})$ for any $B \in \mathcal{B}_l$. Consider any $i \in B \in \mathcal{B}_l$. Note
402 $Z'_i := g_i^l - g_i^{l+1}$ is a function of $Z_i, Z_{[n] \setminus B}$. We now condition on $Z_{[n] \setminus B}$ and then Z'_i is a function of
403 Z_i , which are independent. We can apply Lemma 14 to derive the following inequality

$$\left\| \sum_{i \in B} (g_i^l - g_i^{l+1}) \right\|_p^p(Z_{[n] \setminus B}) \leq \frac{(2\sqrt{p|B|})^p}{|B|} \sum_{i \in B} \left\| g_i^l - g_i^{l+1} \right\|_p^p(Z_{[n] \setminus B}).$$

404 Taking integration over $Z_{[n]\setminus B}$ gives

$$\left\| \left\| \sum_{i \in B} (g_i^l - g_i^{l+1}) \right\| \right\|_p^p \leq \frac{(2\sqrt{p|B|})^p}{|B|} \sum_{i \in B} \left\| \|g_i^l - g_i^{l+1}\| \right\|_p^p,$$

405 which implies

$$\left\| \left\| \sum_{i \in B} (g_i^l - g_i^{l+1}) \right\| \right\|_p \leq 2\sqrt{p|B|} \left(\frac{1}{|B|} \sum_{i \in B} \left\| \|g_i^l - g_i^{l+1}\| \right\|_p^p \right)^{\frac{1}{p}}.$$

406 This together with Eq. (5.3) implies that

$$\left\| \left\| \sum_{i \in B} (g_i^l - g_i^{l+1}) \right\| \right\|_p \leq 2p(\sqrt{2} + 1)|B|^{\frac{1}{2}} \left(\frac{1}{|B|} \sum_{i \in B} \left(\sum_{j \in B^{l+1}(i) \setminus B^l(i)} \beta_j^2 \right)^{\frac{p}{2}} \right)^{\frac{1}{p}}. \quad (5.4)$$

407 The rest of the proof is where we depart from the analysis of [6]. Note for any $i, i' \in B \in \mathcal{B}^l$, we have

$$B^{l+1}(i) \setminus B^l(i) = B^{l+1}(i') \setminus B^l(i').$$

408 Therefore, for any $B \in \mathcal{B}^l$ we have the following well-defined notation

$$\tilde{B} := \{j : j \in B^{l+1}(i) \setminus B^l(i)\}, \quad \forall i \in B,$$

409 which implies

$$\sum_{j \in B^{l+1}(i) \setminus B^l(i)} \beta_j^2 = \sum_{j \in \tilde{B}} \beta_j^2, \quad \text{if } i \in B \in \mathcal{B}^l. \quad (5.5)$$

410 One can interpret \tilde{B} as a sibling set of B in \mathcal{B}^l (they have the same parent set in \mathcal{B}^{l+1}). For example,

411 if $B = \{5, 6\}$, then $\tilde{B} = \{7, 8\}$. If $B = \{7, 8\}$, then $\tilde{B} = \{5, 6\}$. The parent set in \mathcal{B}^2 is $\{5, 6, 7, 8\}$. It

412 then follows from Eq. (5.4) and Eq. (5.5) the following inequality for any $B \in \mathcal{B}^l$ ($|B| = 2^l$)

$$\left\| \left\| \sum_{i \in B} (g_i^l - g_i^{l+1}) \right\| \right\|_p \leq 2p(\sqrt{2} + 1)|B|^{\frac{1}{2}} \left(\left(\sum_{j \in \tilde{B}} \beta_j^2 \right)^{\frac{p}{2}} \right)^{\frac{1}{p}} = 2p(\sqrt{2} + 1)2^{\frac{l}{2}} \left(\sum_{j \in \tilde{B}} \beta_j^2 \right)^{\frac{1}{2}}.$$

413 It then follows from the fact $|\mathcal{B}_l| = 2^{k-l}$ and the Cauchy-Schwartz inequality that

$$\begin{aligned} \left\| \left\| \sum_{i \in [n]} (g_i^l - g_i^{l+1}) \right\| \right\|_p &\leq \sum_{B \in \mathcal{B}_l} \left\| \left\| \sum_{i \in B} (g_i^l - g_i^{l+1}) \right\| \right\|_p \leq 2^{\frac{k-l}{2}} \left(\sum_{B \in \mathcal{B}_l} \left\| \left\| \sum_{i \in B} (g_i^l - g_i^{l+1}) \right\| \right\|_p^2 \right)^{\frac{1}{2}} \\ &\leq 2p(\sqrt{2} + 1)2^{\frac{k}{2}} \left(\sum_{B \in \mathcal{B}_l} \sum_{j \in \tilde{B}} \beta_j^2 \right)^{\frac{1}{2}}. \end{aligned}$$

414 According to our definition of \tilde{B} , one can check $\sum_{B \in \mathcal{B}_l} \sum_{j \in \tilde{B}} \beta_j^2 = \sum_{i=1}^n \beta_i^2$ and therefore

$$\left\| \left\| \sum_{i \in [n]} (g_i^l - g_i^{l+1}) \right\| \right\|_p \leq 2p(\sqrt{2} + 1)\sqrt{n} \left(\sum_{i=1}^n \beta_i^2 \right)^{\frac{1}{2}}.$$

415 This further gives

$$\sum_{l=0}^{k-1} \left\| \left\| \sum_{i=1}^n (g_i^l - g_i^{l+1}) \right\| \right\|_p \leq 2p(\sqrt{2} + 1)k \left(n \sum_{i=1}^n \beta_i^2 \right)^{\frac{1}{2}}.$$

416 We can plug Eq. (5.2) and the above inequality back into Eq. (5.1) to derive

$$\left\| \left\| \sum_{i=1}^n g_i \right\| \right\|_p \leq 2(\sqrt{2} + 1)\sqrt{np}M + 2p(\sqrt{2} + 1)k \left(n \sum_{i=1}^n \beta_i^2 \right)^{\frac{1}{2}}.$$

417 The proof is completed. \square

418 **Remark 11.** We highlight the difference between our proof and the analysis in [6]. We adopt the
 419 analysis in [6] to derive Eq. (5.4), excepting considering vector-valued random variables here. If the
 420 algorithm is β_{unif} -uniformly stable, then the analysis in [6] gives the following inequality similar to Eq.
 421 (5.4)

$$\left\| \left\| \sum_{i \in B} (g_i^l - g_i^{l+1}) \right\| \right\|_p \leq 2p(\sqrt{2} + 1)2^l \beta_{\text{unif}}, \quad \forall B \in \mathcal{B}^l.$$

422 Then one immediately gets

$$\left\| \left\| \sum_{i \in [n]} (g_i^l - g_i^{l+1}) \right\| \right\|_p \leq \sum_{B \in \mathcal{B}_l} \left\| \left\| \sum_{i \in B} (g_i^l - g_i^{l+1}) \right\| \right\|_p \leq 2^{k-l} 2^l 2p(\sqrt{2} + 1) \beta_{\text{unif}} = 2np(\sqrt{2} + 1) \beta_{\text{unif}}. \quad (5.6)$$

423 As a comparison, we get Eq. (5.4) to control $\left\| \left\| \sum_{i \in B} (g_i^l - g_i^{l+1}) \right\| \right\|_p$ in terms of $\sum_{i \in B} \left(\sum_{j \in B^{l+1}(i) \setminus B^l(i)} \beta_j^2 \right)^{\frac{1}{2}}$.
 424 Our observation is that $\sum_{j \in B^{l+1}(i) \setminus B^l(i)} \beta_j^2$ is the same for any $i \in B \in \mathcal{B}^l$, based on which we show

$$\left\| \left\| \sum_{i \in B} (g_i^l - g_i^{l+1}) \right\| \right\|_p^2 \leq 4p^2(\sqrt{2} + 1)^2 2^l \sum_{j \in \tilde{B}} \beta_j^2, \quad (5.7)$$

425 where \tilde{B} is a sibling of B . We then apply the Cauchy-Schwartz inequality to get

$$\left\| \left\| \sum_{i \in [n]} (g_i^l - g_i^{l+1}) \right\| \right\|_p \leq \sum_{B \in \mathcal{B}_l} \left\| \left\| \sum_{i \in B} (g_i^l - g_i^{l+1}) \right\| \right\|_p \leq 2^{\frac{k-l}{2}} \left(\sum_{B \in \mathcal{B}_l} \left\| \left\| \sum_{i \in B} (g_i^l - g_i^{l+1}) \right\| \right\|_p^2 \right)^{\frac{1}{2}}.$$

426 Finally, we can apply Eq. (5.7) to derive a bound similar to Eq. (5.6).

427 5.2. Proof of Theorem 3

428 In this section, we give the proof of Theorem 3.

429 *Proof of Theorem 3.* Let $S' = \{z'_1, \dots, z'_n\}$ be drawn independently from ρ . For any $i \in [n]$, define

$$S^{(i)} = \{z_1, \dots, z_{i-1}, z'_i, z_{i+1}, \dots, z_n\}. \quad (5.8)$$

430 Since $\mathbb{E}_Z[\nabla f(A(S); Z)] = \nabla F(A(S))$, we can decompose $\nabla F(A(S)) - \nabla F_S(A(S))$ as follows

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

431 Since A is β -pointwise uniformly stable in gradients, we know

$$n \|\nabla F(A(S)) - \nabla F_S(A(S))\| \leq 2 \sum_{i=1}^n \beta_i + \left\| \sum_{i=1}^n g_i \right\|, \quad (5.9)$$

432 where $g_i = \mathbb{E}_{z'_i} \left[\mathbb{E}_Z[\nabla f(A(S^{(i)}); Z)] - \nabla f(A(S^{(i)}); z_i) \right]$. According to our assumption, we know $\|g_i\| \leq$
 433 $2M$ and

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

434 where we have used the fact that z_i and Z follow from the same distribution. For any $i \in [n]$, any
 435 $j \neq i$ and any z_j'' , we have

$$\begin{aligned} & \left\| g_i(z_1, \dots, z_{j-1}, z_j, z_{j+1}, \dots, z_n) - g_i(z_1, \dots, z_{j-1}, z_j'', z_{j+1}, \dots, z_n) \right\| \\ &= \left\| \mathbb{E}_{z_i'} [\mathbb{E}_Z [\nabla f(A(S^{(i)}); Z)] - \nabla f(A(S^{(i)}); z_i)] - \mathbb{E}_{z_i'} [\mathbb{E}_Z [\nabla f(A(S_j^{(i)}); Z)] - \nabla f(A(S_j^{(i)}); z_i)] \right\| \\ &\leq \left\| \mathbb{E}_{z_i'} [\mathbb{E}_Z [\nabla f(A(S^{(i)}); Z) - \nabla f(A(S_j^{(i)}); Z)] \right\| + \left\| \mathbb{E}_{z_i'} [\nabla f(A(S^{(i)}); z_i) - \nabla f(A(S_j^{(i)}); z_i)] \right\| \leq 2\beta_j, \end{aligned}$$

436 where

$$S_j^{(i)} = \{z_1, \dots, z_{i-1}, z_i', z_{i+1}, \dots, z_{j-1}, z_j'', z_{j+1}, \dots, z_n\}. \quad (5.10)$$

437 Therefore, all the assumptions in Theorem 1 hold (with M replaced by $2M$ and β_j replaced by $2\beta_j$)
 438 and we can apply Theorem 1 to derive

$$\left\| \left\| \sum_{i=1}^n g_i \right\| \right\|_p \leq 4(\sqrt{2} + 1)\sqrt{np}M + 4p(\sqrt{2} + 1)\lceil \log_2 n \rceil \left(n \sum_{i=1}^n \beta_i^2 \right)^{\frac{1}{2}}.$$

439 We can combine the above inequality and Eq. (5.9) to derive the following inequality

$$n \left\| \left\| \nabla F(A(S)) - \nabla F_S(A(S)) \right\| \right\|_p \leq 2 \sum_{i=1}^n \beta_i + 4(\sqrt{2} + 1)\sqrt{np}M + 4p(\sqrt{2} + 1)\lceil \log_2 n \rceil \left(n \sum_{i=1}^n \beta_i^2 \right)^{\frac{1}{2}}.$$

440 By Lemma A.5, the following inequality holds with probability at least $1 - \delta$

$$n \left\| \left\| \nabla F(A(S)) - \nabla F_S(A(S)) \right\| \right\| \leq 2 \sum_{i=1}^n \beta_i + 4e(\sqrt{2} + 1)\sqrt{n \log(1/\delta)}M + 4e(\sqrt{2} + 1)\lceil \log_2 n \rceil \log(1/\delta) \left(n \sum_{i=1}^n \beta_i^2 \right)^{\frac{1}{2}}$$

441 and therefore

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

442 The proof is completed. \square

443 5.3. Proof of Theorem 7

444 To prove Theorem 7, we require the following lemma on the uniform stability of ERM for strongly
 445 convex problems. It is a direct extension of a similar result in [5] to functions with a structure in
 446 Assumption 1. Since the proof is identical to the classical stability analysis, we omit the proof for
 447 brevity.

448 **Lemma 15.** *Let Assumption 1 hold and ℓ be G -Lipschitz continuous. Then*

$$\max_{i \in [n]} \sup_{S, S^{(i)}} \sup_z [\ell(A_e(S); z) - \ell(A_e(S^{(i)}); z)] \leq 4G^2/(n\lambda),$$

449 where $S^{(i)}$ is defined in Eq. (5.8).

450 *Proof of Theorem 7.* According to Lemma 15, we know that A_e is $4G^2/(n\lambda)$ -uniformly stable in func-
 451 tion values (measured by ℓ). According to Theorem 4, the following inequality holds with probability
 452 at least $1 - \delta$

$$F(A_e(S)) - F_S(A_e(S)) - F(\mathbf{w}^*) + F_S(\mathbf{w}^*) \lesssim (n\lambda)^{-1} \log n \log(1/\delta) + \left(\frac{\sigma_A^2 \log(1/\delta)}{n} \right)^{\frac{1}{2}},$$

453 where (by Lemma 5)

$$\sigma_A^2 \leq C\lambda^{-1} \mathbb{E}_S[F(A_e(S)) - F(\mathbf{w}^*)].$$

454 We know

$$\begin{aligned} \mathbb{E}_S[F(A_e(S)) - F(\mathbf{w}^*)] &= \mathbb{E}_S[F(A_e(S)) - F_S(A_e(S))] + \mathbb{E}_S[F_S(A_e(S)) - F_S(\mathbf{w}^*)] + \mathbb{E}_S[F_S(\mathbf{w}^*) - F(\mathbf{w}^*)] \\ &\leq \mathbb{E}_S[F(A_e(S)) - F_S(A_e(S))] \leq \frac{4G^2}{n\lambda}. \end{aligned}$$

455 We can combine the above inequalities to derive the following inequality with probability at least $1 - \delta$

$$F(A_e(S)) - F_S(A_e(S)) - F(\mathbf{w}^*) + F_S(\mathbf{w}^*) \lesssim (n\lambda)^{-1} \log n \log(1/\delta). \quad (5.11)$$

456 According to Lemma B.1, we know

$$\begin{aligned} \langle A(S) - A_e(S), \nabla F(A_e(S)) \rangle &\leq \|A(S) - A_e(S)\| \|\nabla F(A_e(S))\| \\ &\leq \|A(S) - A_e(S)\| \left(\frac{L_r(1+\alpha)^{\frac{1}{1+\alpha}}}{2L_\alpha^{\frac{1}{1+\alpha}}} \left(F(A_e(S)) - F(\mathbf{w}^*) \right)^{\frac{1}{1+\alpha}} + 2 \left(\frac{L_\alpha}{1+\alpha} \right)^{\frac{1}{1+\alpha}} \left(F(A_e(S)) - F(\mathbf{w}^*) \right)^{\frac{\alpha}{1+\alpha}} \right) \\ &\leq \|A(S) - A_e(S)\| \left(F(A_e(S)) - F(\mathbf{w}^*) \right)^{\frac{\alpha}{1+\alpha}} \left(\frac{\tilde{C}L_r(1+\alpha)^{\frac{1}{1+\alpha}}}{2L_\alpha^{\frac{1}{1+\alpha}}} + 2 \left(\frac{L_\alpha}{1+\alpha} \right)^{\frac{1}{1+\alpha}} \right) \\ &:= C_1 \|A(S) - A_e(S)\| \left(F(A_e(S)) - F(\mathbf{w}^*) \right)^{\frac{\alpha}{1+\alpha}}, \end{aligned} \quad (5.12)$$

457 where we have used the assumption $\left(F(A_e(S)) - F(\mathbf{w}^*) \right)^{\frac{1-\alpha}{1+\alpha}} \leq \tilde{C}$ and introduced C_1 in the last step.

458 Since ℓ has (α, L_α) -Hölder continuous gradients and r is L_r -smooth, Eq. (B.1) implies

$$\begin{aligned} F(A(S)) - F(A_e(S)) &\leq \langle A(S) - A_e(S), \nabla F(A_e(S)) \rangle + \frac{L_\alpha \|A(S) - A_e(S)\|^{1+\alpha}}{1+\alpha} + \frac{L_r \|A(S) - A_e(S)\|^2}{2} \\ &\leq C_1 \|A(S) - A_e(S)\| \left(F(A_e(S)) - F(\mathbf{w}^*) \right)^{\frac{\alpha}{1+\alpha}} + \frac{L_\alpha \|A(S) - A_e(S)\|^{1+\alpha}}{1+\alpha} + \frac{L_r \|A(S) - A_e(S)\|^2}{2}. \end{aligned} \quad (5.13)$$

459 Since $F_S(A_e(S)) \leq F_S(\mathbf{w}^*)$, we can plug Eq. (5.11) to the above inequality and derive the following
 460 inequality with probability at least $1 - \delta$

$$F(A(S)) - F(A_e(S)) \lesssim \|A(S) - A_e(S)\| \left((n\lambda)^{-1} \log n \log(1/\delta) \right)^{\frac{\alpha}{1+\alpha}} + \|A(S) - A_e(S)\|^{1+\alpha} + \|A(S) - A_e(S)\|^2.$$

461 By the following inequality due to the strong convexity of F_S ,

$$F_S(A(S)) - F_S(A_e(S)) \geq \frac{\lambda}{2} \|A(S) - A_e(S)\|^2, \quad (5.14)$$

462 we get the following inequality with probability at least $1 - \delta$

$$F(A(S)) - F(A_e(S)) \lesssim \hat{\Delta}_\lambda^{\frac{1}{2}} \left((n\lambda)^{-1} \log n \log(1/\delta) \right)^{\frac{\alpha}{1+\alpha}} + \hat{\Delta}_\lambda^{\frac{1+\alpha}{2}}.$$

463 We can combine the above inequality and Eq. (5.11) to derive the following inequality with probability
 464 at least $1 - \delta$

$$F(A(S)) - F_S(A_e(S)) - F(\mathbf{w}^*) + F_S(\mathbf{w}^*) \lesssim (n\lambda)^{-1} \log n \log(1/\delta) + \hat{\Delta}_\lambda^{\frac{1}{2}} \left((n\lambda)^{-1} \log n \log(1/\delta) \right)^{\frac{\alpha}{1+\alpha}} + \hat{\Delta}_\lambda^{\frac{1+\alpha}{2}}.$$

465 By the following Young's inequality

$$\hat{\Delta}_\lambda^{\frac{1}{2}} \left((n\lambda)^{-1} \log n \log(1/\delta) \right)^{\frac{\alpha}{1+\alpha}} \leq \frac{\alpha}{1+\alpha} \left((n\lambda)^{-1} \log n \log(1/\delta) \right)^{\frac{\alpha}{1+\alpha} \frac{1+\alpha}{\alpha}} + \frac{1}{1+\alpha} \hat{\Delta}_\lambda^{\frac{1+\alpha}{2}},$$

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

$$F(A(S)) - F_S(A_e(S)) - F(\mathbf{w}^*) + F_S(\mathbf{w}^*) \lesssim (n\lambda)^{-1} \log n \log(1/\delta) + \hat{\Delta}_\lambda^{\frac{1+\alpha}{2}}.$$

467 The stated bound then follows by noting $F_S(\mathbf{w}^*) \geq F_S(A_e(S))$. The proof is completed. \square

468 6. Proofs on Stochastic Gradient Descent

469 In this section, we present the proof on the stability bounds of SGD. Our analysis is based on the
 470 following lemma in [17], which shows that the gradient update $\mathbf{w} \mapsto \mathbf{w} - \eta \nabla f(\mathbf{w}; z)$ is nonexpansive if
 471 f is convex and smooth.

472 **Lemma 16** ([17]). *Suppose the function $\mathbf{w} \mapsto f(\mathbf{w}; z)$ is convex and L -smooth. If $\eta \leq 2/L$, then*

$$\|(\mathbf{w} - \eta \nabla f(\mathbf{w}; z)) - (\mathbf{w}' - \eta \nabla f(\mathbf{w}'; z))\| \leq \|\mathbf{w} - \mathbf{w}'\|.$$

473 Furthermore, if $\mathbf{w} \mapsto f(\mathbf{w}; z)$ is λ -strongly convex and $\eta \leq 1/L$, then

$$\|(\mathbf{w} - \eta \nabla f(\mathbf{w}; z)) - (\mathbf{w}' - \eta \nabla f(\mathbf{w}'; z))\|^2 \leq (1 - \eta\lambda) \|\mathbf{w} - \mathbf{w}'\|^2.$$

474 Let $S^{(i)}$ be defined by Eq. (5.8). Let $\{\mathbf{w}_t^{(i)}\}$ be produced by SGD w.r.t. $S^{(i)}$.

475 *Proof of Theorem 8.* We build a recurrent formula on estimating $\|\mathbf{w}_{t+1} - \mathbf{w}_{t+1}^{(i)}\|$. Consider two cases
 476 at the t -th iteration. If $j_t \neq i$, then Lemma 16 implies that

$$\|\mathbf{w}_{t+1} - \mathbf{w}_{t+1}^{(i)}\| = \|(\mathbf{w}_t - \eta_t \nabla f(\mathbf{w}_t; z_{j_t})) - (\mathbf{w}_t^{(i)} - \eta_t \nabla f(\mathbf{w}_t^{(i)}; z_{j_t}))\| \leq \|\mathbf{w}_t - \mathbf{w}_t^{(i)}\|. \quad (6.1)$$

477 If $j_t = i$, then $\|\mathbf{w}_{t+1} - \mathbf{w}_{t+1}^{(i)}\| \leq \|\mathbf{w}_t - \mathbf{w}_t^{(i)}\| + 2G\eta_t$. We can combine the above two inequalities to get

$$\|\mathbf{w}_{t+1} - \mathbf{w}_{t+1}^{(i)}\| \leq \|\mathbf{w}_t - \mathbf{w}_t^{(i)}\| + 2G\eta_t \mathbb{I}_{[j_t=i]}.$$

478 We apply the above inequality recursively and get

$$\|\mathbf{w}_{t+1} - \mathbf{w}_{t+1}^{(i)}\| \leq 2G\eta \sum_{k=1}^t \mathbb{I}_{[j_k=i]}.$$

479 By the Lipschitz continuity, we know that SGD with T iterations is β -pointwise uniformly stable,
 480 where $\beta_i = 2G^2\eta \sum_{k=1}^T \mathbb{I}_{[j_k=i]}$. It then follows that

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \beta_i^2 &= \frac{4G^4\eta^2}{n} \sum_{i=1}^n \left(\sum_{k=1}^T \mathbb{I}_{[j_k=i]} \right)^2, \\ \frac{1}{n} \sum_{i=1}^n \beta_i &= \frac{2G^2\eta}{n} \sum_{i=1}^n \sum_{k=1}^T \mathbb{I}_{[j_k=i]} = \frac{2G^2\eta}{n} \sum_{k=1}^T \sum_{i=1}^n \mathbb{I}_{[j_k=i]} = \frac{2G^2T\eta}{n}, \end{aligned}$$

481 where we have used $\sum_{i=1}^n \mathbb{I}_{[j_k=i]} = 1$ for any k . The proof is completed. \square

482 *Proof of Theorem 9.* Consider two cases at the t -th iteration. In the first case, assume $j_t \neq i$. Then
 483 the Lipschitz continuity of f implies that

$$\begin{aligned} \|\mathbf{w}_{t+1} - \mathbf{w}_{t+1}^{(i)}\|^2 &= \|(\mathbf{w}_t - \eta_t \nabla f(\mathbf{w}_t; z_{j_t})) - (\mathbf{w}_t^{(i)} - \eta_t \nabla f(\mathbf{w}_t^{(i)}; z_{j_t}))\|^2 \\ &= \|\mathbf{w}_t - \mathbf{w}_t^{(i)}\|^2 - 2\eta_t \langle \mathbf{w}_t - \mathbf{w}_t^{(i)}, \nabla f(\mathbf{w}_t; z_{j_t}) - \nabla f(\mathbf{w}_t^{(i)}; z_{j_t}) \rangle + \eta_t^2 \|\nabla f(\mathbf{w}_t; z_{j_t}) - \nabla f(\mathbf{w}_t^{(i)}; z_{j_t})\|^2 \\ &\leq \|\mathbf{w}_t - \mathbf{w}_t^{(i)}\|^2 + 4G^2\eta_t^2, \end{aligned}$$

484 where we have used the inequality $\langle \mathbf{w}_t - \mathbf{w}_t^{(i)}, \nabla f(\mathbf{w}_t; z_{j_t}) - \nabla f(\mathbf{w}_t^{(i)}; z_{j_t}) \rangle \geq 0$ due to the convexity of
 485 f . For the second case, assume $j_t = i$. Then

$$\begin{aligned} \|\mathbf{w}_{t+1} - \mathbf{w}_{t+1}^{(i)}\|^2 &= \|(\mathbf{w}_t - \eta_t \nabla f(\mathbf{w}_t; z_{j_t})) - (\mathbf{w}_t^{(i)} - \eta_t \nabla f(\mathbf{w}_t^{(i)}; z'_{j_t}))\|^2 \\ &= \|\mathbf{w}_t - \mathbf{w}_t^{(i)}\|^2 + \eta_t^2 \|\nabla f(\mathbf{w}_t; z_{j_t}) - \nabla f(\mathbf{w}_t^{(i)}; z'_{j_t})\|^2 - 2\eta_t \langle \mathbf{w}_t - \mathbf{w}_t^{(i)}, \nabla f(\mathbf{w}_t; z_{j_t}) - \nabla f(\mathbf{w}_t^{(i)}; z'_{j_t}) \rangle \\ &\leq \|\mathbf{w}_t - \mathbf{w}_t^{(i)}\|^2 + 4G^2\eta_t^2 + 4G\eta_t \|\mathbf{w}_t - \mathbf{w}_t^{(i)}\|, \end{aligned}$$

486 where we have used the Lipschitz continuity of f . We can combine the above two cases to get

$$\|\mathbf{w}_{t+1} - \mathbf{w}_{t+1}^{(i)}\|^2 \leq \|\mathbf{w}_t - \mathbf{w}_t^{(i)}\|^2 + 4G^2\eta_t^2 + 4G\eta_t \|\mathbf{w}_t - \mathbf{w}_t^{(i)}\| \mathbb{I}_{[j_t=i]}.$$

487 We apply the above inequality recursively and get

$$\|\mathbf{w}_{T+1} - \mathbf{w}_{T+1}^{(i)}\|^2 \leq 4G^2\eta^2 T + 4G\eta \sum_{t=1}^T \|\mathbf{w}_t - \mathbf{w}_t^{(i)}\| \mathbb{I}_{[j_t=i]}$$

488 and therefore

$$\sup_{S, S^{(i)}} \|\mathbf{w}_{T+1} - \mathbf{w}_{T+1}^{(i)}\|^2 \leq 4G^2\eta^2 T + 4G\eta \sum_{t=1}^T \sup_{S, S^{(i)}} \|\mathbf{w}_t - \mathbf{w}_t^{(i)}\| \mathbb{I}_{[j_t=i]},$$

489 where the supremum is taken over two neighboring datasets differing by the i -th example. Let

$$\beta_{T+1,i} = G \sup_{S, S^{(i)}} \|\mathbf{w}_{T+1} - \mathbf{w}_{T+1}^{(i)}\|. \quad (6.2)$$

490 Then we know SGD with T iterations is β_{T+1} -pointwise uniformly stable in function values, where
 491 $\beta_{T+1,i}$ satisfies the following inequality

$$\beta_{T+1,i}^2 \leq 4G^4\eta^2 T + 4G^2\eta \sum_{t=1}^T \beta_{t,i} \mathbb{I}_{[j_t=i]}. \quad (6.3)$$

492 Let $\Delta_{t,i} = \max_{k \leq t} \beta_{k,i}$. Then the above inequality implies that $\Delta_{T,i}^2 \leq 4G^4\eta^2 T + 4G^2\eta \Delta_{T,i} \sum_{t=1}^T \mathbb{I}_{[j_t=i]}$.

493 Solving this quadratic inequality of $\Delta_{T,i}$ implies that

$$\Delta_{T,i} \leq 2G^2\eta\sqrt{T} + 4G^2\eta \sum_{t=1}^T \mathbb{I}_{[j_t=i]}. \quad (6.4)$$

494 We can plug the above bound back into Eq. (6.3), and get

$$\begin{aligned} \beta_{T+1,i}^2 &\leq 4G^4\eta^2 T + 4G^2\eta \sum_{t=1}^T \mathbb{I}_{[j_t=i]} \left(2G^2\eta\sqrt{t} + 4G^2\eta \sum_{k=1}^t \mathbb{I}_{[j_k=i]} \right) \\ &\leq 4G^4\eta^2 T + 8G^4\eta^2 \sum_{t=1}^T \sqrt{t} \mathbb{I}_{[j_t=i]} + 16G^4\eta^2 \sum_{t=1}^T \sum_{k=1}^t \mathbb{I}_{[j_t=i]} \mathbb{I}_{[j_k=i]}. \end{aligned}$$

495 It then follows that $(\sum_{i=1}^n \mathbb{I}_{[j_i=i]} = 1)$

$$\sum_{i=1}^n \beta_{T+1,i}^2 \leq 4G^4 \eta^2 T n + \frac{16}{3} G^4 \eta^2 (T+1)^{\frac{3}{2}} + 16G^4 \eta^2 \sum_{i=1}^n \sum_{t=1}^T \sum_{k=1}^t \mathbb{I}_{[j_t=i]} \mathbb{I}_{[j_k=i]}.$$

496 The proof is completed. \square

497 *Proof of Theorem 10.* It is clear that for any z , the function $\mathbf{w} \mapsto f(\mathbf{w}; z)$ is L -smooth. We now build
498 a recurrent formula on estimating $\|\mathbf{w}_{t+1} - \mathbf{w}_{t+1}^{(i)}\|$. Consider two cases at the t -th iteration. In the first
499 case, assume $j_t \neq i$. Then Lemma 16 implies that

$$\|\mathbf{w}_{t+1} - \mathbf{w}_{t+1}^{(i)}\| = \|(\mathbf{w}_t - \eta_t \nabla f(\mathbf{w}_t; z_{j_t})) - (\mathbf{w}_t^{(i)} - \eta_t \nabla f(\mathbf{w}_t^{(i)}; z_{j_t}))\| \leq (1 - \eta_t \lambda/2) \|\mathbf{w}_t - \mathbf{w}_t^{(i)}\|. \quad (6.5)$$

500 For the second case, assume $j_t = i$. Then

$$\begin{aligned} \|\mathbf{w}_{t+1} - \mathbf{w}_{t+1}^{(i)}\| &= \|(\mathbf{w}_t - \eta_t \nabla f(\mathbf{w}_t; z_{j_t})) - (\mathbf{w}_t^{(i)} - \eta_t \nabla f(\mathbf{w}_t^{(i)}; z'_{j_t}))\| \\ &\leq \|(\mathbf{w}_t - \eta_t \nabla f(\mathbf{w}_t; z_{j_t})) - (\mathbf{w}_t^{(i)} - \eta_t \nabla f(\mathbf{w}_t^{(i)}; z_{j_t}))\| + \eta_t \|\nabla f(\mathbf{w}_t^{(i)}; z_{j_t}) - \nabla f(\mathbf{w}_t^{(i)}; z'_{j_t})\| \\ &= \|(\mathbf{w}_t - \eta_t \nabla f(\mathbf{w}_t; z_{j_t})) - (\mathbf{w}_t^{(i)} - \eta_t \nabla f(\mathbf{w}_t^{(i)}; z_{j_t}))\| + \eta_t \|\nabla \ell(\mathbf{w}_t^{(i)}; z_{j_t}) - \nabla \ell(\mathbf{w}_t^{(i)}; z'_{j_t})\| \\ &\leq (1 - \eta_t \lambda/2) \|\mathbf{w}_t - \mathbf{w}_t^{(i)}\| + 2G\eta_t, \end{aligned} \quad (6.6)$$

501 where we have used the Lipschitz continuity of ℓ . We can combine Eq. (6.5) and Eq. (6.6) to get

$$\|\mathbf{w}_{t+1} - \mathbf{w}_{t+1}^{(i)}\| \leq (1 - \eta_t \lambda/2) \|\mathbf{w}_t - \mathbf{w}_t^{(i)}\| + 2G\eta_t \mathbb{I}_{[j_t=i]}. \quad (6.7)$$

502 We apply the above inequality recursively, and derive

$$\|\mathbf{w}_{T+1} - \mathbf{w}_{T+1}^{(i)}\| \leq 2G \sum_{t=1}^T \eta_t \mathbb{I}_{[j_t=i]} \prod_{t'=t+1}^T (1 - \eta_{t'} \lambda/2). \quad (6.8)$$

503 Therefore, SGD is β -pointwise uniformly stable in function values, where

$$\beta_i = 2G^2 \sum_{t=1}^T \eta_t \mathbb{I}_{[j_t=i]} \prod_{t'=t+1}^T (1 - \eta_{t'} \lambda/2). \quad (6.9)$$

504 It then follows that

$$\frac{1}{n} \sum_{i=1}^n \beta_i^2 = \frac{4G^4}{n} \sum_{i=1}^n \left(\sum_{t=1}^T \eta_t \mathbb{I}_{[j_t=i]} \prod_{t'=t+1}^T (1 - \eta_{t'} \lambda/2) \right)^2$$

505 and

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \beta_i &= \frac{2G^2}{n} \sum_{i=1}^n \sum_{t=1}^T \eta_t \mathbb{I}_{[j_t=i]} \prod_{t'=t+1}^T (1 - \eta_{t'} \lambda/2) = \frac{2G^2}{n} \sum_{t=1}^T \eta_t \sum_{i=1}^n \mathbb{I}_{[j_t=i]} \prod_{t'=t+1}^T (1 - \eta_{t'} \lambda/2) \\ &= \frac{2G^2}{n} \sum_{t=1}^T \eta_t \prod_{t'=t+1}^T (1 - \eta_{t'} \lambda/2) = \frac{4G^2}{n\lambda} \sum_{t=1}^T (1 - (1 - \eta_t \lambda/2)) \prod_{t'=t+1}^T (1 - \eta_{t'} \lambda/2) \\ &= \frac{4G^2}{n\lambda} \sum_{t=1}^T \left(\prod_{t'=t+1}^T (1 - \eta_{t'} \lambda/2) - \prod_{t'=t}^T (1 - \eta_{t'} \lambda/2) \right) = \frac{4G^2}{n\lambda} \left(1 - \prod_{t=1}^T (1 - \eta_t \lambda/2) \right), \end{aligned}$$

506 where we have used $\sum_{i=1}^n \mathbb{I}_{[j_t=i]} = 1$ for any $t \in [T]$. The proof is completed. \square

507 **7. Conclusion**

508 In this paper, we introduce the pointwise uniform stability to develop high-probability generaliza-
 509 tion bounds. The pointwise uniform stability considers the effect of changing each example in the
 510 dataset, which is weaker than the uniform stability. We first develop a moment bound for a sum-
 511 mation of weakly-dependent vector-valued random variables, and apply it to develop bounds for the
 512 generalization gap as measured by either function values or gradients. We improve the recently fast
 513 high-probability rates in [18] by relaxing the requirement on strong convexity and Lipschitz continu-
 514 ity, and improving the dependency on optimization errors. Finally, we apply our results to develop
 515 improved generalization bounds for SGD.

516 Our generalization bounds involve a factor of $\log(n)$ in front of $(\frac{1}{n} \sum_{i=1}^n \beta_i^2)^{\frac{1}{2}}$. A very interesting
 517 question is to see whether this logarithmic factor can be removed. Indeed, if we can remove this
 518 logarithmic factor, the resulting generalization bound would be optimal up to a constant factor. It is
 519 also interesting to apply the stability analysis to study SGD with functional data [8, 16].

520 **Acknowledgements**

521 The authors are grateful to the associate editor and the anonymous reviewers for their thoughtful
 522 comments and constructive suggestions. The work by Jun Fan is partially supported by the Research
 523 Grants Council of Hong Kong [Project No. HKBU 12302819] and [Project No. HKBU 12303220], and
 524 Hong Kong Baptist University [Project No. RC-FNRA-IG/22-23/SCI/02]. The work by Yunwen Lei
 525 is partially supported by the Research Grants Council of Hong Kong [Project No. 22303723] and URC
 526 Seed Fund for Basic Research for New Staff 2023-24.

527 **Appendix**

528 **A. Useful Inequalities in Probability**

529 *A.1. McDiarmid's Inequality*

530 We first consider McDiarmid's inequality for real-valued functions of random variables, which follows
 531 from the standard tail-bound of McDiarmid's inequality and Proposition 2.5.2 in [40].

532 **Lemma A.1** (McDiarmid's Inequality for Real-Valued Functions). *Let Z_1, \dots, Z_n be independent ran-*
 533 *dom variables, and $f : \mathcal{Z}^n \mapsto \mathbb{R}$ such that the following inequality holds for any i and $z_1, \dots, z_{i-1}, z_{i+1}, \dots, z_n$*

$$\sup_{z_i, z_i'} |f(z_1, \dots, z_{i-1}, z_i, z_{i+1}, \dots, z_n) - f(z_1, \dots, z_{i-1}, z_i', z_{i+1}, \dots, z_n)| \leq \beta_i.$$

534 *Then for any $p \geq 1$ we have*

$$\|f(Z_1, \dots, Z_n) - \mathbb{E}[f(Z_1, \dots, Z_n)]\|_p \leq \left(2p \sum_{i=1}^n \beta_i^2\right)^{\frac{1}{2}}.$$

535 Now we consider **vector-valued** functions of independent random variables. The following lemma
 536 gives the expected distance between $f(Z_1, \dots, Z_n)$ and its expectation.

537 **Lemma A.2** ([31]). *Let Z_1, \dots, Z_n be independent random variables, and $f : \mathcal{Z}^n \mapsto \mathcal{H}$ a function into*
 538 *a Hilbert space \mathcal{H} such that the following inequality holds for any i and $z_1, \dots, z_{i-1}, z_{i+1}, \dots, z_n$*

$$\sup_{z_i, z'_i} \|f(z_1, \dots, z_{i-1}, z_i, z_{i+1}, \dots, z_n) - f(z_1, \dots, z_{i-1}, z'_i, z_{i+1}, \dots, z_n)\| \leq \beta_i. \quad (\text{A.1})$$

539 *Then*

$$\mathbb{E}[\|f(Z_1, \dots, Z_n) - \mathbb{E}[f(Z_1, \dots, Z_n)]\|] \leq \left(\sum_{i=1}^n \beta_i^2 \right)^{\frac{1}{2}}.$$

540 The following lemma controls the p -norm for the vector-valued random variable $f(Z_1, \dots, Z_n) -$
 541 $\mathbb{E}[f(Z_1, \dots, Z_n)]$.

542 **Lemma A.3** (McDiarmid's Inequality for Vector-Valued Functions). *Let assumptions in Lemma A.2*
 543 *hold. Then for any $p \geq 1$ we have*

$$\left\| \|f(Z_1, \dots, Z_n) - \mathbb{E}[f(Z_1, \dots, Z_n)]\| \right\|_p \leq (\sqrt{2} + 1) \left(p \sum_{i=1}^n \beta_i^2 \right)^{\frac{1}{2}}. \quad (\text{A.2})$$

544 *Proof.* We define a real-valued function $g : \mathcal{Z}^n \mapsto \mathbb{R}$ as

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

545 We first show this function satisfies the increment condition. Indeed, for any i and $z_1, \dots, z_{i-1}, z_{i+1}, \dots, z_n$
 546 we have

$$\begin{aligned} & \sup_{z_i, z'_i} |g(z_1, \dots, z_{i-1}, z_i, z_{i+1}, \dots, z_n) - g(z_1, \dots, z_{i-1}, z'_i, z_{i+1}, \dots, z_n)| \\ &= \sup_{z_i, z'_i} \left| \|f(z_1, \dots, z_n) - \mathbb{E}[f(Z_1, \dots, Z_n)]\| - \|f(z_1, \dots, z_{i-1}, z'_i, z_{i+1}, \dots, z_n) - \mathbb{E}[f(Z_1, \dots, Z_n)]\| \right| \\ &\leq \sup_{z_i, z'_i} \left\| f(z_1, \dots, z_n) - f(z_1, \dots, z_{i-1}, z'_i, z_{i+1}, \dots, z_n) \right\| \leq \beta_i. \end{aligned}$$

547 Therefore, we can apply Lemma A.1 to the real-valued function g and derive the following inequality

$$\left\| g(Z_1, \dots, Z_n) - \mathbb{E}[g(Z_1, \dots, Z_n)] \right\|_p \leq \left(2p \sum_{i=1}^n \beta_i^2 \right)^{\frac{1}{2}}.$$

548 According to Lemma A.2, we know the following inequality $\mathbb{E}[g(Z_1, \dots, Z_n)] \leq \left(\sum_{i=1}^n \beta_i^2 \right)^{\frac{1}{2}}$. We can
 549 combine the above two inequalities together and derive the stated inequality. \square

550 A.2. Bernstein Inequality and Tails

551 The following lemma gives a Bernstein inequality to incorporate the variance information in bound-
 552 ing a summation of independent random variables [9].

553 **Lemma A.4** (Bernstein inequality). Let $\{\xi(z_i)\}_{i=1}^n$ be a sequence of independent and identically
554 distributed real-valued random variables and \widetilde{M} be a constant such that $|\xi| \leq \widetilde{M}$ and the variance
555 $\text{Var}(\xi) < \infty$. Then, for any $0 < \delta < 1$ with probability at least $1 - \delta$ there holds

$$\mathbb{E}[\xi] - \frac{1}{n} \sum_{i=1}^n \xi(z_i) \leq \frac{2\widetilde{M} \log \frac{1}{\delta}}{3n} + \sqrt{\frac{2 \text{Var}(\xi) \log \frac{1}{\delta}}{n}}.$$

556 The following lemma shows the relationship between tails and moments [6].

557 **Lemma A.5.** Let Y be a random variable. If $\|Y\|_p \leq \sqrt{pa}$ for any $p \geq 2$, then for any $\delta \in (0, 1)$ with
558 probability at least $1 - \delta$: $|Y| \leq ea\sqrt{\log(e/\delta)}$.

559 B. Self-Bounding Property

560 We present some useful self-bounding properties for functions of a composite structure in Assump-
561 tion 1. The self-bounding property will be very important for our stability and generalization analysis.

562 **Lemma B.1.** Assume $F(\mathbf{w}) = L(\mathbf{w}) + r(\mathbf{w})$, where L has (α, L_α) -Hölder continuous gradients and r
563 is L_r -smooth. Then we have

$$\|\nabla F(\mathbf{w})\| \leq \frac{L_r(1+\alpha)^{\frac{1}{1+\alpha}}}{2L_\alpha^{\frac{1}{1+\alpha}}} \left(F(\mathbf{w}) - F(\mathbf{w}^*)\right)^{\frac{1}{1+\alpha}} + 2\left(\frac{L_\alpha}{1+\alpha}\right)^{\frac{1}{1+\alpha}} \left(F(\mathbf{w}) - F(\mathbf{w}^*)\right)^{\frac{\alpha}{1+\alpha}}.$$

564 **Lemma B.2.** Assume $F(\mathbf{w}) = L(\mathbf{w}) + r(\mathbf{w})$, where L has (α, L_α) -Hölder continuous gradients and r
565 is L_r -smooth. If F is nonnegative, then we have

$$\|\nabla F(\mathbf{w})\| \leq \frac{L_r(1+\alpha)^{\frac{1}{1+\alpha}}}{2L_\alpha^{\frac{1}{1+\alpha}}} F^{\frac{1}{1+\alpha}}(\mathbf{w}) + 2\left(\frac{L_\alpha}{1+\alpha}\right)^{\frac{1}{1+\alpha}} F^{\frac{\alpha}{1+\alpha}}(\mathbf{w}).$$

566 *Proof.* If $\nabla F(\mathbf{w}) = 0$, the inequality holds immediately. Now we only consider the case that $\nabla F(\mathbf{w}) \neq$
567 0. Since L has (α, L_α) -Hölder continuous gradients, we know $L(\mathbf{w}') \leq L(\mathbf{w}) + \langle \mathbf{w}' - \mathbf{w}, \nabla L(\mathbf{w}) \rangle +$
568 $\frac{L_\alpha}{1+\alpha} \|\mathbf{w} - \mathbf{w}'\|^{1+\alpha}$ [43]. Since r is L_r -smooth, we know $r(\mathbf{w}') \leq r(\mathbf{w}) + \langle \mathbf{w}' - \mathbf{w}, \nabla r(\mathbf{w}) \rangle + \frac{L_r}{2} \|\mathbf{w} -$
569 $\mathbf{w}'\|^2$ [26]. It then follows that

$$F(\mathbf{w}') \leq F(\mathbf{w}) + \langle \mathbf{w}' - \mathbf{w}, \nabla F(\mathbf{w}) \rangle + \frac{L_\alpha}{1+\alpha} \|\mathbf{w} - \mathbf{w}'\|^{1+\alpha} + \frac{L_r}{2} \|\mathbf{w} - \mathbf{w}'\|^2. \quad (\text{B.1})$$

570 We choose

$$\mathbf{w}' = \mathbf{w} - A\|\nabla F(\mathbf{w})\|^{-1}\nabla F(\mathbf{w}), \quad A := \left(\frac{(1+\alpha)F(\mathbf{w})}{L_\alpha}\right)^{\frac{1}{1+\alpha}}.$$

571 It then follows that

$$0 \leq F(\mathbf{w}') \leq F(\mathbf{w}) - A\|\nabla F(\mathbf{w})\| + \frac{L_r A^2}{2} + \frac{L_\alpha A^{1+\alpha}}{1+\alpha}.$$

572 That is,

$$\|\nabla F(\mathbf{w})\| \leq \frac{L_r A}{2} + \frac{L_\alpha A^\alpha}{1+\alpha} + \frac{F(\mathbf{w})}{A}.$$

573 According to our construction of A , we further have

$$\begin{aligned}\|\nabla F(\mathbf{w})\| &\leq \frac{L_r}{2} \left(\frac{(1+\alpha)F(\mathbf{w})}{L_\alpha} \right)^{\frac{1}{1+\alpha}} + \frac{L_\alpha}{1+\alpha} \left(\frac{(1+\alpha)F(\mathbf{w})}{L_\alpha} \right)^{\frac{\alpha}{1+\alpha}} + F(\mathbf{w}) \left(\frac{L_\alpha}{(1+\alpha)F(\mathbf{w})} \right)^{\frac{1}{1+\alpha}} \\ &= \frac{L_r}{2} \left(\frac{(1+\alpha)F(\mathbf{w})}{L_\alpha} \right)^{\frac{1}{1+\alpha}} + 2 \left(\frac{L_\alpha}{1+\alpha} \right)^{\frac{1}{1+\alpha}} F^{\frac{\alpha}{1+\alpha}}(\mathbf{w}).\end{aligned}$$

574 The proof is completed. \square

575 We now prove Lemma B.1 as a direct corollary of Lemma B.2.

576 *Proof of Lemma B.1.* Define $\tilde{F} : \mathcal{W} \mapsto \mathbb{R}$ as $\tilde{F}(\mathbf{w}) = F(\mathbf{w}) - F(\mathbf{w}^*)$. It is clear that $\tilde{F}(\mathbf{w}) \geq 0$ and
577 $\tilde{F}(\mathbf{w}) = \tilde{L}(\mathbf{w}) + r(\mathbf{w})$, where $\tilde{L}(\mathbf{w}) = L(\mathbf{w}) - F(\mathbf{w}^*)$ and $\tilde{L}(\mathbf{w})$ has (α, L_α) -Hölder continuous gradients.
578 Therefore, we can apply Lemma B.2 to \tilde{F} and derive

$$\|\nabla \tilde{F}(\mathbf{w})\| \leq \frac{L_r(1+\alpha)^{\frac{1}{1+\alpha}}}{2L_\alpha^{\frac{1}{1+\alpha}}} \tilde{F}^{\frac{1}{1+\alpha}}(\mathbf{w}) + 2 \left(\frac{L_\alpha}{1+\alpha} \right)^{\frac{1}{1+\alpha}} \tilde{F}^{\frac{\alpha}{1+\alpha}}(\mathbf{w}).$$

579 The stated bound then follows directly. The proof is completed. \square

580 The following lemma gives the self-bounding property for a nonnegative function with Hölder
581 continuous gradients [35, 43].

582 **Lemma B.3.** *Assume the map $\mathbf{w} \mapsto g(\mathbf{w})$ is nonnegative, and $\mathbf{w} \mapsto \nabla g(\mathbf{w})$ is (α, L_α) -Hölder contin-*
583 *uous with $\alpha \in [0, 1]$. Let c_α be defined in Eq. (3.9). Then*

$$\|\nabla g(\mathbf{w})\|_2 \leq c_\alpha g^{\frac{\alpha}{1+\alpha}}(\mathbf{w}), \quad \forall \mathbf{w} \in \mathcal{W}. \quad (\text{B.2})$$

584 C. Proof of Theorem 4

585 *Proof of Theorem 4.* Let $S = \{z_1, \dots, z_n\}$, $S' = \{z'_1, \dots, z'_n\}$, $S'' = \{z''_1, \dots, z''_n\}$ be drawn indepen-

586 dently from ρ . For any $i \in [n] := \{1, \dots, n\}$, define

$$g_i(S) = \mathbb{E}_{z'_i} \left[\mathbb{E}_Z [\ell(A(S^{(i)}); Z)] - \ell(A(S^{(i)}); z_i) \right],$$

587 where $S^{(i)}$ is defined in Eq. (5.8). Due to the symmetry between S and S' , we have

$$\begin{aligned}\mathbb{E}_{S \setminus z_i} [g_i(S)] &= \mathbb{E}_{S \setminus z_i} \mathbb{E}_{z'_i} \left[\mathbb{E}_Z [\ell(A(S^{(i)}); Z)] - \ell(A(S^{(i)}); z_i) \right] \\ &= \mathbb{E}_{S^{(i)}} [L(A(S^{(i)}))] - \mathbb{E}_{S^{(i)}} [\ell(A(S^{(i)}); z_i)] = \mathbb{E}_S [L(A(S))] - \mathbb{E}_{S^{(i)}} [\ell(A(S^{(i)}); z_i)].\end{aligned} \quad (\text{C.1})$$

588 According to the definition of pointwise uniform stability, we know

$$\begin{aligned}&\left| \mathbb{E}_Z [\ell(A(S); Z)] - \frac{1}{n} \sum_{i=1}^n \ell(A(S); z_i) - \frac{1}{n} \sum_{i=1}^n g_i(S) \right| \\ &\leq \frac{1}{n} \sum_{i=1}^n \left| \mathbb{E}_Z [\ell(A(S); Z)] - \mathbb{E}_{z'_i, Z} [\ell(A(S^{(i)}); Z)] \right| + \frac{1}{n} \sum_{i=1}^n \left| \ell(A(S); z_i) - \mathbb{E}_{z'_i} [\ell(A(S^{(i)}); z_i)] \right| \\ &\leq \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{z'_i, Z} \left[\left| \ell(A(S); Z) - \ell(A(S^{(i)}); Z) \right| \right] + \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{z'_i} \left[\left| \ell(A(S); z_i) - \ell(A(S^{(i)}); z_i) \right| \right] \leq \frac{2}{n} \sum_{i=1}^n \beta_i.\end{aligned}$$

589 It then follows that

$$\begin{aligned}
& \left| L(A(S)) - L_S(A(S)) - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{S \setminus z_i} [g_i(S)] \right| \\
&= \left| L(A(S)) - L_S(A(S)) - \frac{1}{n} \sum_{i=1}^n \left(\mathbb{E}_{S \setminus z_i} [g_i(S)] - g_i(S) + g_i(S) \right) \right| \\
&\leq \left| L(A(S)) - L_S(A(S)) - \frac{1}{n} \sum_{i=1}^n g_i(S) \right| + \left| \frac{1}{n} \sum_{i=1}^n \left(\mathbb{E}_{S \setminus z_i} [g_i(S)] - g_i(S) \right) \right| \\
&\leq \frac{2}{n} \sum_{i=1}^n \beta_i + \frac{1}{n} \left| \sum_{i=1}^n h_i(S) \right|, \tag{C.2}
\end{aligned}$$

590 where we introduce

$$h_i(S) = g_i(S) - \mathbb{E}_{S \setminus z_i} [g_i(S)], \quad \forall i \in [n].$$

591 We now show that the above h_i satisfies the conditions in Theorem 1. According to the definition of
592 h_i , we know that

$$\mathbb{E}_{S \setminus z_i} [h_i(S)] = \mathbb{E}_{S \setminus z_i} [g_i(S)] - \mathbb{E}_{S \setminus z_i} [g_i(S)] = 0. \tag{C.3}$$

593 It is clear that

$$\mathbb{E}_{z_i} [g_i(S)] = \mathbb{E}_{z_i} \left[\mathbb{E}_{z_i'} \left[\mathbb{E}_Z [\ell(A(S^{(i)}); Z)] - \ell(A(S^{(i)}); z_i) \right] \right] = \mathbb{E}_{z_i'} \left[\mathbb{E}_Z [\ell(A(S^{(i)}); Z)] - \mathbb{E}_{z_i} [\ell(A(S^{(i)}); z_i)] \right] = 0.$$

594 It then follows that

$$\mathbb{E}_{z_i} [h_i(S)] = \mathbb{E}_{z_i} [g_i(S)] - \mathbb{E}_{S \setminus z_i} \mathbb{E}_{z_i} [g_i(S)] = 0. \tag{C.4}$$

595 Finally, for any $j \in [n]$ with $j \neq i$, we have

$$|h_i(S) - h_i(z_1, \dots, z_{j-1}, z_j'', z_{j+1}, \dots, z_n)| \tag{C.5}$$

$$\begin{aligned}
&= \left| \left(g_i(S) - \mathbb{E}_{S \setminus z_i} [g_i(S)] \right) - \left(g_i(S_j'') - \mathbb{E}_{S_j'' \setminus z_i} [g_i(S_j'')] \right) \right| \\
&\leq |g_i(S) - g_i(S_j'')| + \left| \mathbb{E}_{S \setminus z_i} [g_i(S)] - \mathbb{E}_{S_j'' \setminus z_i} [g_i(S_j'')] \right| \\
&\leq |g_i(S) - g_i(S_j'')| + \mathbb{E}_{S \setminus z_i} \mathbb{E}_{S_j'' \setminus z_i} |g_i(S) - g_i(S_j'')|, \tag{C.6}
\end{aligned}$$

596 where

$$S_j'' = \{z_1, \dots, z_{j-1}, z_j'', z_{j+1}, \dots, z_n\}, \quad \forall j \in [n].$$

597 Note that

$$\begin{aligned}
& |g_i(S) - g_i(S_j'')| \\
&= \left| \left(\mathbb{E}_{z_i'} \left[\mathbb{E}_Z [\ell(A(S^{(i)}); Z)] - \ell(A(S^{(i)}); z_i) \right] \right) - \mathbb{E}_{z_i'} \left[\mathbb{E}_Z [\ell(A(S_j^{(i)}); Z)] - \ell(A(S_j^{(i)}); z_i) \right] \right| \\
&\leq \left| \mathbb{E}_{z_i'} \mathbb{E}_Z [\ell(A(S^{(i)}); Z)] - \mathbb{E}_{z_i'} \mathbb{E}_Z [\ell(A(S_j^{(i)}); Z)] \right| + \left| \mathbb{E}_{z_i'} [\ell(A(S^{(i)}); z_i)] - \mathbb{E}_{z_i'} [\ell(A(S_j^{(i)}); z_i)] \right| \leq 2\beta_j,
\end{aligned}$$

598 where $S_j^{(i)}$ is defined in Eq. (5.10). We combine the above inequality and Eq. (C.6) together to get

$$|h_i(S) - h_i(z_1, \dots, z_{j-1}, z_j'', z_{j+1}, \dots, z_n)| \leq 4\beta_j.$$

599 According to Eq. (C.3), (C.4) and the above inequality, the conditions of Theorem 1 hold with $M = 0$.
600 Therefore, we can apply Theorem 1 to derive the following inequality

$$\left| \sum_{i=1}^n h_i(S) \right| \lesssim p \log_2 n \left(n \sum_{i=1}^n \beta_i^2 \right)^{\frac{1}{2}}.$$

601 It then follows the following inequality with probability at least $1 - \delta/2$

$$\left| \sum_{i=1}^n h_i(S) \right| \lesssim \left(n \sum_{i=1}^n \beta_i^2 \right)^{\frac{1}{2}} \log n \log(1/\delta).$$

602 The above inequality together with Eq. (C.1), (C.2) gives the following inequality with probability at
603 least $1 - \delta/2$

$$\begin{aligned} & \left| L(A(S)) - L_S(A(S)) - \mathbb{E}_S[L(A(S))] + \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{S^{(i)}}[\ell(A(S^{(i)}); z_i)] \right| \\ &= \left| L(A(S)) - L_S(A(S)) - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{S \setminus z_i}[g_i(S)] \right| \lesssim \frac{1}{n} \sum_{i=1}^n \beta_i + \left(\frac{1}{n} \sum_{i=1}^n \beta_i^2 \right)^{\frac{1}{2}} \log n \log(1/\delta). \end{aligned} \quad (\text{C.7})$$

604 We have the following identity

$$\begin{aligned} L(A(S)) - L_S(A(S)) - L(\mathbf{w}^*) + L_S(\mathbf{w}^*) &= \left(L(A(S)) - L_S(A(S)) - \mathbb{E}_S[L(A(S))] + \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{S'}[\ell(A(S'); z_i)] \right) \\ &+ \left(\mathbb{E}_S[L(A(S))] - L(\mathbf{w}^*) - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{S'}[\ell(A(S'); z_i)] + L_S(\mathbf{w}^*) \right). \end{aligned} \quad (\text{C.8})$$

605 The first term can be controlled by Eq. (C.7) and the identity $\mathbb{E}_{S^{(i)}}[\ell(A(S^{(i)}); z_i)] = \mathbb{E}_{S'}[\ell(A(S'); z_i)]$.

606 We now control the second term by Bernstein's inequality. To this aim, we introduce $\xi(z) = \mathbb{E}_{S'}[\ell(A(S'); z)] -$
607 $\ell(\mathbf{w}^*; z)$. Due to the symmetry between S and S' , we further get

$$\begin{aligned} & \mathbb{E}_S[L(A(S))] - L(\mathbf{w}^*) - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{S'}[\ell(A(S'); z_i)] + L_S(\mathbf{w}^*) \\ &= \mathbb{E}_{S'}[L(A(S'))] - L(\mathbf{w}^*) - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{S'}[\ell(A(S'); z_i)] + L_S(\mathbf{w}^*) = \mathbb{E}_Z[\xi(Z)] - \frac{1}{n} \sum_{i=1}^n \xi(z_i). \end{aligned}$$

608 We can control the variance of ξ as follows

$$\begin{aligned} \text{Var}(\xi(Z)) &= \mathbb{E}_Z \left[\left(\mathbb{E}_{S'}[\ell(A(S'); Z)] - \ell(\mathbf{w}^*; Z) \right)^2 \right] - \left(\mathbb{E}_Z \left[\mathbb{E}_{S'}[\ell(A(S'); Z)] - \ell(\mathbf{w}^*; Z) \right] \right)^2 \\ &= \mathbb{E}_Z \left[\left(\mathbb{E}_S[\ell(A(S); Z)] - \ell(\mathbf{w}^*; Z) \right)^2 \right] - \left(\mathbb{E}_S[L(A(S))] - L(\mathbf{w}^*) \right)^2, \end{aligned}$$

609 where we have used the symmetry between S and S' . According to Bernstein's inequality (Lemma
610 A.4), the following inequality holds with probability at least $1 - \delta/2$

$$\mathbb{E}_Z[\xi(Z)] - \frac{1}{n} \sum_{i=1}^n \xi(z_i) \leq \frac{2M \log \frac{2}{\delta}}{3n} + \left(\frac{2\sigma_A^2 \log(2/\delta)}{n} \right)^{\frac{1}{2}}.$$

611 We can plug the above inequality and Eq. (C.7) into Eq. (C.8), and derive the following inequality
612 with probability at least $1 - \delta$

$$L(A(S)) - L_S(A(S)) - L(\mathbf{w}^*) + L_S(\mathbf{w}^*) \lesssim \left(\frac{1}{n} \sum_{i=1}^n \beta_i^2 \right)^{\frac{1}{2}} \log n \log(1/\delta) + \frac{M \log \frac{1}{\delta}}{n} + \left(\frac{\sigma_A^2 \log(1/\delta)}{n} \right)^{\frac{1}{2}}.$$

613 The proof is completed by noting the structure of f . □

614 **D. Proof of Lemma 5**

615 *Proof of Lemma 5.* We first prove Eq. (3.10). Since F is λ -strongly convex, we know

$$F(A(S)) - F(\mathbf{w}^*) \geq \lambda \|A(S) - \mathbf{w}^*\|^2/2. \quad (\text{D.1})$$

616 According to the definition of σ_A^2 , we know $\sigma_A^2 \leq \mathbb{E}_Z[(\mathbb{E}_S[\ell(A(S); Z)] - \ell(\mathbf{w}^*; Z))^2]$. Since ℓ is convex,
617 we know

$$|\ell(A(S); Z) - \ell(\mathbf{w}^*; Z)| \leq \|A(S) - \mathbf{w}^*\| \max\{\|\nabla\ell(A(S); Z)\|, \|\nabla\ell(\mathbf{w}^*; Z)\|\}.$$

618 It then follows from Eq. (B.2) that

$$\begin{aligned} \left(\mathbb{E}_S[\ell(A(S); Z)] - \ell(\mathbf{w}^*; Z)\right)^2 &\leq \mathbb{E}_S[\|A(S) - \mathbf{w}^*\|^2] \mathbb{E}_S[\max\{\|\nabla\ell(A(S); Z)\|^2, \|\nabla\ell(\mathbf{w}^*; Z)\|^2\}] \\ &\leq c_\alpha^2 \mathbb{E}_S[\|A(S) - \mathbf{w}^*\|^2] \mathbb{E}_S[\max\{\ell^{\frac{2\alpha}{1+\alpha}}(A(S); Z), \ell^{\frac{2\alpha}{1+\alpha}}(\mathbf{w}^*; Z)\}]. \end{aligned}$$

619 Therefore, we have

$$\sigma_A^2 \leq c_\alpha^2 \mathbb{E}_S[\|A(S) - \mathbf{w}^*\|^2] \mathbb{E}_{S,Z}[\max\{\ell^{\frac{2\alpha}{1+\alpha}}(A(S); Z), \ell^{\frac{2\alpha}{1+\alpha}}(\mathbf{w}^*; Z)\}].$$

620 Eq. (3.10) then follows by combining the above inequality and Eq. (D.1) together.

621 We now turn to Eq. (3.12) on generalization bounds in expectation. We first study the generaliza-
622 tion error for the algorithm A_E . By the definition $F_S, S, S^{(i)}$, we get

$$f(A_E(S^{(i)}); \mathbf{z}_i) = nF_S(A_E(S^{(i)})) - nF_{S^{(i)}}(A_E(S^{(i)})) + f(A_E(S^{(i)}); \mathbf{z}'_i).$$

623 By symmetry on \mathbf{z}_i and \mathbf{z}'_i , we get $\mathbb{E}[f(A_E(S^{(i)}); \mathbf{z}'_i)] = \mathbb{E}[f(A_E(S); \mathbf{z}_i)]$, $\mathbb{E}[F_{S^{(i)}}(A_E(S^{(i)}))] = \mathbb{E}[F_S(A_E(S))]$ and

$$\mathbb{E}[f(A_E(S^{(i)}); \mathbf{z}_i) - f(A_E(S); \mathbf{z}_i)] = n\mathbb{E}[F_S(A_E(S^{(i)})) - F_{S^{(i)}}(A_E(S^{(i)}))] = n\mathbb{E}[F_S(A_E(S^{(i)})) - F_S(A_E(S))].$$

624 Since F_S is λ -strongly convex, we further know $F_S(A_E(S^{(i)})) - F_S(A_E(S)) \leq \|\nabla F_S(A_E(S^{(i)}))\|^2 / (2\lambda)$ [26].

625 We can combine the above two inequalities to get

$$\mathbb{E}[f(A_E(S^{(i)}); \mathbf{z}_i) - f(A_E(S); \mathbf{z}_i)] \leq \frac{n}{2\lambda} \mathbb{E}[\|\nabla F_S(A_E(S^{(i)}))\|^2]. \quad (\text{D.2})$$

626 The definition of $A_E(S^{(i)})$ implies $\nabla F_{S^{(i)}}(A_E(S^{(i)})) = 0$, and

$$\begin{aligned} \mathbb{E}[\|\nabla F_S(A_E(S^{(i)}))\|^2] &= \mathbb{E}\left[\left\|\nabla F_{S^{(i)}}(A_E(S^{(i)})) - \frac{1}{n}\nabla f(A_E(S^{(i)}); \mathbf{z}_i) + \frac{1}{n}\nabla f(A_E(S^{(i)}); \mathbf{z}'_i)\right\|^2\right] \\ &= \frac{1}{n^2} \mathbb{E}\left[\left\|\nabla f(A_E(S^{(i)}); \mathbf{z}'_i) - \nabla f(A_E(S^{(i)}); \mathbf{z}_i)\right\|^2\right] = \frac{1}{n^2} \mathbb{E}\left[\left\|\nabla f(A_E(S); \mathbf{z}_i) - \nabla f(A_E(S); \mathbf{z}'_i)\right\|^2\right], \end{aligned}$$

627 where the last step is due to the symmetry between \mathbf{z}_i and \mathbf{z}'_i . We can combine the above inequality

628 and Eq (D.2) together to derive

$$\begin{aligned} \mathbb{E}[F(A_E(S)) - F_S(A_E(S))] &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[f(A_E(S^{(i)}); \mathbf{z}_i) - f(A_E(S); \mathbf{z}_i)] \\ &\leq \frac{1}{2n^2\lambda} \sum_{i=1}^n \mathbb{E}\left[\left\|\nabla f(A_E(S); \mathbf{z}_i) - \nabla f(A_E(S); \mathbf{z}'_i)\right\|^2\right], \quad (\text{D.3}) \end{aligned}$$

629 where we have used $\mathbb{E}[F(A_E(S))] = \mathbb{E}[F(A_E(S^{(i)}))] = \mathbb{E}[f(A_E(S^{(i)}); \mathbf{z}_i)]$. By the structure of f , we get

$$\begin{aligned} \left\| \nabla f(A_E(S); \mathbf{z}_i) - \nabla f(A_E(S); \mathbf{z}'_i) \right\|^2 &= \left\| \nabla \ell(A_E(S); \mathbf{z}_i) - \nabla \ell(A_E(S); \mathbf{z}'_i) \right\|^2 \\ &\leq 2\|\nabla \ell(A_E(S); \mathbf{z}_i)\|^2 + 2\|\nabla \ell(A_E(S); \mathbf{z}'_i)\|^2 \leq 2c_\alpha^2 \ell^{\frac{2\alpha}{1+\alpha}}(A_E(S); \mathbf{z}_i) + 2c_\alpha^2 \ell^{\frac{2\alpha}{1+\alpha}}(A_E(S); \mathbf{z}'_i). \end{aligned}$$

630 It then follows that

$$\frac{1}{n} \sum_{i=1}^n \left\| \nabla f(A_E(S); \mathbf{z}_i) - \nabla f(A_E(S); \mathbf{z}'_i) \right\|^2 \leq \frac{2c_\alpha^2}{n} \sum_{i=1}^n \ell^{\frac{2\alpha}{1+\alpha}}(A_E(S); \mathbf{z}_i) + \frac{2c_\alpha^2}{n} \sum_{i=1}^n \ell^{\frac{2\alpha}{1+\alpha}}(A_E(S); \mathbf{z}'_i).$$

631 We can use the above inequality, the concavity of the function $x \mapsto x^{\frac{2\alpha}{1+\alpha}}$ and Eq. (D.3) to derive

$$\mathbb{E}[F(A_E(S)) - F_S(A_E(S))] \leq \frac{c_\alpha^2}{n\lambda} \mathbb{E}\left[L_S^{\frac{2\alpha}{1+\alpha}}(A_E(S)) + L_{S'}^{\frac{2\alpha}{1+\alpha}}(A_E(S))\right] \leq \frac{c_\alpha^2}{n\lambda} \mathbb{E}\left[L_S^{\frac{2\alpha}{1+\alpha}}(A_E(S)) + L^{\frac{2\alpha}{1+\alpha}}(A_E(S))\right]. \quad (\text{D.4})$$

632 By Eq. (5.12) and the Cauchy-Schwartz's inequality, we know

$$\mathbb{E}[\langle A(S) - A_E(S), \nabla F(A_E(S)) \rangle] \leq C_1 \left(\mathbb{E}[\|A(S) - A_E(S)\|^2] \right)^{\frac{1}{2}} \left(\mathbb{E}\left[\left(F(A_E(S)) - F(\mathbf{w}^*)\right)^{\frac{2\alpha}{1+\alpha}}\right] \right)^{\frac{1}{2}}.$$

633 Since $\mathbb{E}[F_S(A_E(S))] \leq F(\mathbf{w}^*)$, we can plug the above inequality back into Eq. (5.13), and derive

$$\begin{aligned} \mathbb{E}[F(A(S)) - F(A_E(S))] &\leq C_1 \left(\mathbb{E}[\|A(S) - A_E(S)\|^2] \right)^{\frac{1}{2}} \left(\mathbb{E}\left[F(A_E(S)) - F_S(A_E(S))\right] \right)^{\frac{2\alpha}{1+\alpha}} \frac{1}{2} \\ &\quad + \frac{L_\alpha \mathbb{E}[\|A(S) - A_E(S)\|^{1+\alpha}]}{1+\alpha} + \frac{L_r \mathbb{E}[\|A(S) - A_E(S)\|^2]}{2}, \end{aligned}$$

634 where we have used the concavity of $x \mapsto x^{\frac{2\alpha}{1+\alpha}}$ and the Jensen's inequality. We can plug Eq. (D.4)

635 and Eq. (5.14) into the above inequality to show

$$\begin{aligned} \mathbb{E}[F(A(S)) - F(A_E(S))] &\leq \frac{L_\alpha \mathbb{E}[(2\lambda^{-1}(F_S(A(S)) - F_S(A_E(S))))^{\frac{1+\alpha}{2}}]}{1+\alpha} + \frac{L_r \mathbb{E}[F_S(A(S)) - F_S(A_E(S))]}{\lambda} \\ &\quad + C_1 \left(2\lambda^{-1} \mathbb{E}[F_S(A(S)) - F_S(A_E(S))] \right)^{\frac{1}{2}} \left(\frac{c_\alpha^2}{n\lambda} \mathbb{E}\left[L_S^{\frac{2\alpha}{1+\alpha}}(A_E(S)) + L^{\frac{2\alpha}{1+\alpha}}(A_E(S))\right] \right)^{\frac{\alpha}{1+\alpha}}. \end{aligned}$$

636 According to Eq. (D.4), the concavity of the function $x \mapsto x^{\frac{1+\alpha}{2}}$ and the decomposition

$$\mathbb{E}[F(A(S)) - F_S(A_E(S))] = \mathbb{E}[F(A(S)) - F(A_E(S))] + \mathbb{E}[F(A_E(S)) - F_S(A_E(S))],$$

637 we further get

$$\begin{aligned} \mathbb{E}[F(A(S)) - F_S(A_E(S))] &\leq \frac{2^{\frac{1+\alpha}{2}} L_\alpha}{1+\alpha} \Delta_\lambda^{\frac{1+\alpha}{2}} + L_r \Delta_\lambda + \sqrt{2} C_1 c_\alpha^{\frac{2\alpha}{1+\alpha}} \Delta_\lambda^{\frac{1}{2}} \nabla_\lambda^{\frac{\alpha}{1+\alpha}} + c_\alpha^2 \nabla_\lambda \\ &\leq \left(\frac{2^{\frac{1+\alpha}{2}} L_\alpha}{1+\alpha} + \frac{\sqrt{2} C_1 c_\alpha^{\frac{2\alpha}{1+\alpha}}}{1+\alpha} \right) \Delta_\lambda^{\frac{1+\alpha}{2}} + L_r \Delta_\lambda + \left(\frac{\sqrt{2} C_1 c_\alpha^{\frac{2\alpha}{1+\alpha}} \alpha}{1+\alpha} + c_\alpha^2 \right) \nabla_\lambda, \end{aligned}$$

638 where we have used the Young's inequality $\Delta_\lambda^{\frac{1}{2}} \nabla_\lambda^{\frac{\alpha}{1+\alpha}} \leq \frac{\alpha}{1+\alpha} \nabla_\lambda^{\frac{\alpha}{1+\alpha}} \Delta_\lambda^{\frac{1+\alpha}{2}} + \frac{1}{1+\alpha} \Delta_\lambda^{\frac{\alpha+1}{2}}$. The proof is com-

639 pleted by noting $\mathbb{E}[F_S(A_E(S))] \leq \mathbb{E}[F_S(\mathbf{w}^*)] = F(\mathbf{w}^*)$ and

$$\mathfrak{C} = \max \left\{ \frac{2^{\frac{1+\alpha}{2}} L_\alpha}{1+\alpha} + \frac{\sqrt{2} C_1 c_\alpha^{\frac{2\alpha}{1+\alpha}}}{1+\alpha}, L_r, \frac{\sqrt{2} C_1 c_\alpha^{\frac{2\alpha}{1+\alpha}} \alpha}{1+\alpha} + c_\alpha^2 \right\}. \quad (\text{D.5})$$

640 The proof is completed. \square

641 References

- 642 [1] I. Amir, Y. Carmon, T. Koren, and R. Livni. Never go full batch (in stochastic convex optimization).
643 *Advances in Neural Information Processing Systems*, 34:25033–25043, 2021.
- 644 [2] P. Bartlett and S. Mendelson. Rademacher and Gaussian complexities: Risk bounds and structural results.
645 *Journal of Machine Learning Research*, 3:463–482, 2002.
- 646 [3] R. Bassily, V. Feldman, C. Guzmán, and K. Talwar. Stability of stochastic gradient descent on nonsmooth
647 convex losses. *Advances in Neural Information Processing Systems*, 33, 2020.
- 648 [4] S. Boucheron, G. Lugosi, and P. Massart. *Concentration Inequalities: A Nonasymptotic Theory of Inde-*
649 *pendence*. Oxford university press, 2013.
- 650 [5] O. Bousquet and A. Elisseeff. Stability and generalization. *Journal of Machine Learning Research*, 2
651 (Mar):499–526, 2002.
- 652 [6] O. Bousquet, Y. Klochkov, and N. Zhivotovskiy. Sharper bounds for uniformly stable algorithms. In
653 *Conference on Learning Theory*, pages 610–626, 2020.
- 654 [7] Z. Charles and D. Papailiopoulos. Stability and generalization of learning algorithms that converge to
655 global optima. In *International Conference on Machine Learning*, pages 744–753, 2018.
- 656 [8] X. Chen, B. Tang, J. Fan, and X. Guo. Online gradient descent algorithms for functional data learning.
657 *Journal of Complexity*, 70:101635, 2022.
- 658 [9] F. Cucker and D.-X. Zhou. *Learning Theory: an Approximation Theory Viewpoint*. Cambridge University
659 Press, 2007.
- 660 [10] V. De la Pena and E. Giné. *Decoupling: from dependence to independence*. Springer Science & Business
661 Media, 2012.
- 662 [11] L. Devroye and T. Wagner. Distribution-free inequalities for the deleted and holdout error estimates.
663 *IEEE Transactions on Information Theory*, 25(2):202–207, 1979.
- 664 [12] A. Elisseeff, T. Evgeniou, and M. Pontil. Stability of randomized learning algorithms. *Journal of Machine*
665 *Learning Research*, 6(Jan):55–79, 2005.
- 666 [13] V. Feldman and J. Vondrak. Generalization bounds for uniformly stable algorithms. In *Advances in*
667 *Neural Information Processing Systems*, pages 9747–9757, 2018.
- 668 [14] V. Feldman and J. Vondrak. High probability generalization bounds for uniformly stable algorithms with
669 nearly optimal rate. In *Conference on Learning Theory*, pages 1270–1279, 2019.
- 670 [15] S. Ghadimi and G. Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming.
671 *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.
- 672 [16] X. Guo, Z.-C. Guo, and L. Shi. Capacity dependent analysis for functional online learning algorithms.
673 *Applied and Computational Harmonic Analysis*, 67:101567, 2023.
- 674 [17] M. Hardt, B. Recht, and Y. Singer. Train faster, generalize better: Stability of stochastic gradient descent.
675 In *International Conference on Machine Learning*, pages 1225–1234, 2016.
- 676 [18] Y. Klochkov and N. Zhivotovskiy. Stability and deviation optimal risk bounds with convergence rate
677 $o(1/n)$. *Advances in Neural Information Processing Systems*, 34, 2021.
- 678 [19] T. Koren, R. Livni, Y. Mansour, and U. Sherman. Benign underfitting of stochastic gradient descent. In
679 *Advances in Neural Information Processing Systems*, 2022.

- 680 [20] I. Kuzborskij and C. Lampert. Data-dependent stability of stochastic gradient descent. In *International*
681 *Conference on Machine Learning*, pages 2820–2829, 2018.
- 682 [21] Y. Lei. Stability and generalization of stochastic optimization with nonconvex and nonsmooth problems.
683 In *Annual Conference on Learning Theory*, pages 191–227, 2023.
- 684 [22] Y. Lei and Y. Ying. Fine-grained analysis of stability and generalization for stochastic gradient descent.
685 In *International Conference on Machine Learning*, pages 5809–5819, 2020.
- 686 [23] Y. Lei and Y. Ying. Sharper generalization bounds for learning with gradient-dominated objective func-
687 tions. In *International Conference on Learning Representations*, 2021.
- 688 [24] Y. Lei, R. Jin, and Y. Ying. Stability and generalization analysis of gradient methods for shallow neural
689 networks. *Advances in Neural Information Processing Systems*, 35:38557–38570, 2022.
- 690 [25] W. Mou, L. Wang, X. Zhai, and K. Zheng. Generalization bounds of SGLD for non-convex learning: Two
691 theoretical viewpoints. In *Conference on Learning Theory*, pages 605–638, 2018.
- 692 [26] Y. Nesterov. *Introductory Lectures on Convex Optimization: A Basic Course*, volume 87. Springer Science
693 & Business Media, 2013.
- 694 [27] G. Neu, G. K. Dziugaite, M. Haghifam, and D. M. Roy. Information-theoretic generalization bounds for
695 stochastic gradient descent. In *Conference on Learning Theory*, pages 3526–3545. PMLR, 2021.
- 696 [28] K. Nikolakakis, F. Haddadpour, D. Kalogerias, and A. Karbasi. Black-box generalization: Stability of
697 zeroth-order learning. *Advances in Neural Information Processing Systems*, 35:31525–31541, 2022.
- 698 [29] M. Raab and A. Steger. “balls into bins”—a simple and tight analysis. In *International Workshop on*
699 *Randomization and Approximation Techniques in Computer Science*, pages 159–170. Springer, 1998.
- 700 [30] D. Richards and I. Kuzborskij. Stability & generalisation of gradient descent for shallow neural networks
701 without the neural tangent kernel. *Advances in Neural Information Processing Systems*, 34, 2021.
- 702 [31] O. Rivasplata, E. Parrado-Hernández, J. S. Shawe-Taylor, S. Sun, and C. Szepesvári. Pac-bayes bounds for
703 stable algorithms with instance-dependent priors. In *Advances in Neural Information Processing Systems*,
704 pages 9214–9224, 2018.
- 705 [32] W. H. Rogers and T. J. Wagner. A finite sample distribution-free performance bound for local discrimi-
706 nation rules. *The Annals of Statistics*, pages 506–514, 1978.
- 707 [33] M. Schliserman and T. Koren. Stability vs implicit bias of gradient methods on separable data and beyond.
708 In *Conference on Learning Theory*, pages 3380–3394, 2022.
- 709 [34] S. Shalev-Shwartz, O. Shamir, N. Srebro, and K. Sridharan. Learnability, stability and uniform conver-
710 gence. *Journal of Machine Learning Research*, 11(Oct):2635–2670, 2010.
- 711 [35] N. Srebro, K. Sridharan, and A. Tewari. Smoothness, low noise and fast rates. In *Advances in Neural*
712 *Information Processing Systems*, pages 2199–2207, 2010.
- 713 [36] I. Steinwart and A. Christmann. *Support Vector Machines*. Springer Science & Business Media, 2008.
- 714 [37] H. Taheri and C. Thrampoulidis. Generalization and stability of interpolating neural networks with
715 minimal width. *arXiv preprint arXiv:2302.09235*, 2023.
- 716 [38] E. Ullah, T. Mai, A. Rao, R. A. Rossi, and R. Arora. Machine unlearning via algorithmic stability. In
717 *Conference on Learning Theory*, pages 4126–4142. PMLR, 2021.
- 718 [39] V. Vapnik. *The nature of statistical learning theory*. Springer, 2013.
- 719 [40] R. Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47.

- 720 Cambridge university press, 2018.
- 721 [41] S. Villa, S. Matet, B. C. Vũ, and L. Rosasco. Implicit regularization with strongly convex bias: Stability
722 and acceleration. *Analysis and Applications*, pages 1–27, 2022.
- 723 [42] P. Wang, Y. Lei, Y. Ying, and H. Zhang. Differentially private sgd with non-smooth losses. *Applied and*
724 *Computational Harmonic Analysis*, 56:306–336, 2022.
- 725 [43] Y. Ying and D.-X. Zhou. Unregularized online learning algorithms with general loss functions. *Applied*
726 *and Computational Harmonic Analysis*, 42(2):224–244, 2017.
- 727 [44] Y. Zhang, W. Zhang, S. Bald, V. Pingali, C. Chen, and M. Goswami. Stability of sgd: Tightness analysis
728 and improved bounds. In *Uncertainty in artificial intelligence*, pages 2364–2373. PMLR, 2022.
- 729 [45] D.-X. Zhou. The covering number in learning theory. *Journal of Complexity*, 18(3):739–767, 2002.
- 730 [46] L. Zhu, M. Gurbuzbalaban, A. Raj, and U. Simsekli. Uniform-in-time wasserstein stability bounds for
731 (noisy) stochastic gradient descent. In *Advances in Neural Information Processing Systems*, 2023.