
INTERACTIVE LEARNING OF SINGLE-INDEX MODELS VIA STOCHASTIC GRADIENT DESCENT

005 **Anonymous authors**

006 Paper under double-blind review

ABSTRACT

011 Stochastic gradient descent (SGD) is a cornerstone algorithm for high-
012 dimensional optimization, renowned for its empirical successes. Recent theoretical
013 advances have provided a deep understanding of how SGD enables feature
014 learning in high-dimensional nonlinear models, most notably the *single-index*
015 *model* with i.i.d. data. In this work, we study the sequential learning problem
016 for single-index models, also known as generalized linear bandits or ridge bandits,
017 where SGD is a simple and natural solution, yet its learning dynamics remain
018 largely unexplored. We show that, similar to the optimal interactive learner, SGD
019 undergoes a distinct “burn-in” phase before entering the “learning” phase in this
020 setting. Moreover, with an appropriately chosen learning rate schedule, a single
021 SGD procedure simultaneously achieves near-optimal (or best-known) sample
022 complexity and regret guarantees across both phases, for a broad class of link functions.
023 Our results demonstrate that SGD remains highly competitive for learning
024 single-index models under adaptive data.

1 INTRODUCTION

028 Stochastic gradient descent (SGD) and its many variants have achieved remarkable empirical suc-
029 cess in solving high-dimensional optimization problems in machine learning. Recent theoretical
030 advances have provided rigorous analyses of SGD in high-dimensional, non-convex settings for a
031 range of statistical and machine learning tasks, such as tensor decomposition (Ge et al., 2015), PCA
032 (Wang et al., 2017), phase retrieval (Chen et al., 2019; Tan & Vershynin, 2023), to name a few. A
033 particularly intriguing setting is that of single-index models (Dudeja & Hsu, 2018; Ben Arous et al.,
034 2021) (and generalizations to multi-index models (Abbe et al., 2022; 2023; Damian et al., 2022;
035 Arnaboldi et al., 2023; Bietti et al., 2025)) with Gaussian data. In this framework, each observation
(x_t, y_t) consists of a Gaussian feature $x_t \sim \mathcal{N}(0, I_d)$ and a noisy outcome

$$y_t = f(\langle \theta^*, x_t \rangle) + \varepsilon_t,$$

036 where $f : \mathbb{R} \rightarrow \mathbb{R}$ is a known link function, $\theta^* \in \mathbb{S}^{d-1}$ is an unknown parameter vector on the
037 unit sphere in \mathbb{R}^d , and ε_t denotes the unobserved noise. With a learning rate $\eta_t > 0$ and a random
038 initialization $\theta_1 \sim \text{Unif}(\mathbb{S}^{d-1})$, the SGD update for learning single-index models is given by

$$\theta_{t+1/2} = \theta_t - \eta_t(f(\langle \theta_t, x_t \rangle) - y_t)f'(\langle \theta_t, x_t \rangle) \cdot (I - \theta_t \theta_t^\top)x_t, \quad \theta_{t+1} = \frac{\theta_{t+1/2}}{\|\theta_{t+1/2}\|}. \quad (1)$$

044 Here, the first update is a descent step of the population loss $\theta \mapsto \frac{1}{2}\mathbb{E}(f(\langle \theta, x \rangle) - y)^2$ at $\theta = \theta_t$,
045 whose spherical gradient¹ is estimated based on the current sample (x_t, y_t) . It is well known (cf. e.g.
046 (Ben Arous et al., 2021)) that the evolution of SGD in this context exhibits two distinct phases: an
047 initial “search” phase, during which the *correlation* $\langle \theta_t, \theta^* \rangle$ gradually improves from $O(d^{-1/2})$ to
048 $\Omega(1)$, followed by a “descent” phase in which the iterates θ_t converge rapidly to the global optimum
049 θ^* , driving $\langle \theta_t, \theta^* \rangle$ arbitrarily close to 1.

050 Beyond statistical learning, single-index models have found applications in interactive decision-
051 making problems, including bandits and reinforcement learning, where the reward is a nonlinear

052 ¹Recall that the spherical gradient of a function $f : \mathbb{S}^{d-1} \rightarrow \mathbb{R}$ is defined as $\nabla f = Df - \frac{\partial f}{\partial r} \frac{\partial}{\partial r}$, where Df
053 is the Euclidean gradient, and $\frac{\partial}{\partial r}$ is the derivative in the radial direction.

function of the action. An example is manipulation with object interaction, which represents one of the largest open problems in robotics (Billard & Kragic, 2019) and requires designing good sequential decision rules that can deal with sparse and non-linear reward functions and continuous action spaces (Zhu et al., 2019). This setting is known as the *generalized linear bandit* or *ridge bandit* in the bandit literature, where the mean reward satisfies $\mathbb{E}[r_t|a_t] = f(\langle \theta^*, a_t \rangle)$ with a known link function f . Classical results (Filippi et al., 2010; Russo & Van Roy, 2014) show that when $0 < c_1 \leq f'(x) \leq c_2$ everywhere, both the optimal regret and the optimal learner are essentially the same as in the linear bandit case (where $f(x) = x$). Recent studies (Lattimore & Hao, 2021; Huang et al., 2021; Rajaraman et al., 2024) have considered challenging settings where $f'(x)$ could be small around $x = 0$. This line of work yields two main insights:

1. While the final “learning” phase has the same regret as linear bandits, there could be a long “burn-in” period until the learner can identify some action a_t with $\langle \theta^*, a_t \rangle = \Omega(1)$;
2. New exploration algorithms are necessary during this burn-in period, as classical methods such as UCB are provably suboptimal for minimizing the initial exploration cost.

In response to the second point, this line of research has proposed various exploration strategies for the burn-in phase that are often tailored to the specific link function f and rely on noisy gradient estimates via zeroth-order optimization. In contrast, SGD offers a natural and straightforward alternative, as its intrinsic “search” and “descent” phases align well with the “burn-in” and “learning” phases encountered in interactive decision-making.

This paper is devoted to a systematic study of SGD for learning single-index models, including the aforementioned challenging setting where $f'(x)$ could be small around $x \approx 0$, within interactive decision-making settings. In these scenarios, the actions a_t are no longer Gaussian, prompting us to adopt the following exploration strategy:

$$a_t = \sqrt{1 - \sigma_t^2} \theta_t + \sigma_t Z_t, \quad Z_t \sim \text{Unif}(\{a \in \mathbb{S}^{d-1} : \langle a, \theta_t \rangle = 0\}), \quad (2)$$

where an additional hyperparameter $\sigma_t \in [0, 1]$ governs the exploration-exploitation tradeoff. After playing the action a_t and observing the reward r_t , we update the parameter θ_t via the same SGD as Equation (1):

$$\theta_{t+1/2} = \theta_t - \eta_t (f(\langle \theta_t, a_t \rangle) - r_t) f'(\langle \theta_t, a_t \rangle) \cdot (I - \theta_t \theta_t^\top) a_t, \quad \theta_{t+1} = \frac{\theta_{t+1/2}}{\|\theta_{t+1/2}\|}. \quad (3)$$

By simple algebra, the stochastic gradient in Equation (3) is also an unbiased estimator of the population (spherical) gradient of $\theta \mapsto \frac{1}{2} \mathbb{E}(f(\langle \theta, a \rangle) - r)^2$ at $\theta = \theta_t$, with the distribution of a given by Equation (2) and the reward $r = f(\langle \theta^*, a \rangle) + \varepsilon$. Our main result will establish that, for a broad class of link functions, this SGD procedure, with appropriately chosen hyperparameters (η_t, σ_t) , achieves near-optimal performance in both the burn-in and learning phases.

Notation. For $x \in \mathbb{R}^d$, let $\|x\|$ be its ℓ_2 norm. For $x, y \in \mathbb{R}^d$, let $\langle x, y \rangle$ be their inner product. Let \mathbb{S}^{d-1} be the unit sphere in \mathbb{R}^d . Throughout this paper we will use $\theta^* \in \mathbb{S}^{d-1}$ to denote the true parameter, θ_t to denote the current estimate, and $m_t = \langle \theta^*, \theta_t \rangle \in [-1, 1]$ to denote the *correlation*. The standard asymptotic notations o, O, Ω , etc. are used throughout the paper, and we also use $\tilde{O}, \tilde{\Omega}$, etc. to denote the respective meanings with hidden poly-logarithmic factors.

1.1 MAIN RESULTS

First we give a formal formulation of the single-index model in the interactive setting. Let $\theta^* \in \mathbb{S}^{d-1}$ be an unknown parameter vector, and $\mathcal{A} = \mathbb{S}^{d-1}$ be the action space. Upon choosing an action $a_t \in \mathcal{A}$, the learner receives a reward $r_t = f(\langle \theta^*, a_t \rangle) + \varepsilon_t$ for a known link function $f : [-1, 1] \rightarrow \mathbb{R}$ and an unobserved noise ε_t which is assumed to be zero-mean and 1-subGaussian.

Remark 1. *The scaling considered here differs crucially from the prior study on learning single-index models under non-interactive environments (such as (Dudeja & Hsu, 2018; Ben Arous et al., 2021) with Gaussian or i.i.d. features). In the non-interactive setting, it is usually assumed that $x_t \sim \mathcal{N}(0, I_d)$, so that $\|x_t\| \asymp \sqrt{d}$. In the interactive setting, we stick to the convention that actions belong to the unit ℓ_2 ball, in line with settings considered in the bandit literature (Filippi et al., 2010; Russo & Van Roy, 2014; Rajaraman et al., 2024). As a consequence, sample complexity comparisons between the interactive and non-interactive settings must be made with care. We discuss*

108 this in more detail in Section 5 and compare with results established for online SGD with Gaussian
109 features (Ben Arous et al., 2021) after normalizing for the difference in scaling.
110

111 Throughout the paper we make the following mild assumptions on the link function f .

112 **Assumption 1.** *The following conditions hold for the link function f :*

- 113 1. (monotonicity) $f : [-1, 1] \rightarrow [-1, 1]$ is non-decreasing, with $\|f\|_\infty \leq 1$;
- 114 2. (locally linear near $x = 1$) $0 < \gamma_1 \leq f'(x) \leq \gamma_2$ for all $x \in [1 - \gamma_0, 1]$, with absolute
115 constants $\gamma_0, \gamma_1, \gamma_2 > 0$. Without loss of generality we assume that $\gamma_0 \leq 0.1$.

116 In Assumption 1, the monotonicity condition is taken from (Rajaraman et al., 2024) to ensure that
117 reward maximization is aligned with parameter estimation, where improving the alignment $\langle \theta^*, a_t \rangle$
118 directly increases the learner’s reward. In addition, when it comes to SGD, we will show in Sec-
119 tion 5 that the population loss associated with the SGD dynamics in Equation (3) is decreasing in
120 the correlation $m_t = \langle \theta^*, \theta_t \rangle$ only if f is increasing. Without monotonicity, there also exists a coun-
121 terexample where the SGD can never make meaningful progress (cf. Proposition 1). Similar to
122 (Rajaraman et al., 2024), this condition can be generalized to f being even and non-decreasing on
123 $[0, 1]$, which covers, for example, $f(x) = |x|^p$ for all $p > 0$. The second condition in Assumption 1
124 is very mild, satisfied by many natural functions, and ensures that the problem locally resembles a
125 linear bandit near the global optimum $a_t \approx \theta^*$. Finally, we emphasize that this local linearity does
126 not exclude the nontrivial scenario where $f'(x)$ is very small when $x \approx 0$.

127 Our first result is the SGD dynamics in the learning phase, under Assumption 1.

128 **Theorem 1** (Learning Phase). *Let $\varepsilon, \delta > 0$. Under Assumption 1, let $(a_t, \theta_t)_{t \geq 1}$ be given by the SGD
129 evolution in Equation (2) and Equation (3), with an initialization θ_1 such that $\langle \theta_1, \theta^* \rangle \geq 1 - \gamma_0/4$.*

- 130 1. (Pure exploration) By choosing $\eta_t = \tilde{\Theta}(\frac{d}{t} \wedge \frac{1}{d})$ and $\sigma_t^2 = \Theta(1)$, it holds that $m_T \geq 1 - \varepsilon$
131 with probability at least $1 - \delta T$, with $T = \tilde{O}(\frac{d^2}{\varepsilon})$.
- 133 2. (Regret minimization) By choosing $\eta_t = \tilde{\Theta}(\frac{1}{\sqrt{t}} \wedge \frac{1}{d})$ and $\sigma_t^2 = \tilde{\Theta}(\frac{d}{\sqrt{t}} \wedge 1)$, with probability
134 at least $1 - \delta T$ it holds that $\sum_{t=1}^T (f(1) - f(m_t)) = \tilde{O}(d\sqrt{T})$.

135 Both upper bounds in Theorem 1 are near-optimal and match the lower bounds $\Omega(\frac{d^2}{\varepsilon})$ and $\Omega(d\sqrt{T})$
136 for the respective tasks, shown in Theorem 1.7 of (Rajaraman et al., 2024). In other words, SGD
137 with proper learning rate and exploration schedules achieves an optimal learning performance in the
138 learning phase, given a “warm start” θ_1 with $\langle \theta_1, \theta^* \rangle \geq 1 - \gamma_0/2$. To search for this “warm start”
139 through the burn-in phase, we additionally make one of the following assumptions.

140 **Assumption 2.** *There is an absolute constant $c_0 > 0$ such that $f'(x) \geq c_0$ for all $x \in [0, 1]$.*

141 **Assumption 3.** *The link function f is convex on $[0, 1]$.*

143 Specifically, Assumption 2 and 3 cover two different regimes of generalized linear bandits: Assump-
144 tion 2 corresponds to the classical “linear bandit” regime studied in (Filippi et al., 2010; Russo &
145 Van Roy, 2014), and Assumption 3 covers the case with a long burn-in period where $f'(x)$ is small
146 at the beginning, e.g. in (Lattimore & Hao, 2021; Huang et al., 2021). We will discuss the challenges
147 in dropping the convexity assumption for the SGD analysis in Section 5.

148 Under Assumption 2 or 3, our next result characterizes the SGD dynamics in the burn-in phase.

149 **Theorem 2** (Burn-in Phase). *Let $\delta > 0$, and Assumption 1 hold. Let $(a_t, \theta_t)_{t \geq 1}$ be given by the
150 SGD evolution in Equation (2) and Equation (3), with an initialization θ_1 such that $\langle \theta_1, \theta^* \rangle \geq \frac{1}{\sqrt{d}}$.*

- 152 1. Under Assumption 2, by choosing $\eta_t = \tilde{\Theta}(\frac{1}{d^2})$ and $\sigma_t^2 = \Theta(1)$, it holds that $m_T \geq 1 - \gamma_0/4$
153 with probability at least $1 - \delta T$, where $T = \tilde{O}(d^2)$.
- 154 2. Under Assumption 3, by choosing an appropriate learning rate schedule (cf. Lemma 8) and
155 $\sigma_t^2 = \Theta(1)$, it holds that $m_T \geq 1 - \gamma_0/4$ with probability at least $1 - \delta T$, where

$$156 \quad T = \tilde{O}\left(d^2 \int_{1/(2\sqrt{d})}^{1-\gamma_0/4} \frac{m}{f'(m)^2} dm\right).$$

157 Note that for $\theta_1 \sim \text{Unif}(\mathbb{S}^{d-1})$, the condition $\langle \theta^*, \theta_1 \rangle \geq 1/\sqrt{d}$ is fulfilled with a constant proba-
158 bility. A simple hypothesis testing subroutine in (Rajaraman et al., 2024, Lemma 3.1) could further
159 certify it using $\tilde{O}((f(1/\sqrt{d}) - f(0))^{-2})$ samples. Therefore, combining Theorem 1 and 2, we have
160 the following corollary on the overall complexity of SGD.

162 **Corollary 1** (Overall sample complexity and regret). *Under Assumption 1 and Assumption 2 or 3,
163 the SGD evolution in Equation (2) and Equation (3) with proper $(\eta_t, \sigma_t)_{t \geq 1}$ and a hypothesis testing
164 subroutine for initialization satisfies the following:*

165 1. (Pure exploration) For $\varepsilon, \delta > 0$, $m_T \geq 1 - \varepsilon$ with probability at least $1 - \delta T$, where

$$167 \quad 168 \quad 169 \quad T = \tilde{O} \left(d^2 \int_{1/(2\sqrt{d})}^{1-\gamma_0/4} \frac{m}{f'(m)^2} dm + \frac{d^2}{\varepsilon} \right).$$

170 2. (Regret minimization) For $\delta > 0$, with probability at least $1 - \delta T$, the cumulative regret
171 satisfies

$$173 \quad 174 \quad 175 \quad \sum_{t=1}^T (f(1) - f(m_t)) = \tilde{O} \left(\min \left\{ T, d^2 \int_{1/(2\sqrt{d})}^{1-\gamma_0/2} \frac{m}{f'(m)^2} dm + d\sqrt{T} \right\} \right).$$

176 Under Assumption 2, Corollary 1 yields an overall sample complexity bound $\tilde{O}(d^2/\varepsilon)$ and a regret
177 bound $\tilde{O}(\min\{T, d\sqrt{T}\})$, both of which are known to be near-optimal, e.g., in the case of linear
178 bandits (Lattimore & Szepesvári, 2020; Wagenmaker et al., 2022). Under Assumption 3, the upper
179 bounds in Corollary 1 also match the best known guarantees in (Rajaraman et al., 2024), using a
180 different algorithm based on successive hypothesis testing. In the special case $f(x) = x^p$ with odd
181 $p \geq 3$, SGD achieves a regret bound $\tilde{O}(\min\{T, d^p + d\sqrt{T}\})$, which is near-optimal (Huang et al.,
182 2021; Rajaraman et al., 2024). By contrast, many other approaches, including all non-interactive
183 algorithms (in particular, non-interactive SGD) and UCB-based methods, provably incur a larger
184 burn-in cost of $\tilde{\Omega}(d^{p+1})$ (Rajaraman et al., 2024). Therefore, it is striking that SGD attains optimal
185 performance even in the burn-in phase, while simultaneously staying optimal in the learning phase.
186 Taken together, these results highlight SGD as a natural, efficient, and highly competitive algorithm
187 with near-optimal statistical guarantees for learning single-index models in the interactive setting.

188 1.2 RELATED WORK

189 **Single-index models.** Analyzing feature learning in non-linear functions of low-dimensional fea-
190 tures has a long history. The approximation and statistical aspects are well understood in (Barron,
191 2002; Bach, 2017); by contrast, the computational aspects remain more challenging, and positive
192 results typically require additional assumptions on the link function and/or the data distribution. Fo-
193 cusing on the link function f , a rich line of work (Kalai & Sastry, 2009; Shalev-Shwartz et al., 2010;
194 Kakade et al., 2011; Soltanolkotabi, 2017; Frei et al., 2020; Yehudai & Ohad, 2020; Wu, 2022) has
195 exploited its monotonicity or invertibility to obtain efficient learning guarantees under broad distri-
196 butional assumptions. At the other end of the spectrum, the seminal works (Dudeja & Hsu, 2018;
197 Ben Arous et al., 2021) developed a harmonic-analysis framework for studying SGD under Gaus-
198 sian data, sparking extensive follow-up research (Abbe et al., 2022; Bietti et al., 2022; Damian et al.,
199 2022; Ben Arous et al., 2022; Abbe et al., 2023; Zweig et al., 2023; Damian et al., 2024; Ben Arous
200 et al., 2024; Bietti et al., 2025).

201 A representative finding for the single-index model is that the sample complexity of SGD is gov-
202 erned by the *information exponent* of the link function, i.e., the index of its first non-zero Hermite
203 coefficient. In the interactive setting, however, where the data distribution is no longer i.i.d., the
204 information exponent ceases to be an informative measure of SGD’s performance. We defer more
205 discussions to Section 5.

206 **Generalized linear bandits.** The most canonical examples of generalized linear bandits are linear
207 bandits (Dani et al., 2008; Rusmevichientong & Tsitsiklis, 2010; Chu et al., 2011) and ridge bandits
208 with $0 < c_1 \leq |f'(\cdot)| \leq c_2$ everywhere. In both cases, the minimax regret is $\tilde{\Theta}(d\sqrt{T})$ (Filippi
209 et al., 2010; Abbasi-Yadkori et al., 2011; Russo & Van Roy, 2014), attained by algorithms such
210 as LinUCB and information-directed sampling. For more challenging convex link functions, the
211 special cases $f(x) = x^2$ and $f(x) = x^p$ with $p \geq 2$ were analyzed in (Lattimore & Hao, 2021;
212 Huang et al., 2021), using either successive searching algorithms or noisy power methods. These
213 results were substantially generalized by (Rajaraman et al., 2024), which identified the existence of
214 a general burn-in period and established tight upper and lower bounds on the optimal burn-in cost
215 via differential equations. In particular, their upper bound strengthens Corollary 1 in the absence
of convexity, using a refined algorithm of (Lattimore & Hao, 2021) during the burn-in phase and

216 an ETC (explore-then-commit) algorithm in the learning phase. By contrast, we show that a single,
 217 much simpler SGD algorithm achieves the same upper bound for convex f .
 218

219 A related line of work (Fan et al., 2023; Kang et al., 2025) studied the single-index model with
 220 an unknown link function, where the central idea is to estimate the score function. Their resulting
 221 algorithms are of the ETC type, and the regret guarantees rely on a positive lower bound for f' .
 222

223 **Gradient descent in online learning and bandits.** Gradient and mirror descent are classical al-
 224 gorithms in online settings (including online learning and online convex optimization (Cesa-Bianchi
 225 & Lugosi, 2006; Hazan et al., 2016; Orabona, 2019)), as well as in bandit problems with gradient es-
 226 timation, such as EXP3 for adversarial multi-armed bandits and FTRL for adversarial linear bandits
 227 (Lattimore & Szepesvári, 2020). A distinct feature of our work is that our SGD remains a first-order
 228 method even in this bandit problem, in contrast to the zeroth-order stochastic optimization usually
 229 used for single-index models such as (Huang et al., 2021). Moreover, the SGD dynamics for single-
 230 index models demand a more fine-grained analysis than that required by standard online learning
 231 guarantees. Further details are provided in Section 5.
 232

233 1.3 ORGANIZATION

234 The rest of this paper is organized as follows. In Section 2 we present a general analysis of the SGD
 235 update, including bounds on the mean drift, stochastic term, and normalization error. In Section 3
 236 and 4, we analyze the learning and burn-in phases, respectively. Additional discussion is provided
 237 in Section 5, and detailed proofs are deferred to the appendix.
 238

2 ANALYSIS OF THE SGD UPDATE

240 To establish our main results Theorem 1 and 2, we first understand the properties of each SGD
 241 update in Equation (2) and Equation (3). At each time step t , the improvement on the correlation
 242 from $m_t := \langle \theta^*, \theta_t \rangle$ to $m_{t+1} := \langle \theta^*, \theta_{t+1} \rangle$ consists of three parts:

- 243 *1. Drift:* the mean improvement $\mathbb{E}[m_{t+1/2} | \mathcal{F}_t] - m_t$ of the descent step in Equation (3), where
 244 $m_{t+1/2} := \langle \theta^*, \theta_{t+1/2} \rangle$, and \mathcal{F}_t denotes all historic observations up to the end of time t .
 245
- 246 *2. Martingale difference:* the stochastic term $m_{t+1/2} - \mathbb{E}[m_{t+1/2} | \mathcal{F}_t]$ with zero mean.
 247
- 248 *3. Normalization error:* the difference $m_{t+1} - m_{t+1/2}$ due to the normalization step in Equa-
 249 tion (3).
 250

251 We will present generic lemmas to bound each of the above terms in this section, and use them to
 252 analyze the learning and burn-in phases in the next two sections. We start from the drift.
 253

254 **Lemma 1** (Drift). *Let $d \geq 3$. The following identity holds for the drift:*

$$255 \mathbb{E}[m_{t+1/2} | \mathcal{F}_t] - m_t \\ 256 = \frac{\eta_t \sigma_t^2}{d-2} f'(\sqrt{1-\sigma_t^2})(1-m_t^2) \cdot \mathbb{E} \left[f' \left(\sqrt{1-\sigma_t^2} m_t + \sigma_t \sqrt{1-m_t^2} X \right) (1-X^2) \middle| \mathcal{F}_t \right]. \\ 257$$

258 where X follows the one-dimensional marginal of the uniform distribution over \mathbb{S}^{d-2} . In particular,
 259 if $m_t > 0$,
 260

$$261 \mathbb{E}[m_{t+1/2} | \mathcal{F}_t] - m_t \geq c_{dr} \begin{cases} \frac{\eta_t \sigma_t^2}{d} f'(\sqrt{1-\sigma_t^2})(1-m_t^2) & \text{under Assumption 2} \\ \frac{\eta_t \sigma_t^2}{d} f'(\sqrt{1-\sigma_t^2})(1-m_t^2) f'(\sqrt{1-\sigma_t^2} m_t) & \text{under Assumption 3} \end{cases}, \\ 262$$

263 for a universal constant $c_{dr} > 0$.
 264

265 The next result concerns the subexponential concentration property of the martingale difference.
 266

267 **Lemma 2** (Martingale difference). *The Ψ_1 -Orlicz norm (i.e., the subexponential norm) of the mar-
 268 tingale difference has the following upper bound, conditioned on \mathcal{F}_t :*
 269

$$270 \|m_{t+1/2} - \mathbb{E}[m_{t+1/2} | \mathcal{F}_t]\|_{\Psi_1} \leq K_t \triangleq C_{se} \sqrt{\frac{1-m_t^2}{d}} \eta_t \sigma_t f'(\sqrt{1-\sigma_t^2}), \\ 271$$

272 where $C_{se} > 0$ is a universal constant.
 273

270 Based on Lemma 2, we proceed to consider the (discounted) sum of martingale differences. For
 271 $t_0 \geq 0$ and $\beta > 0$, let
 272

$$273 \quad S_t^{t_0, \beta} := \sum_{s=t_0}^{t-1} \beta^{s-t_0} (m_{s+1/2} - \mathbb{E}[m_{s+1/2} | \mathcal{F}_s])$$

$$274$$

$$275$$

276 be a martingale adapted to $\{\mathcal{F}_t\}_{t \geq t_0}$, and $V_t^{t_0, \beta} := \sum_{s=t_0}^{t-1} \beta^{2(s-t_0)} K_s^2$ be a proxy for its predictable
 277 quadratic variation. The following result is a self-normalized concentration inequality for such pro-
 278 cesses established in (Whitehouse et al., 2023, Theorem 3.1):

279 **Lemma 3** (Sum of martingale differences). *Let $\eta_t \leq (C_{se}\gamma_2)^{-1}$ and $\sigma_t^2 \leq \gamma_0$ for all $t \geq 1$, and
 280 $\beta \geq 0$. For $\delta > 0$, it holds that*

$$282 \quad \mathbb{P} \left(\exists t \geq t_0 : |S_t^{t_0, \beta}| \geq C_{mt} \sqrt{V_t^{t_0, \beta} \vee 1} \log \left(\frac{1 + \log(V_t^{t_0, \beta} \vee 1)}{\delta} \right) \text{ and } \beta^{t-t_0} \leq 2 \middle| \mathcal{F}_{t_0} \right) \leq \delta,$$

$$283$$

$$284$$

285 for some universal constant $C_{mt} > 0$.

286 Note that when $\beta \in [0, 1]$, the condition $\beta^{t-t_0} \leq 2$ is vacuous. For $\beta > 1$, this condition results in
 287 a smaller range of $t \in [t_0, t_0 + \log_\beta 2]$. Finally, we bound the normalization error $m_{t+1} - m_{t+1/2}$.

288 **Lemma 4** (Normalization error). *With probability at least $1 - \delta$, it holds that*

$$290 \quad m_{t+1} \geq m_{t+1/2} - C_{nm} \cdot \eta_t^2 \sigma_t^2 (f'(\sqrt{1 - \sigma_t^2}))^2 \log(1/\delta),$$

$$291$$

292 for some universal constant $C_{nm} > 0$. In addition, if $m_{t+1/2} \geq 0$, then with probability $1 - \delta$,

$$293 \quad m_{t+1/2} \geq m_{t+1} \geq m_{t+1/2} \left(1 - C_{nm} \cdot \eta_t^2 \sigma_t^2 (f'(\sqrt{1 - \sigma_t^2}))^2 \log(1/\delta) \right).$$

$$294$$

$$295$$

296 3 ANALYSIS OF THE LEARNING PHASE

297 In this section we analyze the SGD dynamics in the learning phase, given a “warm start” θ_1 with
 298 $m_1 = \langle \theta^*, \theta_1 \rangle \geq 1 - \gamma_0/4$.

301 3.1 PURE EXPLORATION

303 The crux of the proof of Theorem 1 lies in the following lemma, which shows that starting from a
 304 correlation $m_t \geq 1 - \varepsilon$, SGD will improve it to $1 - \varepsilon/2$ after $\tilde{O}(d^2/\varepsilon)$ steps.

305 **Lemma 5** (Local improvement for pure exploration). *Suppose $m_t \geq 1 - \varepsilon$ for some $\varepsilon \leq \gamma_0/4$. Let
 306 $\iota := \log^2(d/\varepsilon\delta)$, and for $s \geq t$, set*

$$308 \quad \eta_s \equiv \eta := \frac{c\varepsilon}{d\iota}, \quad \sigma_s^2 \equiv \sigma^2 := \gamma_0,$$

$$309$$

310 where $c > 0$ is a small absolute constant. Then for $\Delta := Cd/\eta$ and a large absolute constant $C > 0$
 311 independent of c , we have $m_{t+\Delta} \geq 1 - \varepsilon/2$ with probability at least $1 - \Delta\delta$.

312 We call the time interval $[t, t + \Delta]$ an “epoch”, and choose the learning rate based on the epoch.
 313 Lemma 5 shows that, as long as the correlation is large at the beginning of an epoch, then it must
 314 be improved in a linear rate at the end of the epoch. Therefore, by induction and a geometric
 315 series calculation, it is clear that the learning rate schedule given by Lemma 5 corresponds to $\eta_t =$
 316 $\tilde{\Theta}(\frac{d}{t} \wedge \frac{1}{d})$, and Lemma 5 gives an overall sample complexity $\tilde{O}(\frac{d^2}{\varepsilon})$ for pure exploration.

318 In the sequel we prove Lemma 5. We first show that by induction on s that with probability at least
 319 $1 - \Delta\delta/3$, $m_s \geq 1 - 2\varepsilon$ for all $s \in [t, t + \Delta]$. The base case $s = t$ is ensured by the assumption
 320 $m_t \geq 1 - \varepsilon$. For the inductive step, suppose $m_t, \dots, m_{s-1} \geq 1 - 2\varepsilon$. Then

$$321 \quad m_s - m_t = \sum_{r=t}^{s-1} \left[\underbrace{(\mathbb{E}[m_{r+1/2} | \mathcal{F}_r] - m_r)}_{\geq 0 \text{ by Lemma 1}} + \underbrace{(m_{r+1/2} - \mathbb{E}[m_{r+1/2} | \mathcal{F}_r])}_{=: A_r} + \underbrace{(m_{r+1} - m_{r+1/2})}_{=: B_r} \right].$$

$$322$$

$$323$$

324 Thanks to the inductive hypothesis, $K_r = O(\eta\sqrt{\frac{\varepsilon}{d}})$ for all $r \in [t, s-1]$ in Lemma 2, so Lemma 3
 325 (with $t_0 = t, \beta = 1$) gives $|\sum_{r=t}^{s-1} A_r| = O(\eta\sqrt{\frac{\Delta\varepsilon}{d}} \log(\frac{\Delta}{\delta})) = O(\sqrt{\eta\varepsilon} \log(\frac{\Delta}{\delta})) < \frac{\varepsilon}{8}$ with prob-
 326 ability $1 - \frac{\delta}{6}$, by choosing $c > 0$ small enough. Similarly, $\sum_{r=t}^{s-1} |B_r| = O(\Delta\eta^2 \log(\frac{\Delta}{\delta})) =$
 327 $O(d\eta \log(\frac{\Delta}{\delta})) < \frac{\varepsilon}{8}$ with probability $1 - \frac{\delta}{6}$, by Lemma 4 and choosing $c > 0$ small enough. This
 328 implies that $m_s \geq m_t - \frac{\varepsilon}{4} > 1 - 2\varepsilon$ with probability $1 - \frac{\delta}{3}$, completing the induction.
 329

330 Conditioned on the event $m_s \geq 1 - 2\varepsilon$ for all $s \in [t, t + \Delta]$, we distinguish into two regimes in this
 331 epoch. Let $T_0 \geq t$ be the stopping time when $m_s > 1 - \varepsilon/4$ for the first time.
 332

333 **Regime I:** $t \leq s < T_0$. In this regime $m_s \in [1 - 2\varepsilon, 1 - \varepsilon/4]$. We show that $T_0 \leq t + \Delta$ with
 334 probability $1 - \Delta\delta/3$. If $T_0 > t + \Delta$, using the same high-probability bounds, we have
 335

$$336 \quad m_{t+\Delta} - m_t \geq \sum_{s=t}^{t+\Delta-1} (\mathbb{E}[m_{s+1/2} | \mathcal{F}_s] - m_s) - \frac{\varepsilon}{4}$$

$$337$$

$$338$$

339 with probability $1 - \Delta\delta/3$. By Lemma 1 with $1 - m_s^2 = \Omega(\varepsilon)$ and $\sqrt{1 - \sigma_s^2} m_s \geq 1 - \gamma_0$ for $s < T_0$,
 340 the total drift is $\Omega(\frac{\Delta\eta\varepsilon}{d}) = \Omega(C\varepsilon)$. Therefore, for a large absolute constant $C > 0$, we would have
 341 $m_{t+\Delta} \geq 1 - \varepsilon/4$, a contradiction to the assumption $T_0 > t + \Delta$.
 342

343 **Regime II:** $s \geq T_0$. As shown above, this regime is non-empty with high probability. The same
 344 induction starting from $s = T_0$ shows that, with probability $1 - \Delta\delta/3$, $m_s \geq m_{T_0} - \varepsilon/4$ holds for
 345 all $s \in [T_0, t + \Delta]$. In particular, choosing $s = t + \Delta$ gives the desired result $m_{t+\Delta} \geq 1 - \varepsilon/2$.
 346

3.2 REGRET MINIMIZATION

348 The proof of Theorem 1 for regret minimization follows similarly from the following lemma.
 349

350 **Lemma 6** (Local improvement for regret minimization). *Suppose $m_t \geq 1 - \varepsilon$ for some $\varepsilon \leq \gamma_0/4$.
 351 Let $\iota := \log^2(d/\varepsilon\delta)$, and for $s \geq t$, set*

$$352 \quad \eta_s \equiv \eta := \frac{c\varepsilon}{dt}, \quad \sigma_s^2 \equiv \sigma^2 := \varepsilon,$$

$$353$$

$$354$$

355 where $c > 0$ is a small absolute constant. Then for $\Delta := Cd/(\eta\varepsilon)$ and a large absolute constant
 356 $C > 0$ independent of c , with probability at least $1 - \Delta\delta$, we have $\langle \theta^*, a_s \rangle \geq 1 - 4\varepsilon$ for all
 357 $s \in [t, t + \Delta]$, and $m_{t+\Delta} \geq 1 - \varepsilon/2$.
 358

359 The main distinction in Lemma 6 is the choice of a smaller σ_s^2 to encourage exploitation for a small
 360 regret: using the local linearity assumption in Assumption 1, the total regret in the epoch is
 361

$$362 \quad \sum_{s=t}^{t+\Delta} (f(1) - f(\langle \theta^*, a_s \rangle)) \leq (\Delta + 1) \cdot 4\gamma_2\varepsilon = \tilde{O}\left(\frac{d^2}{\varepsilon}\right) \quad \text{with probability } 1 - \Delta\delta.$$

$$363$$

$$364$$

365 In addition, the duration of each epoch becomes longer, with a correspondence $\varepsilon = \tilde{\Theta}(\frac{d}{\sqrt{t}} \wedge 1)$.
 366 This correspondence gives the learning rate and exploration schedule in Theorem 1, as well as the
 367 $\tilde{O}(d\sqrt{T})$ regret bound. The proof of Lemma 6 is deferred to the appendix.
 368

369 4 ANALYSIS OF THE BURN-IN PHASE

$$370$$

371 The analysis of the SGD dynamics in the burn-in phase relies on similar induction ideas, with a more
 372 complicated tradeoff among the three components in the correlation improvement $m_{t+1} - m_t$.
 373

374 4.1 LINK FUNCTION WITH DERIVATIVE LOWER BOUND

$$375$$

376 We first investigate the simpler scenario in Assumption 2, i.e., $f'(x) \geq c_0$ for all $x \in [0, 1]$. In this
 377 case, Theorem 2 is a direct consequence of the following lemma:

378 **Lemma 7** (Burn-in phase under Assumption 2). Suppose $m_1 \geq \frac{1}{\sqrt{d}}$. Let $\iota := \log^2(d/\delta)$, and set

$$380 \quad \eta_t \equiv \eta := \frac{c}{dt}, \quad \sigma_t^2 \equiv \sigma^2 := \gamma_0, \\ 381$$

382 where $c > 0$ is a universal constant. Then for $T := Cd/\eta$ and a large absolute constant $C > 0$
383 independent of c , we have $m_T \geq 1 - \gamma_0/4$ with probability at least $1 - T\delta$.

384 In the sequel we present the proof of Lemma 7. Again we consider the stopping time $T_0 = \min\{t \geq 385 : m_t \geq 1 - \gamma_0/8\}$ and splits into two regimes.

387 **Regime I:** $t \leq T_0$. If $T_0 > T$, we prove by induction that $m_t \geq \frac{1}{2\sqrt{d}} + c_1 \frac{\eta(t-1)}{d}$ for all $t \in [1, T]$
388 with probability at least $1 - T\delta$, for some absolute constant $c' > 0$ independent of c . The base case
389 $t = 1$ is our assumption. Now suppose this lower bound holds for m_1, \dots, m_{t-1} , then by Lemma 1
390 and 4, with probability at least $1 - \frac{\delta}{4}$, for each $s = 1, \dots, t-1$,

$$392 \quad (\mathbb{E}[m_{s+1/2} | \mathcal{F}_s] - m_s) + (m_{s+1} - m_{s+1/2}) = \Omega\left(\frac{\eta}{d}\right) - O\left(\eta^2 \log\left(\frac{2}{\delta}\right)\right) = \Omega\left(\frac{\eta}{d}\right)$$

393 by our choice of η . Here we have critically used the condition $m_s = 1 - \Omega(1)$ for $s < T_0$ when
394 applying Lemma 1, and the inductive hypothesis to ensure $m_s > 0$. By Lemma 2 and 3, with
395 probability $1 - \frac{\delta}{4}$, the sum of martingale difference is at most $O(\eta \sqrt{\frac{T}{d}} \log(\frac{T}{\delta})) = O(\sqrt{\eta} \log(\frac{T}{\delta})) \leq$
396 $\frac{1}{2\sqrt{d}}$ for $c > 0$ small enough. Therefore,

$$399 \quad m_t \geq m_1 - \frac{1}{2\sqrt{d}} + \sum_{s=1}^{t-1} \Omega\left(\frac{\eta}{d}\right) \geq \frac{1}{2\sqrt{d}} + \Omega\left(\frac{\eta(t-1)}{d}\right),$$

401 completing the induction step. Now choosing $t = T$ with $C > 0$ large enough shows the opposite
402 result $m_T \geq 1 - \gamma_0/8$, implying that the event $T_0 > T$ only occurs with probability at most $T\delta/2$.

403 **Regime II:** $T_0 \leq t \leq T$. Under the high-probability event $T_0 \leq T$ and starting from $t = T_0$,

$$405 \quad m_T - m_{T_0} = \sum_{t=T_0}^{T-1} \left[\underbrace{(\mathbb{E}[m_{t+1/2} | \mathcal{F}_t] - m_t)}_{\geq 0 \text{ by Lemma 1}} + \underbrace{(m_{t+1/2} - \mathbb{E}[m_{t+1/2} | \mathcal{F}_t])}_{=: A_t} + \underbrace{(m_{t+1} - m_{t+1/2})}_{=: B_t} \right].$$

409 By Lemma 2 and 3, $|\sum_{t=T_0}^{T-1} A_t| = O(\eta \sqrt{\frac{T}{d}} \log(\frac{T}{\delta})) = O(\sqrt{\eta} \log(\frac{T}{\delta})) < \frac{\gamma_0}{16}$ with probability
410 $1 - T\delta/2$, for $c > 0$ small enough. In addition, Lemma 4 gives $\sum_{t=T_0}^{T-1} |B_t| = O(T\eta^2 \log(\frac{T}{\delta})) =$
411 $O(d\eta \log(\frac{T}{\delta})) < \frac{\gamma_0}{16}$ with probability $1 - T\delta/2$, again for $c > 0$ small enough. Therefore, at the end
412 of this regime, $m_T \geq m_{T_0} - \gamma_0/8 \geq 1 - \gamma_0/4$ with probability $1 - T\delta$, as desired.

4.2 CONVEX LINK FUNCTION

416 When f is convex in Assumption 3, we establish the following lemma.

417 **Lemma 8** (Local improvement for convex link function). For $1 \leq k \leq d-1$, let $\underline{m}_k := (1 -$
418 $\gamma_0)^2 \sqrt{k/d}$, and $\bar{m}_k := (1 - \gamma_0/4) \sqrt{k/d}$. Suppose that $m_t \geq \bar{m}_k$ at the beginning of the k -th
419 epoch. Let $\iota := \log^2(d/\delta)$, and for $s \geq t$, set

$$421 \quad \eta_s \equiv \eta := \frac{cf'(\underline{m}_k)}{\iota d \underline{m}_k}, \quad \sigma_s^2 \equiv \sigma^2 := \gamma_0,$$

423 where $c > 0$ is a small absolute constant. Then for $\Delta := Cd(\underline{m}_{k+1} - \underline{m}_k)/(\eta f'(\underline{m}_k))$ and a large
424 absolute constant $C > 0$ independent of c , we have $m_{t+\Delta} \geq \bar{m}_{k+1}$ with probability at least $1 - \Delta\delta$.

425 Since $m_1 \geq \sqrt{1/d} \geq \bar{m}_1$, a recursive application of Lemma 8 for $k = 1, \dots, d-1$ leads to
426 $m_T \geq 1 - \gamma_0/4$ with probability at least $1 - T\delta$, with (recall that $\gamma_0 \leq 0.1$)

$$428 \quad T = O\left(\log^2\left(\frac{d}{\delta}\right) \cdot d^2 \sum_{k=1}^{d-1} \frac{\underline{m}_k(\underline{m}_{k+1} - \underline{m}_k)}{f'(\underline{m}_k)^2}\right) = \tilde{O}\left(d^2 \int_{\frac{1}{2\sqrt{d}}}^{1-\gamma_0/4} \frac{x dx}{f'(x)^2}\right).$$

431 This completes the proof of Theorem 2. The proof of Lemma 8 is more involved, and we defer the
432 details to the appendix.

432 5 DISCUSSION

434 **Comparison with other descent algorithms.** Our SGD update in Equation (3) is an online
435 gradient descent applied to the loss $\ell_t(\theta) := \frac{1}{2}(r_t - f(\langle \theta, a_t \rangle))^2$, with a_t chosen according to Equation
436 (2). A typical guarantee in online learning takes the form (e.g., via the sequential Rademacher
437 complexity (Rakhlin et al., 2015))

$$438 \quad 439 \quad 440 \quad 441 \quad \sum_{t=1}^T (f(\langle \theta_t, a_t \rangle) - f(\langle \theta^*, a_t \rangle))^2 = \tilde{O}(d).$$

442 which is known as an online regression oracle (Foster & Rakhlin, 2020; Foster et al., 2021). However,
443 this oracle guarantee alone does not yield the optimal regret of θ_t in single-index models; see
444 Theorem 1.5 of (Rajaraman et al., 2024) for a general negative result. This motivates us to move
445 beyond standard online learning guarantees and directly analyze the SGD dynamics.

446 A different descent algorithm for single-index models is also in (Huang et al., 2021), using zeroth-
447 order stochastic optimization to approximate the gradient and implement a noisy power method. In
448 contrast, our SGD is *not* a zeroth-order method: rather than performing gradient descent on the link
449 function $\theta \mapsto f(\langle \theta^*, \theta \rangle)$ where only a zeroth-order oracle is available, we apply gradient descent
450 to the *population loss* $\theta \mapsto \frac{1}{2}\mathbb{E}(r - f(\langle \theta, a \rangle))^2$ for which an unbiased gradient estimator exists for
451 every θ . This change of objective makes SGD a natural yet novel solution to nonlinear ridge bandits.

452 **Necessity of monotonicity.** Throughout this paper we assume that the link function f is mono-
453 tone, an assumption that is not needed in the non-interactive setting (see, e.g., (Ben Arous et al.,
454 2021)). This condition, however, turns out to be essentially necessary for SGD to succeed under our
455 exploration strategy equation 2. Indeed, when $\sigma_t \equiv \sigma$, SGD is performed on the population loss

$$456 \quad 457 \quad \mathbb{E}[(r_t - f(\langle \theta_t, a_t \rangle))^2] = \mathbb{E}[(f(\langle \theta^*, a_t \rangle) - f(\langle \theta_t, a_t \rangle))^2] + \text{Var}(r_t) \\ 458 \quad 459 \quad = \mathbb{E}\left[\left(f\left(\sqrt{1-\sigma^2}\right) - f\left(\sqrt{1-\sigma^2}\langle \theta^*, \theta_t \rangle + \sigma\langle \theta^*, Z_t \rangle\right)\right)^2\right] + \text{Var}(r_t) \\ 460 \quad 461 \quad \approx \left(f\left(\sqrt{1-\sigma^2}\right) - f\left(\sqrt{1-\sigma^2}\langle \theta^*, \theta_t \rangle\right)\right)^2 + \text{Var}(r_t),$$

462 where the last approximation uses that $\langle \theta^*, Z_t \rangle$ is typically of order $\tilde{O}(1/\sqrt{d})$ and thus often neg-
463 ligible. Recall that for SGD to succeed at the population level, the population loss must decrease
464 with the alignment $\langle \theta^*, \theta_t \rangle$ (stated as Assumption A in (Ben Arous et al., 2021)). Treating $\text{Var}(r_t)$
465 as a constant, this requires f to be increasing on $[0, \sqrt{1-\sigma^2}]$ in the interactive setting (assuming
466 $f'(0) > 0$). Hence, whenever σ is bounded away from 1, a monotonicity assumption on f is indis-
467 pensable in the interactive setting. By contrast, when $\sigma = 1$ the monotonicity condition is unneces-
468 sary: in this case Equation (2) reduces to pure exploration, and the problem essentially collapses to
469 the non-interactive setting. However, this would eliminate the statistical benefits of interaction.

470 We also provide an explicit counterexample to formally support the above intuition.

471 **Proposition 1.** Consider the SGD dynamics in Equation (3) applied to the link function

$$472 \quad 473 \quad 474 \quad 475 \quad 476 \quad 477 \quad 478 \quad f(m) = \begin{cases} 0 & \text{if } m \leq 0 \\ -m & \text{if } 0 < m \leq \frac{1}{3}, \\ m - \frac{2}{3} & \text{if } \frac{1}{3} < m \leq 1 \end{cases}$$

479 with any initialization $m_1 = \langle \theta^*, \theta_1 \rangle \leq 0.1$, any exploration schedule $\sigma_t \leq 0.1$, and any learning
480 rate $\eta_t \leq \frac{c}{\log(T/\delta)}$ for some small absolute constant $c > 0$. Then $\mathbb{P}(\max_{t \in [T]} m_t \leq 0.2) \geq 1 - \delta$.

481 Note that the above link function f violates the monotonicity condition: it first decreases and then
482 increases on $[0, 1]$. By choosing $\delta = T^{-2}$, Proposition 1 shows that with any practical initialization,
483 any exploration schedule that does not essentially correspond to a non-interactive exploration, and
484 any learning rate that is not too large to escape the local optima, with high probability the resulting
485 SGD cannot achieve an alignment better than a small constant (say 0.2).

486 **Comparison with information exponent.** In the non-interactive case with $a_t \sim \mathcal{N}(0, I_d)$, it
 487 is known that the *information exponent* of f determines the sample complexity of SGD. In the
 488 interactive case, however, the monotonicity of f ensures that the information exponent is always 1.
 489 Indeed, for the first Hermite polynomial $H_1(x) = x$, Chebyshev's sum inequality yields

$$490 \quad 491 \quad \mathbb{E}_{Z \sim \mathcal{N}(0,1)}[f(Z)H_1(Z)] \geq \mathbb{E}_{Z \sim \mathcal{N}(0,1)}[f(Z)] \cdot \mathbb{E}_{Z \sim \mathcal{N}(0,1)}[H_1(Z)] = 0,$$

492 with equality iff $f \equiv c$ is a constant. Moreover, the sample complexity predicted by the information
 493 exponent is no longer tight. For instance, when $f(x) = x^p$ with an odd $p \geq 3$, the sample complexity
 494 of SGD with $a_t \sim \mathcal{N}(0, I_d/d)$ is $\tilde{O}(d^{p+1})$ (see remark below), which is strictly worse than the $\tilde{O}(d^p)$
 495 guarantee obtained by our interactive SGD. These observations show that the information exponent
 496 ceases to be an informative measure for SGD in the interactive case, for the actions a_t are no longer
 497 Gaussian.

498 **Remark 2.** For $f(x) = x^p$ with odd $p \geq 3$, the population square loss has information exponent
 499 equal to 1. Let c_1 be the coefficient of the linear term $\langle \theta^*, \theta_t \rangle$ in

$$500 \quad 501 \quad \mathbb{E}_{X \sim \mathcal{N}(0, I_d)} \left[(f(\langle \theta^*, X \rangle) - f(\langle \theta_t, X \rangle))^2 \right],$$

502 then $c_1 = -2u_1(f)^2$ with $u_1(f)$ being the first Hermite coefficient of f . When we scale down the
 503 input features into $X \sim \mathcal{N}(0, I_d/d)$, we effectively changes f to $\tilde{f}(x) = (x/\sqrt{d})^p$, so c_1 becomes
 504 $d^{-p}c_1$. Therefore, the SNR effectively worsens by a factor of d^p .

505 **Dropping the convexity assumption.** The convexity assumption in Assumption 3 is not required
 506 in the statistical complexity framework developed for ridge bandits in (Rajaraman et al., 2024).
 507 Relying only on the monotonicity of f , they establish the upper bound

$$508 \quad 509 \quad \tilde{O} \left(d^2 \int_{1/\sqrt{d}}^{1/2} \frac{d[x^2]}{\max_{\frac{1}{\sqrt{d}} \leq y \leq x} f'(y)^2} \right)$$

510 on the sample complexity of finding an action a_t with $\langle \theta^*, a_t \rangle \geq 1/2$. In comparison, under our con-
 511 vevity assumption the denominator simplifies to $f'(x)^2$. There are two main obstacles to recovering
 512 this sharper bound. First, our analysis in Lemma 4 requires a conservative choice of the learning
 513 rate η_t , which in turn depends on having a lower bound for $f'(m_t)$ at the current correlation m_t .
 514 Obtaining such a bound is challenging without further conditions on f . In this paper we handle this
 515 by using $f'(m_t) \geq c$ in the generalized linear case, and $f'(m_t) \geq f'(\underline{m}_t)$ in the convex case, where
 516 $\underline{m}_t \leq m_t$ is known. Second, achieving the factor $\max_{1/\sqrt{d} \leq y \leq x} f'(y)^2$ requires a careful tuning of
 517 σ_t to target the maximizer of f' , which in turn relies on knowledge of the current correlation m_t . In
 518 (Rajaraman et al., 2024), this is accomplished by running a separate hypothesis test. However, such
 519 an additional testing step is not compatible with the dynamics of SGD.

520 REFERENCES

- 521 Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic
 522 bandits. *Advances in neural information processing systems*, 24, 2011.
- 523 Emmanuel Abbe, Enric Boix Adsera, and Theodor Misiakiewicz. The merged-staircase property: a
 524 necessary and nearly sufficient condition for sgd learning of sparse functions on two-layer neural
 525 networks. In *Conference on Learning Theory*, pp. 4782–4887. PMLR, 2022.
- 526 Emmanuel Abbe, Enric Boix Adsera, and Theodor Misiakiewicz. Sgd learning on neural networks:
 527 leap complexity and saddle-to-saddle dynamics. In *The Thirty Sixth Annual Conference on Learn-
 528 ing Theory*, pp. 2552–2623. PMLR, 2023.
- 529 Luca Arnaboldi, Ludovic Stephan, Florent Krzakala, and Bruno Loureiro. From high-dimensional &
 530 mean-field dynamics to dimensionless odes: A unifying approach to sgd in two-layers networks.
 531 In *The Thirty Sixth Annual Conference on Learning Theory*, pp. 1199–1227. PMLR, 2023.
- 532 Francis Bach. Breaking the curse of dimensionality with convex neural networks. *Journal of Ma-
 533 chine Learning Research*, 18(19):1–53, 2017.

-
- 540 Andrew R Barron. Universal approximation bounds for superpositions of a sigmoidal function.
541 *IEEE Transactions on Information theory*, 39(3):930–945, 2002.
542
- 543 Gerard Ben Arous, Reza Gheissari, and Aukosh Jagannath. Online stochastic gradient descent on
544 non-convex losses from high-dimensional inference. *Journal of Machine Learning Research*, 22
545 (106):1–51, 2021.
- 546 Gerard Ben Arous, Reza Gheissari, and Aukosh Jagannath. High-dimensional limit theorems for
547 sgd: Effective dynamics and critical scaling. *Advances in neural information processing systems*,
548 35:25349–25362, 2022.
- 549 Gérard Ben Arous, Cédric Gerbelot, and Vanessa Piccolo. High-dimensional optimization for multi-
550 spiked tensor pca. *arXiv preprint arXiv:2408.06401*, 2024.
551
- 552 Alberto Bietti, Joan Bruna, Clayton Sanford, and Min Jae Song. Learning single-index models with
553 shallow neural networks. *Advances in neural information processing systems*, 35:9768–9783,
554 2022.
- 555 Alberto Bietti, Joan Bruna, and Lucas Pillaud-Vivien. On learning gaussian multi-index models
556 with gradient flow part i: General properties and two-timescale learning. *Communications on
557 Pure and Applied Mathematics*, 2025.
- 558 Aude Billard and Danica Kragic. Trends and challenges in robot manipulation. *Science*, 364(6446):
559 eaat8414, 2019.
560
- 561 Nicolo Cesa-Bianchi and Gábor Lugosi. *Prediction, learning, and games*. Cambridge university
562 press, 2006.
563
- 564 Yuxin Chen, Yuejie Chi, Jianqing Fan, and Cong Ma. Gradient descent with random initialization:
565 Fast global convergence for nonconvex phase retrieval. *Mathematical Programming*, 176:5–37,
566 2019.
- 567 Wei Chu, Lihong Li, Lev Reyzin, and Robert Schapire. Contextual bandits with linear payoff func-
568 tions. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and
569 Statistics*, pp. 208–214. JMLR Workshop and Conference Proceedings, 2011.
- 570 Alex Damian, Lucas Pillaud-Vivien, Jason Lee, and Joan Bruna. Computational-statistical gaps in
571 gaussian single-index models. In *The Thirty Seventh Annual Conference on Learning Theory*, pp.
572 1262–1262. PMLR, 2024.
573
- 574 Alexandru Damian, Jason Lee, and Mahdi Soltanolkotabi. Neural networks can learn representations
575 with gradient descent. In *Conference on Learning Theory*, pp. 5413–5452. PMLR, 2022.
- 576 Varsha Dani, Thomas P Hayes, and Sham M Kakade. Stochastic linear optimization under bandit
577 feedback. *Conference on Learning Theory*, pp. 355–366, 2008.
578
- 579 Rishabh Dudeja and Daniel Hsu. Learning single-index models in gaussian space. In *Conference
580 On Learning Theory*, pp. 1887–1930. PMLR, 2018.
- 581 Jianqing Fan, Zhuoran Yang, and Mengxin Yu. Understanding implicit regularization in over-
582 parameterized single index model. *Journal of the American Statistical Association*, 118(544):
583 2315–2328, 2023.
- 584 Sarah Filippi, Olivier Cappe, Aurélien Garivier, and Csaba Szepesvári. Parametric bandits: The
585 generalized linear case. *Advances in neural information processing systems*, 23, 2010.
586
- 587 Dylan Foster and Alexander Rakhlin. Beyond ucb: Optimal and efficient contextual bandits with
588 regression oracles. In *International conference on machine learning*, pp. 3199–3210. PMLR,
589 2020.
- 590 Dylan J Foster, Sham M Kakade, Jian Qian, and Alexander Rakhlin. The statistical complexity of
591 interactive decision making. *arXiv preprint arXiv:2112.13487*, 2021.
592
- 593 Spencer Frei, Yuan Cao, and Quanquan Gu. Agnostic learning of a single neuron with gradient
descent. *Advances in Neural Information Processing Systems*, 33:5417–5428, 2020.

-
- 594 Rong Ge, Furong Huang, Chi Jin, and Yang Yuan. Escaping from saddle points—online stochastic
595 gradient for tensor decomposition. In *Conference on learning theory*, pp. 797–842. PMLR, 2015.
596
- 597 Elad Hazan et al. Introduction to online convex optimization. *Foundations and Trends® in Opti-*
598 *mization*, 2(3-4):157–325, 2016.
- 599 Baihe Huang, Kaixuan Huang, Sham Kakade, Jason D Lee, Qi Lei, Runzhe Wang, and Jiaqi Yang.
600 Optimal gradient-based algorithms for non-concave bandit optimization. *Advances in Neural*
601 *Information Processing Systems*, 34:29101–29115, 2021.
- 602
- 603 Sham M Kakade, Varun Kanade, Ohad Shamir, and Adam Kalai. Efficient learning of general-
604 ized linear and single index models with isotonic regression. *Advances in Neural Information*
605 *Processing Systems*, 24, 2011.
- 606 Adam Tauman Kalai and Ravi Sastry. The isotron algorithm: High-dimensional isotonic regression.
607 In *COLT*, volume 1, pp. 9, 2009.
- 608
- 609 Yue Kang, Mingshuo Liu, Bongsoo Yi, Jing Lyu, Zhi Zhang, Doudou Zhou, and Yao Li. Single index
610 bandits: Generalized linear contextual bandits with unknown reward functions. *arXiv preprint*
611 *arXiv:2506.12751*, 2025.
- 612 Tor Lattimore and Botao Hao. Bandit phase retrieval. *Advances in Neural Information Processing*
613 *Systems*, 34:18801–18811, 2021.
- 614
- 615 Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- 616 Francesco Orabona. A modern introduction to online learning. *arXiv preprint arXiv:1912.13213*,
617 2019.
- 618
- 619 Nived Rajaraman, Yanjun Han, Jiantao Jiao, and Kannan Ramchandran. Statistical complexity and
620 optimal algorithms for nonlinear ridge bandits. *The Annals of Statistics*, 52(6):2557–2582, 2024.
- 621 Alexander Rakhlin, Karthik Sridharan, and Ambuj Tewari. Online learning via sequential complex-
622 ities. *J. Mach. Learn. Res.*, 16(1):155–186, 2015.
- 623
- 624 Paat Rusmevichientong and John N Tsitsiklis. Linearly parameterized bandits. *Mathematics of*
625 *Operations Research*, 35(2):395–411, 2010.
- 626
- 627 Daniel Russo and Benjamin Van Roy. Learning to optimize via information-directed sampling.
628 *Advances in neural information processing systems*, 27, 2014.
- 629
- 630 Shai Shalev-Shwartz, Ohad Shamir, and Karthik Sridharan. Learning kernel-based halfspaces with
631 the zero-one loss. *arXiv preprint arXiv:1005.3681*, 2010.
- 632
- 633 Mahdi Soltanolkotabi. Learning relus via gradient descent. *Advances in neural information pro-*
634 *cessing systems*, 30, 2017.
- 635
- 636 Yan Shuo Tan and Roman Vershynin. Online stochastic gradient descent with arbitrary initialization
637 solves non-smooth, non-convex phase retrieval. *Journal of Machine Learning Research*, 24(58):
638 1–47, 2023.
- 639
- 640 Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*,
641 volume 47. Cambridge university press, 2018.
- 642
- 643 Andrew J Wagenmaker, Yifang Chen, Max Simchowitz, Simon Du, and Kevin Jamieson. Reward-
644 free RL is no harder than reward-aware RL in linear Markov decision processes. In *International*
645 *Conference on Machine Learning*, pp. 22430–22456. PMLR, 2022.
- 646
- 647 Chuang Wang, Jonathan Mattingly, and Yue M Lu. Scaling limit: Exact and tractable analysis of
648 online learning algorithms with applications to regularized regression and pca. *arXiv preprint*
649 *arXiv:1712.04332*, 2017.
- 650 Justin Whitehouse, Zhiwei Steven Wu, and Aaditya Ramdas. Time-uniform self-normalized con-
651 centration for vector-valued processes. *arXiv preprint arXiv:2310.09100*, 2023.

-
- 648 Lei Wu. Learning a single neuron for non-monotonic activation functions. In *International conference*
649 *on artificial intelligence and statistics*, pp. 4178–4197. PMLR, 2022.
650
- 651 Gilad Yehudai and Shamir Ohad. Learning a single neuron with gradient methods. In *Conference*
652 *on Learning Theory*, pp. 3756–3786. PMLR, 2020.
- 653 Henry Zhu, Abhishek Gupta, Aravind Rajeswaran, Sergey Levine, and Vikash Kumar. Dexterous
654 manipulation with deep reinforcement learning: Efficient, general, and low-cost. In *2019 Inter-*
655 *national Conference on Robotics and Automation (ICRA)*, pp. 3651–3657. IEEE, 2019.
- 656
- 657 Aaron Zweig, Loucas Pillaud-Vivien, and Joan Bruna. On single-index models beyond gaussian
658 data. *Advances in Neural Information Processing Systems*, 36:10210–10222, 2023.
- 659
- 660
- 661
- 662
- 663
- 664
- 665
- 666
- 667
- 668
- 669
- 670
- 671
- 672
- 673
- 674
- 675
- 676
- 677
- 678
- 679
- 680
- 681
- 682
- 683
- 684
- 685
- 686
- 687
- 688
- 689
- 690
- 691
- 692
- 693
- 694
- 695
- 696
- 697
- 698
- 699
- 700
- 701

A PROOFS OF MAIN LEMMAS

A.1 PROOF OF COROLLARY 1

By Theorem 1 and 2, it remains to show that both the initialization cost $\tilde{O}((f(1/\sqrt{d}) - f(0))^{-2})$ and the burn-in cost $\tilde{O}(d^2)$ under Assumption 2 are dominated by the integral.

For the initialization cost, we have

$$\begin{aligned} \frac{1}{(f(\frac{1}{\sqrt{d}}) - f(0))^2} &\stackrel{(a)}{\leq} \frac{1}{(f(\frac{1}{\sqrt{d}}) - f(\frac{1}{2\sqrt{d}}))^2} = \frac{4d}{\left(2\sqrt{d} \int_{1/(2\sqrt{d})}^{1/\sqrt{d}} f'(m) dm\right)^2} \\ &\stackrel{(b)}{\leq} 4d \cdot 2\sqrt{d} \int_{1/(2\sqrt{d})}^{1/\sqrt{d}} \frac{1}{f'(m)^2} dm \leq 16d^2 \int_{1/(2\sqrt{d})}^{1/\sqrt{d}} \frac{m}{f'(m)^2} dm, \end{aligned}$$

where (a) follows from the monotonicity of f , and (b) applies Jensen's inequality.

For the burn-in cost $\tilde{O}(d^2)$ under Assumption 2, we simply note that $f'(x) \leq \gamma_2$ when $x \in [1-\gamma_0, 1]$ by Assumption 1, so that

$$d^2 \int_{1-\gamma_0}^{1-\gamma_0/4} \frac{m}{f'(m)^2} dm \geq d^2 \cdot \frac{3\gamma_0}{4} \frac{1-\gamma_0}{\gamma_2^2} = \Omega(d^2).$$

These complete the proof.

A.2 PROOF OF LEMMA 1

Observe that

$$\begin{aligned} \mathbb{E}[\theta_{t+1/2} | \mathcal{F}_t] &= \mathbb{E}[\theta_t - \eta_t \sigma_t [(f(\langle a_t, \theta_t \rangle) - f(\langle a_t, \theta^* \rangle)) - N_t) f'(\langle a_t, \theta_t \rangle)] \cdot Z_t | \mathcal{F}_t] \\ &= \theta_t - \eta_t \sigma_t \mathbb{E}[(f(\langle a_t, \theta_t \rangle) - f(\langle a_t, \theta^* \rangle)) f'(\langle a_t, \theta_t \rangle)] \cdot Z_t | \mathcal{F}_t]. \end{aligned}$$

Recall that $a_t = \sqrt{1 - \sigma_t^2} \theta_t + \sigma_t Z_t$ in Equation (2). Since $Z_t \perp \theta_t$ almost surely, $\langle a_t, \theta_t \rangle = \sqrt{1 - \sigma_t^2}$. Taking an inner product with θ^* on both sides,

$$\mathbb{E}[m_{t+1/2} | \mathcal{F}_t] - m_t = \eta_t \sigma_t f'(\sqrt{1 - \sigma_t^2}) \cdot \mathbb{E}\left[f\left(\sqrt{1 - \sigma_t^2} \langle \theta_t, \theta^* \rangle + \sigma_t \langle Z_t, \theta^* \rangle\right) \langle Z_t, \theta^* \rangle \middle| \mathcal{F}_t\right].$$

Since $Z_t \sim \text{Unif}(\{x \in \mathbb{S}^{d-1} : x \perp \theta_t\})$, the random variable $(1 - m_t^2)^{-1/2} \langle Z_t, \theta^* \rangle$ is distributed as the one-dimensional marginal of a uniform random vector on \mathbb{S}^{d-2} ; denote by X a random variable following this distribution. Consequently, for

$$g(x) = f\left(\sqrt{1 - \sigma_t^2} \langle \theta_t, \theta^* \rangle + \sigma_t \sqrt{1 - m_t^2} x\right),$$

an application of the spherical Stein's lemma (cf. Lemma 10) gives

$$\begin{aligned} \mathbb{E}[m_{t+1/2} | \mathcal{F}_t] - m_t &= \eta_t \sigma_t f'(\sqrt{1 - \sigma_t^2}) \sqrt{1 - m_t^2} \cdot \mathbb{E}[g(X) X | \mathcal{F}_t] \\ &= \frac{\eta_t \sigma_t}{d-2} f'(\sqrt{1 - \sigma_t^2}) \sqrt{1 - m_t^2} \cdot \mathbb{E}[g'(X) (1 - X^2) | \mathcal{F}_t] \\ &= \frac{\eta_t \sigma_t^2}{d-2} f'(\sqrt{1 - \sigma_t^2}) (1 - m_t^2) \cdot \mathbb{E}\left[f'\left(\sqrt{1 - \sigma_t^2} m_t + \sigma_t \sqrt{1 - m_t^2} X\right) (1 - X^2) \middle| \mathcal{F}_t\right]. \end{aligned}$$

This is the desired identity. For the other inequalities, under Assumption 2 and $m_t \geq 0$, for

$$h(x) = f'\left(\sqrt{1 - \sigma_t^2} \langle \theta_t, \theta^* \rangle + \sigma_t \sqrt{1 - m_t^2} x\right) \geq 0,$$

756 we have

$$\begin{aligned}
758 \quad \mathbb{E}[h(X)(1-X^2)] &\geq \mathbb{E}[h(X)(1-X^2)\mathbb{1}(X \geq 0)] \\
759 \quad &\geq c_0 \mathbb{E}[(1-X^2)\mathbb{1}(X \geq 0)] = c_0 \cdot \frac{d-2}{2(d-1)} = \Omega(1)
760
\end{aligned}$$

761 for $d \geq 3$. Under Assumption 3 and $m_t \geq 0$, we then write

$$\begin{aligned}
763 \quad \mathbb{E}[h(X)(1-X^2)] &\geq \mathbb{E}[h(X)(1-X^2)\mathbb{1}(X \geq 0)] \\
764 \quad &\geq f'(\sqrt{1-\sigma_t^2}m_t) \cdot \mathbb{E}[(1-X^2)\mathbb{1}(X \geq 0)] \\
765 \quad &= \Omega(f'(\sqrt{1-\sigma_t^2}m_t)).
766
\end{aligned}$$

768 A.3 PROOF OF LEMMA 2

769 By definition,

$$771 \quad m_{t+1/2} - m_t = \eta_t \sigma_t (f(\langle a_t, \theta^* \rangle) + N_t - f(\langle a_t, \theta_t \rangle)) f'(\langle a_t, \theta_t \rangle) \cdot \langle Z_t, \theta^* \rangle$$

772 Define two new random variables:

$$\begin{aligned}
774 \quad \xi^{(1)} &= \eta_t \sigma_t (f(\langle a_t, \theta^* \rangle) - f(\langle a_t, \theta_t \rangle)) f'(\langle a_t, \theta_t \rangle) \cdot \langle Z_t, \theta^* \rangle, \\
775 \quad \xi^{(2)} &= \eta_t \sigma_t N_t f'(\langle a_t, \theta_t \rangle) \cdot \langle Z_t, \theta^* \rangle,
776
\end{aligned}$$

777 such that $m_{t+1/2} - m_t = \xi^{(1)} + \xi^{(2)}$. We will show that each of these random variables is subex-
778 ponential with a bounded Ψ_1 -Orlicz norm.

779 For $\xi^{(1)}$, note that $|f(\langle a_t, \theta^* \rangle) - f(\langle a_t, \theta_t \rangle)| \leq 2\|f\|_\infty$ and $\langle a_t, \theta_t \rangle = \sqrt{1-\sigma_t^2}$. In addition,

$$781 \quad \langle Z_t, \theta^* \rangle \stackrel{d}{=} \sqrt{1-m_t^2} X,$$

783 where X follows the one-dimensional marginal of a uniform random vector on \mathbb{S}^{d-2} . By Lemma 11,
784 it holds that $\|X\|_{\Psi_2} \leq \|\mathcal{N}(0, d^{-1})\|_{\Psi_2} = O(d^{-1/2})$. Therefore,

$$786 \quad \|\xi^{(1)}\|_{\Psi_1} \stackrel{(a)}{=} O(\|\xi^{(1)}\|_{\Psi_2}) = O\left(\frac{\eta_t \sigma_t f'(\sqrt{1-\sigma_t^2}) \sqrt{1-m_t^2}}{\sqrt{d}}\right),$$

788 where (a) follows from (Vershynin, 2018, Remark 2.8.8).

790 For $\xi^{(2)}$, note that $\|N_t\|_{\Psi_2} \leq 1$ by the 1-subGaussian assumption on the noise. Therefore, by
791 independence of Z_t and N_t , (Vershynin, 2018, Lemma 2.8.6) gives

$$793 \quad \|\xi^{(2)}\|_{\Psi_1} \leq \eta_t \sigma_t f'(\sqrt{1-\sigma_t^2}) \|N_t\|_{\Psi_2} \|\langle Z_t, \theta^* \rangle\|_{\Psi_2} = O\left(\frac{\eta_t \sigma_t f'(\sqrt{1-\sigma_t^2}) \sqrt{1-m_t^2}}{\sqrt{d}}\right).$$

796 Finally, the triangle inequality of the Ψ_1 norm gives

$$798 \quad \|m_{t+1/2} - \mathbb{E}[m_{t+1/2} | \mathcal{F}_t]\|_{\Psi_1} \leq \|\xi^{(1)}\|_{\Psi_1} + \|\xi^{(2)}\|_{\Psi_1} = O\left(\frac{\eta_t \sigma_t f'(\sqrt{1-\sigma_t^2}) \sqrt{1-m_t^2}}{\sqrt{d}}\right).$$

801 A.4 PROOF OF LEMMA 3

803 For notational simplicity we write $S_t := S_t^{t_0, \beta}$. By Lemma 2,

$$804 \quad \log \mathbb{E}[\exp(\lambda(S_{t+1} - S_t)) | \mathcal{F}_t] \leq C \beta^{2(t-t_0)} K_t^2 \lambda^2, \quad \text{for all } |\lambda| \leq \frac{1}{C \beta^{t-t_0} K_t}.$$

807 Here $C > 0$ is a universal constant. We show that $K_t \leq 1$ almost surely. In fact, $f'(\sqrt{1-\sigma_t^2}) \leq \gamma_2$
808 by Assumption 1 when $\sigma_t^2 \leq \gamma_0$, and

$$809 \quad K_t \leq C_{\text{se}} \eta_t \gamma_2 \leq 1$$

810 by the choice of η_t . Consequently, for $V_t = C \sum_{s=t_0}^{t-1} \beta^{2(s-t_0)} K_s^2$, $\lambda_{\max} := \frac{1}{2C}$, and

$$812 \quad 813 \quad \psi(\lambda) = \frac{\lambda^2}{1 - \lambda/\lambda_{\max}}, \quad \lambda \in [0, \lambda_{\max}),$$

814 it holds that

$$815 \quad 816 \quad \mathbb{E}[\exp(\lambda S_{t+1} - \psi(\lambda)V_{t+1})|\mathcal{F}_t] \leq \exp(\lambda S_t - \psi(\lambda)V_t), \quad \lambda \in [0, \lambda_{\max}), \beta^{t-t_0} \leq 2.$$

817 Therefore, the conditions of Lemma 9 are fulfilled, and the claimed upper tail of S_t follows from
818 choosing $\omega = 1$. Replacing S_t by $-S_t$ in the above analysis gives the lower tail of S_t .

819 A.5 PROOF OF LEMMA 4

820 Since $\theta_t \perp Z_t$, the iterate $\theta_{t+1/2}$ in Equation (3) satisfies

$$821 \quad 822 \quad \|\theta_{t+1/2}\|^2 = 1 + \eta_t^2 \sigma_t^2 f'(\sqrt{1 - \sigma_t^2})^2 (f(\langle \theta_t, a_t \rangle) - r_t)^2.$$

823 Therefore, $\|\theta_{t+1/2}\| \geq 1$, it is clear that

$$824 \quad 825 \quad m_{t+1} = \frac{m_{t+1/2}}{\|\theta_{t+1/2}\|} = m_{t+1/2} - \frac{m_{t+1/2}}{\|\theta_{t+1/2}\|} (\|\theta_{t+1/2}\| - 1) \\ 826 \quad 827 \quad \geq m_{t+1/2} - \frac{1}{2} \eta_t^2 \sigma_t^2 f'(\sqrt{1 - \sigma_t^2})^2 (f(\langle \theta_t, a_t \rangle) - r_t)^2,$$

830 using $\sqrt{1+x} - 1 \leq \frac{x}{2}$ for $x \geq 0$, and $|m_{t+1/2}|/\|\theta_{t+1/2}\| \leq 1$. The first statement now follows
831 from the sub-Gaussian concentration of r_t , which implies $(f(\langle \theta_t, a_t \rangle) - r_t)^2 = O(\log(1/\delta))$ with
832 probability at least $1 - \delta$.

833 For the second statement, $m_{t+1} \leq m_{t+1/2}$ follows from $\|\theta_{t+1/2}\| \geq 1$. The other direction follows
834 from the same high-probability upper bound of $\|\theta_{t+1/2}\| - 1$, and the simple inequality $\frac{1}{1+x} \geq 1 - x$
835 for $x \geq 0$.

836 A.6 PROOF OF LEMMA 6

837 As we showed in the proof of Lemma 5, we will show by induction on s that with probability at least
838 $1 - \Delta\delta/3$, $m_s \geq m_t - \frac{\varepsilon}{4}$ for all $s \in [t, t + \Delta]$. The base case $s = t$ is ensured by the assumption
839 $m_t \geq 1 - \varepsilon$. For the inductive step, the induction hypothesis implies that $m_t, \dots, m_{s-1} \geq 1 - 2\varepsilon$.
840 Then

$$841 \quad 842 \quad m_s - m_t = \sum_{r=t}^{s-1} \left[\underbrace{(\mathbb{E}[m_{r+1/2}|\mathcal{F}_r] - m_r)}_{\geq 0 \text{ by Lemma 1}} + \underbrace{(m_{r+1/2} - \mathbb{E}[m_{r+1/2}|\mathcal{F}_r])}_{=:A_r} + \underbrace{(m_{r+1} - m_{r+1/2})}_{=:B_r} \right].$$

843 Thanks to the inductive hypothesis, $K_r = O(\eta \frac{\varepsilon}{\sqrt{d}})$ for all $r \in [t, s-1]$ in Lemma 2, so Lemma 3
844 (with $t_0 = t$, $\beta = 1$) gives $|\sum_{r=t}^{s-1} A_r| = O(\eta \sqrt{\frac{\Delta\varepsilon^2}{d}} \log(\frac{\Delta}{\delta})) = O(\sqrt{\eta\varepsilon} \log(\frac{\Delta}{\delta})) < \frac{\varepsilon}{8}$ with prob-
845 ability $1 - \frac{\delta}{6}$, by choosing $c > 0$ small enough. Similarly, $|\sum_{r=t}^{s-1} B_r| = O(\Delta\eta^2\varepsilon \log(\frac{\Delta}{\delta})) =$
846 $O(d\eta \log(\frac{\Delta}{\delta})) < \frac{\varepsilon}{8}$ with probability $1 - \frac{\delta}{6}$, by Lemma 4 and choosing $c > 0$ small enough. This
847 implies that $m_s \geq m_t - \frac{\varepsilon}{4}$ with probability $1 - \frac{\delta}{3}$, completing the induction.

848 Conditioned on the event $m_s \geq 1 - 2\varepsilon$ for all $s \in [t, t + \Delta]$, we distinguish into two regimes in this
849 epoch. Let $T_0 \geq t$ be the stopping time where $m_s > 1 - \varepsilon/4$ for the first time.

850 **Regime I:** $t \leq s < T_0$. In this regime $m_s \in [1 - 2\varepsilon, 1 - \varepsilon/4]$. We show that $T_0 \leq t + \Delta$ with
851 probability $1 - \Delta\delta/3$. If $T_0 > t + \Delta$, using the same high-probability bounds, we have

$$852 \quad 853 \quad m_{t+\Delta} - m_t \geq \sum_{s=t}^{t+\Delta-1} (\mathbb{E}[m_{s+1/2}|\mathcal{F}_s] - m_s) - \frac{\varepsilon}{4}$$

854 with probability $1 - \Delta\delta/3$. By Lemma 1 with $1 - m_s^2 = \Omega(\varepsilon)$ and $\sqrt{1 - \sigma_s^2} m_s \geq 1 - \gamma_0$ for $s < T_0$,
855 the total drift is $\Omega(\frac{\Delta\eta\varepsilon^2}{d}) = \Omega(C\varepsilon)$. Therefore, for a large absolute constant $C > 0$, we would have
856 $m_{t+\Delta} \geq 1 - \varepsilon/2$, a contradiction to the assumption $T_0 > t + \Delta$.

864 **Regime II:** $s \geq T_0$. As shown above, this regime is non-empty with high probability. The same
865 induction starting from $s = T_0$ shows that, with probability $1 - \Delta\delta/3$, $m_s \geq m_{T_0} - \varepsilon/4$ holds for
866 all $s \in [T_0, t + \Delta]$. In particular, choosing $s = t + \Delta$ gives the desired result $m_{t+\Delta} \geq 1 - \varepsilon/2$.
867

868 Finally, to lower bound $\langle \theta^*, a_s \rangle$ during this epoch, we simply note that

$$\begin{aligned} 869 \langle \theta^*, a_s \rangle &= \sqrt{1 - \sigma_s^2} m_s + \sigma_s \langle \theta^*, Z_s \rangle \\ 870 &= \sqrt{1 - \sigma_s^2} m_s + \sigma_s \langle \theta^* - m_s \theta_s, Z_s \rangle \\ 871 &\geq \sqrt{1 - \sigma_s^2} m_s - \sigma_s \|\theta^* - m_s \theta_s\| \\ 872 &= \sqrt{1 - \sigma_s^2} m_s - \sigma_s \sqrt{1 - m_s^2}. \\ 873 \end{aligned}$$

874 Under the good event $m_s \geq m_t - \frac{\varepsilon}{4} \geq 1 - \frac{3\varepsilon}{2}$, by $\sigma_s \equiv \sqrt{\varepsilon}$ we have $\langle \theta^*, a_s \rangle \geq 1 - 4\varepsilon$, as desired.
875

877 A.7 PROOF OF LEMMA 8

878 Let

$$879 \beta := 1 - C_{\text{nm}} \gamma_2^2 \eta^2 \sigma^2 \log \left(\frac{4\Delta}{\delta} \right), \quad (4)$$

880 with C_{nm} given in Lemma 4. By the choice of η , when the constant $c > 0$ is small enough, we have
881 $\beta \in (1/2, 1)$. In addition, let

$$882 T_0 = \min \left\{ s \geq t : m_s \geq \left(1 - \frac{\gamma_0}{8} \right) \sqrt{\frac{k+1}{d}} \right\} \quad (5)$$

883 be the stopping time when the correlation m_s first hits a given threshold. Unlike the other proofs,
884 the event $T_0 \leq t + \Delta$ no longer occurs with high probability, and our proof will discuss both cases.
885

886 **Case I:** $T_0 > t + \Delta$. Define the following event:
887

$$888 \mathcal{E}_s := \left\{ m_s \geq \bar{m}_k - \frac{\gamma_0}{d} + \frac{c' \eta f'(\bar{m}_k)}{d} (s - t) \right\}, \quad (6)$$

889 where $c' > 0$ is a small absolute constant (to be chosen later) independent of c . We will prove by
890 induction that

$$891 \mathbb{P}((\cup_{r=t}^s \mathcal{E}_r^c) \cap \{T_0 > t + \Delta\}) \leq (s - t) \frac{\delta}{2}, \quad \text{for all } s = t, t + 1, \dots, t + \Delta. \quad (7)$$

892 The base case follows from the assumption $m_t \geq \bar{m}_k$, so that $\mathbb{P}(\mathcal{E}_t^c) = 0$. For the inductive step,
893 suppose that Equation (7) holds for $s - 1$. Since $\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(A^c \cap B)$, it suffices to prove
894 that

$$895 \mathbb{P}(\mathcal{E}_s^c \cap (\cap_{r=t}^{s-1} \mathcal{E}_r) \cap \{T_0 > t + \Delta\}) \leq \frac{\delta}{2}. \quad (8)$$

896 To this end, we introduce some additional events. First, applying Lemma 3 with $t_0 = t$ and $\beta^{-1} \leq 2$
897 in Equation (4) gives

$$898 \mathbb{P}(\mathcal{E}_{s,1}) := \mathbb{P} \left(\left| \sum_{r=t}^s \frac{m_{r+1/2} - \mathbb{E}[m_{r+1/2} | \mathcal{F}_r]}{\beta^{r-t}} \right| \leq C \eta \sqrt{\frac{\Delta}{d} \log \left(\frac{d}{\delta} \right)} \right) \geq 1 - \frac{\delta}{4\Delta}, \quad (9)$$

899 for some absolute constant $C > 0$. To see Equation (9), note that

$$900 \beta^{-\Delta} = \exp(O((1 - \beta)\Delta)) = \exp \left(O \left(\eta^2 \Delta \log \frac{d}{\delta} \right) \right) = \exp \left(O \left(\frac{cC}{\eta d} \right) \right) = 1 + \frac{o_c(1)}{d}, \quad (10)$$

901 so that the condition $\beta^{-\Delta} \leq 2$ holds for small $c > 0$, and $\sum_{r=t}^s \beta^{-2(r-t)} = O(s - t + 1) = O(\Delta)$.
902 In addition, let $\mathcal{E}_{s,2}$ be the good event that the lower bound in Lemma 4 holds for m_{s+1} , with $\delta/(4\Delta)$
903 in place of δ .

918 Note that $\mathcal{E}_r \cap \mathcal{E}_{r,2} \cap \{T_0 > t + \Delta\}$ implies that
919

$$\begin{aligned} 920 \quad m_{r+1} &\geq \beta m_{r+1/2} \\ 921 \quad &= \beta (m_{r+1/2} - \mathbb{E}[m_{r+1/2} | \mathcal{F}_t] + \mathbb{E}[m_{r+1/2} | \mathcal{F}_t] - m_r + m_r) \\ 922 \quad &\geq \beta \left(m_{r+1/2} - \mathbb{E}[m_{r+1/2} | \mathcal{F}_t] + c_1 \frac{\eta f'(\underline{m}_k)}{d} + m_r \right), \\ 924 \end{aligned}$$

925 where $c_1 > 0$ is an absolute constant, and the last step invokes Lemma 1, uses $m_r \leq 1 - \Omega(1)$ since
926 $r \leq t + \Delta < T_0$, and
927

$$928 \quad \sqrt{1 - \sigma^2} m_r \geq \sqrt{1 - \gamma_0} \left(\bar{m}_k - \frac{\gamma_0}{d} \right) \geq (1 - \gamma_0)^2 \sqrt{\frac{k}{d}} = \underline{m}_k \\ 929$$

930 by Equation (6) and the definitions of \bar{m}_k , \underline{m}_k . Summing over $r = t, \dots, s-1$, the event $\cap_{r=t}^{s-1} (\mathcal{E}_r \cap \mathcal{E}_{r,2}) \cap \{T_0 > t + \Delta\}$ implies that
931
932

$$933 \quad m_s \geq \beta^{s-t} \left(m_t + \sum_{r=t}^{s-1} \frac{m_{r+1/2} - \mathbb{E}[m_{r+1/2} | \mathcal{F}_t]}{\beta^{r-t}} \right) + c_1 \frac{\eta f'(\underline{m}_k)}{d} \sum_{r=t}^{s-1} \beta^{r+1-t}. \\ 934 \\ 935$$

936 In view of Equation (9) and Equation (10), a further intersection with \mathcal{E}_{s-1} implies that
937

$$\begin{aligned} 938 \quad m_s &\geq \left(1 - \frac{o_c(1)}{d} \right) \bar{m}_k - C \eta \sqrt{\frac{\Delta}{d}} \log \left(\frac{d}{\delta} \right) + \frac{c' \eta f'(\underline{m}_k)}{d} (s-t) \\ 939 \quad &= \left(1 - \frac{o_c(1)}{d} \right) \bar{m}_k - O \left(\frac{cC}{d} \right) + \frac{c' \eta f'(\underline{m}_k)}{d} (s-t) \\ 940 \quad &\geq \bar{m}_k - \frac{\gamma_0}{d} + \frac{c' \eta f'(\underline{m}_k)}{d} (s-t) \\ 941 \\ 942 \\ 943 \\ 944 \end{aligned}$$

945 for $c > 0$ small enough; this is precisely the event \mathcal{E}_s . In other words, we have shown that
946

$$947 \quad \mathcal{E}_s^c \cap (\cap_{r=t}^{s-1} (\mathcal{E}_r \cap \mathcal{E}_{r,1} \cap \mathcal{E}_{r,2})) \cap \{T_0 > t + \Delta\} = \emptyset. \quad (11)$$

948 By Equation (11), we have

$$\begin{aligned} 949 \quad \mathbb{P}(\mathcal{E}_s^c \cap (\cap_{r=t}^{s-1} \mathcal{E}_r) \cap \{T_0 > t + \Delta\}) \\ 950 \quad \leq \mathbb{P}(\cup_{r=t}^{s-1} \mathcal{E}_{r,1}^c) + \mathbb{P}((\cup_{r=t}^{s-1} \mathcal{E}_{r,2}^c) \cap (\cap_{r=t}^{s-1} (\mathcal{E}_r \cap \mathcal{E}_{r,1}))) \cap \{T_0 > t + \Delta\}). \\ 951 \end{aligned}$$

952 By Equation (9) and the union bound, the first probability is at most $\frac{\delta}{4}$. For the second probability, the
953 same program above shows that $(\cap_{i=t}^{s-1} \mathcal{E}_{i,2}) \cap (\cap_{i=t}^r (\mathcal{E}_r \cap \mathcal{E}_{r,1})) \cap \{T_0 > t + \Delta\}$ implies $m_{r+1/2} \geq 0$, which is the prerequisite of Lemma 4. Therefore, the conditional probability of $\mathcal{E}_{r,2}$ is at least
954 $1 - \frac{\delta}{4\Delta}$, and by a union bound the second probability is at most $\frac{\delta}{4}$. This proves Equation (8) and
955 completes the induction.
956

957 Finally, note that $\mathcal{E}_{t+\Delta}$ implies that
958

$$\begin{aligned} 959 \quad m_{t+\Delta} &\geq \bar{m}_k - \frac{\gamma_0}{d} + \frac{c' \eta f'(\underline{m}_k)}{d} \Delta \\ 960 \quad &= \bar{m}_k - \frac{\gamma_0}{d} + c' C (\underline{m}_{k+1} - \underline{m}_k) \geq \bar{m}_{k+1}, \\ 961 \\ 962 \\ 963 \end{aligned}$$

964 by choosing $C > 0$ large enough. Therefore, Equation (7) with $s = t + \Delta$ implies that
965

$$966 \quad \mathbb{P}(\{m_{t+\Delta} < \bar{m}_{k+1}\} \cap \{T_0 > t + \Delta\}) \leq \frac{\Delta \delta}{2}. \quad (12)$$

967 **Case II:** $T_0 \leq t + \Delta$. We apply our usual program to this case: if $T_0 \leq t + \Delta$, then
968

$$969 \quad m_{t+\Delta} - m_{T_0} = \sum_{s=T_0}^{t+\Delta-1} \left[\underbrace{(\mathbb{E}[m_{s+1/2} | \mathcal{F}_s] - m_t)}_{\geq 0 \text{ by Lemma 1}} + \underbrace{(\mathbb{E}[m_{s+1/2} | \mathcal{F}_s] - m_{s+1/2})}_{=: A_s} + \underbrace{(m_{s+1} - m_{s+1/2})}_{=: B_s} \right]. \\ 970 \\ 971$$

972 By Lemma 3, with probability at least $1 - \frac{\Delta\delta}{4}$,

$$974 \quad \left| \sum_{s=T_0}^{t+\Delta-1} A_s \right| = O\left(\eta \sqrt{\frac{\Delta}{d}} \log\left(\frac{d}{\delta}\right)\right) = O\left(\frac{cC}{d}\right).$$

977 By Lemma 4, with probability at least $1 - \frac{\Delta\delta}{4}$,

$$979 \quad \sum_{s=T_0}^{t+\Delta-1} |B_s| = O\left(\Delta \cdot \eta^2 \log\left(\frac{d}{\delta}\right)\right) = O\left(\frac{cC}{d}\right).$$

982 Therefore, conditioned on $T_0 \leq t + \Delta$, with probability at least $1 - \frac{\Delta\delta}{2}$,

$$984 \quad m_{t+\Delta} \geq \left(1 - \frac{\gamma_0}{8}\right) \sqrt{\frac{k}{d}} - O\left(\frac{cC}{d}\right) \geq \left(1 - \frac{\gamma_0}{4}\right) \sqrt{\frac{k}{d}} = \bar{m}_k$$

986 for a small enough constant $c > 0$. In other words,

$$988 \quad \mathbb{P}(\{m_{t+\Delta} < \bar{m}_{k+1}\} \cap \{T_0 \leq t + \Delta\}) \leq \frac{\Delta\delta}{2}. \quad (13)$$

991 Finally, a combination of Equation (12) and Equation (13) gives $\mathbb{P}(m_{t+\Delta} < \bar{m}_{k+1}) \leq \Delta\delta$, which
992 is the desired result.

993 A.8 PROOF OF PROPOSITION 1

995 Let T_0 be the first time $t \geq 1$ such that $m_t \geq 0.1$. If $T_0 > T$, the target claim $\max_{t \in [T]} m_t \leq 0.2$
996 is clearly true. Hence in the sequel we condition on the event $T_0 \leq T$. In addition, by Gaussian tail
997 bounds, we have $\max_{t \in [T]} |r_t| = O(\sqrt{\log(T/\delta)})$ with probability at least $1 - \delta/4$. By Equation (3),
998 we then have a deterministic inequality

$$1000 \quad m_{T_0-1/2} \leq m_{T_0-1} + C\eta_{T_0-1} \sqrt{\log(T/\delta)} \leq m_{T_0-1} + 0.05 \leq 0.15,$$

1001 by assumption of $\eta_t \leq \frac{c}{\log(T/\delta)}$ for a sufficiently small constant $c > 0$, and the definition of T_0 that
1002 $m_{T_0-1} \leq 0.1$. By Lemma 4, this implies that $m_{T_0} \leq 0.15$.

1004 In the sequel, we start from $m_{T_0} \in [0.1, 0.15]$, and for notational simplicity we redefine m_{T_0} to be
1005 our starting point, i.e. $T_0 = 1$. Next we consider the time interval $[1, T_1]$ with

$$1007 \quad T_1 = \min \left\{ t \geq 1 : \sum_{s \leq t} \frac{\eta_s^2 \sigma_s^2}{d} \geq \frac{c_1}{\log^2(T/\delta)} \right\},$$

1009 for some absolute constant $c_1 > 0$ to be chosen later. We prove the following claims.

1011 **Claim I:** $\max_{t \in [T_1]} m_t \leq 0.2$ with probability at least $1 - \delta T_1 / (4T)$. To prove this claim, we
1012 first show that when $m_t \leq 0.2$, then

$$1014 \quad \mathbb{E}[m_{t+1/2} | \mathcal{F}_t] \leq m_t. \quad (14)$$

1015 Indeed, by Lemma 1,

$$1017 \quad \mathbb{E}[m_{t+1/2} | \mathcal{F}_t] - m_t \\ 1018 \quad = \frac{\eta_t \sigma_t^2}{d-2} f'(\sqrt{1 - \sigma_t^2})(1 - m_t^2) \cdot \mathbb{E} \left[f' \left(\sqrt{1 - \sigma_t^2} m_t + \sigma_t \sqrt{1 - m_t^2} X \right) (1 - X^2) \middle| \mathcal{F}_t \right].$$

1021 Since $\sigma_t \leq 0.1$, $m_t \leq 0.2$, and $|X| \leq 1$ almost surely, we have

$$1022 \quad \sqrt{1 - \sigma_t^2} m_t + \sigma_t \sqrt{1 - m_t^2} X \leq m_t + \sigma_t \leq 0.3 < \frac{1}{3}.$$

1024 Since $f'(m) \leq 0$ for all $m \leq 1/3$ in our construction, and $f'(\sqrt{1 - \sigma_t^2}) > 0$, we obtain Equation (14).

1026 Next, without loss of generality we assume that $m_t \geq 0$ for all $t \in [T_1]$, since a negative m_t only
1027 makes the target claim simpler. For every $t \in [T_1]$,

$$1029 \quad m_t - m_1 = \sum_{r=1}^{t-1} \left[\underbrace{(\mathbb{E}[m_{r+1/2}|\mathcal{F}_r] - m_r)}_{\leq 0 \text{ by Equation (14)}} + \underbrace{(m_{r+1/2} - \mathbb{E}[m_{r+1/2}|\mathcal{F}_r])}_{=:A_r} + \underbrace{(m_{r+1} - m_{r+1/2})}_{\leq 0 \text{ by Lemma 4}} \right].$$

1032 By Lemma 3 with $\beta = 1$, we get

$$1034 \quad \left| \sum_{r=1}^{t-1} A_r \right| \leq C \log \left(\frac{T}{\delta} \right) \sqrt{\sum_{r=1}^{t-1} \frac{\sigma_r^2 \eta_r^2}{d}}$$

1037 with probability at least $1 - \delta/(4T)$, for some absolute constant $C > 0$. By the definition of T_1 , we
1038 obtain $|\sum_{r=1}^{t-1} A_r| \leq 0.05$ for a sufficiently small $c_1 > 0$. Therefore, $m_t \leq m_1 + 0.05 \leq 0.2$ with
1039 probability at least $1 - \delta/(4T)$, and an induction on t with a union bound gives the target claim.

1041 **Claim II:** $\min_{t \in [T_1]} m_t \leq 0.1$ with probability at least $1 - \delta T_1/(4T)$. In the sequel, we condition
1042 on the good event in Claim I. Let T_2 be the first time $t \geq 1$ such that $m_t \leq 0.1$; note that it is possible
1043 to have $T_2 > T_1$ or even $T_2 = \infty$. We first show that if $m_t \geq 0.1$, then

$$1045 \quad \mathbb{E}[m_{t+1/2}|\mathcal{F}_t] - m_t \leq -\frac{c_2 \eta_t \sigma_t^2}{d} \quad (15)$$

1047 for some absolute constant $c_2 > 0$. Indeed, for $\sigma_t \leq 0.1$, $m_t \in [0.1, 0.2]$, and $|X| \leq 1$, we have

$$1049 \quad 0 \leq \sqrt{0.99}m_t - \sqrt{0.99}\sigma_t \leq \sqrt{1 - \sigma_t^2}m_t + \sigma_t \sqrt{1 - m_t^2}X \leq m_t + \sigma_t < \frac{1}{3}.$$

1051 Since $f'(m) = -1$ for all $m \in [0, 1/3]$ in our construction, Equation (15) follows from Lemma 1.

1052 Next, for every $t \leq \min\{T_2, T_1\}$, we write

$$1054 \quad m_t - m_1 = \sum_{r=1}^{t-1} \left[\underbrace{(\mathbb{E}[m_{r+1/2}|\mathcal{F}_r] - m_r)}_{\leq -\frac{c_2 \eta_t \sigma_t^2}{d} \text{ by Equation (15)}} + \underbrace{(m_{r+1/2} - \mathbb{E}[m_{r+1/2}|\mathcal{F}_r])}_{=:A_r} + \underbrace{(m_{r+1} - m_{r+1/2})}_{\leq 0 \text{ by Lemma 4}} \right].$$

1058 Similar to Claim I, we have $|\sum_{r=1}^{t-1} A_r| \leq 0.05$ with probability at least $1 - \delta/(4T)$. On the other
1059 hand, the total drift is

$$1060 \quad \sum_{r=1}^{T_1-1} (\mathbb{E}[m_{r+1/2}|\mathcal{F}_r] - m_r) \leq -\frac{c_2}{d} \sum_{r=1}^{T_1-1} \eta_r \sigma_r^2 \stackrel{(a)}{\leq} -\frac{c_2 \log^2(T/\delta)}{cd} \sum_{r=1}^{T_1-1} \eta_r^2 \sigma_r^2 \stackrel{(b)}{\leq} -\frac{c_1 c_2}{2c},$$

1063 where (a) uses the upper bound of η_r , and (b) uses the definition of T_1 . By choosing $c > 0$ small
1064 enough, the total drift can be made smaller than -0.1 , so that if $T_2 > T_1$, then $m_{T_1} \leq m_1 - 0.1 +$
1065 $0.05 \leq 0.1$, which in turn means that $T_2 \leq T_1$, a contradiction. Therefore, with probability at least
1066 $1 - \delta T_1/(4T)$, we have $T_2 \leq T_1$, or equivalently $\min_{t \in [T_1]} m_t \leq 0.1$.

1068 Finally, it is clear that a repeated application of Claim I and II implies Proposition 1: starting from
1069 the first time T_0 with $m_{T_0} \geq 0.1$, the above claims show that with high probability, future alignment
1070 m_t will fall below 0.1 before it rises above 0.2. Once m_t falls below 0.1, we repeat the entire process
1071 again and wait for the next time it falls below 0.1. Since the failure probability at each step of the
1072 analysis is at most δ/T , a union bound gives the total failure probability of δ .

1073 B AUXILIARY RESULTS

1076 Below we state a self-normalized concentration inequality for martingales (Whitehouse et al., 2023,
1077 Theorem 3.1) adapted to our setting.

1078 **Definition 1** (CGF-like function). A function $\psi : [0, \lambda_{\max}] \rightarrow \mathbb{R}_{\geq 0}$ is said to be CGF-like if it is (a)
1079 twice continuously-differentiable on its domain, (b) strictly convex, (c) satisfies $\psi(0) = \psi'(0) = 0$,
and (d) $\psi''(0) > 0$.

1080 **Definition 2** (Sub- ψ). Let $\psi : [0, \lambda_{\max}] \rightarrow \mathbb{R}_{\geq 0}$ be a CGF-like function. Let $\{S_t\}_{t \geq 0}$ and $\{V_t\}_{t \geq 0}$
1081 be respectively \mathbb{R} -valued and $\mathbb{R}_{\geq 0}$ -valued processes adapted to some filtration $\{\mathcal{F}_t\}_{t \geq 0}$. We say that
1082 $\{S_t, V_t\}_{t \geq 0}$ is sub- ψ if for every $\lambda \in [0, \lambda_{\max}]$,

1083
$$M_t^\lambda := \exp(\lambda S_t - \psi(\lambda) V_t) \leq L_t^\lambda,$$

1084 where $\{L_t^\lambda\}_{t \geq 0}$ is a non-negative supermartingale adapted to $\{\mathcal{F}_t\}_{t \geq 0}$.

1085 The following result is a corollary of (Whitehouse et al., 2023, Theorem 3.1) with the choice $h(k) =$
1086 $(1+k)^2$ for $k \geq 1$.

1087 **Lemma 9** (Self-normalized concentration inequality). Suppose $\{S_t, V_t\}_{t \geq 0}$ is a real-valued sub- ψ
1088 process for $\psi : [0, \lambda_{\max}] \rightarrow \mathbb{R}_{\geq 0}$ satisfying

1089
$$\psi(\lambda) = \frac{\lambda^2}{1 - \lambda/\lambda_{\max}}$$

1090 on its domain. Let $\delta \in (0, 1)$ denote the error probability. Define the function $\ell : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}$ by

1091
$$\ell_\omega(v) = 2 \log(1 + \log(v\omega \vee 1)) + \log\left(\frac{1}{\delta}\right),$$

1092 then there exists a universal constant $C > 0$ such that,

1093
$$\Pr\left(\exists t \geq 1 : S_t \geq C \left(\sqrt{(V_t \vee \omega^{-1}) \ell_\omega(V_t)} + \lambda_{\max}^{-1} \ell_\omega(V_t)\right)\right) \leq \delta.$$

1094 *Proof.* By simple algebra, the convex conjugate ψ^* of ψ satisfies $(\psi^*)^{-1}(u) = 2\sqrt{u} + \lambda_{\max}^{-1}u$. The
1095 rest follows from (Whitehouse et al., 2023, Theorem 3.1). \square

1096 **Lemma 10** (Spherical Stein's Lemma). Suppose $Z \sim \text{Unif}(\mathbb{S}^{d-1})$ and consider a fixed $\alpha \in \mathbb{R}^d$ and
1097 let $X = \langle \alpha, Z \rangle$. For any bounded function f ,

1098
$$\mathbb{E}[Xf(X)] = \frac{1}{d-1} \mathbb{E}[f'(X)(1-X^2)].$$

1100 *Proof.* The density of X is given by

1101
$$P_d(x) \triangleq \frac{2(1-x^2)^{\frac{d-1}{2}-1}}{\text{Beta}(\frac{1}{2}, \frac{d-1}{2})} \mathbb{I}(|x| \leq 1).$$

1102 Consequently,

1103
$$\begin{aligned} \mathbb{E}[Xf(X)] &= \int_{-1}^1 xf(x) \cdot \frac{2(1-x^2)^{\frac{d-1}{2}-1}}{\text{Beta}(\frac{1}{2}, \frac{d-1}{2})} dx \\ 1104 &\stackrel{(a)}{=} \frac{2}{d-1} \int_{-1}^1 f'(x) \cdot \frac{(1-x^2)^{\frac{d-1}{2}}}{\text{Beta}(\frac{1}{2}, \frac{d-1}{2})} dx \\ 1105 &= \frac{1}{d-1} \int_{-1}^1 f'(x)(1-x^2) \cdot \frac{2(1-x^2)^{\frac{d-1}{2}-1}}{\text{Beta}(\frac{1}{2}, \frac{d-1}{2})} dx \\ 1106 &= \frac{1}{d-1} \mathbb{E}[f'(X)(1-X^2)], \end{aligned}$$

1107 where (a) follows from integration by parts. \square

1108 **Lemma 11.** Suppose $X \sim \mathcal{N}(0, I/d)$ and $X' \sim \text{Unif}(\mathbb{S}^{d-1})$. For any fixed $\alpha \in \mathbb{R}^d$, $\langle \alpha, X \rangle^2$
1109 dominates $\langle \alpha, X' \rangle^2$ in the convex order. Namely, for every convex function $g : \mathbb{R} \rightarrow \mathbb{R}$,

1110
$$\mathbb{E}[g(\langle \alpha, X' \rangle^2)] \leq \mathbb{E}[g(\langle \alpha, X \rangle^2)].$$

1134 *Proof.* Observe that X follows the same distribution as NX' , where N and X' are independent,
1135 and N is a scaled chi-squared random variable such that $\mathbb{E}[N^2] = 1$. Therefore,
1136

$$\begin{aligned} 1137 \quad \mathbb{E}[g(\langle \alpha, X \rangle^2)] &= \mathbb{E}[g(N^2 \langle \alpha, X' \rangle^2)] \\ 1138 &= \mathbb{E}[\mathbb{E}[g(N^2 \langle \alpha, X' \rangle^2) | X']] \\ 1139 &\geq \mathbb{E}[g(\mathbb{E}[N^2] \langle \alpha, X' \rangle^2)] \\ 1140 &= \mathbb{E}[g(\langle \alpha, X' \rangle^2)]. \\ 1141 \end{aligned}$$

1142 \square
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187