
On the Generalization Error of Stochastic Mirror Descent for Quadratically-Bounded Losses: an Improved Analysis

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In this work, we revisit the generalization error of stochastic mirror descent for
2 quadratically bounded losses studied in Telgarsky (2022). Quadratically bounded
3 losses is a broad class of loss functions, capturing both Lipschitz and smooth
4 functions, for both regression and classification problems. We study the high
5 probability generalization for this class of losses on linear predictors in both
6 realizable and non-realizable cases when the data are sampled IID or from a
7 Markov chain. The prior work relies on an intricate coupling argument between
8 the iterates of the original problem and those projected onto a bounded domain.
9 This approach enables blackbox application of concentration inequalities, but
10 also leads to suboptimal guarantees due in part to the use of a union bound
11 across all iterations. In this work, we depart significantly from the prior work of
12 Telgarsky (2022), and introduce a novel approach for establishing high probability
13 generalization guarantees. In contrast to the prior work, our work directly analyzes
14 the moment generating function of a novel supermartingale sequence and leverages
15 the structure of stochastic mirror descent. As a result, we obtain improved bounds
16 in all aforementioned settings. Specifically, in the realizable case and non-realizable
17 case with light-tailed sub-Gaussian data, we improve the bounds by a $\log T$ factor,
18 matching the correct rates of $1/T$ and $1/\sqrt{T}$, respectively. In the more challenging
19 case of heavy-tailed polynomial data, we improve the existing bound by a poly T
20 factor.

21 1 Introduction

22 Along with convergence analysis of optimization methods, understanding the generalization of models
23 trained by these methods on unseen data is an important question in machine learning. However,
24 despite the number of works attempting to answer it, the problem has not been fully understood, even
25 in the simplest setting of linear predictors constructed with the standard stochastic gradient/mirror
26 descent. A great part of prior works [28, 10, 25, 26, 27] focus only on the generalization on linearly
27 separable data and/or of models trained with specific losses with exponentially decaying tails such as
28 logistic loss. The question of what we can guarantee beyond these settings remains open.

29 Recently, [30] proposes a new approach to analyze the generalization error with *high probability* of
30 stochastic mirror descent for a broad class of quadratically bounded losses, beyond the realizable
31 setting. This class of losses encapsulates both Lipschitz and smooth functions, for both regression
32 and classification problems. The obtained bounds complement existing in-expectation bounds [7]
33 and nearly match the counterpart of convergence rates in optimization. While this result pushes
34 forward the state of the art, the obtained guarantees do not completely resolve the problem. The

35 central piece of the proposed approach is a “coupling” technique between the iterates of the original
36 problem and those projected onto a bounded domain. In this technique, one first constrains the
37 problem in a bounded domain with a well chosen diameter. The bounded domain diameter allows to
38 apply concentration inequalities as a blackbox and obtain bounds in high probability. Then using
39 an inductive argument and a union bound across all iterations, one can show that the iterates in
40 the original problem coincide with the ones in the constrained problem. Due to the union bound,
41 the success probability decreases from $1 - \delta$ to $1 - T\delta$, where T is the number of iterations in the
42 algorithm. This loss translates to a milder $\log T$ factor loss in the guarantee in the case of realizable
43 data, and a more significant poly T factor loss in the non-realizable setting when the data has
44 polynomial tails. Thus a natural question arises of whether we can obtain a stronger analysis that
45 closes these remaining gaps.

46 In this paper, we revisit these generalization bounds for quadratically bounded losses by [30]. We
47 introduce a novel approach to analyze the generalization errors of stochastic mirror descent in both
48 realizable and non-realizable cases when the data are sampled IID or from a Markov chain. In all
49 these cases, we remove the need to use the union bound argument, thus preventing the loss in the
50 success probability. This translates to the following improvements:

51 – In the realizable, and the non-realizable cases with sub-gaussian tailed data and Markovian data,
52 we improve the bounds by a $\log T$ factor. This improvement comes from analyzing the moment
53 generating function of a martingale difference sequence with well-chosen coefficients. In these cases,
54 we also remove the necessity of using the coupling-based argument used in the same work by [30].
55 Instead, by solely making use of the problem structure, we arrive at the same conclusion that with
56 high probability, the iterates of stochastic mirror descent for quadratically bounded losses behave as
57 if the problem domain is bounded.

58 – In the non-realizable case with polynomial tailed data, we improve the existing bound by a poly T
59 factor. Due to the polynomial dependency on $\frac{1}{\delta}$, being able to maintain the same success probability
60 through all iterations is crucial in this case. Unlike the previous work, we rely on a truncation
61 technique. Using a more refined analysis of the truncated random variables, in combination with
62 suitable concentration inequalities and the coupling technique, we improve the existing bounds
63 significantly.

64 1.1 Related Work

65 Broadly speaking, there is a rich body of works in optimization and generalization that provide
66 convergence guarantees and generalization bounds for stochastic methods. Earlier works often focus
67 on in-expectation bounds [3, 19, 21, 13, 7], and bounds in high probability [11, 23, 9, 8] for problems
68 with bounded domains or under various additional assumptions such as strong convexity, noise with
69 light tails. Recent developments for optimization [20, 5, 15, 18, 6, 12, 4, 14, 24, 17, 16] are able to
70 handle unconstrained problems and relax these assumptions, but also require changes to the algorithm
71 such as gradient clipping. In generalization error analysis, specifically, a number of prior works,
72 including [28, 10, 25, 26, 27], focus only on linearly separable data. Among these, [28, 10, 27] only
73 deal with exponentially tailed losses while [25, 26] show generalization bounds for general smooth
74 convex losses. Our work, similarly to [30], goes beyond the realizable setting and specific losses. We
75 show high probability generalization bounds in both realizable and non-realizable settings for the
76 broad class of quadratically bounded losses, for both regression and classification problems.

77 The main point of reference for this paper is the work by [30]. This work develops a “coupling”
78 technique to bound the generalization error of stochastic mirror descent for quadratically bounded
79 losses. This technique has been employed in prior works [5, 6, 4, 24, 22, 17] to obtain high probability
80 convergence bounds of stochastic methods in optimization. Our work improves their results by using
81 a different approach that takes a closer look at the mechanism of the concentration inequalities and
82 leverages the problem structure. When the data are bounded or have sub-gaussian tails, analyzing
83 the moment generating function of a novel martingale difference sequence allows us to maintain the
84 same success probability, without using either the coupling technique or the union bound. This new
85 analysis, however, does not change the observation by [30] that the iterates of the unconstrained and
86 the constrained problems coincide with high probability. When the data have a polynomial tail, we
87 rely on a truncation technique. In this case, the coupling technique is necessary but not the union
88 bound, and we are still able to significantly improve the success rate.

89 In terms of techniques, the work by [16] for optimization is the closest to ours. In this work, the
90 authors develop the whitebox approach to analyzing stochastic methods for optimization with light-
91 tailed noise. In this work, we study generalization errors. Moreover, in all settings, our choice of
92 martingale difference sequences and coefficients are a significant departure from the prior work. In
93 particular, in [16] the choice of coefficients only depends on the problem parameters whereas in the
94 realizable case, our coefficients depend also on the historical data. Our approach also allows for a
95 flexible use of an induction argument without decreasing the success probability, while in [16] the
96 bounds are simpler and can be easily achieved in a single step.

97 2 Preliminaries

98 In this section, we provide the general set up and necessary notations before analyzing stochastic
99 mirror descent in the subsequent sections. Overall, we closely follow notations used in [30].

100 **Domain and norms.** In this work, we consider \mathcal{X} —the domain of the problem—to be a closed
101 convex set or \mathbb{R}^d . We will use $\|\cdot\|$ to denote an arbitrary norm on \mathcal{X} and let $\|\cdot\|_*$ be its dual norm. We
102 define the Bregman divergence as $\mathbf{D}_\psi(w; v) = \psi(w) - \psi(v) - \langle \nabla \psi(v), w - v \rangle$ where $\psi : \mathbb{R}^d \rightarrow \mathbb{R}$
103 is a differentiable function that is 1-strongly convex with respect to the norm $\|\cdot\|$.

104 **Loss functions.** Each loss function $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}_{\geq 0}$ in our consideration can be written using
105 a convex scalar function $\tilde{\ell}$ in one of the two following forms: 1) $\ell(y, \hat{y}) = \tilde{\ell}(\text{sign}(y)\hat{y})$ where
106 $\text{sign}(y) = 1$ if $y \geq 0$ and $= -1$ otherwise; and 2) $\ell(y, \hat{y}) = \tilde{\ell}(y - \hat{y})$. The first form captures
107 classification losses and the second regression losses. We will assume that subgradients $\partial \ell$ of ℓ in
108 the second argument always exist, and let ℓ' denote a subgradient in $\partial \ell$. For a function f , we also
109 use $\|\partial f(w)\| := \sup \{\|g\| : g \in \partial f(w)\}$. We further make the following assumptions, introduced in
110 [30] as quadratic boundedness and self-boundedness.

111 **Assumption 1.** We assume that ℓ is (C_1, C_2) -quadratically-bounded, for some constants $C_1, C_2 \geq 0$,
112 i.e., for all y, \hat{y}

$$|\ell'(y, \hat{y})| \leq C_1 + C_2 (|y| + |\hat{y}|).$$

113 This condition captures both classes of Lipschitz and smooth functions. Indeed, Lemma 1.2 from
114 [30] shows that α -Lipschitz functions are $(\alpha, 0)$ -quadratically-bounded while β -smooth functions
115 are $(\|\partial \tilde{\ell}(0)\|, \beta)$ -quadratically-bounded.

116 **Assumption 2.** In the realizable setting, we assume that ℓ is ρ -self-bounding, i.e., $\tilde{\ell}$ satisfies
117 $\tilde{\ell}'(z)^2 \leq 2\rho \tilde{\ell}(z)$ for all $z \in \mathbb{R}$.

118 The second assumption is a generalization of smoothness. This assumption is satisfied by smooth
119 losses but also certain non-smooth losses such as the exponential loss. This condition is necessary in
120 the current analysis to prove $1/T$ rates in the realizable setting. The readers can refer to [29, 30] for
121 more detailed discussion on this assumption.

122 Assumptions 1 and 2 are satisfied by commonly used loss functions in machine learning. These
123 include the logistic loss $\ell(y, \hat{y}) = \ln(1 + \exp(-y\hat{y}))$ and the squared loss $\ell(y, \hat{y}) = \frac{1}{2}(y - \hat{y})^2$ (see
124 Lemma 1.4 in [30]).

125 For the loss function ℓ and the configuration w , and sample (x, y) where x denotes the attribute and
126 y the label, we will write $\ell_{x,y} = \ell(y, w^T x)$. We state the following crucial lemma which is the same
127 as Lemma A.1 in [30], whose proof will be omitted.

128 **Lemma 1** (Lemma A.1 in ([30])). *Suppose ℓ is (C_1, C_2) -quadratically-bounded and $B_x \geq 0$ is
129 given. Given (x, y) such that $\max \{\|x\|_*, |y|\} \leq B_x$ and any u, v ,*

$$\begin{aligned} \|\partial \ell_{x,y}(u)\|_* &\leq B_x (C_1 + C_2 B_x (1 + \|u\|)) \\ |\ell_{x,y}(u) - \ell_{x,y}(v)| &\leq B_x \|u - v\| (C_1 + C_2 B_x (1 + \|u\|)). \end{aligned}$$

130 **Risk, IID and Markovian data.** When sample (x_i, y_i) arrives in iteration i of an algorithm,
131 we will use the notation $\ell_i(w) = \ell(y_i, w^T x_i)$. For an algorithm of T iterations, we use $\mathcal{F}_t =$
132 $\sigma((x_1, y_1), \dots, (x_t, y_t))$ to denote the natural filtration up to and including time t . When the data
133 are IID and generated from a distribution π , we define the risk

$$\mathcal{R}(w) = \mathbb{E}_{(x,y) \sim \pi} [\ell(y, w^T x)];$$

Algorithm 1 Stochastic Mirror Descent

Input w_0 , step size η For t in $1 \dots T$ $g_t \in \partial \ell_t(w_{t-1})$ $w_t = \arg \min_{w \in \mathcal{X}} \{\langle \eta g_t, w \rangle + \mathbf{D}_\psi(w; w_{t-1})\}$

134 In contrast to IID data, Markovian data come from a stochastic process. This setting has also been
135 considered in [1]. We let P_s^t be the distribution of (x_t, y_t) at iteration t conditioned on \mathcal{F}_s . We make
136 the following assumption regarding the uniform mixing time of the stochastic process. Note that
137 similar assumptions have also appeared in [30, 1].

138 **Assumption 3.** We assume that for some $\epsilon, \tau \geq 0$ of our choice, there is a distribution π such that

$$\sup_{t \in \mathbb{Z}_{\geq 0}} \sup_{\mathcal{F}_t} \text{TV}(P_t^{t+\tau}, \pi) \leq \epsilon.$$

139 We refer to the triple (π, τ, ϵ) as an approximate stationarity witness. We then define the risk according
140 to the approximate stationary distribution π : $\mathcal{R}(w) = \mathbb{E}_{(x,y) \sim \pi} [\ell(y, w^T x)]$.

141 **Algorithm.** Stochastic Mirror Descent is given in Algorithm 1. In this algorithm, for the simplicity of
142 the analysis, we consider a fixed step size η . In each iteration, we pick a subgradient $g_t \in \partial \ell_t(w_{t-1})$
143 and perform the update step.

144 We finally introduce a standard lemma used in the analysis of Stochastic Mirror Descent.

145 **Lemma 2.** For $t \geq 0$ and $w_{\text{ref}} \in \mathcal{X}$, we have

$$\mathbf{D}_\psi(w_{\text{ref}}; w_{t+1}) - \mathbf{D}_\psi(w_{\text{ref}}; w_t) \leq \eta(\ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(w_t)) + \frac{\eta^2}{2} \|g_{t+1}\|_*^2.$$

146 **Other notations.** We will use w_{ref} to refer to a comparator of interest. For the simplicity of the
147 exposition, we let $D_0 = \mathbf{D}_\psi(w_{\text{ref}}; w_0)$, and $\mathcal{R}^* = \inf_{v \in \mathcal{X}} \mathcal{R}(v)$. For a loss function ℓ that is
148 (C_1, C_2) -quadratically-bounded, we let $C_4 = C_1 + C_2(1 + \|w_{\text{ref}}\|)$.

149 3 Generalization bounds of SMD for IID data

150 In this section, we distinguish between two cases: the realizable case and the non-realizable case. In
151 the realizable case, there exists an optimal solution $w^* \in \mathcal{X}$ such that $\mathcal{R}(w^*) = 0$. We will show that
152 under mild assumptions, the risks of the solutions output by Algorithm 1 are bounded by $O(1/T)$. In
153 the non-realizable case, we will show, on the other hand, a weaker statement that the excess risks of
154 the solutions are bounded by $O(1/\sqrt{T})$.

155 3.1 Realizable case

156 In the realizable case, the comparator w_{ref} is not necessarily the global minimizer. To show the $1/T$
157 rate, we will assume w_{ref} satisfies $\mathcal{R}(w_{\text{ref}}) \leq \rho \mathbf{D}_\psi(w_{\text{ref}}; w_0)/T$ and that the loss at w_{ref} is bounded.
158 The guarantee for the iterates of Algorithm 1 is provided in Theorem 3.

159 **Theorem 3.** Suppose ℓ is convex, (C_1, C_2) -quadratically-bounded, and ρ -self-bounding. Given T ,
160 $((x_t, y_t))_{t \leq T}$ are IID samples with $\max\{\|x_t\|_*, |y_t|\} \leq 1$ almost surely, w_{ref} satisfies $\mathcal{R}(w_{\text{ref}}) \leq$
161 $\rho \mathbf{D}_\psi(w_{\text{ref}}; w_0)/T$, and $\max_{t < T} \ell_{t+1}(w_{\text{ref}}) \leq C_3$ almost surely. Then for $\eta \leq \frac{1}{2\rho}$, with probability
162 at least $1 - 2\delta$, for every $0 \leq k \leq T - 1$

$$\frac{1}{k+1} \sum_{t=0}^k \mathcal{R}(w_t) + \frac{16 \mathbf{D}_\psi(w_{\text{ref}}; w_{k+1})}{5(k+1)\eta} \leq \frac{C}{k+1} + 3\mathcal{R}(w_{\text{ref}}).$$

163 where $C = \frac{16C_4}{5} \log \frac{1}{\delta} \sqrt{\frac{15}{4} D_0 + 4\eta\gamma C_3} + \left(\frac{6}{\eta} D_0 + \frac{32}{5} \gamma C_3\right)$ with $\gamma = \max\{1, \log \frac{1}{\delta}\}$.

164 The analysis of Theorem 3 relies on the use of concentration inequalities. In contrast to existing
165 works that utilize concentration inequalities as a blackbox, we will make use of the mechanism for

166 proving concentration inequalities in order to obtain stronger guarantees. The type of concentration
 167 inequalities we consider are shown by analyzing the moment generating function of suitably chosen
 168 martingale sequences. We will use Lemma 14 (Appendix) which gives a basic inequality that bounds
 169 the moment generating function of a bounded random variable. To start the analysis, we use Lemma
 170 2 and Assumption 2 to obtain

171 **Lemma 4.** *For all $t \geq 0$, we have*

$$\mathbf{D}_\psi(w_{\text{ref}}; w_{t+1}) - \mathbf{D}_\psi(w_{\text{ref}}; w_t) \leq \eta \ell_{t+1}(w_{\text{ref}}) - \frac{\eta}{2} \ell_{t+1}(w_t),$$

$$\text{and hence, } \mathbf{D}_\psi(w_{\text{ref}}; w_t) \leq \mathbf{D}_\psi(w_{\text{ref}}; w_0) + \eta \sum_{i=1}^t \ell_i(w_{\text{ref}}) = D_0 + \eta \sum_{i=1}^t \ell_i(w_{\text{ref}}).$$

172 First, let us pay attention to the term $\sum_{i=1}^t \ell_i(w_{\text{ref}})$. Recall that the terms $\ell_i(w_{\text{ref}})$ are non-negative
 173 and bounded by a constant C_3 almost surely. We can analyze the term $\sum_{i=1}^T \ell_i(w_{\text{ref}})$ which upper
 174 bounds all sums $\sum_{i=1}^t \ell_i(w_{\text{ref}})$ by studying its moment generating function (or via a concentration
 175 inequality). We state this bound in the next lemma and defer the proof to the appendix.

176 **Lemma 5.** *With probability at least $1 - \delta$, $\sum_{i=1}^T \ell_i(w_{\text{ref}}) \leq \frac{7}{4} T \mathcal{R}(w_{\text{ref}}) + C_3 \log \frac{1}{\delta}$.*

177 Lemma 4 and lemma 5 and the assumption that $\mathcal{R}(w_{\text{ref}}) = O(1/T)$ imply that with probability at
 178 least $1 - \delta$, $\mathbf{D}_\psi(w_{\text{ref}}; w_t)$ is bounded. In other words, with probability at least $1 - \delta$, the iterates w_t
 179 all lie in a bounded region. One could therefore proceed to assume that the problem domain is simply
 180 this bounded ball around w_{ref} . This is the basic idea behind the ‘‘coupling’’ technique demonstrated
 181 in [30]. However, the important question is how to obtain a bound for the risk of all iterates even
 182 when we are working with a problem with unbounded domain. Here, not paying close attention to
 183 the structure of the problem and the blackbox use of concentration inequalities lead to suboptimal
 184 bounds. On the other hand, as discussed above, a crucial novelty in our analysis is the choice of a
 185 supermartingale difference sequence, defined in the proof below. By working from first principles
 186 using moment generating function of this sequence, we derive two conclusions: 1) an improved risk
 187 bound can be obtained, and 2) the coupling technique is not necessary.

188 *Proof Sketch.* Towards bounding the risk $\sum_{t=0}^k \mathcal{R}(w_t)$, we define random variables

$$Z_t = \frac{1}{2} z_t \eta (\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}}) - \ell_{t+1}(w_{\text{ref}})) + z_t (\mathbf{D}_\psi(w_{\text{ref}}; w_{t+1}) - \mathbf{D}_\psi(w_{\text{ref}}; w_t)) \\ - \frac{3}{16} z_t \eta (\mathcal{R}(w_{\text{ref}}) + \mathcal{R}(w_t)), \quad \forall 0 \leq t \leq T - 1$$

$$\text{where } z_t = \frac{1}{\eta C_4 \sqrt{2\eta\gamma C_3 + 2D_0 + 2\eta \sum_{i=1}^t \ell_i(w_{\text{ref}})}}; \quad \gamma = \max \left\{ 1, \log \frac{1}{\delta} \right\}$$

189 and we let $S_t = \sum_{i=0}^t Z_i$; $\forall 0 \leq t \leq T - 1$. Using Lemma 4, we can relate Z_t and
 190 $z_t \eta (\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}}) + \ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(w_t))$, which is a random variable with expectation 0.
 191 By Lemma 4, we can show $\mathbb{E}[\exp(Z_t) | \mathcal{F}_t] \leq 1$ and hence $(\exp(S_t))_{t \geq 0}$ is a supermartingale. By
 192 Ville’s inequality, we have with probability at least $1 - \delta$, for all $0 \leq k \leq T - 1$

$$\sum_{t=0}^k Z_t \leq \log \frac{1}{\delta}$$

193 Expanding this inequality, in combination with Lemma 5, we obtain the conclusion. \square

194 *Remark 6.* The new analysis does not change the conclusion observed in [30]—that is, with high
 195 probability, the iterate sequence $(w_t)_{t \geq 0}$ behaves as if the domain of the problem is bounded. We
 196 improve the probability that this event happens.

197 3.2 Non-realizable case

198 In the non-realizable case, we do not aim for $1/T$ but only $1/\sqrt{T}$ rates. Hence we do not assume that
 199 the comparator w_{ref} satisfies $\mathcal{R}(w_{\text{ref}}) \leq \rho \mathbf{D}_\psi(w_{\text{ref}}; w_0)/T$ but rather the following assumption on
 200 the excess risk:

201 **Assumption 4.** Let $\mathcal{R}^* = \inf_{v \in \mathcal{X}} \mathcal{R}(v)$, assume that $\mathcal{R}(w_{\text{ref}}) - \mathcal{R}^* \leq \frac{\mathbf{D}_\psi(w_{\text{ref}}; w_0)}{\sqrt{T}}$.

202 We also relax the assumption on the data samples. In the previous case, the data are bounded, i.e.
 203 $\{\|x\|_*, |y|\} \leq 1$ a.s. We will consider in this section two settings, one when the data come from a
 204 sub-Gaussian distribution and one when the data distribution has a polynomial tail.

205 3.2.1 IID data with sub-Gaussian tails

206 We will show the following guarantee:

207 **Theorem 7.** Suppose ℓ is convex, (C_1, C_2) -quadratically-bounded. Given $T, ((x_t, y_t))_{t \leq T}$ are IID
 208 samples with $Q_t = \max\{1, \|x_t\|_*^2, |y_t|^2\}$ and there exists $\sigma \geq 0$ such that for all λ

$$\max\{\mathbb{E}[\exp(\lambda(Q_t^2 - \mathbb{E}[Q_t^2]))], \mathbb{E}[\exp(\lambda(Q_t - \mathbb{E}[Q_t]))]\} \leq \exp(\lambda^2 \sigma^2)$$

209 Let $\mu_1 = \mathbb{E}[Q_t]$ and $\mu_2 = \mathbb{E}[Q_t^2]$. Suppose that w_{ref} satisfies Assumption 4. Then for $\eta \leq$
 210 $\frac{1}{4C_2\sqrt{T\mu_2+2\sigma}\sqrt{T\log\frac{1}{\delta}}}$, with probability at least $1 - 2\delta$, for every $0 \leq k \leq T - 1$

$$\frac{1}{k+1} \sum_{t=0}^k (\mathcal{R}(w_t) - \mathcal{R}^*) + \frac{\mathbf{D}_\psi(w_{\text{ref}}; w_{k+1})}{\eta(k+1)} \leq \frac{R^2}{\eta(k+1)}$$

211 where $R^2 = 16C_4^2(\sigma^2 + 4\mu_1^2) \log\frac{1}{\delta} \eta^2 T + 4D_0(1 + \eta\sqrt{T}) + 4\eta^2 C_4^2(T\mu_2 + 2\sigma\sqrt{T\log\frac{1}{\delta}}) = O(1)$.

212 **Remark 8.** For zero-mean sub-Gaussian variable X , the definition $\mathbb{E}[\exp(\lambda X)] \leq \exp(\lambda^2 \sigma^2)$ for
 213 all λ is equivalent to $\mathbb{E}[\exp(\lambda^2 X^2)] \leq \exp(\lambda^2 \sigma^2)$ for all $0 \leq \lambda \leq \frac{1}{\sigma}$ (see [32]). The lemma
 214 below shows a property of sub-Gaussian variables under scaling and translating. First let us consider
 215 $\sum_{t=1}^T Q_t^2$. Similar to Lemma 5, by bounding the moment generating function of this term, we have
 216 the following (see also Section B4 in [30]).

217 **Lemma 9.** With probability at least $1 - \delta$, $\sum_{t=1}^T Q_t^2 \leq T\mu_2 + 2\sigma\sqrt{T\log\frac{1}{\delta}}$.

218 *Proof of Theorem 7.* The proof of this Theorem uses the technique developed in [16]. We will also
 219 analyze the moment generating function of a suitable martingale sequence. However, the choice
 220 of the coefficients will differ significantly from the previous proof. In this case the structure of the
 221 problem is deeply integrated into the analysis of the martingale. We define

$$Z_t = z_t \eta (\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}})) + z_t (\mathbf{D}_\psi(w_{\text{ref}}; w_{t+1}) - \mathbf{D}_\psi(w_{\text{ref}}; w_t)) \\ - \frac{1}{2} z_t \eta^2 \|g_{t+1}\|_*^2 - 4z_t^2 \eta^2 C_4^2 (\sigma^2 + 4\mu_1^2) \mathbf{D}_\psi(w_{\text{ref}}; w_t) \quad \forall 0 \leq t \leq T - 1$$

$$\text{where } z_t = \frac{1}{4\eta^2 C_4^2 (\sigma^2 + 4\mu_1^2) (T + t + 1)} \quad \forall -1 \leq t \leq T - 1$$

222 and let $S_t = \sum_{i=0}^t Z_i$; $\forall 0 \leq t \leq T - 1$. By Lemma 2, we have

$$Z_t + 4z_t^2 \eta^2 C_4^2 (\sigma^2 + 4\mu_1^2) \mathbf{D}_\psi(w_{\text{ref}}; w_t) \leq z_t \eta (\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}})) + \ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(w_t)$$

223 where we have $\mathbb{E}[(\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}})) + \ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(w_t)] = 0$, and using the same notation
 224 $C_4 = C_1 + C_2(1 + \|w_{\text{ref}}\|)$, by Lemma 1,

$$|(\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}})) + \ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(w_t)| \\ \leq |\ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(w_t)| + |\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}})| \\ \leq |\ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(w_t)| + \mathbb{E}[|\ell_{x,y}(w_t) - \ell_{x,y}(w_{\text{ref}})|] \\ \leq (Q_t + \mu_1) \|w_{\text{ref}} - w_t\| C_4 = ((Q_t - \mu_1) + 2\mu_1) \|w_{\text{ref}} - w_t\| C_4$$

225 Hence applying Lemma 15, we have

$$\mathbb{E}[\exp(Z_t) | \mathcal{F}_t] \exp(4z_t^2 \eta^2 C_4^2 (\sigma^2 + 4\mu_1^2) \mathbf{D}_\psi(w_{\text{ref}}; w_t))$$

$$\begin{aligned}
&= \mathbb{E} [\exp (z_t \eta_t (\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}}) + \ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(w_t))) \mid \mathcal{F}_t] \\
&\leq \exp \left(2z_t^2 \eta^2 C_4^2 \|w_{\text{ref}} - w_t\|^2 (\sigma^2 + 4\mu_1^2) \right) \\
&\leq \exp \left(4z_t^2 \eta^2 C_4^2 (\sigma^2 + 4\mu_1^2) \mathbf{D}_\psi (w_{\text{ref}}; w_t) \right)
\end{aligned}$$

226 Therefore $\mathbb{E} [\exp (Z_t) \mid \mathcal{F}_t] \leq 1$ and hence $(\exp (S_t))_{t \geq 0}$ is a supermartingale. By Ville's inequality,
227 we have with probability at least $1 - \delta$, for all $0 \leq k \leq T - 1$

$$\sum_{t=0}^k Z_t \leq \log \frac{1}{\delta}$$

228 Expanding this inequality we have

$$\begin{aligned}
&\sum_{t=0}^k z_t \eta \mathcal{R}(w_t) + z_k \mathbf{D}_\psi (w_{\text{ref}}; w_{k+1}) \\
&\leq \log \frac{1}{\delta} + z_{-1} D_0 + \eta \mathcal{R}(w_{\text{ref}}) \sum_{t=0}^k z_t + \frac{1}{2} \sum_{t=0}^k z_t \eta^2 \|g_{t+1}\|_*^2 \\
&\quad + \sum_{t=0}^k \underbrace{(z_t + 4z_t^2 \eta^2 C_4^2 (\sigma^2 + 4\mu_1^2) - z_{t-1})}_{\leq 0} \mathbf{D}_\psi (w_{\text{ref}}; w_t) \\
&\stackrel{(a)}{\leq} \log \frac{1}{\delta} + z_{-1} D_0 + \eta \mathcal{R}(w_{\text{ref}}) \sum_{t=0}^k z_t + \frac{1}{2} \sum_{t=0}^k z_t \eta^2 \|g_{t+1}\|_*^2
\end{aligned}$$

229 where for (a), by the choice of $z_t = \frac{1}{4\eta^2 C_4^2 (\sigma^2 + 4\mu_1^2) (T+1+t)}$ we have $z_{t-1} - z_t \geq$
230 $4z_t^2 \eta^2 C_4^2 (\sigma^2 + 4\mu_1^2)$. We highlight that this is where the structure of the problem comes into
231 play. That is, by setting appropriate coefficients, we can leverage gain in the distance in the
232 martingale difference sequence $((z_t - z_{t-1}) \mathbf{D}_\psi (w_{\text{ref}}; w_t))$ to cancel out the loss from bounding the
233 moment generating function $(4z_t^2 \eta^2 C_4^2 (\sigma^2 + 4\mu_1^2) \mathbf{D}_\psi (w_{\text{ref}}; w_t))$. Another important property of
234 the sequence (z_t) is that it is a decreasing sequence and $\frac{z_t}{z_k} \leq 2$ for all t, k . Hence we have

$$\begin{aligned}
&\eta \sum_{t=0}^k (\mathcal{R}(w_t) - \mathcal{R}^*) + \mathbf{D}_\psi (w_{\text{ref}}; w_{k+1}) \\
&\leq 4C_4^2 (\sigma^2 + 4\mu_1^2) \log \frac{1}{\delta} \eta^2 (T+1+k) + 2D_0 + 2(\mathcal{R}(w_{\text{ref}}) - \mathcal{R}^*) \eta (k+1) + \eta^2 \sum_{t=0}^k \|g_{t+1}\|_*^2.
\end{aligned}$$

235 Combined with Lemma 9, with probability at least $1 - 2\delta$, for all $0 \leq k \leq T - 1$

$$\begin{aligned}
&\eta \sum_{t=0}^k (\mathcal{R}(w_t) - \mathcal{R}^*) + \mathbf{D}_\psi (w_{\text{ref}}; w_{k+1}) \\
&\leq 4C_4^2 (\sigma^2 + 4\mu_1^2) \log \frac{1}{\delta} \eta^2 (T+1+k) + 2D_0 + 2(\mathcal{R}(w_{\text{ref}}) - \mathcal{R}^*) \eta (k+1) + \eta^2 \sum_{t=0}^k \|g_{t+1}\|_*^2; \\
&\text{and } \sum_{t=1}^{k+1} Q_t^2 \leq T\mu_2 + 2\sigma \sqrt{T \log \frac{1}{\delta}}
\end{aligned}$$

236 Conditioned on this event, we will prove by induction that

$$\mathbf{D}_\psi (w_{\text{ref}}; w_k) \leq R^2 := 16C_4^2 (\sigma^2 + 4\mu_1^2) \log \frac{1}{\delta} \eta^2 T + 4D_0 + 4D_0 \eta \sqrt{T} + 4\eta^2 C_4^2 \left(T\mu_2 + 2\sigma \sqrt{T \log \frac{1}{\delta}} \right)$$

237 For the base case $k = 0$, it is trivial. Suppose for all $t \leq k$ we have $\mathbf{D}_\psi (w_{\text{ref}}; w_t) \leq R^2$, now we
238 prove for $t = k + 1$. By Lemma 1,

$$\eta^2 \sum_{t=0}^k \|g_{t+1}\|_*^2 \leq \eta^2 \sum_{t=0}^k Q_{t+1}^2 (C_1 + C_2 (1 + \|w_t\|))^2 \leq \eta^2 \sum_{t=0}^k Q_{t+1}^2 (C_4 + C_2 \|w_t - w_{\text{ref}}\|)^2$$

$$\begin{aligned}
&\leq 2\eta^2 C_4^2 \sum_{t=1}^{k+1} Q_t^2 + 2\eta^2 C_2^2 \sum_{t=0}^k Q_{t+1}^2 \|w_t - w_{\text{ref}}\|^2 \\
&\leq \eta^2 (2C_4^2 + 4C_2^2 R^2) \left(T\mu_2 + 2\sigma \sqrt{T \log \frac{1}{\delta}} \right)
\end{aligned}$$

239 Therefore

$$\begin{aligned}
\mathbf{D}_\psi(w_{\text{ref}}; w_{k+1}) &\leq 8C_4^2 (\sigma^2 + 4\mu_1^2) \log \frac{1}{\delta} \eta^2 T + 2D_0 + 2(\mathcal{R}(w_{\text{ref}}) - \mathcal{R}^*) \eta(k+1) \\
&\quad + \eta^2 (2C_4^2 + 4C_2^2 R^2) \left(T\mu_2 + 2\sigma \sqrt{T \log \frac{1}{\delta}} \right) \\
&\leq \frac{R^2}{2} + 4\eta^2 C_2^2 \left(T\mu_2 + 2\sigma \sqrt{T \log \frac{1}{\delta}} \right) R^2 \leq R^2.
\end{aligned}$$

240 Finally we obtain, $\eta \sum_{t=0}^k (\mathcal{R}(w_t) - \mathcal{R}^*) + \mathbf{D}_\psi(w_{\text{ref}}; w_{k+1}) \leq R^2$, as needed. \square

241 3.2.2 IID data with polynomial tails

242 **Theorem 10.** Suppose ℓ is convex, (C_1, C_2) -quadratically bounded. Given T , $((x_t, y_t))_{t \leq T}$ are IID
243 samples with $Q_t = \max \left\{ 1, \|x_t\|_*^2, |y_t|^2 \right\}$ and for some $p \geq 2$ there exists $M \geq \frac{p}{e}$ such that for all
244 λ

$$\max \left\{ \sup_{2 \leq r \leq 2p} \left\{ \mathbb{E} [|Q_t - \mathbb{E}[Q_t]|^r] \right\}, \sup_{2 \leq r \leq p} \left\{ \mathbb{E} [|Q_t^2 - \mathbb{E}[Q_t^2]|^r] \right\} \right\} \leq M$$

245 Let $\mu_1 = \mathbb{E}[Q_t]$ and $\mu_2 = \mathbb{E}[Q_t^2]$. Suppose that w_{ref} satisfies Assumption 4. Then for $\eta \leq$
246 $\frac{1}{C_2 \sqrt{6 \left(T\mu_2 + 2M\sqrt{T} \left(\frac{2}{\delta} \right)^{\frac{1}{p}} \right)}}$, with probability at least $1 - 3\delta$, for every $0 \leq k \leq T - 1$

$$\frac{1}{k+1} \sum_{t=0}^k (\mathcal{R}(w_t) - \mathcal{R}^*) + \frac{\mathbf{D}_\psi(w_{\text{ref}}; w_{k+1})}{\eta(k+1)} \leq \frac{R^2}{2\eta(k+1)}$$

247 where $R = \max \left\{ \sqrt{6 \left(D_0 \left(1 + \eta\sqrt{T} \right) + \eta^2 C_4^2 \left(T\mu_2 + 2M\sqrt{T} \left(\frac{2}{\delta} \right)^{\frac{1}{p}} \right) \right)}, 6 \left(\frac{2}{3} \gamma \left(7 \left(\frac{MT}{\delta} \right)^{1/2p} + \right. \right. \right.$
248 $\left. \left. 2\mu_1 \right) + \sqrt{\log \frac{2}{\delta} T\mu_2 \eta C_4} \right\} = O(1)$, $\gamma = \max \left\{ 1, \log \frac{2}{\delta} \right\}$.

249 **Remark 11.** Since $p \geq 2$, the rate is $O\left(\frac{1}{T^{1/2}} \log \frac{1}{\delta} + \frac{1}{T^{3/4}} \left(\frac{1}{\delta}\right)^{\frac{1}{2p}}\right)$. This rate improves over
250 the $O\left(\left(\frac{1}{T^{1/2}} + \frac{1}{T^{3/4}} \left(\frac{T}{\delta}\right)^{\frac{1}{2p}}\right) \log \frac{T}{\delta}\right)$ rate by [30] by a polynomial factor $T^{\frac{1}{2p}} \log \frac{T}{\delta}$ in the high
251 probability regime where $\delta = \frac{1}{\text{poly}(T)}$.

252 We will give a proof sketch for this theorem. The full proof is deferred to the appendix.

253 *Proof Sketch.* The heavy tailed distribution of the data does not allow us to analyze the moment
254 generating function. In this case, we rely on the coupling technique as in [30]. Since it is not
255 possible to apply Azuma's inequality due to the bounds on the variables being not measurable, and
256 the variables are heavy tailed, we use truncation technique. We define,

$$v_t = \arg \min_{\|w - w_{\text{ref}}\| \leq R} \left\{ \langle \eta_t g_t(v_{t-1}), w \rangle + \mathbf{D}_\psi(w; v_{t-1}) \right\}$$

257 where we use $g_t(v_{t-1})$ to denote the gradient at v_{t-1} using the same data point (x_t, y_t) when
258 computing w_t and we define

$$U_t = (\mathcal{R}(v_t) - \mathcal{R}(w_{\text{ref}})) + \ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(v_t)$$

$$P_t = \begin{cases} U_t & \text{if } |U_t| \leq (A + 2\mu_1) RC_4 \\ (A + 2\mu_1) RC_4 \text{sign}(U_t) & \text{otherwise} \end{cases}$$

$$\text{where } A = \left(\frac{MT}{\delta}\right)^{1/2p} \text{ and } B_t = U_t - P_t.$$

259 We can write

$$\sum_{t=0}^k U_t = \sum_{t=0}^k (P_t - \mathbb{E}[P_t | \mathcal{F}_t]) + \sum_{t=0}^k \mathbb{E}[P_t | \mathcal{F}_t] + \sum_{t=0}^k B_t$$

260 We bound $\sum_{t=0}^k (P_t - \mathbb{E}[P_t | \mathcal{F}_t])$ by applying Freedman's inequality. The terms $\sum_{t=0}^k \mathbb{E}[P_t | \mathcal{F}_t]$
 261 and $\sum_{t=0}^k B_t$ are both the bias terms can be bounded by analyzing the tail of the distribution and
 262 Markov's inequality. We also use Lemma 12 to bound $\sum_{t=0}^k \|g_{t+1}(v_t)\|_*^2$. Finally, using the induction
 263 technique, we can prove that $w_t = v_t$ with high probability and obtain the desired result. \square

264 **Lemma 12** (Lemma A.5 from [30]). *With probability $\geq 1 - \delta$, $\sum_{t=1}^T Q_t^2 \leq T\mu_2 + 2M\sqrt{T} \left(\frac{2}{\delta}\right)^{\frac{1}{p}}$.*

265 4 Generalization bounds of SMD for Markovian data

266 The final result we present in this paper is the following theorem for the case when the data are
 267 sampled from a Markov chain.

268 **Theorem 13.** *Suppose ℓ is convex, (C_1, C_2) -quadratically bounded. Given T , $((x_t, y_t))_{t \leq T}$ are
 269 sampled from a Markov chain with $\max\{\|x_t\|_*^2, |y_t|^2\} \leq 1$ and $(\pi, \tau, \epsilon = \frac{1}{\sqrt{T}})$ is an approximate
 270 stationarity witness. Suppose that w_{ref} satisfies Assumption 4. Then for $\eta \leq \frac{1}{2C_2\sqrt{T(1+2\tau)}}$, with
 271 probability at least $1 - \tau\delta$, for every $0 \leq k \leq T - 1$*

$$\frac{1}{k+1} \sum_{t=0}^k (\mathcal{R}(w_t) - \mathcal{R}^*) + \frac{\mathbf{D}_\psi(w_{\text{ref}}; w_{k+1})}{\eta(k+1)} \leq \frac{R^2}{2\eta(k+1)}.$$

272 where $R = \max\left\{\sqrt{6\left(2D_0 + 2\eta D_0\sqrt{T} + 16\eta^2 C_4^2 T\tau \log \frac{1}{\delta} + 2T\eta^2 C_4^2 + 4\eta^2\tau T C_4^2\right)}, 6(2\eta\tau C_4 +\right.$
 273 $\left.2\eta C_4\epsilon T + 4\eta\tau C_4)\right\} = O(1)$ and $C_4 = C_1 + C_2(1 + \|w_{\text{ref}}\|)$.

274 We will give a proof sketch for this theorem.

275 *Proof Sketch.* The proof of this Theorem follow similarly to that of Theorem 7. The difference here
 276 is we need to bound τ different martingale difference sequences in the form of

$$\mathbb{E}[\ell_{\tau(i+1)+j}(w_{\text{ref}}) | \mathcal{F}_{\tau i+j}] - \mathbb{E}[\ell_{\tau(i+1)+j}(w_{\tau i+j}) | \mathcal{F}_{\tau i+j}] + \ell_{\tau(i+1)+j}(w_{\text{ref}}) - \ell_{\tau(i+1)+j}(w_{\tau i+j})$$

277 for $0 \leq j \leq \tau - 1$, $0 \leq i \leq \frac{T-1-j}{\tau}$. We also need the assumption on the approximate stationarity
 278 witness to see that

$$|\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}}) - \mathbb{E}[\ell_{t+\tau}(w_{\text{ref}}) | \mathcal{F}_t] + \mathbb{E}[\ell_{t+\tau}(w_t) | \mathcal{F}_t]| \leq C_4 R \epsilon.$$

279 Now we only need the union bound over τ sequences, instead of all iterations. The success probability
 280 will decrease from $1 - \delta$ to $1 - \tau\delta$. \square

281 5 Conclusion

282 In this paper, we show a new approach to analyze the generalization error of SMD for quadratically
 283 bounded losses. Our approach improves a logarithmic factor for the realizable setting and non-
 284 realizable setting with light tailed data and a poly T factor for the non-realizable setting with
 285 polynomial tailed data from the prior work by [30]. An inherent limitation of the current approach is
 286 the assumption that we can obtain a fresh sample in each iteration, whereas the setting with finite
 287 training data is still not well understood. In the realizable setting, we require that the data is bounded,
 288 as opposed to more relaxed assumptions in the non-realizable settings. We leave the question of
 289 resolving these issues for future works.

References

- 290
- 291 [1] John C Duchi, Alekh Agarwal, Mikael Johansson, and Michael I Jordan. Ergodic mirror descent.
292 *SIAM Journal on Optimization*, 22(4):1549–1578, 2012.
- 293 [2] David A Freedman. On tail probabilities for martingales. *the Annals of Probability*, pages
294 100–118, 1975.
- 295 [3] Saeed Ghadimi and Guanghui Lan. Optimal stochastic approximation algorithms for strongly
296 convex stochastic composite optimization i: A generic algorithmic framework. *SIAM Journal
297 on Optimization*, 22(4):1469–1492, 2012.
- 298 [4] Eduard Gorbunov, Marina Danilova, David Dobre, Pavel Dvurechenskii, Alexander Gasnikov,
299 and Gauthier Gidel. Clipped stochastic methods for variational inequalities with heavy-tailed
300 noise. *Advances in Neural Information Processing Systems*, 35:31319–31332, 2022.
- 301 [5] Eduard Gorbunov, Marina Danilova, and Alexander Gasnikov. Stochastic optimization with
302 heavy-tailed noise via accelerated gradient clipping. *Advances in Neural Information Processing
303 Systems*, 33:15042–15053, 2020.
- 304 [6] Eduard Gorbunov, Marina Danilova, Innokentiy Shibaev, Pavel Dvurechensky, and Alexander
305 Gasnikov. Near-optimal high probability complexity bounds for non-smooth stochastic
306 optimization with heavy-tailed noise. *arXiv preprint arXiv:2106.05958*, 2021.
- 307 [7] Moritz Hardt, Ben Recht, and Yoram Singer. Train faster, generalize better: Stability of
308 stochastic gradient descent. In *International conference on machine learning*, pages 1225–1234.
309 PMLR, 2016.
- 310 [8] Nicholas JA Harvey, Christopher Liaw, Yaniv Plan, and Sikander Randhawa. Tight analyses for
311 non-smooth stochastic gradient descent. In *Conference on Learning Theory*, pages 1579–1613.
312 PMLR, 2019.
- 313 [9] Elad Hazan and Satyen Kale. Beyond the regret minimization barrier: optimal algorithms
314 for stochastic strongly-convex optimization. *The Journal of Machine Learning Research*,
315 15(1):2489–2512, 2014.
- 316 [10] Ziwei Ji and Matus Telgarsky. Risk and parameter convergence of logistic regression. *arXiv
317 preprint arXiv:1803.07300*, 2018.
- 318 [11] Sham M Kakade and Ambuj Tewari. On the generalization ability of online strongly convex
319 programming algorithms. *Advances in Neural Information Processing Systems*, 21, 2008.
- 320 [12] Ali Kavis, Kfir Yehuda Levy, and Volkan Cevher. High probability bounds for a class
321 of nonconvex algorithms with adagrad stepsize. In *International Conference on Learning
322 Representations*, 2021.
- 323 [13] Ahmed Khaled and Peter Richtárik. Better theory for sgd in the nonconvex world. *arXiv
324 preprint arXiv:2002.03329*, 2020.
- 325 [14] Shaojie Li and Yong Liu. High probability guarantees for nonconvex stochastic gradient descent
326 with heavy tails. In *International Conference on Machine Learning*, pages 12931–12963.
327 PMLR, 2022.
- 328 [15] Xiaoyu Li and Francesco Orabona. A high probability analysis of adaptive sgd with momentum.
329 *arXiv preprint arXiv:2007.14294*, 2020.
- 330 [16] Zijian Liu, Ta Duy Nguyen, Thien Hang Nguyen, Alina Ene, and Huy Lê Nguyen. High
331 probability convergence of stochastic gradient methods. *arXiv preprint arXiv:2302.14843*,
332 2023.
- 333 [17] Zijian Liu, Jiawei Zhang, and Zhengyuan Zhou. Breaking the lower bound with (little) structure:
334 Acceleration in non-convex stochastic optimization with heavy-tailed noise. *arXiv preprint
335 arXiv:2302.06763*, 2023.

- 336 [18] Liam Madden, Emiliano Dall’Anese, and Stephen Becker. High probability convergence
337 and uniform stability bounds for nonconvex stochastic gradient descent. *arXiv preprint*
338 *arXiv:2006.05610*, 2020.
- 339 [19] Eric Moulines and Francis Bach. Non-asymptotic analysis of stochastic approximation
340 algorithms for machine learning. *Advances in neural information processing systems*, 24,
341 2011.
- 342 [20] Alexander V Nazin, Arkadi S Nemirovsky, Alexandre B Tsybakov, and Anatoli B Juditsky.
343 Algorithms of robust stochastic optimization based on mirror descent method. *Automation and*
344 *Remote Control*, 80(9):1607–1627, 2019.
- 345 [21] Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust stochastic
346 approximation approach to stochastic programming. *SIAM Journal on optimization*, 19(4):1574–
347 1609, 2009.
- 348 [22] Ta Duy Nguyen, Thien Hang Nguyen, Alina Ene, and Huy Le Nguyen. High probability
349 convergence of clipped-sgd under heavy-tailed noise. *arXiv preprint arXiv:2302.05437*, 2023.
- 350 [23] Alexander Rakhlin, Ohad Shamir, and Karthik Sridharan. Making gradient descent optimal for
351 strongly convex stochastic optimization. *arXiv preprint arXiv:1109.5647*, 2011.
- 352 [24] Abdurakhmon Sadiev, Marina Danilova, Eduard Gorbunov, Samuel Horváth, Gauthier Gidel,
353 Pavel Dvurechensky, Alexander Gasnikov, and Peter Richtárik. High-probability bounds for
354 stochastic optimization and variational inequalities: the case of unbounded variance. *arXiv*
355 *preprint arXiv:2302.00999*, 2023.
- 356 [25] Matan Schliserman and Tomer Koren. Stability vs implicit bias of gradient methods on separable
357 data and beyond. In *Conference on Learning Theory*, pages 3380–3394. PMLR, 2022.
- 358 [26] Matan Schliserman and Tomer Koren. Tight risk bounds for gradient descent on separable data.
359 *arXiv preprint arXiv:2303.01135*, 2023.
- 360 [27] Ohad Shamir. Gradient methods never overfit on separable data. *The Journal of Machine*
361 *Learning Research*, 22(1):3847–3866, 2021.
- 362 [28] Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The
363 implicit bias of gradient descent on separable data. *The Journal of Machine Learning Research*,
364 19(1):2822–2878, 2018.
- 365 [29] Matus Telgarsky. Margins, shrinkage, and boosting. In *International Conference on Machine*
366 *Learning*, pages 307–315. PMLR, 2013.
- 367 [30] Matus Telgarsky. Stochastic linear optimization never overfits with quadratically-bounded
368 losses on general data. In *Conference on Learning Theory*, pages 5453–5488. PMLR, 2022.
- 369 [31] Joel Tropp. Freedman’s inequality for matrix martingales. 2011.
- 370 [32] Roman Vershynin. *High-dimensional probability: An introduction with applications in data*
371 *science*, volume 47. Cambridge university press, 2018.

372 **A Concentration Inequalities**

373 **Lemma 14.** *Let X be a random variable such that $\mathbb{E}[X] = 0$ and $|X| \leq R$ almost surely. Then for*
 374 $0 \leq \lambda \leq \frac{1}{R}$

$$\mathbb{E}[\exp(\lambda X)] \leq \exp\left(\frac{3}{4}\lambda^2\mathbb{E}[X^2]\right).$$

375 The following lemma is similar to Lemma 2.2 in [16].

376 **Lemma 15.** *Suppose that Q satisfies for all $0 \leq \lambda \leq \frac{1}{\sigma}$, $\mathbb{E}[\exp(\lambda^2 Q^2)] \leq \exp(\lambda^2 \sigma^2)$. Then for*
 377 *variable X such that $\mathbb{E}[X] = 0$ and $|X| \leq a(Q + b)$ for some $a \geq 0$ then for all $\lambda \geq 0$*

$$\mathbb{E}[\exp(\lambda X)] \leq \exp(2\lambda^2 a^2 (\sigma^2 + b^2)).$$

378 *In particular, if $b = 0$ we can have a tighter constant: $\mathbb{E}[\exp(\lambda X)] \leq \exp(\lambda^2 a^2 \sigma^2)$.*

379 *Proof.* We consider $\mathbb{E}[\exp(\lambda X)]$

380 If $0 \leq \lambda \leq \frac{1}{\sqrt{2a\sigma}}$ then using $\exp(x) \leq x + \exp(x^2)$

$$\begin{aligned} \mathbb{E}[\exp(\lambda X)] &\leq \mathbb{E}[\exp(\lambda^2 X^2)] \\ &\leq \mathbb{E}\left[\exp\left(\lambda^2 a^2 (Q + b)^2\right)\right] \\ &\leq \mathbb{E}\left[\exp\left(2\lambda^2 a^2 Q^2 + 2\lambda^2 a^2 b^2\right)\right] \\ &\leq \exp(2\lambda^2 a^2 b^2) \mathbb{E}[\exp(2\lambda^2 a^2 Q^2)] \\ &\leq \exp(2\lambda^2 a^2 (\sigma^2 + b^2)) \end{aligned}$$

381 Otherwise $\frac{1}{\sigma} \leq \lambda\sqrt{2a}$

$$\begin{aligned} \mathbb{E}[\exp(\lambda X)] &\leq \mathbb{E}\left[\exp\left(\lambda^2 a^2 \sigma^2 + \frac{X^2}{4a^2 \sigma^2}\right)\right] \\ &\leq \exp(\lambda^2 a^2 \sigma^2) \mathbb{E}\left[\exp\left(\frac{(Q + b)^2}{4\sigma^2}\right)\right] \\ &\leq \exp(\lambda^2 a^2 \sigma^2) \mathbb{E}\left[\exp\left(\frac{Q^2 + b^2}{2\sigma^2}\right)\right] \\ &\leq \exp(\lambda^2 a^2 \sigma^2) \exp\left(\frac{b^2}{2\sigma^2}\right) \exp\left(\frac{1}{2}\right) \\ &\leq \exp(\lambda^2 a^2 \sigma^2) \exp(\lambda^2 a^2 b^2) \exp(\lambda^2 a^2 \sigma^2) \\ &\leq \exp(2\lambda^2 a^2 (\sigma^2 + b^2)). \end{aligned}$$

382

□

383 **Theorem 16** (Freedman's inequality [2, 31]). *Let $(X_t)_{t \geq 1}$ be a martingale difference sequence.*
 384 *Assume that there exists a constant c such that $|X_t| \leq c$ almost surely for all $t \geq 1$ and define*
 385 $\sigma_t^2 = \mathbb{E}[X_t^2 | X_{t-1}, \dots, X_1]$. *Then for all $b > 0$, $F > 0$ and $T \geq 1$*

$$\Pr\left[\exists T \geq 1 : \left|\sum_{t=1}^T X_t\right| > b \text{ and } \sum_{t=1}^T \sigma_t^2 \leq F\right] \leq 2 \exp\left(-\frac{b^2}{2F + 2cb/3}\right).$$

386 **B Missing Proofs**

387 *Proof of Lemma 2.* Using the optimality condition

$$\langle \eta g_{t+1} + \nabla \psi(w_{t+1}) - \nabla \psi(w_t), w_{\text{ref}} - w_{t+1} \rangle \geq 0$$

388 we have

$$\begin{aligned}
\langle \eta g_{t+1}, w_t - w_{\text{ref}} \rangle &= \langle \eta g_{t+1}, w_{t+1} - w_{\text{ref}} \rangle + \langle \eta g_{t+1}, w_t - w_{t+1} \rangle \\
&\leq \langle \nabla \psi(w_{t+1}) - \nabla \psi(w_t), w_{\text{ref}} - w_{t+1} \rangle + \langle \eta g_{t+1}, w_t - w_{t+1} \rangle \\
&= \mathbf{D}_\psi(w_{\text{ref}}; w_t) - \mathbf{D}_\psi(w_{\text{ref}}; w_{t+1}) - \mathbf{D}_\psi(w_{t+1}; w_t) + \langle \eta g_{t+1}, w_t - w_{t+1} \rangle \\
&\leq \mathbf{D}_\psi(w_{\text{ref}}; w_t) - \mathbf{D}_\psi(w_{\text{ref}}; w_{t+1}) - \frac{1}{2} \|w_t - w_{t+1}\|^2 + \langle \eta g_{t+1}, w_t - w_{t+1} \rangle \\
&\leq \mathbf{D}_\psi(w_{\text{ref}}; w_t) - \mathbf{D}_\psi(w_{\text{ref}}; w_{t+1}) + \frac{\eta^2}{2} \|g_{t+1}\|_*^2
\end{aligned}$$

389 Hence

$$\begin{aligned}
\mathbf{D}_\psi(w_{\text{ref}}; w_{t+1}) - \mathbf{D}_\psi(w_{\text{ref}}; w_t) &\leq \langle \eta g_{t+1}, w_{\text{ref}} - w_t \rangle + \frac{\eta^2}{2} \|g_{t+1}\|_*^2 \\
&\leq \eta (\ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(w_t)) + \frac{\eta^2}{2} \|g_{t+1}\|_*^2
\end{aligned}$$

390 as needed. \square

391 *Proof of Lemma 4.* We have

$$\begin{aligned}
\mathbf{D}_\psi(w_{\text{ref}}; w_{t+1}) - \mathbf{D}_\psi(w_{\text{ref}}; w_t) &\leq \eta (\ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(w_t)) + \frac{\eta^2}{2} \|g_{t+1}\|_*^2 \\
&\leq \eta (\ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(w_t)) + \frac{\eta^2}{2} \ell'_{t+1}(w_t)^2 \\
&\leq \eta (\ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(w_t)) + \eta^2 \rho \ell_{t+1}(w_t) \\
&= \eta \ell_{t+1}(w_{\text{ref}}) - \frac{\eta}{2} \ell_{t+1}(w_t) \leq \eta \ell_{t+1}(w_{\text{ref}}).
\end{aligned}$$

392 Summing up, we have, for any $0 \leq t \leq T$

$$\mathbf{D}_\psi(w_{\text{ref}}; w_t) \leq \mathbf{D}_\psi(w_{\text{ref}}; w_0) + \eta \sum_{i=1}^t \ell_i(w_{\text{ref}}) = D_0 + \eta \sum_{i=1}^t \ell_i(w_{\text{ref}}).$$

393 \square

394 *Proof of Lemma 5.* We have $|\ell_i(w_{\text{ref}}) - \mathcal{R}(w_{\text{ref}})| \leq \max\{\ell_i(w_{\text{ref}}), \mathcal{R}(w_{\text{ref}})\} \leq C_3$ thus by
395 lemma 14, for $\lambda \leq \frac{1}{C_3}$

$$\begin{aligned}
&\mathbb{E}[\exp(\lambda (\ell_i(w_{\text{ref}}) - \mathcal{R}(w_{\text{ref}})))] \\
&\leq \exp\left(\frac{3}{4}\lambda^2 \mathbb{E}[(\ell_i(w_{\text{ref}}) - \mathcal{R}(w_{\text{ref}}))^2]\right) \\
&\stackrel{(a)}{\leq} \exp\left(\frac{3}{4}\lambda^2 \mathbb{E}[\ell_i(w_{\text{ref}})^2]\right) \\
&\stackrel{(b)}{\leq} \exp\left(\frac{3}{4}\lambda^2 C_3 \mathcal{R}(w_{\text{ref}})\right) \leq \exp\left(\frac{3}{4}\lambda \mathcal{R}(w_{\text{ref}})\right),
\end{aligned}$$

396 where for (a) we use $\mathbb{E}[(X - \mathbb{E}[X])^2] \leq \mathbb{E}[X^2]$ and for (b) we use $\ell_i(w_{\text{ref}}) \leq C_3$. Since $\ell_i(w_{\text{ref}})$
397 are independent random variables, we have

$$\begin{aligned}
&\mathbb{E}\left[\exp\left(\lambda \sum_{i=1}^T (\ell_i(w_{\text{ref}}) - \mathcal{R}(w_{\text{ref}}))\right)\right] = \mathbb{E}\left[\prod_{i=1}^T \exp(\lambda (\ell_i(w_{\text{ref}}) - \mathcal{R}(w_{\text{ref}})))\right] \\
&= \prod_{i=1}^T \mathbb{E}[\exp(\lambda (\ell_i(w_{\text{ref}}) - \mathcal{R}(w_{\text{ref}})))] \leq \prod_{i=1}^T \exp\left(\frac{3}{4}\lambda \mathcal{R}(w_{\text{ref}})\right) = \exp\left(\frac{3}{4}\lambda T \mathcal{R}(w_{\text{ref}})\right).
\end{aligned}$$

398 Hence by Markov's inequality

$$\begin{aligned} & \Pr \left[\lambda \sum_{i=1}^T (\ell_i(w_{\text{ref}}) - \mathcal{R}(w_{\text{ref}})) \geq \frac{3}{4} \lambda T \mathcal{R}(w_{\text{ref}}) + \log \frac{1}{\delta} \right] \\ &= \Pr \left[\exp \left(\lambda \sum_{i=1}^T (\ell_i(w_{\text{ref}}) - \mathcal{R}(w_{\text{ref}})) \right) \geq \frac{1}{\delta} \exp \left(\frac{3}{4} \lambda T \mathcal{R}(w_{\text{ref}}) \right) \right] \\ &\leq \frac{\mathbb{E} \left[\exp \left(\lambda \sum_{i=1}^T (\ell_i(w_{\text{ref}}) - \mathcal{R}(w_{\text{ref}})) \right) \right]}{\frac{1}{\delta} \exp \left(\frac{3}{4} \lambda T \mathcal{R}(w_{\text{ref}}) \right)} \leq \delta \end{aligned}$$

399 Choose $\lambda = \frac{1}{C_3}$ we have with probability at least $1 - \delta$

$$\sum_{i=1}^T (\ell_i(w_{\text{ref}}) - \mathcal{R}(w_{\text{ref}})) \leq \frac{3}{4} T \mathcal{R}(w_{\text{ref}}) + C_3 \log \frac{1}{\delta}.$$

400

□

401 *Proof of Theorem 3.* Towards bounding the risk $\sum_{t=0}^k \mathcal{R}(w_t)$, we define random variables

$$\begin{aligned} Z_t &= \frac{1}{2} z_t \eta (\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}}) - \ell_{t+1}(w_{\text{ref}})) + z_t (\mathbf{D}_\psi(w_{\text{ref}}; w_{t+1}) - \mathbf{D}_\psi(w_{\text{ref}}; w_t)) \\ &\quad - \frac{3}{16} z_t \eta (\mathcal{R}(w_{\text{ref}}) + \mathcal{R}(w_t)); \quad \forall 0 \leq t \leq T-1 \end{aligned}$$

$$\text{where } z_t = \frac{1}{\eta C_4 \sqrt{2\eta\gamma C_3 + 2D_0 + 2\eta \sum_{i=1}^t \ell_i(w_{\text{ref}})}}; \quad \gamma = \max \left\{ 1, \log \frac{1}{\delta} \right\}$$

$$\text{and } S_t = \sum_{i=0}^t Z_i; \quad \forall 0 \leq t \leq T-1$$

402 The reason to define these variables is because from Lemma 4, we can bound

$$\begin{aligned} & \mathbb{E} [\exp(Z_t) \mid \mathcal{F}_t] \times \exp \left(\frac{3}{16} z_t \eta (\mathcal{R}(w_{\text{ref}}) + \mathcal{R}(w_t)) \right) \\ &\leq \mathbb{E} \left[\exp \left(\frac{1}{2} z_t \eta (\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}}) - \ell_{t+1}(w_{\text{ref}})) + z_t \left(\eta \ell_{t+1}(w_{\text{ref}}) - \frac{\eta}{2} \ell_{t+1}(w_t) \right) \right) \mid \mathcal{F}_t \right] \\ &= \mathbb{E} \left[\exp \left(\frac{1}{2} z_t \eta (\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}}) + \ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(w_t)) \right) \mid \mathcal{F}_t \right] \end{aligned}$$

403 where now inside the expectation, we have the term $\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}}) + \ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(w_t)$
404 which has expectation 0. This reminds us of Lemma 14. To use this lemma, we notice that, by the
405 assumption that the samples are IID with $\max \{\|x\|_*, |y|\} \leq 1$ and Lemma 1,

$$|\ell_{x,y}(w_{\text{ref}}) - \ell_{x,y}(w_t)| \leq \|w_{\text{ref}} - w_t\| \underbrace{(C_1 + C_2(1 + \|w_{\text{ref}}\|))}_{C_4}$$

406 We also have

$$|\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}})| = |\mathbb{E} [\ell_{x,y}(w_{\text{ref}}) - \ell_{x,y}(w_t)]| \leq C_4 \|w_{\text{ref}} - w_t\|$$

407 Therefore

$$\left| \frac{\eta}{2} (\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}}) + \ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(w_t)) \right| \leq \eta C_4 \|w_{\text{ref}} - w_t\|$$

408 By the choice of z_t we have

$$z_t \leq \frac{1}{\eta C_4 \sqrt{2\eta C_3 + 2D_0 + 2\eta \sum_{i=1}^t \ell_i(w_{\text{ref}})}} \leq \frac{1}{\eta C_4 \sqrt{2\mathbf{D}_\psi(w_{\text{ref}}; w_t)}} \leq \frac{1}{\eta C_4 \|w_{\text{ref}} - w_t\|}$$

409 Now we can apply Lemma 14 to bound

$$\begin{aligned}
& \mathbb{E} [\exp(Z_t) \mid \mathcal{F}_t] \times \exp\left(\frac{3}{16} z_t \eta (\mathcal{R}(w_{\text{ref}}) + \mathcal{R}(w_t))\right) \\
& \leq \exp\left(\frac{3}{4} \frac{1}{4} z_t^2 \eta^2 \mathbb{E}\left[(\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}}) + \ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(w_t))^2 \mid \mathcal{F}_t\right]\right) \\
& \leq \exp\left(\frac{3}{16} z_t^2 \eta^2 \mathbb{E}\left[(\ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(w_t))^2 \mid \mathcal{F}_t\right]\right) \\
& \leq \exp\left(\frac{3}{16} z_t^2 \eta^2 C_4 \|w_{\text{ref}} - w_t\| \mathbb{E}[\ell_{t+1}(w_{\text{ref}}) + \ell_{t+1}(w_t) \mid \mathcal{F}_t]\right) \\
& \leq \exp\left(\frac{3}{16} z_t \eta (\mathcal{R}(w_{\text{ref}}) + \mathcal{R}(w_t))\right)
\end{aligned}$$

410 Therefore $\mathbb{E} [\exp(Z_t) \mid \mathcal{F}_t] \leq 1$ and hence $(\exp(S_t))_{t \geq 0}$ is a supermartingale. By Ville's inequality,
411 we have with probability at least $1 - \delta$, for all $0 \leq k \leq T - 1$

$$\sum_{t=0}^k Z_t \leq \log \frac{1}{\delta}$$

412 Expanding this inequality, we obtain

$$\begin{aligned}
& \sum_{t=0}^k \frac{5}{16} z_t \eta \mathcal{R}(w_t) + z_k \mathbf{D}_\psi(w_{\text{ref}}; w_{k+1}) \\
& \leq \log \frac{1}{\delta} + z_0 \mathbf{D}_\psi(w_{\text{ref}}; w_0) + \frac{11}{16} \eta \mathcal{R}(w_{\text{ref}}) \sum_{t=0}^k z_t + \frac{1}{2} \sum_{t=0}^k z_t \eta \ell_{t+1}(w_{\text{ref}}) + \sum_{t=1}^k \underbrace{(z_t - z_{t-1})}_{\leq 0} \mathbf{D}_\psi(w_{\text{ref}}; w_t) \\
& \stackrel{(a)}{\leq} \log \frac{1}{\delta} + z_0 D_0 + \frac{11}{16} \eta \mathcal{R}(w_{\text{ref}}) (k+1) z_0 + \frac{1}{2} \sum_{t=0}^k \frac{\eta \ell_{t+1}(w_{\text{ref}})}{\eta C_4 \sqrt{2\eta C_3 + 2D_0 + 2\eta \sum_{i=1}^t \ell_i(w_{\text{ref}})}} \\
& \stackrel{(b)}{\leq} \log \frac{1}{\delta} + z_0 D_0 + \frac{11}{16} \eta \mathcal{R}(w_{\text{ref}}) (k+1) z_0 + \frac{1}{2\eta C_4} \sum_{t=0}^k \frac{\eta \ell_{t+1}(w_{\text{ref}})}{\sqrt{2D_0 + 2\eta \sum_{i=1}^{t+1} \ell_i(w_{\text{ref}})}} \tag{1}
\end{aligned}$$

413 For (a) we use the fact that (z_t) is a decreasing sequence and $\mathcal{R}(w_{\text{ref}}) \leq \frac{\rho D_0}{T}$. For (b) we
414 use the assumption $\ell_{t+1}(w_{\text{ref}}) \leq C_3$. Now notice that we can write $\frac{\eta \ell_{t+1}(w_{\text{ref}})}{\sqrt{2D_0 + 2\eta \sum_{i=1}^{t+1} \ell_i(w_{\text{ref}})}} \leq$

415 $\sqrt{2D_0 + 2\eta \sum_{i=1}^{t+1} \ell_i(w_{\text{ref}})} - \sqrt{2D_0 + 2\eta \sum_{i=1}^t \ell_i(w_{\text{ref}})}$ and sum over t we obtain

$$\frac{5}{16} z_k \eta \sum_{t=0}^k \mathcal{R}(w_t) + z_k \mathbf{D}_\psi(w_{\text{ref}}; w_{k+1}) \leq \log \frac{1}{\delta} + \frac{11(k+1)\mathcal{R}(w_{\text{ref}})}{16C_4\sqrt{2\eta\gamma C_3 + 2D_0}} + \frac{1}{\sqrt{2\eta}C_4} \sqrt{D_0 + \eta \sum_{i=1}^{k+1} \ell_i(w_{\text{ref}})}$$

416 Hence

$$\begin{aligned}
& \sum_{t=0}^k \mathcal{R}(w_t) + \frac{16}{5\eta} \mathbf{D}_\psi(w_{\text{ref}}; w_{k+1}) \\
& \leq \frac{16C_4}{5} \left(\log \frac{1}{\delta} + \frac{11(k+1)\mathcal{R}(w_{\text{ref}})}{16C_4\sqrt{2\eta\gamma C_3 + 2D_0}} + \frac{1}{\sqrt{2\eta}C_4} \sqrt{D_0 + \eta \sum_{i=1}^T \ell_i(w_{\text{ref}})} \right) \sqrt{2\eta\gamma C_3 + 2D_0 + 2\eta \sum_{i=1}^T \ell_i(w_{\text{ref}})}
\end{aligned}$$

417 By Lemma 5, with probability at least $1 - \delta$ we have

$$\sum_{i=1}^T \ell_i(w_{\text{ref}}) \leq \frac{7}{4} T \mathcal{R}(w_{\text{ref}}) + C_3 \log \frac{1}{\delta} \leq \frac{7}{4} \rho D_0 + C_3 \gamma$$

418 Therefore with probability at least $1 - 2\delta$

$$\begin{aligned}
& \sum_{t=0}^k \mathcal{R}(w_t) + \frac{16}{5\eta} \mathbf{D}_\psi(w_{\text{ref}}; w_{k+1}) \\
& \leq \left(\frac{16C_4}{5} \log \frac{1}{\delta} + \frac{11(k+1)}{5\sqrt{2\eta\gamma C_3 + 2D_0}} \mathcal{R}(w_{\text{ref}}) + \frac{8}{5\eta} \sqrt{\frac{15}{4}D_0 + 2\eta\gamma C_3} \right) \sqrt{\frac{15}{4}D_0 + 4\eta\gamma C_3} \\
& \leq \frac{16C_4}{5} \log \frac{1}{\delta} \sqrt{\frac{15}{4}D_0 + 4\eta\gamma C_3} + \left(\frac{6}{\eta}D_0 + \frac{32}{5}\gamma C_3 \right) + 3(k+1)\mathcal{R}(w_{\text{ref}}).
\end{aligned}$$

419 which gives us the conclusion. \square

420 *Proof of Theorem 10.* First we consider the bounded domain case. Let

$$v_t = \arg \min_{\|w - w_{\text{ref}}\| \leq R} \{ \langle \eta_t g_t(v_{t-1}), w \rangle + \mathbf{D}_\psi(w; v_{t-1}) \}$$

421 where we use $g_t(v_{t-1})$ to denote the gradient at v_{t-1} using the same data point (x_t, y_t) when
422 computing w_t and we choose

$$R = \max \left\{ \sqrt{6 \left(D_0 + \eta^2 C_4^2 \left(T\mu_2 + 2M\sqrt{T} \left(\frac{2}{\delta} \right)^{\frac{1}{p}} \right) \right)}, 6 \left(\frac{2}{3}\gamma \left(7 \left(\frac{MT}{\delta} \right)^{1/2p} + 2\mu_1 \right) + \sqrt{\log \frac{2}{\delta} T\mu_2} \right) \eta C_4 \right\}$$

423 We have

$$\begin{aligned}
& |(\mathcal{R}(v_t) - \mathcal{R}(w_{\text{ref}}) + \ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(v_t))| \\
& \leq |\ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(v_t)| + |\mathcal{R}(v_t) - \mathcal{R}(w_{\text{ref}})| \\
& \leq (Q_t + \mu_1) \|w_{\text{ref}} - v_t\| C_4 \leq (Q_t + \mu_1) RC_4
\end{aligned} \tag{2}$$

424 Let us define the following variables

$$\begin{aligned}
U_t &= (\mathcal{R}(v_t) - \mathcal{R}(w_{\text{ref}}) + \ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(v_t)) \\
P_t &= \begin{cases} U_t & \text{if } |U_t| \leq (A + 2\mu_1) RC_4 \\ (A + 2\mu_1) RC_4 \text{sign}(U_t) & \text{otherwise} \end{cases} \\
\text{where } A &= \left(\frac{MT}{\delta} \right)^{1/2p} \text{ and } B_t = U_t - P_t.
\end{aligned}$$

425 In words, U_t is the variable of our interest and P_t is the truncated version of U_t and B_t is the bias.

426 We would want to control these terms in order to bound $\sum_{t=0}^k U_t$. We start with the following
427 decomposition

$$\sum_{t=0}^k U_t = \sum_{t=0}^k (P_t - \mathbb{E}[P_t | \mathcal{F}_t]) + \sum_{t=0}^k \mathbb{E}[P_t | \mathcal{F}_t] + \sum_{t=0}^k B_t$$

428 First, we consider the term $\sum_{t=0}^k \mathbb{E}[P_t | \mathcal{F}_t]$.

$$\begin{aligned}
& \mathbb{E}[P_t | \mathcal{F}_t] = \mathbb{E}[P_t - U_t | \mathcal{F}_t] \leq \mathbb{E}[|P_t - U_t| | \mathcal{F}_t] \\
& = \mathbb{E} \left[|P_t - U_t| \left(\mathbf{1}[|U_t| \leq (A + 2\mu_1)RC_4] + \sum_{k=2}^{\infty} \mathbf{1}[(k-1)ARC_4 + 2\mu_1RC_4 \leq |U_t| \leq kARC_4 + 2\mu_1RC_4] \right) \right] \\
& = \mathbb{E} \left[\sum_{k=2}^{\infty} |P_t - U_t| \mathbf{1}[(k-1)ARC_4 + 2\mu_1RC_4 \leq |U_t| \leq kARC_4 + 2\mu_1RC_4] \right] \\
& \leq \sum_{k=2}^{\infty} (kARC_4 + 2\mu_1RC_4 - (A + 2\mu_1)RC_4) RC_4 \mathbb{E}[\mathbf{1}[|U_t| \geq (k-1)ARC_4 + 2\mu_1RC_4]]
\end{aligned}$$

$$\begin{aligned}
&\leq \sum_{k=1}^{\infty} kARC_4 \mathbb{E} [\mathbf{1} [(Q_t + \mu_1) RC_4 \geq kARC_4 + 2\mu_1 RC_4]] \quad (\text{due to 2}) \\
&= \sum_{k=1}^{\infty} kARC_4 \Pr [Q_t \geq kA + \mu_1] \leq ARC_4 \sum_{k=1}^{\infty} k \Pr \left[|Q_t - \mu_1|^{2p} \geq (kA)^{2p} \right] \\
&\leq ARC_4 \sum_{k=1}^{\infty} \frac{Mk}{k^{2p}A^{2p}} = A^{1-2p} RC_4 M \sum_{k=1}^{\infty} k^{1-2p} \leq 2A^{1-2p} RC_4 M
\end{aligned}$$

429 where the last inequality is because $p \geq 2$. We obtain

$$\sum_{t=0}^k \mathbb{E} [P_t | \mathcal{F}_t] \leq 2A^{1-2p} RC_4 MT$$

430 The term $\sum_{t=0}^k B_t \leq \sum_{t=0}^k |B_t| \leq \sum_{t=0}^{T-1} |B_t|$ will be bounded by Markov inequality. From the
431 above deduction,

$$\mathbb{E} \left[\sum_{t=0}^{T-1} |B_t| \right] = \sum_{t=0}^{T-1} \mathbb{E} [|B_t|] = \sum_{t=0}^{T-1} \mathbb{E} [\mathbb{E} [|U_t - P_t| | \mathcal{F}_t]] \leq 2A^{1-2p} RC_4 MT$$

432 With probability at least $1 - \delta$,

$$\sum_{t=0}^{T-1} |B_t| \leq 2TA^{1-2p} RC_4 M \frac{1}{\delta} = 2RC_4 A^{1-2p} \left(\frac{MT}{\delta} \right)$$

433 Finally, we will use Freedman's inequality to bound the remaining term $\sum_{t=0}^k (P_t - \mathbb{E} [P_t | \mathcal{F}_t])$.
434 First, notice that

$$\begin{aligned}
&\mathbb{E} \left[|P_t - \mathbb{E} [P_t | \mathcal{F}_t]|^2 | \mathcal{F}_t \right] \leq \mathbb{E} [P_t^2 | \mathcal{F}_t] \\
&\leq \mathbb{E} [U_t^2 | \mathcal{F}_t] \leq \mathbb{E} \left[(\ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(v_t))^2 | \mathcal{F}_t \right] \\
&\leq R^2 C_4^2 \mathbb{E} [Q_t^2] \leq R^2 C_4^2 \mu_2.
\end{aligned}$$

435 We have $(P_t - \mathbb{E} [P_t | \mathcal{F}_t])$ is a martingale difference sequence with $|P_t - \mathbb{E} [P_t | \mathcal{F}_t]| \leq$
436 $2(A + 2\mu_1) RC_4$. We can apply Freedman's inequality,

$$\begin{aligned}
&\Pr \left[\exists k \geq 0 : \left| \sum_{t=0}^k P_t - \mathbb{E} [P_t | \mathcal{F}_t] \right| > a \text{ and } \sum_{t=0}^k \mathbb{E} \left[|P_t - \mathbb{E} [P_t | \mathcal{F}_t]|^2 | \mathcal{F}_t \right] \leq F \right] \\
&\leq 2 \exp \left(\frac{-2a^2}{2F + 4(A + 2\mu_1) RC_4 a / 3} \right)
\end{aligned}$$

437 If we select

$$\begin{aligned}
F &= T\mu_2 R^2 C_4^2 \\
\text{and } a &= \frac{2}{3} \log \frac{2}{\delta} (A + 2\mu_1) RC_4 + RC_4 \sqrt{\log \frac{2}{\delta} T\mu_2}
\end{aligned}$$

438 we obtain with probability at least $1 - \delta$, for all $k \geq 0$

$$\sum_{t=0}^k P_t - \mathbb{E} [P_t | \mathcal{F}_t] \leq \frac{2}{3} \log \frac{2}{\delta} (A + 2\mu_1) RC_4 + RC_4 \sqrt{\log \frac{2}{\delta} T\mu_2}$$

439 Therefore with probability at least $1 - 3\delta$ we have the following event E : for all $k \geq 0$

$$\begin{aligned}
\sum_{t=0}^k U_t &\leq \frac{2}{3} \log \frac{2}{\delta} (A + 2\mu_1) RC_4 + RC_4 \sqrt{\log \frac{2}{\delta} T\mu_2} + 4RC_4 A^{1-2p} \left(\frac{MT}{\delta} \right) \\
&\leq \frac{2}{3} \gamma \left(7 \left(\frac{MT}{\delta} \right)^{1/2p} + 2\mu_1 \right) RC_4 + RC_4 \sqrt{\log \frac{2}{\delta} T\mu_2}
\end{aligned}$$

$$\text{and } \sum_{t=1}^{k+1} Q_t^2 \leq T\mu_2 + 2M\sqrt{T} \left(\frac{2}{\delta}\right)^{\frac{1}{p}}.$$

440 where we denote $\gamma = \max\{1, \log \frac{2}{\delta}\}$. Furthermore

$$\begin{aligned} \frac{\eta^2}{2} \sum_{t=0}^k \|g_{t+1}(v_t)\|_*^2 &\leq \frac{\eta^2}{2} \sum_{t=0}^k Q_{t+1}^2 (C_1 + C_2(1 + \|v_t\|))^2 \\ &\leq \frac{\eta^2}{2} \sum_{t=0}^k Q_{t+1}^2 (C_4 + C_2\|v_t - w_{\text{ref}}\|)^2 \\ &\leq \eta^2 C_4^2 \sum_{t=1}^{k+1} Q_t^2 + \eta^2 C_2^2 \sum_{t=0}^k Q_{t+1}^2 \|v_t - w_{\text{ref}}\|^2 \\ &\leq \eta^2 (C_4^2 + C_2^2 R^2) \left(T\mu_2 + 2M\sqrt{T} \left(\frac{2}{\delta}\right)^{\frac{1}{p}}\right) \end{aligned}$$

441 Now we will proceed by induction to show that conditioned on the event E , $w_t = v_t$. For the base
442 case, we have $w_0 = v_0$. Suppose that we have $w_t = v_t$ for all $t \leq k$. We will show that $w_{k+1} = v_{k+1}$.
443 From Lemma 2, we have

$$\begin{aligned} &\sum_{t=0}^k \eta (\mathcal{R}(w_t) - \mathcal{R}^*) + \mathbf{D}_\psi(w_{\text{ref}}; w_{k+1}) \\ &\leq D_0 + \sum_{t=0}^k \eta (\mathcal{R}(w_{\text{ref}}) - \mathcal{R}^*) + \eta \sum_{t=0}^k (\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}}) + \ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(w_t)) + \frac{\eta^2}{2} \sum_{t=0}^k \|g_{t+1}\|_*^2 \\ &\leq D_0 + \eta\sqrt{T}D_0 + \eta \sum_{t=0}^k U_t + \frac{\eta^2}{2} \sum_{t=0}^k \|g_{t+1}(v_t)\|_*^2 \\ &\leq D_0 (1 + \eta\sqrt{T}) \\ &\quad + \left(\frac{2}{3}\gamma \left(7 \left(\frac{MT}{\delta}\right)^{1/2p} + 2\mu_1\right) + \sqrt{\log \frac{2}{\delta} T\mu_2}\right) \eta RC_4 + \eta^2 (C_4^2 + C_2^2 R^2) \left(T\mu_2 + 2M\sqrt{T} \left(\frac{2}{\delta}\right)^{\frac{1}{p}}\right) \\ &\leq \frac{R^2}{2} \end{aligned}$$

444 Thus $\|w_{k+1} - w_{\text{ref}}\| \leq R$. And thus $w_{k+1} = v_{k+1}$. Finally, we can conclude that with probability
445 at least $1 - 3\delta$, for all $0 \leq k \leq T - 1$

$$\frac{1}{k+1} \sum_{t=0}^k (\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}})) + \frac{\mathbf{D}_\psi(w_{\text{ref}}; w_{k+1})}{\eta(k+1)} \leq \frac{R^2}{2\eta(k+1)}.$$

446

□

447 *Proof of Theorem 13.* For $0 \leq j \leq \tau - 1$, we define

$$\begin{aligned} Z_i^j &= z_{\tau i+j} \eta (\mathbb{E}[\ell_{\tau(i+1)+j}(w_{\text{ref}}) | \mathcal{F}_{\tau i+j}] - \mathbb{E}[\ell_{\tau(i+1)+j}(w_{\tau i+j}) | \mathcal{F}_{\tau i+j}]) \\ &\quad + z_{\tau i+j} \eta (\ell_{\tau(i+1)+j}(w_{\text{ref}}) - \ell_{\tau(i+1)+j}(w_{\tau i+j})) - 8(z_{\tau i+j} \eta)^2 C_4^2 \mathbf{D}_\psi(w_{\text{ref}}; w_{\tau i+j}) \quad \forall 0 \leq i \leq \frac{T-1-j}{\tau} \\ S_k^j &= \sum_{i=0}^k Z_i^j \quad \forall 0 \leq k \leq \frac{T-1-j}{\tau} \end{aligned}$$

448 where

$$z_t = \frac{1}{8\eta^2 C_4^2 (T+1+t)} \quad \forall -1 \leq t \leq T-1$$

449 We bound

$$\begin{aligned} & \left| \mathbb{E} [\ell_{\tau(i+1)+j}(w_{\text{ref}}) \mid \mathcal{F}_{\tau i+j}] - \mathbb{E} [\ell_{\tau(i+1)+j}(w_{\tau i+j}) \mid \mathcal{F}_{\tau i+j}] + \ell_{\tau(i+1)+j}(w_{\text{ref}}) - \ell_{\tau(i+1)+j}(w_{\tau i+j}) \right| \\ & \leq 2C_4 \|w_{\text{ref}} - w_{\tau i+j}\| \end{aligned}$$

450 By Lemma 15

$$\begin{aligned} & \mathbb{E} \left[\exp \left(Z_i^j \right) \mid \mathcal{F}_{\tau i+j} \right] \\ & = \exp \left(-8(z_{\tau i+j}\eta)^2 C_4^2 \mathbf{D}_\psi(w_{\text{ref}}; w_{\tau i+j}) \right) \\ & \quad \times \mathbb{E} \left[\exp \left(z_{\tau i+j}\eta \left(\mathbb{E} [\ell_{\tau(i+1)+j}(w_{\text{ref}}) \mid \mathcal{F}_{\tau i+j}] - \mathbb{E} [\ell_{\tau(i+1)+j}(w_{\tau i+j}) \mid \mathcal{F}_{\tau i+j}] \right. \right. \right. \\ & \quad \left. \left. \left. + \ell_{\tau(i+1)+j}(w_{\text{ref}}) - \ell_{\tau(i+1)+j}(w_{\tau i+j}) \right) \right) \mid \mathcal{F}_{\tau i+j} \right] \\ & \leq \exp \left(-8(z_{\tau i+j}\eta)^2 C_4^2 \mathbf{D}_\psi(w_{\text{ref}}; w_{\tau i+j}) \right) \exp \left(4(z_{\tau i+j}\eta)^2 C_4^2 \|w_{\text{ref}} - w_{\tau i+j}\|^2 \right) \leq 1 \end{aligned}$$

451 Therefore $\mathbb{E} \left[\exp \left(Z_i^j \right) \mid \mathcal{F}_{\tau i+j} \right] \leq 1$ and hence $\left(\exp \left(S_k^j \right) \right)_{k \geq 0}$ is a supermartingale. By Ville's
452 inequality, we have with probability at least $1 - \delta$, for all $0 \leq k \leq \kappa$

$$\sum_{i=0}^k Z_i^j \leq \log \frac{1}{\delta}$$

453 By union bound over $j = 0, \dots, \tau - 1$, and with probability at least $1 - \tau\delta$ we have for all
454 $0 \leq k \leq T - \tau - 1$

$$\begin{aligned} & \sum_{t=0}^k z_t \eta \left(\mathbb{E} [\ell_{t+\tau}(w_{\text{ref}}) \mid \mathcal{F}_t] - \mathbb{E} [\ell_{t+\tau}(w_t) \mid \mathcal{F}_t] + \ell_{t+\tau}(w_{\text{ref}}) - \ell_{t+\tau}(w_t) \right) \\ & \leq \sum_{t=0}^k 8(z_t \eta)^2 C_4^2 \mathbf{D}_\psi(w_{\text{ref}}; w_t) + \tau \log \frac{1}{\delta} \end{aligned}$$

455 We will proceed to prove by induction that $\mathbf{D}_\psi(w_{\text{ref}}; w_t) \leq \frac{1}{2}R^2$

456 For the base case $t = 0$, this holds trivially. Suppose that this is true for all $0 \leq t \leq k$, we now show
457 for $t = k + 1$.

458 If $k \leq \tau - 1$,

$$\begin{aligned} & \sum_{t=0}^k \eta (\mathcal{R}(w_t) - \mathcal{R}^*) + \mathbf{D}_\psi(w_{\text{ref}}; w_{k+1}) \\ & \leq D_0 + \sum_{t=0}^k \eta (\mathcal{R}(w_{\text{ref}}) - \mathcal{R}^*) + \sum_{t=0}^k \eta (\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}}) + \ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(w_t)) + \frac{\eta^2}{2} \sum_{t=0}^k \|g_{t+1}\|_*^2 \end{aligned}$$

459 We have

$$\sum_{t=0}^k \eta |\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}}) + \ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(w_t)| \leq \sum_{t=0}^k 2\eta C_4 \|w_{\text{ref}} - w_t\| \leq 2\eta C_4 R(k+1) \leq 2\eta C_4 R\tau$$

460 and

$$\begin{aligned} \frac{\eta^2}{2} \sum_{t=0}^k \|g_{t+1}\|_*^2 & \leq \frac{\eta^2}{2} \sum_{t=0}^k (C_1 + C_2(1 + \|w_t\|))^2 \\ & \leq \frac{\eta^2}{2} \sum_{t=0}^k (C_4 + C_2 \|w_t - w_{\text{ref}}\|)^2 \end{aligned}$$

$$\begin{aligned} &\leq \eta^2 C_4^2 (k+1) + \eta^2 C_2^2 R^2 (k+1) \\ &\leq \tau \eta^2 (C_4^2 + C_2^2 R^2) \end{aligned}$$

461 Therefore

$$\mathbf{D}_\psi(w_{\text{ref}}; w_{k+1}) \leq D_0 + \eta D_0 \sqrt{T} + 2\eta C_4 R \tau + \tau \eta^2 (C_4^2 + C_2^2 R^2) \leq \frac{R^2}{2}.$$

462 If $k \geq \tau$,

$$\begin{aligned} &\sum_{t=0}^k z_t \eta (\mathcal{R}(w_t) - \mathcal{R}^*) + z_k \mathbf{D}_\psi(w_{\text{ref}}; w_{k+1}) - z_{-1} \mathbf{D}_\psi(w_{\text{ref}}; w_0) \\ &\leq \sum_{t=0}^k z_t \eta (\mathcal{R}(w_{\text{ref}}) - \mathcal{R}^*) + \sum_{t=0}^k z_t \eta (\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}})) \\ &\quad + \sum_{t=0}^k z_t \eta (\ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(w_t)) + \sum_{t=0}^k \frac{z_t \eta^2}{2} \|g_{t+1}\|_*^2 + \sum_{t=0}^k (z_t - z_{t-1}) \mathbf{D}_\psi(w_{\text{ref}}; w_t) \\ &\leq \sum_{t=0}^k z_t \eta (\mathcal{R}(w_{\text{ref}}) - \mathcal{R}^*) + \sum_{t=k-\tau+1}^k z_t \eta (\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}})) \\ &\quad + \sum_{t=0}^{k-\tau} z_t \eta (\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}}) - \mathbb{E}[\ell_{t+\tau}(w_{\text{ref}}) | \mathcal{F}_t] + \mathbb{E}[\ell_{t+\tau}(w_t) | \mathcal{F}_t]) \\ &\quad + \sum_{t=0}^{k-\tau} z_t \eta (\mathbb{E}[\ell_{t+\tau}(w_{\text{ref}}) | \mathcal{F}_t] - \mathbb{E}[\ell_{t+\tau}(w_t) | \mathcal{F}_t] + \ell_{t+\tau}(w_{\text{ref}}) - \ell_{t+\tau}(w_t)) \\ &\quad + \sum_{t=0}^{k-\tau} z_t \eta (\ell_{t+\tau}(w_t) - \ell_{t+\tau}(w_{\text{ref}})) + \sum_{t=0}^k z_t \eta (\ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(w_t)) \\ &\quad + \sum_{t=0}^{k-\tau} (z_t - z_{t-1}) \mathbf{D}_\psi(w_{\text{ref}}; w_t) + \sum_{t=0}^k \frac{z_t \eta^2}{2} \|g_{t+1}\|_*^2 \\ &\leq \sum_{t=0}^k z_t \eta (\mathcal{R}(w_{\text{ref}}) - \mathcal{R}^*) + \sum_{t=k-\tau+1}^k z_t \eta (\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}})) \\ &\quad + \sum_{t=0}^{k-\tau} z_t \eta (\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}}) - \mathbb{E}[\ell_{t+\tau}(w_{\text{ref}}) | \mathcal{F}_t] + \mathbb{E}[\ell_{t+\tau}(w_t) | \mathcal{F}_t]) \\ &\quad + \tau \log \frac{1}{\delta} + \sum_{t=0}^{k-\tau} 8(z_t \eta)^2 C_4^2 \mathbf{D}_\psi(w_{\text{ref}}; w_t) + \sum_{t=0}^{k-\tau} (z_t - z_{t-1}) \mathbf{D}_\psi(w_{\text{ref}}; w_t) \\ &\quad + \sum_{t=0}^k \frac{z_t \eta^2}{2} \|g_{t+1}\|_*^2 + \sum_{t=0}^{k-\tau} z_t \eta (\ell_{t+\tau}(w_t) - \ell_{t+\tau}(w_{t+\tau-1})) \\ &\quad + \sum_{t=\tau-1}^{k-1} z_{t-\tau+1} \eta (\ell_{t+1}(w_t) - \ell_{t+1}(w_{\text{ref}})) + \sum_{t=0}^k z_t \eta (\ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(w_t)) \\ &\leq \sum_{t=0}^k z_t \eta (\mathcal{R}(w_{\text{ref}}) - \mathcal{R}^*) + \sum_{t=k-\tau+1}^k z_t \eta (\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}})) \\ &\quad + \sum_{t=0}^{k-\tau} z_t \eta (\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}}) - \mathbb{E}[\ell_{t+\tau}(w_{\text{ref}}) | \mathcal{F}_t] + \mathbb{E}[\ell_{t+\tau}(w_t) | \mathcal{F}_t]) + \tau \log \frac{1}{\delta} \\ &\quad + \sum_{t=0}^k \frac{z_t \eta^2}{2} \|g_{t+1}\|_*^2 + \sum_{t=0}^{k-\tau} z_t \eta (\ell_{t+\tau}(w_t) - \ell_{t+\tau}(w_{t+\tau-1})) \end{aligned}$$

$$\begin{aligned}
& + \sum_{t=\tau-1}^{k-1} \eta (z_{t-\tau+1} - z_t) (\ell_{t+1}(w_t) - \ell_{t+1}(w_{\text{ref}})) \\
& + \sum_{t=0}^{\tau-2} z_t \eta (\ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(w_t)) + z_k \eta (\ell_{k+1}(w_{\text{ref}}) - \ell_{k+1}(w_k))
\end{aligned}$$

463 where in the last inequality we use $z_t = \frac{1}{8\eta^2 C_4^2 (T+1+t)}$ to see that $z_t + 8(z_t \eta)^2 C_4^2 \leq z_{t-1}$. Notice
464 that, $\frac{z_t}{z_k} \leq 2$, and $(\mathcal{R}(w_{\text{ref}}) - \mathcal{R}^*) \leq \frac{D_0}{\sqrt{T}}$ we have

$$\begin{aligned}
& \sum_{t=0}^k \eta (\mathcal{R}(w_t) - \mathcal{R}^*) + \mathbf{D}_\psi(w_{\text{ref}}; w_{k+1}) \\
& \leq 2D_0 + 2\eta D_0 \sqrt{T} + 16\eta^2 C_4^2 T \tau \log \frac{1}{\delta} + 2\eta \underbrace{\sum_{t=k-\tau+1}^k |\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}})|}_A \\
& + 2\eta \underbrace{\sum_{t=0}^{k-\tau} |\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}}) - \mathbb{E}[\ell_{t+\tau}(w_{\text{ref}}) | \mathcal{F}_t] + \mathbb{E}[\ell_{t+\tau}(w_t) | \mathcal{F}_t]|}_B \\
& + \underbrace{\eta^2 \sum_{t=0}^k \|g_{t+1}\|_*^2}_C + \underbrace{2\eta \sum_{t=0}^{k-\tau} |\ell_{t+\tau}(w_t) - \ell_{t+\tau}(w_{t+\tau-1})|}_D \\
& + \underbrace{\frac{2(\tau-1)\eta}{T} \sum_{t=\tau-1}^{k-1} |\ell_{t+1}(w_t) - \ell_{t+1}(w_{\text{ref}})| + 2\eta \sum_{t=0}^{\tau-2} |\ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(w_t)| + 2\eta |\ell_{k+1}(w_{\text{ref}}) - \ell_{k+1}(w_k)|}_E
\end{aligned}$$

465 Now we bound each term. For A

$$A = 2\eta \sum_{t=k-\tau+1}^k |\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}})| \leq 2\eta \sum_{t=k-\tau+1}^k C_4 \|w_{\text{ref}} - w_t\| \leq 2\eta \tau C_4 R$$

466 For B , by Assumption 3, $\sup_{t \in \mathbb{Z}_{\geq 0}} \sup_{\mathcal{F}_t} \text{TV}(P_t^{t+\tau}, \pi) \leq \epsilon$,

$$2\eta |\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}}) - \mathbb{E}[\ell_{t+\tau}(w_{\text{ref}}) | \mathcal{F}_t] + \mathbb{E}[\ell_{t+\tau}(w_t) | \mathcal{F}_t]| \leq 2\eta C_4 R \epsilon$$

467 Thus

$$B = 2\eta \sum_{t=0}^{k-\tau} |\mathcal{R}(w_t) - \mathcal{R}(w_{\text{ref}}) - \mathbb{E}[\ell_{t+\tau}(w_{\text{ref}}) | \mathcal{F}_t] + \mathbb{E}[\ell_{t+\tau}(w_t) | \mathcal{F}_t]| \leq 2\eta C_4 R \epsilon T$$

468 For C , similarly to before

$$C = \eta^2 \sum_{t=0}^k \|g_{t+1}\|_*^2 \leq 2T\eta^2 (C_4^2 + C_2^2 R^2)$$

469 For D , we have

$$\begin{aligned}
& |\ell_{t+\tau}(w_t) - \ell_{t+\tau}(w_{t+\tau-1})| \\
& \leq \sum_{i=t+1}^{t+\tau-1} |\ell_{t+\tau}(w_i) - \ell_{t+\tau}(w_{i-1})| \\
& \leq \sum_{i=t+1}^{t+\tau-1} \|w_i - w_{i-1}\| (C_1 + C_2 (1 + \|w_i\|))
\end{aligned}$$

$$\begin{aligned}
&\leq \sum_{i=t+1}^{t+\tau-1} \eta \|\nabla \ell_i(w_{i-1})\| (C_4 + C_2 \|w_i - w_{\text{ref}}\|) \\
&\leq \eta (C_4 + C_2 R) \sum_{i=t+1}^{t+\tau-1} (C_4 + C_2 \|w_{i-1} - w_{\text{ref}}\|) \\
&\leq \eta (C_4 + C_2 R)^2 \tau \leq \eta \tau (2C_4^2 + 2C_2^2 R^2)
\end{aligned}$$

470 We obtain

$$D = 2\eta \sum_{t=0}^{k-\tau} |\ell_{t+\tau}(w_t) - \ell_{t+\tau}(w_{t+\tau-1})| \leq 2\eta^2 \tau T (2C_4^2 + 2C_2^2 R^2)$$

471 For E , since

$$|\ell_{t+1}(w_t) - \ell_{t+1}(w_{\text{ref}})| \leq C_4 R$$

472 Hence

$$\begin{aligned}
E &= \frac{2(\tau-1)\eta}{T} \sum_{t=\tau-1}^{k-1} |\ell_{t+1}(w_t) - \ell_{t+1}(w_{\text{ref}})| \\
&\quad + 2\eta \sum_{t=0}^{\tau-2} |\ell_{t+1}(w_{\text{ref}}) - \ell_{t+1}(w_t)| + 2\eta |\ell_{k+1}(w_{\text{ref}}) - \ell_{k+1}(w_k)| \\
&\leq 2\eta C_4 R \left(\frac{(\tau-1)(k-\tau+1)}{T} + \tau \right) \leq 4\eta \tau C_4 R
\end{aligned}$$

473 Sum up we have

$$\begin{aligned}
&\sum_{t=0}^k \eta (\mathcal{R}(w_t) - \mathcal{R}^*) + \mathbf{D}_\psi(w_{\text{ref}}; w_{k+1}) \\
&\leq 2D_0 + 2\eta D_0 \sqrt{T} + 16\eta^2 C_4^2 T \tau \log \frac{1}{\delta} \\
&\quad + 2\eta \tau C_4 R + 2\eta C_4 R \epsilon T + 2T\eta^2 (C_4^2 + C_2^2 R^2) \\
&\quad + 2\eta^2 \tau T (2C_4^2 + 2C_2^2 R^2) + 4\eta \tau C_4 R \\
&\leq \frac{R^2}{2}
\end{aligned}$$

474 as needed. Finally we have

$$\frac{1}{k+1} \sum_{t=0}^k (\mathcal{R}(w_t) - \mathcal{R}^*) + \frac{\mathbf{D}_\psi(w_{\text{ref}}; w_{k+1})}{\eta(k+1)} \leq \frac{R^2}{2\eta(k+1)}.$$

475

□