
OPT-AMSGrad: An Optimistic Acceleration of AMSGrad for Nonconvex Optimization

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In this paper, we propose a new variant of AMSGrad [31], a popular adaptive gra-
2 dent based optimization algorithm widely used in training deep neural networks.
3 Our algorithm adds prior knowledge about the sequence of consecutive mini-batch
4 gradients leveraging an underlying structure which makes the gradients sequen-
5 tially predictable. By exploiting the predictability and ideas from Optimistic On-
6 line Learning, the proposed algorithm can accelerate the convergence and increase
7 sample efficiency. After establishing a tighter upper bound under some convexity
8 conditions on the regret, we offer a complimentary view of our algorithm which
9 generalizes the offline and stochastic versions of nonconvex optimization. In the
10 nonconvex case, we establish a $\mathcal{O}\left(\sqrt{d/T} + d/T\right)$ non-asymptotic bound inde-
11 pendent of the initialization of the method. We illustrate the practical speedup on
12 several deep learning models through numerical experiments.

13 1 Introduction

14 Deep learning models have been successful in several applications, from robotics (e.g. [20]), com-
15 puter vision (e.g. [17, 14]), reinforcement learning (e.g. [25]), to natural language processing (e.g.
16 [15]). With the sheer size of modern data sets and the dimension of neural networks, speeding up
17 training is of utmost importance. To do so, several algorithms have been proposed in recent years,
18 such as AMSGRAD [31], ADAM [18], RMSPROP [35], ADADELTA [41], and NADAM [10].

19 All the prevalent algorithms for training deep networks mentioned above combine two ideas: the
20 idea of adaptivity from ADAGRAD [11, 23] and the idea of momentum from NESTEROV’S METHOD
21 [27] or HEAVY BALL method [28]. ADAGRAD is an online learning algorithm that works well
22 compared to the standard online gradient descent when the gradient is sparse. Its update has a
23 notable feature: it leverages an anisotropic learning rate depending on the magnitude of gradient in
24 each dimension which helps in exploiting the geometry of data. On the other hand, NESTEROV’S
25 METHOD or HEAVY BALL Method [28] is an accelerated optimization algorithm which update not
26 only depends on the current iterate and current gradient but also depends on the past gradients (i.e.
27 momentum). State-of-the-art algorithms like AMSGRAD [31] and ADAM [18] leverage these ideas
28 to accelerate the training of nonconvex objective functions such as deep neural networks losses.

29 In this paper, we propose an algorithm that goes further than the hybrid of the adaptivity and mo-
30 mentum approach. Our algorithm is inspired by OPTIMISTIC ONLINE LEARNING [7, 29, 34, 1, 24],
31 which assumes that a good *predictable process* of the gradient of the loss function in each round of
32 online learning is available, and plays an action by exploiting these predictors. By exploiting this
33 (possibly) arbitrary process, algorithms in OPTIMISTIC ONLINE LEARNING enjoy smaller regret
34 than the ones without. We combine the OPTIMISTIC ONLINE LEARNING idea with the adaptivity
35 and the momentum ideas to design a new algorithm — OPT-AMSGRAD.

A single work along that direction stands out. [8] develops OPTIMISTIC-ADAM in their paper leveraging optimistic online mirror descent [30]. Yet, OPTIMISTIC-ADAM is specifically designed to optimize two-player games (e.g. GANs [14]). GANs is a two-player zero-sum game. There have been some related works in OPTIMISTIC ONLINE LEARNING like [7, 30, 34] showing that if both players use some kinds of OPTIMISTIC-update, then accelerating the convergence to the equilibrium of the game is possible. [8] was inspired by these related works and showed that OPTIMISTIC-MIRROR-DESCENT can avoid the cycle behavior in a bilinear zero-sum game, which accelerates the convergence. In contrast, in this paper, the proposed algorithm is designed to accelerate nonconvex optimization (e.g. empirical risk minimization). To the best of our knowledge, this is the first work exploring towards this direction and bridging the unfilled *theoretical* gap at the crossroads of online learning and stochastic optimization. The contributions of this paper are as follows:

- We derive an optimistic variant of AMSGRAD borrowing techniques from online learning procedures. Our method relies on (I) the addition of *prior knowledge* in the sequence of the model parameter estimations alleviating a predictable process able to provide good guesses of gradients of the loss functions through the iterations and (II) the construction of a *double update* algorithm done sequentially. We interpret this two-projection step as the learning of both an underlying scheme which makes the gradients sequentially predictable and the global parameter learning.
- We focus on the *theoretical* justifications of our method by establishing novel *non-asymptotic* and *global* convergence rates in both the convex and nonconvex case. Based on both *convex regret minimization* and *nonconvex stochastic optimization* views, we prove, respectively, that our algorithm suffers regret of $\mathcal{O}(\sqrt{\sum_{t=1}^T \|g_t - m_t\|_{\psi_{t-1}}^2})$ and achieves a rate of convergence $\mathcal{O}(\sqrt{d/T} + d/T)$.

The proposed algorithm not only adapts to the informative dimensions, exhibits momentum, but also exploits a good guess of the next gradient to facilitate acceleration. Besides the global analysis of OPT-AMSGRAD, we conduct experiments and show that the proposed algorithm not only accelerates convergence of loss function, but also leads to better empirical generalization performance.

Section 2 is devoted to introductory notions on online learning for regret minimization and adaptive learning methods for nonconvex stochastic optimization. We introduce in Section 3 our new algorithm called OPT-AMSGRAD and provide a comprehensive global analysis in both *convex/online* and *nonconvex/offline* settings in Section 4. We illustrate the benefits of our method on several finite-sum nonconvex optimization problem in Section 5. The supplementary material of this paper is devoted to the proofs of our theoretical results.

Notations: We follow the notations in related adaptive optimization papers [18, 31]. For any vector $u, v \in \mathbb{R}^d$, u/v represents element-wise division, u^2 represents element-wise square, \sqrt{u} represents element-wise square-root. We denote $g_{1:T}[i]$ as the sum of the i_{th} element of T vectors $g_1, g_2, \dots, g_T \in \mathbb{R}^d$.

2 Preliminaries

Optimistic Online learning. The standard setup of ONLINE LEARNING is that, in each round t , an online learner selects an action $w_t \in \Theta \subseteq \mathbb{R}^d$, then the learner observes $\ell_t(\cdot)$ and suffers loss $\ell_t(w_t)$ after the action is committed. The goal of the learner is to minimize the regret,

$$\mathcal{R}_T(\{w_t\}) := \sum_{t=1}^T \ell_t(w_t) - \sum_{t=1}^T \ell_t(w^*),$$

which is the cumulative loss of the learner minus the cumulative loss of some benchmark $w^* \in \Theta$. The idea of OPTIMISTIC ONLINE LEARNING (e.g. [7, 29, 34, 1]) is as follows. In each round t , the learner exploits a guess $m_t(\cdot)$ of the gradient $\nabla \ell_t(\cdot)$ of the loss function to choose an action w_t ¹. Consider the FOLLOW-THE-REGULARIZED-LEADER (FTRL, [16]) online learning algorithm

¹Imagine that if the learner would had been known $\nabla \ell_t(\cdot)$ (i.e., exact guess) before committing its action, then it would exploit the knowledge to determine its action and consequently minimizes the regret.

78 which update reads

$$w_t = \arg \min_{w \in \Theta} \langle w, L_{t-1} \rangle + \frac{1}{\eta} \mathbf{R}(w), \quad (1)$$

79 where η is a parameter, $\mathbf{R}(\cdot)$ is a 1-strongly convex function with respect to a norm ($\|\cdot\|$) on
 80 the constraint set Θ , and $L_{t-1} := \sum_{s=1}^{t-1} g_s$ is the cumulative sum of gradient vectors of the loss
 81 functions up to round $t-1$. It has been shown that FTRL has regret at most $O(\sqrt{\sum_{t=1}^T \|g_t\|_*^2})$. The
 82 update of its optimistic variant, noted OPTIMISTIC-FTRL and developed in [34] reads

$$w_t = \arg \min_{w \in \Theta} \langle w, L_{t-1} + m_t \rangle + \frac{1}{\eta} \mathbf{R}(w), \quad (2)$$

83 where $\{m_t\}_{t>0}$ is a predictable process incorporating (possibly arbitrarily) knowledge about the
 84 sequence of gradients $\{g_t := \nabla \ell_t(w_t)\}_{t>0}$. Under the assumption that loss functions are convex,
 85 the regret of OPTIMISTIC-FTRL is at most $O(\sqrt{\sum_{t=1}^T \|g_t - m_t\|_*^2})$.

86 *Remark:* Note that the usual worst-case bound is preserved even when the predictors $\{m_t\}_{t>0}$ do
 87 not predict well the gradients. Indeed, if we take the example of OPTIMISTIC-FTRL, the bound
 88 reads $\sqrt{\sum_{t=1}^T \|g_t - m_t\|_*^2} \leq 2 \max_{w \in \Theta} \|\nabla \ell_t(w)\| \sqrt{T}$ which is equal to the usual bounds up to a factor
 89 2. Yet, when the predictions are well constructed, the regret will be lower. We will have a similar
 90 argument when we compare OPT-AMSGRAD and AMSGRAD.

91 We emphasize in Section 3 the importance of leveraging a good guess m_t for updating w_t in order
 92 to get a fast convergence rate (or equivalently, small regret) and present Section 5 a simple, yet
 93 effective, predictable process $\{m_t\}_{t>0}$ leading to empirical acceleration.

94 **Adaptive optimization methods.** Recently, adaptive optimization has been popular in
 95 various deep learning applications due to their superior empirical performance. ADAM
 96 [18] is a very popular adaptive algorithm for training deep neural networks. It com-
 97 bines the momentum idea [28] with the idea of ADAGRAD [11], which has different learn-
 98 ing rates for different dimensions, adaptive to the learning process. More specifically,
 99 the learning rate of ADAGRAD in iteration t for dimension j is proportional to the in-
 100 verse of $\sqrt{\sum_{s=1}^t g_s[j]^2}$, where $g_s[j]$ is the j -th element of the gradient vector g_s at time s .
 101 This adaptive learning rate helps accelerating
 102 the convergence when the gradient vector is
 103 sparse [11] but, when applying ADAGRAD to
 104 train deep networks, it is observed that the
 105 learning rate might decay too fast [18]. There-
 106 fore, [18] proposes ADAM that uses a moving
 107 average of gradients divided by the square root
 108 of the second moment of the moving average
 109 (element-wise multiplication), for updating the
 110 model parameter w . A variant, called AMS-
 111 GRAD and detailed in Algorithm 1, has been
 112 developed in [31] to fix ADAM failures at some online convex optimization problems. The dif-
 113 ference between ADAM and AMSGRAD lies in line 7 of Algorithm 1. AMSGRAD [31] adds the
 114 max operation to guarantee a non-increasing learning rate, $\frac{\eta_t}{\sqrt{\hat{v}_t}}$, which helps for the convergence (i.e.
 115 average regret $\mathcal{R}_T/T \rightarrow 0$).

Algorithm 1 AMSGRAD [31]

- 1: Required: parameter β_1, β_2 , and η_t .
 - 2: Init: $w_1 \in \Theta \subseteq \mathbb{R}^d$ and $v_0 = \epsilon \mathbf{1} \in \mathbb{R}^d$.
 - 3: **for** $t = 1$ to T **do**
 - 4: Get mini-batch stochastic gradient g_t at w_t .
 - 5: $\theta_t = \beta_1 \theta_{t-1} + (1 - \beta_1) g_t$.
 - 6: $v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$.
 - 7: $\hat{v}_t = \max(\hat{v}_{t-1}, v_t)$.
 - 8: $w_{t+1} = w_t - \eta_t \frac{\theta_t}{\sqrt{\hat{v}_t}}$. (element-wise division)
 - 9: **end for**
-

116 3 OPT-AMSGRAD Algorithm

117 We formulate in this section the proposed optimistic acceleration of AMSGrad, noted OPT-
 118 AMSGRAD, and detailed in Algorithm 2. It combines the idea of adaptive optimization with opti-
 119 mistic learning. At each iteration, the learner computes a gradient vector $g_t := \nabla \ell_t(w_t)$ at w_t (line
 120 4), then it maintains an exponential moving average of $\theta_t \in \mathbb{R}^d$ (line 5) and $v_t \in \mathbb{R}^d$ (line 6), which
 121 is followed by the max operation to get $\hat{v}_t \in \mathbb{R}^d$ (line 7). The learner first updates an auxiliary
 122 variable $\tilde{w}_{t+1} \in \Theta$ (line 8) and then computes the next model parameter (line 9).

Algorithm 2 OPT-AMSGRAD

1: Required: parameter $\beta_1, \beta_2, \epsilon$, and η_t .
 2: Init: $w_1 = w_{-1/2} \in \Theta \subseteq \mathbb{R}^d$ and $v_0 = \epsilon 1 \in \mathbb{R}^d$.
 3: **for** $t = 1$ to T **do**
 4: Get mini-batch stochastic gradient g_t at w_t .
 5: $\theta_t = \beta_1 \theta_{t-1} + (1 - \beta_1) g_t$.
 6: $v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$.
 7: $\hat{v}_t = \max(\hat{v}_{t-1}, v_t)$.
 8: $\tilde{w}_{t+1} = \tilde{w}_t - \eta_t \frac{\theta_t}{\sqrt{\hat{v}_t}}$.
 9: $w_{t+1} = \tilde{w}_{t+1} - \eta_t \frac{h_{t+1}}{\sqrt{\hat{v}_t}}$,
 where $h_{t+1} := \beta_1 \theta_{t-1} + (1 - \beta_1) m_{t+1}$ and m_{t+1}
 is the guess of g_{t+1} .
 10: **end for**

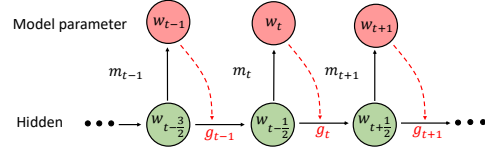


Figure 1: OPT-AMSGRAD UNDERLYING STRUCTURE.

124 Observe that the proposed algorithm does not reduce to AMSGRAD when $m_t = 0$, contrary to the
 125 optimistic variant of FTRL. Furthermore, combining line 8 and line 9 yields the following single
 126 update $w_{t+1} = \tilde{w}_t - \eta_t \frac{\theta_t}{\sqrt{\hat{v}_t}} - \eta_t \frac{h_{t+1}}{\sqrt{\hat{v}_t}}$.

127 Compared to AMSGRAD, the algorithm is characterized by a *two-level* update that interlinks some
 128 *auxiliary state* \tilde{w}_t and the model parameter state, w_t , similarly to the OPTIMISTIC MIRROR DE-
 129 SCENT algorithm developed in [29]. It leverages the auxiliary variable (hidden model) to update
 130 and commit w_{t+1} , which exploits the guess m_{t+1} , see Figure 1 for a schematic illustration. In the
 131 following analysis, we show that the interleaving actually leads to some cancellation in the regret
 132 bound. Such two-levels method where the guess m_t is equal to the last known gradient g_{t-1} has
 133 been exhibited recently in [7]. The gradient prediction procedure plays naturally an important role
 134 and will be tackled Section 5. The proposed OPT-AMSGRAD inherits three properties:

- 135 • Adaptive learning rate of each dimension as ADAGRAD [11]. (line 6, line 8 and line 9)
- 136 • Exponential moving average of the past gradients as NESTEROV'S METHOD [27] and the
 137 HEAVY-BALL method [28]. (line 5)
- 138 • Optimistic update that exploits *prior knowledge* of the next gradient vector as optimistic
 139 online learning algorithms [7, 29, 34]. (line 9)

140 The first property helps for acceleration when the gradient has a sparse structure. The second one
 141 is from the well-recognized idea of momentum which can also help for acceleration. The last one
 142 can actually lead to an acceleration when the prediction of the next gradient is good as mentioned
 143 above when introducing the regret bound for the OPTIMISTIC-FTRL algorithm. The latter property
 144 will be elaborated in the following subsection in which we provide the theoretical analysis of OPT-
 145 AMSGRAD. Observe that the proposed algorithm does not reduce to AMSGRAD when $m_t = 0$.

146 4 Global Convergence of OPT-AMSGRAD

147 For conciseness, we place all the proofs of the following results in the supplementary material.

148 **Notations.** We denote the Mahalanobis norm $\|\cdot\|_H := \sqrt{\langle \cdot, H \cdot \rangle}$ for some positive semidefinite
 149 (PSD) matrix H . We let $\psi_t(x) := \langle x, \text{diag}\{\hat{v}_t\}^{1/2} x \rangle$ for a PSD matrix $H_t^{1/2} := \text{diag}\{\hat{v}_t\}^{1/2}$, where
 150 $\text{diag}\{\hat{v}_t\}$ represents the diagonal matrix which i_{th} diagonal element is $\hat{v}_t[i]$ defined in Algorithm 2.
 151 We define its corresponding Mahalanobis norm $\|\cdot\|_{\psi_t} := \sqrt{\langle \cdot, \text{diag}\{\hat{v}_t\}^{1/2} \cdot \rangle}$, where we abuse
 152 the notation ψ_t to represent the PSD matrix $H_t^{1/2} := \text{diag}\{\hat{v}_t\}^{1/2}$. Note that $\psi_t(\cdot)$ is 1-strongly
 153 convex with respect to the norm $\|\cdot\|_{\psi_t}$. Namely, $\psi_t(\cdot)$ satisfies $\psi_t(u) \geq \psi_t(v) + \langle \psi_t(v), u - v \rangle + \frac{1}{2} \|u - v\|_{\psi_t}^2$ for any point $(u, v) \in \Theta^2$. A consequence of 1-strongly convexity of $\psi_t(\cdot)$ is
 154 that $B_{\psi_t}(u, v) \geq \frac{1}{2} \|u - v\|_{\psi_t}^2$, where the Bregman divergence $B_{\psi_t}(u, v)$ is defined as $B_{\psi_t}(u, v) :=$
 155 $\psi_t(u) - \psi_t(v) - \langle \psi_t(v), u - v \rangle$ with $\psi_t(\cdot)$ as the distance generating function. We can also define
 156 the corresponding dual norm $\|\cdot\|_{\psi_t^*} := \sqrt{\langle \cdot, \text{diag}\{\hat{v}_t\}^{-1/2} \cdot \rangle}$.
 157

158 4.1 Convex Regret Analysis

159 We prove the following result regarding the regret in the convex optimization setting. That is, we
 160 assume that the loss functions $\{\ell_t\}_{t>0}$ are convex. We also assume that Θ has bounded diameter
 161 D_∞ , which is a standard assumption in previous works [31, 18] on adaptive methods. It is necessary
 162 in regret analysis since if the boundedness assumption is lifted, one might construct a scenario such
 163 that the benchmark is $w^* = \infty$ and the learner's regret is infinite.

164 **Theorem 1.** *Suppose the learner incurs a sequence of convex loss functions $\{\ell_t(\cdot)\}$. Then,*
 165 *OPTIMISTIC-AMSGRAD (Algorithm 2) has regret*

$$\mathcal{R}_T \leq \frac{B_{\psi_1}(w^*, \tilde{w}_1)}{\eta_1} + \sum_{t=1}^T \frac{\eta_t}{2} \|g_t - \tilde{m}_t\|_{\psi_{t-1}^*}^2 + \frac{D_\infty^2}{\eta_{\min}} \sum_{i=1}^d \hat{v}_T^{1/2}[i] + D_\infty^2 \beta_1^2 \sum_{t=1}^T \|g_t - \theta_{t-1}\|_{\psi_{t-1}^*}, \quad (3)$$

166 where $\tilde{m}_{t+1} = \beta_1 \theta_{t-1} + (1 - \beta_1) m_{t+1}$, $g_t := \nabla \ell_t(w_t)$, $\eta_{\min} := \min_t \eta_t$ and D_∞^2 is the diameter of
 167 the bounded set Θ . The result holds for any benchmark $w^* \in \Theta$ and any step size sequence $\{\eta_t\}_{t>0}$.

168 **Corollary 1.** *Suppose $\beta_1 = 0$ and $\{v_t\}_{t>0}$ is an increasing monotone sequence, then we obtain the*
 169 *following regret bound for any $w^* \in \Theta$ and sequence $\{\eta_t = \eta/\sqrt{t}\}_{t>0}$:*

$$\mathcal{R}_T \leq \frac{B_{\psi_1}(w^*, \tilde{w}_1)}{\eta_1} + \frac{\eta \sqrt{1 + \log T}}{\sqrt{1 - \beta_2}} \sum_{i=1}^d \|(g - m)_{1:T}[i]\|_2 + \frac{D_\infty^2}{\eta_{\min}} \sum_{i=1}^d \left[(1 - \beta_2) \sum_{s=1}^T \beta_2^{T-s} g_s[i]^2 \right]^{1/2},$$

170 where $g_t := \nabla \ell_t(w_t)$ and $\eta_{\min} := \min_t \eta_t$.

171 We can compare the bound of Corollary 1 with that of AMSGRAD [31] with $\eta_t = \eta/\sqrt{t}$:

$$\mathcal{R}_T \leq \frac{\eta \sqrt{1 + \log T}}{\sqrt{1 - \beta_2}} \sum_{i=1}^d \|g_{1:T}[i]\|_2 + \frac{\sqrt{T}}{2\eta} D_\infty^2 \sum_{i=1}^d \hat{v}_T[i]^2. \quad (4)$$

172 For convex regret minimization, the results above yields that the learner suffers regret of
 173 $\mathcal{O}(\sqrt{\sum_{t=1}^T \|g_t - m_t\|_{\psi_{t-1}}^2})$ with an access to an arbitrary predictable process $\{m_t\}_{t>0}$ of the mini-
 174 batch gradient. The better the predictors, the lower the regret, see the second term in Corollary 1
 175 compared to the first term in eq:boundAMS. One can thus wonder how the learner can build those
 176 good gradients predictions $\{m_t\}_{t>0}$. Is this process can be learned through the iterations? Those
 177 questions are interesting research questions and will not be dealt in this paper for the sake of page
 178 limit. Though, for implementation purposes, we derive a simple, yet effective, gradient prediction
 179 algorithm, see Algorithm 3 in Section 5 to embed to our new OPT-AMSGRAD.

180 4.2 Nonconvex Analysis (Finite-Time Upper Bound)

181 In this section, we discuss the offline and stochastic non-convex optimization properties of our online
 182 framework. In the stochastic optimization literature, the problem we are tackling reads as follows:

$$\min_{w \in \Theta} f(w) := \mathbb{E}[f(w, \xi)], \quad (5)$$

183 where ξ is some random noise and only noisy versions of the objective function are accessible in
 184 this work. The objective function $f(w)$ is (potentially) nonconvex and has Lipschitz gradients.

185 Set the terminating iteration number, $T \in \{0, \dots, T_{\max} - 1\}$, as a discrete r.v. with:

$$P(T = \ell) = \frac{\eta_\ell}{\sum_{j=0}^{T_{\max}-1} \eta_j}, \quad (6)$$

186 where T_{\max} is the maximum number of iteration. The random termination number (6) is inspired
 187 by [13] which enables to derive novel non-asymptotic convergence results to stationary point for
 188 non-convex optimization. We make the following mild assumptions necessary to our analysis:

189 **H1.** *For any $t > 0$, the estimated weight w_t stays within a ℓ_∞ -ball. There exists a constant $W > 0$*
 190 *such that $\|w_t\| \leq W$ almost surely.*

191 **H2.** The function $f(w)$ is L -smooth (has L -Lipschitz gradients) w.r.t. the parameter w . There exist
 192 some constant $L > 0$ such that for $(w, \vartheta) \in \Theta^2$:

$$f(w) - f(\vartheta) - \nabla f(\vartheta)^\top (w - \vartheta) \leq \frac{L}{2} \|w - \vartheta\|^2 .$$

193 We assume that the optimistic guess m_t at iteration k and the true gradient g_t are correlated:

194 **H3.** There exists a constant $a \in \mathbb{R}$ such that for any $t > 0$, $\langle m_t | g_t \rangle \leq a \|g_t\|^2$.

195 Classically in nonconvex optimization [13] we make an assumption on the magnitude of the gradient:

196 **H4.** There exists a constant $M > 0$ such that for any w and ξ , it holds $\|\nabla f(w, \xi)\| < M$.

197 We begin with some auxiliary Lemmas important for the analysis. The first one ensures bounded
 198 norms of various quantities of interests (resulting from the classical stochastic gradient boundedness
 199 assumption):

200 **Lemma 1.** Assume assumption H4, then the quantities defined in Algorithm 2 satisfy for any $w \in \Theta$
 201 and $t > 0$, $\|\nabla f(w_t)\| < M$, $\|\theta_t\| < M$ and $\|\hat{v}_t\| < M^2$.

202 Then, following [39] and their study of the SGD with Momentum we denote for any $t > 0$:

$$\bar{w}_t = w_t + \frac{\beta_1}{1 - \beta_1} (w_t - \tilde{w}_{t-1}) = \frac{1}{1 - \beta_1} w_t - \frac{\beta_1}{1 - \beta_1} \tilde{w}_{t-1} , \quad (7)$$

203 and derive an important Lemma:

204 **Lemma 2.** Assume a strictly positive and non increasing sequence of stepsizes $\{\eta_t\}_{t>0}$, $\beta_1 < \beta_2 \in$
 205 $[0, 1)$, then the following holds:

$$\bar{w}_{t+1} - \bar{w}_t \leq \frac{\beta_1}{1 - \beta_1} \tilde{\theta}_{t-1} \left[\eta_{t-1} \hat{v}_{t-1}^{-1/2} - \eta_t \hat{v}_t^{-1/2} \right] - \eta_t \hat{v}_t^{-1/2} \tilde{g}_t ,$$

206 where $\tilde{\theta}_t = \theta_t + \beta_1 \theta_{t-1}$ and $\tilde{g}_t = g_t - \beta_1 m_t + \beta_1 g_{t-1} + m_{t+1}$.

207 **Lemma 3.** Assume H4, a strictly positive and a sequence of constant stepsizes $\{\eta_t\}_{t>0}$, $\beta_1 < \beta_2 \in$
 208 $[0, 1)$, denote $\gamma = \beta_1 / \beta_2$, then the following holds:

$$\sum_{k=1}^{T_{\max}} \eta_t^2 \mathbb{E} \left[\left\| \hat{v}_t^{-1/2} \theta_t \right\|_2^2 \right] \leq \frac{\eta^2 d T_{\max} (1 - \beta_1)}{(1 - \beta_2)(1 - \gamma)} .$$

209 We now formulate the main result of our paper giving a finite-time upper bound of the quantity
 210 $\mathbb{E} [\|\nabla f(w_T)\|^2]$ where T is a random termination number distributed according to 6, see [13].

211 **Theorem 2.** Assume H1-H4, $\beta_1 < \beta_2 \in [0, 1)$ and a sequence of decreasing stepsizes $\{\eta_t\}_{t>0}$, then
 212 the following result holds:

$$\mathbb{E} [\|\nabla f(w_T)\|^2] \leq \tilde{C}_1 \sqrt{\frac{d}{T_{\max}}} + \tilde{C}_2 \frac{1}{T_{\max}} , \quad (8)$$

213 where K is a random termination number distributed according to (6).

214 We remark that the bound for our OPT-AMSGrad method matches the complexity bound of
 215 $\mathcal{O} \left(\sqrt{d/T_{\max}} + 1/T_{\max} \right)$ of [13] for SGD and [43] for AMSGrad method.

216 4.3 Checking H1 for a Deep Neural Network

217 We show in this section that the weights satisfy assumption H1 and stay in a bounded set when
 218 the model we are fitting, using our method, is a fully connected feed forward neural network. The
 219 activation function for this section will be sigmoid function and we add a ℓ_2 regularization.

220 We consider a fully connected feed forward neural network with L layers modeled by the function
 221 $\text{MLN}(w, \xi) : \mathbb{R}^l \rightarrow \mathbb{R}$:

$$\text{MLN}(w, \xi) = \sigma \left(w^{(L)} \sigma \left(w^{(L-1)} \dots \sigma \left(w^{(1)} \xi \right) \right) \right) \quad (9)$$

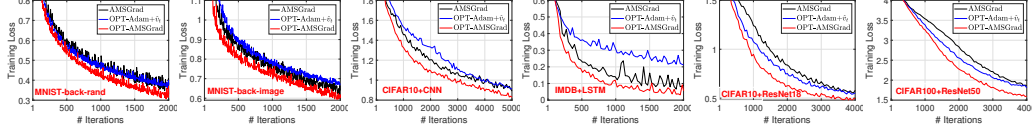


Figure 2: Training loss vs. Number of iterations. The first row are results with fully-connected NN.

where $w = [w^{(1)}, w^{(2)}, \dots, w^{(L)}]$ is the vector of parameters, $\xi \in \mathbb{R}^l$ is the input data and σ is the sigmoid activation function. We assume a l dimension input data and a scalar output for simplicity. The stochastic objective function (5) reads:

$$f(w, \xi) = \mathcal{L}(\text{MLN}(w, \xi), y) + \frac{\lambda}{2} \|w\|^2 \quad (10)$$

where $\mathcal{L}(\cdot, y)$ is the loss function (can be Huber loss or cross entropy), y are the true labels and $\lambda > 0$ is the regularization parameter. For any layer index $\ell \in [1, L]$ we denote the output of layer ℓ by $h^{(\ell)}(w, \xi) = \sigma(w^{(\ell)} \sigma(w^{(\ell-1)} \dots \sigma(w^{(1)} \xi)))$. The following Lemma proves that assumption H1 is satisfied with a feed forward neural net (9):

Lemma 4. *Given the multilayer model (9), assume the boundedness of the input data and of the loss function, i.e., for any $\xi \in \mathbb{R}^l$ and $y \in \mathbb{R}$ there is a constant $T > 0$ such that $\|\xi\| \leq 1$ a.s. and $|\mathcal{L}'(\cdot, y)| \leq T$ where $\mathcal{L}'(\cdot, y)$ denotes its derivative w.r.t. the parameter. Then for each layer $\ell \in [1, L]$, there exist a constant $A_{(\ell)}$ such that $\|w^{(\ell)}\| \leq A_{(\ell)}$*

5 Numerical Experiments

5.1 Gradient Estimation

From the analysis in the previous section, we know that whether OPT-AMSGRAD converges faster than its counterpart depends on how m_t is chosen. In OPTIMISTIC-ONLINE LEARNING, m_t is usually set to $m_t = g_{t-1}$, which means that it uses the previous gradient as a guess of the next one. The choice can accelerate the convergence to equilibrium in some two-player zero-sum games [29, 34, 8], in which each player uses an optimistic online learning algorithm against its opponent.

However, this paper is about solving optimization problems, as in (5), instead of solving zero-sum games. In most classical deep learning tasks, as we will develop in the numerical section, (5) even reads $\min_{w \in \Theta} f(w) = \sum_{i=1}^n f(w, \xi_i)$ for a fixed batch of n samples $\{\xi_i\}_{i=1}^n$. We propose to use the extrapolation algorithm of [32]. Extrapolation studies estimating the limit of sequence using the last few iterates [3]. Some classical works include Anderson acceleration [37], minimal polynomial extrapolation [4], reduced rank extrapolation [12]. These methods typically assume that the sequence $\{x_t\} \in \mathbb{R}^d$ has a linear relation $x_t = A(x_{t-1} - x^*) + x^*$ and $A \in \mathbb{R}^{d \times d}$ is an unknown, not necessarily symmetric, matrix. The goal is to find the fixed point of x^* . [32] relaxes the assumption to certain degrees. It assumes that the sequence $\{x_t\} \in \mathbb{R}^d$ satisfies

$$x_t - x^* = A(x_{t-1} - x^*) + e_t, \quad (11)$$

where e_t is a second order term satisfying $\|e_t\|_2 = O(\|x_{t-1} - x^*\|_2^2)$ and $A \in \mathbb{R}^{d \times d}$ is an unknown matrix. The extrapolation algorithm we used is shown in Algorithm 3. Some theoretical guarantees regarding the distance between the output and x^* are provided in [32]. For our numerical experi-

Algorithm 3 REGULARIZED APPROXIMATE MINIMAL POLYNOMIAL EXTRAPOLATION (RMPE) [32]

- 1: **Input:** sequence $\{x_s \in \mathbb{R}^d\}_{s=0}^{s=r}$, parameter $\lambda > 0$.
 - 2: Compute matrix $U = [x_1 - x_0, \dots, x_r - x_{r-1}] \in \mathbb{R}^{d \times r}$.
 - 3: Obtain z by solving $(U^\top U + \lambda I)z = \mathbf{1}$.
 - 4: Get $c = z / (z^\top \mathbf{1})$.
 - 5: **Output:** $\sum_{i=0}^{r-1} c_i x_i$, the approximation of the fixed point x^* .
-

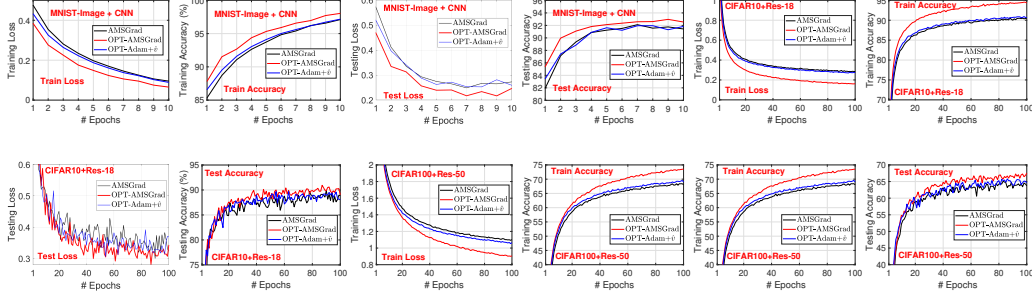


Figure 3: *MNIST-back-image* + CNN, *CIFAR10* + Res-18 and *CIFAR100* + Res-50 . We compare three methods in terms of training (cross-entropy) loss and accuracy, testing loss and accuracy.

ments in the next section, we run OPT-AMSGRAD using Algorithm 3 to get m_t . Specifically, m_t is obtained by (a) calling Algorithm 3 with input being a sequence of some past $r + 1$ gradients, $\{g_t, g_{t-1}, g_{t-2}, \dots, g_{t-r}\}$, where r is a parameter and (b) setting $m_t := \sum_{i=0}^{r-1} c_i g_{t-r+i}$ from the output of Algorithm 3. To see why the output from the extrapolation method may be a reasonable estimation, assume that the update converges to a stationary point (i.e. $g^* := \nabla f(w^*) = 0$ for the underlying function f). Then, we might rewrite (11) as

$$g_t = Ag_{t-1} + O(\|g_{t-1}\|_2^2)u_{t-1}, \quad (12)$$

for some vector u_{t-1} with a unit norm. The equation suggests that the next gradient vector g_t is a linear transform of g_{t-1} plus an error vector that may not be in the span of A whose length is $O(\|g_{t-1}\|_2^2)$. If the algorithm is guaranteed to converge to a stationary point, the magnitude of the error component will eventually go to zero. We remark that the choice of algorithm for gradient prediction is surely not unique. We propose to use the recent result among various related works. Indeed, one can use any method that can provide reasonable guess of gradient in next iteration.

5.2 Classification Experiments

In this section, we provide experiments on classification tasks with various neural network architectures and datasets to demonstrate the effectiveness of OPTIMISTIC-AMSGRAD.

Methods. We consider two baselines. The first one is the original AMSGRAD. The hyperparameters are set to be β_1 and β_2 to be 0.9 and 0.999 respectively, as recommended by [31]. We tune the learning rate η over a fine grid and report the best result. The other competing method is the aforementioned OPTIMISTIC-ADAM+ \hat{v}_t method, see [8]. The key difference is that it uses previous gradient as the gradient prediction of the next iteration. We also report the best result achieved by tuning the step size η for OPTIMISTIC-ADAM+ \hat{v}_t . For OPTIMISTIC-AMSGRAD, we use the same β_1 , β_2 and the best step size η of AMSGRAD for a fair evaluation of the improvement brought by the extra optimistic step. Yet, OPTIMISTIC-AMSGRAD has an additional parameter r that controls the number of previous gradients used for gradient prediction. Fortunately, we observe similar performance of OPTIMISTIC-AMSGRAD with different values of r . Hence, we report $r = 5$ for now when comparing with other baselines. We will address on the choice of r at the end of this section. In all experiments, all the optimization algorithms are initialized at the same point. We report the results averaged over 5 repetitions.

Datasets. Following [31] and [18], we compare different algorithms on *MNIST*, *CIFAR10*, *CIFAR100*, and *IMDB* datasets. For *MNIST*, we use two noisy variants named as 1.65*MNIST-back-rand* and 1.65*MNIST-back-image* from [19]. They both have 12000 training samples and 50000 test samples, where random background is inserted to the original *MNIST* hand written digit images. For *MNIST-back-rand*, each image is inserted with a random background, whose pixel values generated uniformly from 0 to 255, while *MNIST-back-image* takes random patches from a black and white as noisy background. The input dimension is 784 (28×28) and the number of classes is 10. *CIFAR10* and *CIFAR100* are popular computer-vision datasets consisting of 50000 training images and 10000 test images, of size 32×32 . The number of classes are 10 and 100, respectively. The *IMDB* movie

review dataset is a binary classification dataset with 25000 training and testing samples respectively. It is a popular datasets for text classification.

Network architecture. We adopt a multi-layer fully-connected neural network with input layer followed by a hidden layer with 200 nodes, which is connected to another layer with 100 nodes before the output layer. The activation function is ReLU for hidden layers, and softmax for the output layer. This network is tested on *MNIST* variants. Since convolutional networks are popular for image classification tasks, we consider an ALL-CNN architecture proposed by [33], which is constructed with several convolutional blocks and dropout layers. In addition, we also apply residual networks, Resnet-18 and Resnet-50 [17], which have achieved many state-of-the-art results. For the texture *IMDB* dataset, we consider training a Long-Short Term Memory (LSTM) network. The network includes a word embedding layer with 5000 input entries representing most frequent words in the dataset, and each word is embedded into a 32 dimensional space. The output of the embedding layer is passed to 100 LSTM units, which is then connected to 100 fully connected ReLU's before the output layer. For all the models, we use cross-entropy loss. A mini-batch size of 128 is used to compute the stochastic gradients.

Results. Firstly, to illustrate the acceleration effect of OPTIMISTIC-AMSGRAD at early stage, we provide the training loss against number of iterations in Figure 2. We clearly observe that on all datasets, the proposed OPTIMISTIC-AMSGRAD converges faster than the other competing methods, right after the training begins. In other words, we need fewer iterations (samples) to achieve the same training loss. This validates one of the main advantages of OPTIMISTIC-AMSGRAD, which is a higher sample efficiency. We are also curious about the long-term performance and generalization of the proposed method in test phase. In Figure 3, we plot the corresponding results when the model is trained to the state with stable test accuracy. We observe: 1) In the long term, OPTIMISTIC-AMSGRAD algorithm may converge to a better point with smaller objective function value, and 2) In this three applications, the proposed OPTIMISTIC-AMSGRAD also outperforms the competing methods in terms of test accuracy. These are also important benefits of OPTIMISTIC-AMSGRAD.

5.3 Choice of parameter r

Recall that our proposed algorithm has the parameter r that governs the use of past information. Figure 4 compares the performance under different values of $r = 3, 5, 10$ on two datasets. From the result we see that the choice of r does not have significant impact on learning performance. Taking into consideration both quality of gradient prediction and computational cost, it appears that $r = 5$ is a good choice. We remark that empirically, the performance comparison among $r = 3, 5, 10$ is not absolutely consistent (i.e. more means better) in all cases. One possible reason is that for deep neural nets (with highly non-convex loss), using gradient information from too long ago may not be helpful in accurate gradient prediction. Nevertheless, $r = 5$ seems to be good for most applications.

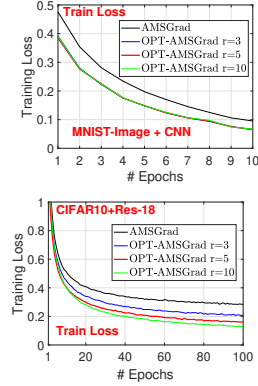


Figure 4: Training loss with different r .

5.4 Some Remarks on the Experiments

Discussion on the iteration cost: We observe that the iteration cost (i.e., actual running time per iteration) of our implementation of OPTIMISTIC-AMSGRAD with $r = 5$ is roughly two times larger than the standard AMSGRAD. When $r = 3$, the cost is roughly 0.7 times longer. Nevertheless, OPTIMISTIC-AMSGRAD may still be beneficial in terms of training efficiency, since fewer iterations are typically needed. For example, in Figure 3, to reach the training loss of AMSGRAD at 100 epochs, the proposed method only needs roughly 20 and 40 epochs, respectively. That said, OPTIMISTIC-AMSGRAD needs 40% and 80% time to achieve same training loss as AMSGRAD, in this two problems.

The computational overhead mostly comes from the gradient extrapolation step. More specifically, recall that the extrapolation step consists of: (a) The step of constructing the linear system ($U^T U$). The cost of this step can be optimized and reduced to $\mathcal{O}(d)$, since the matrix U only changes one column at a time. (b) The step of solving the linear system. The cost of this step is $\mathcal{O}(r^3)$, which is negligible as the linear system is very small (5-by-5 if $r = 5$). (c) The step that outputs an estimated gradient as a weighted average of previous gradients. The cost of this step is $\mathcal{O}(r \times d)$. Thus, the

343 computational overhead is $\mathcal{O}((r+1)d + r^3)$. Yet, we notice that step (a) and (c) is parallelizable,
344 so they can be accelerated in practice.

345 **Memory usage:** Our algorithm needs a storage of past r gradients for each coordinate, in addition
346 to the estimated second moments and the moving average. Though it seems demanding compared
347 to the standard AMSGrad, it is relatively cheap compared to Natural gradient method (e.g., [22]), as
348 Natural gradient method needs to store some matrix inverse.

349 6 Conclusion

350 In this paper, we propose OPTIMISTIC-AMSGRAD, which combines optimistic learning and AMS-
351 GRAD to improve sampling efficiency and accelerate the process of training, in particular for deep
352 neural networks. With a good gradient prediction, the regret can be smaller than that of standard
353 AMSGRAD. Experiments on various deep learning problems demonstrate the effectiveness of the
354 proposed method in improving the training efficiency.

References

- [1] J. Abernethy, K. A. Lai, K. Y. Levy, and J.-K. Wang. Faster rates for convex-concave games. *COLT*, 2018.
- [2] N. Agarwal, B. Bullins, X. Chen, E. Hazan, K. Singh, C. Zhang, and Y. Zhang. Efficient full-matrix adaptive regularization. *ICML*, 2019.
- [3] C. Brezinski and M. R. Zaglia. Extrapolation methods: theory and practice. *Elsevier*, 2013.
- [4] S. Cabay and L. Jackson. A polynomial extrapolation method for finding limits and antilimits of vector sequences. *SIAM Journal on Numerical Analysis*, 1976.
- [5] X. Chen, S. Liu, R. Sun, and M. Hong. On the convergence of a class of adam-type algorithms for non-convex optimization. *ICLR*, 2019.
- [6] Z. Chen, Z. Yuan, J. Yi, B. Zhou, E. Chen, and T. Yang. Universal stagewise learning for non-convex problems with convergence on averaged solutions. *ICLR*, 2019.
- [7] C.-K. Chiang, T. Yang, C.-J. Lee, M. Mahdavi, C.-J. Lu, R. Jin, and S. Zhu. Online optimization with gradual variations. *COLT*, 2012.
- [8] C. Daskalakis, A. Ilyas, V. Syrgkanis, and H. Zeng. Training gans with optimism. *ICLR*, 2018.
- [9] A. Défossez, L. Bottou, F. Bach, and N. Usunier. On the convergence of adam and adagrad. *arXiv preprint arXiv:2003.02395*, 2020.
- [10] T. Dozat. Incorporating nesterov momentum into adam. *ICLR (Workshop Track)*, 2016.
- [11] J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research (JMLR)*, 2011.
- [12] R. Eddy. Extrapolating to the limit of a vector sequence. *Information linkage between applied mathematics and industry*, Elsevier, 1979.
- [13] S. Ghadimi and G. Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.
- [14] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. *NIPS*, 2014.
- [15] A. Graves, A. rahman Mohamed, and G. Hinton. Speech recognition with deep recurrent neural networks. *ICASSP*, 2013.
- [16] E. Hazan. Introduction to online convex optimization. *Foundations and Trends in Optimization*, 2016.
- [17] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *CVPR*, 2016.
- [18] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *ICLR*, 2015.
- [19] H. Larochelle, D. Erhan, A. Courville, J. Bergstra, and Y. Bengio. An empirical evaluation of deep architectures on problems with many factors of variation. *ICML*, 2007.
- [20] S. Levine, C. Finn, T. Darrell, and P. Abbeel. End-to-end training of deep visuomotor policies. *NIPS*, 2017.
- [21] X. Li and F. Orabona. On the convergence of stochastic gradient descent with adaptive step-sizes. *AISTAT*, 2019.
- [22] J. Martens and R. Grosse. Optimizing neural networks with kronecker-factored approximate curvature. *ICML*, 2015.

- [23] H. B. McMahan and M. J. Streeter. Adaptive bound optimization for online convex optimization. *COLT*, 2010.
- [24] P. Mertikopoulos, B. Lecouat, H. Zenati, C.-S. Foo, V. Chandrasekhar, and G. Piliouras. Optimistic mirror descent in saddle-point problems: Going the extra (gradient) mile. *arXiv preprint arXiv:1807.02629*, 2018.
- [25] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller. Playing atari with deep reinforcement learning. *NIPS (Deep Learning Workshop)*, 2013.
- [26] M. Mohri and S. Yang. Accelerating optimization via adaptive prediction. *AISTATS*, 2016.
- [27] Y. Nesterov. Introductory lectures on convex optimization: A basic course. *Springer*, 2004.
- [28] B. T. Polyak. Some methods of speeding up the convergence of iteration methods. *Mathematics and Mathematical Physics*, 1964.
- [29] A. Rakhlin and K. Sridharan. Optimization, learning, and games with predictable sequences. *NIPS*, 2013.
- [30] S. Rakhlin and K. Sridharan. Optimization, learning, and games with predictable sequences. In *Advances in Neural Information Processing Systems*, pages 3066–3074, 2013.
- [31] S. J. Reddi, S. Kale, and S. Kumar. On the convergence of adam and beyond. *ICLR*, 2018.
- [32] D. Scieur, A. d’Aspremont, and F. Bach. Regularized nonlinear acceleration. *NIPS*, 2016.
- [33] J. Springenberg, A. Dosovitskiy, T. Brox, and M. Riedmiller. Striving for simplicity: The all convolutional net. *ICLR*, 2015.
- [34] V. Syrgkanis, A. Agarwal, H. Luo, and R. E. Schapire. Fast convergence of regularized learning in games. *NIPS*, 2015.
- [35] T. Tieleman and G. Hinton. Rmsprop: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural Networks for Machine Learning*, 2012.
- [36] P. Tseng. On accelerated proximal gradient methods for convex-concave optimization. 2008.
- [37] H. F. Walker and P. Ni. Anderson acceleration for fixed-point iterations. *SIAM Journal on Numerical Analysis*, 2011.
- [38] R. Ward, X. Wu, and L. Bottou. Adagrad stepsizes: Sharp convergence over nonconvex landscapes, from any initialization. *ICML*, 2019.
- [39] Y. Yan, T. Yang, Z. Li, Q. Lin, and Y. Yang. A unified analysis of stochastic momentum methods for deep learning. *arXiv preprint arXiv:1808.10396*, 2018.
- [40] M. Zaheer, S. Reddi, D. Sachan, S. Kale, and S. Kumar. Adaptive methods for nonconvex optimization. *NeurIPS*, 2018.
- [41] M. D. Zeiler. Adadelta: An adaptive learning rate method. *arXiv:1212.5701*, 2012.
- [42] D. Zhou, Y. Tang, Z. Yang, Y. Cao, and Q. Gu. On the convergence of adaptive gradient methods for nonconvex optimization. *arXiv:1808.05671*, 2018.
- [43] D. Zhou, Y. Tang, Z. Yang, Y. Cao, and Q. Gu. On the convergence of adaptive gradient methods for nonconvex optimization. *arXiv preprint arXiv:1808.05671*, 2018.
- [44] F. Zou and L. Shen. On the convergence of adagrad with momentum for training deep neural networks. *arXiv:1808.03408*, 2018.

435 A Proof of Theorem 1

436 **Theorem.** Suppose the learner incurs a sequence of convex loss functions $\{\ell_t(\cdot)\}$. Then,
 437 OPTIMISTIC-AMSGRAD (Algorithm 2) has regret

$$\mathcal{R}_T \leq \frac{B_{\psi_1}(w^*, \tilde{w}_1)}{\eta_1} + \sum_{t=1}^T \frac{\eta_t}{2} \|g_t - \tilde{m}_t\|_{\psi_{t-1}^*}^2 + \frac{D_\infty^2}{\eta_{\min}} \sum_{i=1}^d \hat{v}_T^{1/2}[i] + D_\infty^2 \beta_1^2 \sum_{t=1}^T \|g_t - \theta_{t-1}\|_{\psi_{t-1}^*}, \quad (13)$$

438 where $\tilde{m}_{t+1} = \beta_1 \theta_{t-1} + (1 - \beta_1) m_{t+1}$, $g_t := \nabla \ell_t(w_t)$, $\eta_{\min} := \min_t \eta_t$ and D_∞^2 is the diameter of
 439 the bounded set Θ . The result holds for any benchmark $w^* \in \Theta$ and any step size sequence $\{\eta_t\}_{t>0}$.

440 **Proof** Beforehand, note:

$$\begin{aligned} \tilde{g}_t &= \beta_1 \theta_{t-1} + (1 - \beta_1) g_t \\ \tilde{m}_{t+1} &= \beta_1 \theta_{t-1} + (1 - \beta_1) m_{t+1} \end{aligned} \quad (14)$$

441 where we recall that g_t and m_{t+1} are respectively the gradient $\nabla \ell_t(w_t)$ and the predictable guess.
 442 By regret decomposition, we have that

$$\begin{aligned} \text{Regret}_T &:= \sum_{t=1}^T \ell_t(w_t) - \min_{w \in \Theta} \sum_{t=1}^T \ell_t(w) \\ &\leq \sum_{t=1}^T \langle w_t - w^*, \nabla \ell_t(w_t) \rangle \\ &= \sum_{t=1}^T \langle w_t - \tilde{w}_{t+1}, g_t - \tilde{m}_t \rangle + \langle w_t - \tilde{w}_{t+1}, \tilde{m}_t \rangle + \langle \tilde{w}_{t+1} - w^*, \tilde{g}_t \rangle + \langle \tilde{w}_{t+1} - w^*, g_t - \tilde{g}_t \rangle. \end{aligned} \quad (15)$$

443 Recall the notation $\psi_t(x)$ and the Bregman divergence $B_{\psi_t}(u, v)$ we defined in the beginning of this
 444 section. Now we are going to exploit a useful inequality (which appears in e.g., [36]); for any update
 445 of the form $\hat{w} = \arg \min_{w \in \Theta} \langle w, \theta \rangle + B_\psi(w, v)$, it holds that

$$\langle \hat{w} - u, \theta \rangle \leq B_\psi(u, v) - B_\psi(u, \hat{w}) - B_\psi(\hat{w}, v) \quad \text{for any } u \in \Theta. \quad (16)$$

446 For $\beta_1 = 0$, we can rewrite the update on line 8 of (Algorithm 2) as

$$\tilde{w}_{t+1} = \arg \min_{w \in \Theta} \eta_t \langle w, \tilde{g}_t \rangle + B_{\psi_t}(w, \tilde{w}_t), \quad (17)$$

447 By using (16) for (17) with $\hat{w} = \tilde{w}_{t+1}$ (the output of the minimization problem), $u = w^*$ and
 448 $v = \tilde{w}_t$, we have

$$\langle \tilde{w}_{t+1} - w^*, \tilde{g}_t \rangle \leq \frac{1}{\eta_t} [B_{\psi_t}(w^*, \tilde{w}_t) - B_{\psi_t}(w^*, \tilde{w}_{t+1}) - B_{\psi_t}(\tilde{w}_{t+1}, \tilde{w}_t)]. \quad (18)$$

449 We can also rewrite the update on line 9 of (Algorithm 2) at time t as

$$w_{t+1} = \arg \min_{w \in \Theta} \eta_{t+1} \langle w, \tilde{m}_{t+1} \rangle + B_{\psi_t}(w, \tilde{w}_{t+1}). \quad (19)$$

450 and, by using (16) for (19) (written at iteration t), with $\hat{w} = w_{t+1}$ (the output of the minimization
 451 problem), $u = \tilde{w}_{t+1}$ and $v = \tilde{w}_t$, we have

$$\langle w_t - \tilde{w}_{t+1}, \tilde{m}_t \rangle \leq \frac{1}{\eta_t} [B_{\psi_{t-1}}(\tilde{w}_{t+1}, \tilde{w}_t) - B_{\psi_{t-1}}(\tilde{w}_{t+1}, w_t) - B_{\psi_{t-1}}(w_t, \tilde{w}_t)], \quad (20)$$

452 By (15), (18), and (20), we obtain

$$\begin{aligned} \mathcal{R}_T &\stackrel{(15)}{\leq} \sum_{t=1}^T \langle w_t - \tilde{w}_{t+1}, g_t - \tilde{m}_t \rangle + \langle w_t - \tilde{w}_{t+1}, \tilde{m}_t \rangle + \langle \tilde{w}_{t+1} - w^*, \tilde{g}_t \rangle + \langle \tilde{w}_{t+1} - w^*, g_t - \tilde{g}_t \rangle \\ &\stackrel{(18), (20)}{\leq} \sum_{t=1}^T \|w_t - \tilde{w}_{t+1}\|_{\psi_{t-1}} \|g_t - \tilde{m}_t\|_{\psi_{t-1}^*} + \|\tilde{w}_{t+1} - w^*\|_{\psi_{t-1}} \|g_t - \tilde{g}_t\|_{\psi_{t-1}^*} \\ &\quad + \frac{1}{\eta_t} [B_{\psi_{t-1}}(\tilde{w}_{t+1}, \tilde{w}_t) - B_{\psi_{t-1}}(\tilde{w}_{t+1}, w_t) - B_{\psi_{t-1}}(w_t, \tilde{w}_t) + B_{\psi_t}(w^*, \tilde{w}_t) - B_{\psi_t}(w^*, \tilde{w}_{t+1}) - B_{\psi_t}(\tilde{w}_{t+1}, \tilde{w}_t)], \end{aligned} \quad (21)$$

453 which is further bounded by

$$\begin{aligned} \mathcal{R}_T \leq & \sum_{t=1}^T \left\{ \frac{1}{2\eta_t} \|w_t - \tilde{w}_{t+1}\|_{\psi_{t-1}}^2 + \frac{\eta_t}{2} \|g_t - m_t\|_{\psi_{t-1}^*}^2 + \|\tilde{w}_{t+1} - w^*\|_{\psi_{t-1}} \|g_t - \tilde{g}_t\|_{\psi_{t-1}^*} \right. \\ & \left. + \frac{1}{\eta_t} \underbrace{(B_{\psi_{t-1}}(\tilde{w}_{t+1}, \tilde{w}_t) - B_{\psi_t}(\tilde{w}_{t+1}, \tilde{w}_t))}_{A_1} - \frac{1}{2} \|\tilde{w}_{t+1} - w_t\|_{\psi_{t-1}}^2 + \underbrace{(B_{\psi_t}(w^*, \tilde{w}_t) - B_{\psi_t}(w^*, \tilde{w}_{t+1}))}_{A_2} \right\}, \end{aligned} \quad (22)$$

454 where the inequality is due to $\|w_t - \tilde{w}_{t+1}\|_{\psi_{t-1}} \|g_t - m_t\|_{\psi_{t-1}^*} = \inf_{\beta > 0} \frac{1}{2\beta} \|w_t - \tilde{w}_{t+1}\|_{\psi_{t-1}}^2 +$
 455 $\frac{\beta}{2} \|g_t - m_t\|_{\psi_{t-1}^*}^2$ by Young's inequality and the 1-strongly convex of $\psi_{t-1}(\cdot)$ with respect to $\|\cdot\|_{\psi_{t-1}}$
 456 which yields that $B_{\psi_{t-1}}(\tilde{w}_{t+1}, w_t) \geq \frac{1}{2} \|\tilde{w}_{t+1} - w_t\|_{\psi_t}^2 \geq 0$.

457 To proceed, notice that

$$A_1 = B_{\psi_{t-1}}(\tilde{w}_{t+1}, \tilde{w}_t) - B_{\psi_t}(\tilde{w}_{t+1}, \tilde{w}_t) = \langle \tilde{w}_{t+1} - \tilde{w}_t, \text{diag}(\hat{v}_{t-1}^{1/2} - \hat{v}_t^{1/2})(\tilde{w}_{t+1} - \tilde{w}_t) \rangle \leq 0, \quad (23)$$

458 as the sequence $\{\hat{v}_t\}$ is non-decreasing. And that

$$\begin{aligned} A_2 &= B_{\psi_t}(w^*, \tilde{w}_t) - B_{\psi_t}(w^*, \tilde{w}_{t+1}) = \langle w^* - \tilde{w}_{t+1}, \text{diag}(\hat{v}_{t+1}^{1/2} - \hat{v}_t^{1/2})(w^* - \tilde{w}_{t+1}) \rangle \\ &\leq (\max_i (w^*[i] - \tilde{w}_{t+1}[i])^2) \cdot \left(\sum_{i=1}^d \hat{v}_{t+1}^{1/2}[i] - \hat{v}_t^{1/2}[i] \right) \end{aligned} \quad (24)$$

459 Therefore, by (22),(24),(23), we have

$$\begin{aligned} \mathcal{R}_T \leq & \frac{D_\infty^2}{\eta_{\min}} \sum_{i=1}^d \hat{v}_T^{1/2}[i] + \frac{B_{\psi_1}(w^*, \tilde{w}_1)}{\eta_1} + \sum_{t=1}^T \frac{\eta_t}{2} \|g_t - \tilde{m}_t\|_{\psi_{t-1}^*}^2 \\ & + D_\infty^2 \beta_1^2 \sum_{t=1}^T \|g_t - \theta_{t-1}\|_{\psi_{t-1}^*}. \end{aligned}$$

460 since $\|g_t - \tilde{g}_t\|_{\psi_{t-1}^*} = \|g_t - \beta_1 \theta_{t-1} - (1 - \beta_1)g_t\|_{\psi_{t-1}^*} = \beta^2 \|g_t - \theta_{t-1}\|_{\psi_{t-1}^*}$. This completes the
 461 proof.

462 □

463 B Proof of Corollary 1

464 **Corollary.** Suppose $\beta_1 = 0$ and $\{v_t\}_{t>0}$ is an increasing monotone sequence, then we obtain the
 465 following regret bound for any $w^* \in \Theta$ and sequence $\{\eta_t = \eta/\sqrt{t}\}_{t>0}$:

$$\mathcal{R}_T \leq \frac{B_{\psi_1}(w^*, \tilde{w}_1)}{\eta_1} + \frac{\eta\sqrt{1+\log T}}{\sqrt{1-\beta_2}} \sum_{i=1}^d \|(g - m)_{1:T}[i]\|_2 + \frac{D_\infty^2}{\eta_{\min}} \sum_{i=1}^d \left[(1 - \beta_2) \sum_{s=1}^T \beta_2^{T-s} g_s[i]^2 \right]^{1/2},$$

466 where $g_t := \nabla \ell_t(w_t)$ and $\eta_{\min} := \min_t \eta_t$.

467 **Proof** Recall the bound in Theorem 1:

$$\mathcal{R}_T \leq \frac{B_{\psi_1}(w^*, \tilde{w}_1)}{\eta_1} + \sum_{t=1}^T \frac{\eta_t}{2} \|g_t - \tilde{m}_t\|_{\psi_{t-1}^*}^2 + \frac{D_\infty^2}{\eta_{\min}} \sum_{i=1}^d \hat{v}_T^{1/2}[i] + D_\infty^2 \beta_1^2 \sum_{t=1}^T \|g_t - \theta_{t-1}\|_{\psi_{t-1}^*},$$

468 The second term reads:

$$\begin{aligned}
\sum_{t=1}^T \frac{\eta_t}{2} \|g_t - m_t\|_{\psi_{t-1}^*}^2 &= \sum_{t=1}^{T-1} \frac{\eta_t}{2} \|g_t - m_t\|_{\psi_{t-1}^*}^2 + \eta_T \sum_{i=1}^d \frac{(g_T[i] - m_T[i])^2}{\sqrt{v_{T-1}[i]}} \\
&= \sum_{t=1}^{T-1} \frac{\eta_t}{2} \|g_t - m_t\|_{\psi_{t-1}^*}^2 + \eta \sum_{i=1}^d \frac{(g_T[i] - m_T[i])^2}{\sqrt{T((1-\beta_2) \sum_{s=1}^{T-1} \beta_2^{T-1-s} (g_s[i] - m_s[i])^2)}} \\
&\leq \eta \sum_{i=1}^d \sum_{t=1}^T \frac{(g_t[i] - m_t[i])^2}{\sqrt{t((1-\beta_2) \sum_{s=1}^{t-1} \beta_2^{t-1-s} (g_s[i] - m_s[i])^2)}}.
\end{aligned}$$

469 To interpret the bound, let us make a rough approximation such that $\sum_{s=1}^{t-1} \beta_2^{t-1-s} (g_s[i] - m_s[i])^2 \simeq$
470 $(g_t[i] - m_t[i])^2$. Then, we can further get an upper-bound as

$$\sum_{t=1}^T \frac{\eta_t}{2} \|g_t - m_t\|_{\psi_{t-1}^*}^2 \leq \frac{\eta}{\sqrt{1-\beta_2}} \sum_{i=1}^d \sum_{t=1}^T \frac{|g_t[i] - m_t[i]|}{\sqrt{t}} \leq \frac{\eta \sqrt{1+\log T}}{\sqrt{1-\beta_2}} \sum_{i=1}^d \|(g - m)_{1:T}[i]\|_2,$$

471 where the last inequality is due to Cauchy-Schwarz.

472

□

473 C Proofs of Auxiliary Lemmas

474 C.1 Proof of Lemma 1

Lemma. Assume assumption H4, then the quantities defined in Algorithm 2 satisfy for any $w \in \Theta$ and $t > 0$:

$$\|\nabla f(w_t)\| < M, \quad \|\theta_t\| < M, \quad \|\hat{v}_t\| < M^2.$$

Proof Assume assumption H4 we have:

$$\|\nabla f(w)\| = \|\mathbb{E}[\nabla f(w, \xi)]\| \leq \mathbb{E}[\|\nabla f(w, \xi)\|] \leq M$$

475 By induction reasoning, since $\|\theta_0\| = 0 \leq M$ and suppose that for $\|\theta_t\| \leq M$ then we have

$$\|\theta_{t+1}\| = \|\beta_1 \theta_t + (1 - \beta_1) g_{t+1}\| \leq \beta_1 \|\theta_t\| + (1 - \beta_1) \|g_{t+1}\| \leq M \quad (25)$$

476 Using the same induction reasoning we prove that

$$\|\hat{v}_{t+1}\| = \|\beta_2 \hat{v}_t + (1 - \beta_2) g_{t+1}^2\| \leq \beta_2 \|\hat{v}_t\| + (1 - \beta_1) \|g_{t+1}^2\| \leq M^2 \quad (26)$$

477

□

478 C.2 Proof of Lemma 2

479 **Lemma.** Assume a strictly positive and non increasing sequence of stepsizes $\{\eta_t\}_{t>0}$, $\beta \in [0, 1]$, then
480 the following holds:

$$\bar{w}_{t+1} - \bar{w}_t \leq \frac{\beta_1}{1 - \beta_1} \tilde{\theta}_{t-1} \left[\eta_{t-1} \hat{v}_{t-1}^{-1/2} - \eta_t \hat{v}_t^{-1/2} \right] - \eta_t \hat{v}_t^{-1/2} \tilde{g}_t, \quad (27)$$

481 where $\tilde{\theta}_t = \theta_t + \beta_1 \theta_{t-1}$ and $\tilde{g}_t = g_t - \beta_1 m_t + \beta_1 g_{t-1} + m_{t+1}$.

482 **Proof** By definition (7) and using the Algorithm updates, we have:

$$\begin{aligned} \bar{w}_{t+1} - \bar{w}_t &= \frac{1}{1 - \beta_1} (w_{t+1} - \tilde{w}_t) - \frac{\beta_1}{1 - \beta_1} (w_t - \tilde{w}_{t-1}) \\ &= -\frac{1}{1 - \beta_1} \eta_t \hat{v}_t^{-1/2} (\theta_t + h_{t+1}) + \frac{\beta_1}{1 - \beta_1} \eta_{t-1} \hat{v}_{t-1}^{-1/2} (\theta_{t-1} + h_t) \\ &= -\frac{1}{1 - \beta_1} \eta_t \hat{v}_t^{-1/2} (\theta_t + \beta_1 \theta_{t-1}) - \frac{1}{1 - \beta_1} \eta_t \hat{v}_t^{-1/2} (1 - \beta_1) m_{t+1} \\ &\quad + \frac{\beta_1}{1 - \beta_1} \eta_{t-1} \hat{v}_{t-1}^{-1/2} (\theta_{t-1} + \beta_1 \theta_{t-2}) + \frac{\beta_1}{1 - \beta_1} \eta_{t-1} \hat{v}_{t-1}^{-1/2} (1 - \beta_1) m_t \end{aligned} \quad (28)$$

483 Denote $\tilde{\theta}_t = \theta_t + \beta_1 \theta_{t-1}$ and $\tilde{g}_t = g_t - \beta_1 m_t + \beta_1 g_{t-1} + m_{t+1}$. Notice that $\tilde{\theta}_t = \beta_1 \tilde{\theta}_{t-1} + (1 -$
484 $\beta_1)(g_t + \beta_1 g_{t-1})$.

$$\bar{w}_{t+1} - \bar{w}_t \leq \frac{\beta_1}{1 - \beta_1} \tilde{\theta}_{t-1} \left[\eta_{t-1} \hat{v}_{t-1}^{-1/2} - \eta_t \hat{v}_t^{-1/2} \right] - \eta_t \hat{v}_t^{-1/2} \tilde{g}_t \quad (29)$$

485

□

486 C.3 Proof of Lemma 3

487 **Lemma.** Assume H4, a strictly positive and a sequence of constant stepsizes $\{\eta_t\}_{t>0}$, $\beta \in [0, 1]$, then
488 the following holds:

$$\sum_{t=1}^{T_{\max}} \eta_t^2 \mathbb{E} \left[\left\| \hat{v}_t^{-1/2} \theta_t \right\|_2^2 \right] \leq \frac{\eta^2 d T_{\max} (1 - \beta_1)}{(1 - \beta_2)(1 - \gamma)} \quad (30)$$

489 **Proof** We denote by index $p \in [1, d]$ the dimension of each component of vectors of interest. Noting
 490 that for any $t > 0$ and dimension p we have $\hat{v}_{t,p} \geq v_{t,p}$, then:

$$\begin{aligned} \eta_t^2 \mathbb{E} \left[\left\| \hat{v}_t^{-1/2} \theta_t \right\|_2^2 \right] &= \eta_t^2 \mathbb{E} \left[\sum_{p=1}^d \frac{\theta_{t,p}^2}{\hat{v}_{t,p}} \right] \\ &\leq \eta_t^2 \mathbb{E} \left[\sum_{i=1}^d \frac{\theta_{t,p}^2}{v_{t,p}} \right] \\ &\leq \eta_t^2 \mathbb{E} \left[\sum_{i=1}^d \frac{(\sum_{r=1}^t (1 - \beta_1) \beta_1^{t-r} g_{r,p})^2}{\sum_{r=1}^t (1 - \beta_2) \beta_2^{t-r} g_{r,p}^2} \right] \end{aligned} \quad (31)$$

491 where the last inequality is due to initializations. Denote $\gamma = \frac{\beta_1}{\beta_2}$. Then,

$$\begin{aligned} \eta_t^2 \mathbb{E} \left[\left\| \hat{v}_t^{-1/2} \theta_t \right\|_2^2 \right] &\leq \frac{\eta_t^2 (1 - \beta_1)^2}{1 - \beta_2} \mathbb{E} \left[\sum_{i=1}^d \frac{(\sum_{r=1}^t \beta_1^{t-r} g_{r,p})^2}{\sum_{r=1}^t \beta_2^{t-r} g_{r,p}^2} \right] \\ &\stackrel{(a)}{\leq} \frac{\eta_t^2 (1 - \beta_1)}{1 - \beta_2} \mathbb{E} \left[\sum_{i=1}^d \frac{\sum_{r=1}^t \beta_1^{t-r} g_{r,p}^2}{\sum_{r=1}^t \beta_2^{t-r} g_{r,p}^2} \right] \\ &\leq \frac{\eta_t^2 (1 - \beta_1)}{1 - \beta_2} \mathbb{E} \left[\sum_{i=1}^d \sum_{r=1}^t \gamma^{t-r} \right] = \frac{\eta_t^2 d (1 - \beta_1)}{1 - \beta_2} \mathbb{E} \left[\sum_{r=1}^t \gamma^{t-r} \right] \end{aligned} \quad (32)$$

492 where (a) is due to $\sum_{r=1}^t \beta_1^{t-r} \leq \frac{1}{1 - \beta_1}$. Summing from $t = 1$ to $t = T_{\max}$ on both sides yields:

$$\begin{aligned} \sum_{t=1}^{T_{\max}} \eta_t^2 \mathbb{E} \left[\left\| \hat{v}_t^{-1/2} \theta_t \right\|_2^2 \right] &\leq \frac{\eta_t^2 d (1 - \beta_1)}{1 - \beta_2} \mathbb{E} \left[\sum_{t=1}^{T_{\max}} \sum_{r=1}^t \gamma^{t-r} \right] \\ &\leq \frac{\eta^2 d T (1 - \beta_1)}{1 - \beta_2} \mathbb{E} \left[\sum_{t=1}^T \gamma^{t-r} \right] \\ &\leq \frac{\eta^2 d T (1 - \beta_1)}{(1 - \beta_2)(1 - \gamma)} \end{aligned} \quad (33)$$

493 where the last inequality is due to $\sum_{r=1}^t \gamma^{t-r} \leq \frac{1}{1 - \gamma}$ by definition of γ . \square

494 **D Proof of Theorem 2**

495 **Theorem.** Assume H2-H4, $(\beta_1, \beta_2) \in [0, 1]$ and a sequence of decreasing stepsizes $\{\eta_t\}_{t \geq 0}$, then
 496 the following result holds:

$$\mathbb{E} [\|\nabla f(w_T)\|^2] \leq \tilde{C}_1 \sqrt{\frac{d}{T_{\max}}} + \tilde{C}_2 \frac{1}{T_{\max}} \quad (34)$$

497 where T is a random termination number distributed according (6) and the constants are defined as
 498 follows:

$$\begin{aligned} \tilde{C}_1 &= C_1 + \frac{M}{(1 - a\beta_1) + (\beta_1 + a)} \left[\frac{a(1 - \beta_1)^2}{1 - \beta_2} + 2L \frac{1}{1 - \beta_2} \right] \\ C_1 &= \frac{M}{(1 - a\beta_1) + (\beta_1 + a)} \Delta f + \frac{4L \left(\frac{\beta_1}{1 - \beta_1} \right)^2 M}{(1 - a\beta_1) + (\beta_1 + a)} \frac{(1 + \beta_1^2)(1 - \beta_1)}{(1 - \beta_2)(1 - \gamma)} \\ \tilde{C}_2 &= \frac{M}{(1 - \beta_1)((1 - a\beta_1) + (\beta_1 + a))} \tilde{M}^2 \mathbb{E} \left[\left\| \hat{v}_0^{-1/2} \right\| \right] \end{aligned} \quad (35)$$

499 **Proof** Using H2 and the iterate \bar{w}_t we have:

$$\begin{aligned} f(\bar{w}_{t+1}) &\leq f(\bar{w}_t) + \nabla f(\bar{w}_t)^\top (\bar{w}_{t+1} - \bar{w}_t) + \frac{L}{2} \|\bar{w}_{t+1} - \bar{w}_t\|^2 \\ &\leq f(\bar{w}_t) + \underbrace{\nabla f(w_t)^\top (\bar{w}_{t+1} - \bar{w}_t)}_A + \underbrace{(\nabla f(\bar{w}_t) - \nabla f(w_t))^\top (\bar{w}_{t+1} - \bar{w}_t)}_B + \frac{L}{2} \|\bar{w}_{t+1} - \bar{w}_t\| \end{aligned} \quad (36)$$

500 **Term A.** Using Lemma 2, we have that:

$$\begin{aligned} \nabla f(w_t)^\top (\bar{w}_{t+1} - \bar{w}_t) &\leq \nabla f(w_t)^\top \left[\frac{\beta_1}{1 - \beta_1} \tilde{\theta}_{t-1} \left[\eta_{t-1} \hat{v}_{t-1}^{-1/2} - \eta_t \hat{v}_t^{-1/2} \right] - \eta_t \hat{v}_t^{-1/2} \tilde{g}_t \right] \\ &\leq \frac{\beta_1}{1 - \beta_1} \|\nabla f(w_t)\| \left\| \eta_{t-1} \hat{v}_{t-1}^{-1/2} - \eta_t \hat{v}_t^{-1/2} \right\| \left\| \tilde{\theta}_{t-1} \right\| - \nabla f(w_t)^\top \eta_t \hat{v}_t^{-1/2} \tilde{g}_t \end{aligned} \quad (37)$$

501 where the inequality is due to trivial inequality for positive diagonal matrix. Using Lemma 1 and
502 assumption H3 we obtain:

$$\nabla f(w_t)^\top (\bar{w}_{t+1} - \bar{w}_t) \leq \frac{\beta_1(1 + \beta_1)}{1 - \beta_1} \mathbf{M}^2 \left[\left\| \eta_{t-1} \hat{v}_{t-1}^{-1/2} \right\| - \left\| \eta_t \hat{v}_t^{-1/2} \right\| \right] - \nabla f(w_t)^\top \eta_t \hat{v}_t^{-1/2} \tilde{g}_t \quad (38)$$

503 where we have used the fact that $\eta_t \hat{v}_t^{-1/2}$ is a diagonal matrix such that $\eta_{t-1} \hat{v}_{t-1}^{-1/2} \succcurlyeq \eta_t \hat{v}_t^{-1/2} \succcurlyeq 0$
504 (decreasing stepsize and max operator). Also note that:

$$\begin{aligned} -\nabla f(w_t)^\top \eta_t \hat{v}_t^{-1/2} \tilde{g}_t &= -\nabla f(w_t)^\top \eta_{t-1} \hat{v}_{t-1}^{-1/2} \bar{g}_t - \nabla f(w_t)^\top \left[\eta_t \hat{v}_t^{-1/2} - \eta_{t-1} \hat{v}_{t-1}^{-1/2} \right] \bar{g}_t \\ &\quad - \nabla f(w_t)^\top \eta_{t-1} \hat{v}_{t-1}^{-1/2} (\beta_1 g_{t-1} + m_{t+1}) \\ &\leq -\nabla f(w_t)^\top \eta_{t-1} \hat{v}_{t-1}^{-1/2} \bar{g}_t + (1 - a\beta_1) \mathbf{M}^2 \left[\left\| \eta_{t-1} \hat{v}_{t-1}^{-1/2} \right\| - \left\| \eta_t \hat{v}_t^{-1/2} \right\| \right] \\ &\quad - \nabla f(w_t)^\top \eta_t \hat{v}_t^{-1/2} (\beta_1 g_{t-1} + m_{t+1}) \end{aligned} \quad (39)$$

505 using Lemma 1 on $\|g_t\|$ and where that $\tilde{g}_t = \bar{g}_t + \beta_1 g_{t-1} + m_{t+1} = g_t - \beta_1 m_t + \beta_1 g_{t-1} + m_{t+1}$.
506 Plugging (39) into (38) yields:

$$\begin{aligned} \nabla f(w_t)^\top (\bar{w}_{t+1} - \bar{w}_t) &\leq -\nabla f(w_t)^\top \eta_{t-1} \hat{v}_{t-1}^{-1/2} \bar{g}_t + \frac{1}{1 - \beta_1} (a\beta_1^2 - 2a\beta_1 + \beta_1) \mathbf{M}^2 \left[\left\| \eta_{t-1} \hat{v}_{t-1}^{-1/2} \right\| - \left\| \eta_t \hat{v}_t^{-1/2} \right\| \right] \\ &\quad - \nabla f(w_t)^\top \eta_t \hat{v}_t^{-1/2} (\beta_1 g_{t-1} + m_{t+1}) \end{aligned} \quad (40)$$

507 **Term B.** By Cauchy-Schwarz (CS) inequality we have:

$$(\nabla f(\bar{w}_t) - \nabla f(w_t))^\top (\bar{w}_{t+1} - \bar{w}_t) \leq \|\nabla f(\bar{w}_t) - \nabla f(w_t)\| \|\bar{w}_{t+1} - \bar{w}_t\| \quad (41)$$

508 Using smoothness assumption H2:

$$\begin{aligned} \|\nabla f(\bar{w}_t) - \nabla f(w_t)\| &\leq L \|\bar{w}_t - w_t\| \\ &\leq L \frac{\beta_1}{1 - \beta_1} \|w_t - \tilde{w}_{t-1}\| \end{aligned} \quad (42)$$

509 By Lemma 2 we also have:

$$\begin{aligned} \bar{w}_{t+1} - \bar{w}_t &= \frac{\beta_1}{1 - \beta_1} \tilde{\theta}_{t-1} \left[\eta_{t-1} \hat{v}_{t-1}^{-1/2} - \eta_t \hat{v}_t^{-1/2} \right] - \eta_t \hat{v}_t^{-1/2} \tilde{g}_t \\ &= \frac{\beta_1}{1 - \beta_1} \tilde{\theta}_{t-1} \eta_{t-1} \hat{v}_{t-1}^{-1/2} \left[I - (\eta_t \hat{v}_t^{-1/2})(\eta_{t-1} \hat{v}_{t-1}^{-1/2})^{-1} \right] - \eta_t \hat{v}_t^{-1/2} \tilde{g}_t \\ &= \frac{\beta_1}{1 - \beta_1} \left[I - (\eta_t \hat{v}_t^{-1/2})(\eta_{t-1} \hat{v}_{t-1}^{-1/2})^{-1} \right] (\tilde{w}_{t-1} - w_t) - \eta_t \hat{v}_t^{-1/2} \tilde{g}_t \end{aligned} \quad (43)$$

510 where the last equality is due to $\tilde{\theta}_{t-1}\eta_{t-1}\hat{v}_{t-1}^{-1/2} = \tilde{w}_{t-1} - w_t$ by construction of $\tilde{\theta}_t$. Taking the
 511 norms on both sides, observing $\left\|I - (\eta_t\hat{v}_t^{-1/2})(\eta_{t-1}\hat{v}_{t-1}^{-1/2})^{-1}\right\| \leq 1$ due to the decreasing stepsize
 512 and the construction of \hat{v}_t and using CS inequality yield:

$$\|\bar{w}_{t+1} - \bar{w}_t\| \leq \frac{\beta_1}{1 - \beta_1} \|\tilde{w}_{t-1} - w_t\| + \left\|\eta_t\hat{v}_t^{-1/2}\tilde{g}_t\right\| \quad (44)$$

We recall Young's inequality with a constant $\delta \in (0, 1)$ as follows:

$$\langle X | Y \rangle \leq \frac{1}{\delta} \|X\|^2 + \delta \|Y\|^2$$

513 Plugging (42) and (44) into (41) returns:

$$\begin{aligned} (\nabla f(\bar{w}_t) - \nabla f(w_t))^\top (\bar{w}_{t+1} - \bar{w}_t) &\leq L \frac{\beta_1}{1 - \beta_1} \left\|\eta_t\hat{v}_t^{-1/2}\tilde{g}_t\right\| \|w_t - \tilde{w}_{t-1}\| \\ &\quad + L \left(\frac{\beta_1}{1 - \beta_1}\right)^2 \|\tilde{w}_{t-1} - w_t\|^2 \end{aligned} \quad (45)$$

514 Applying Young's inequality with $\delta \rightarrow \frac{\beta_1}{1 - \beta_1}$ on the product $\left\|\eta_t\hat{v}_t^{-1/2}\tilde{g}_t\right\| \|w_t - \tilde{w}_{t-1}\|$ yields:

$$(\nabla f(\bar{w}_t) - \nabla f(w_t))^\top (\bar{w}_{t+1} - \bar{w}_t) \leq L \left\|\eta_t\hat{v}_t^{-1/2}\tilde{g}_t\right\|^2 + 2L \left(\frac{\beta_1}{1 - \beta_1}\right)^2 \|\tilde{w}_{t-1} - w_t\|^2 \quad (46)$$

515 The last term $\frac{L}{2} \|\bar{w}_{t+1} - \bar{w}_t\|^2$ can be upper bounded using (44):

$$\begin{aligned} \frac{L}{2} \|\bar{w}_{t+1} - \bar{w}_t\|^2 &\leq \frac{L}{2} \left[\frac{\beta_1}{1 - \beta_1} \|\tilde{w}_{t-1} - w_t\| + \left\|\eta_t\hat{v}_t^{-1/2}\tilde{g}_t\right\| \right]^2 \\ &\leq L \left\|\eta_t\hat{v}_t^{-1/2}\tilde{g}_t\right\|^2 + 2L \left(\frac{\beta_1}{1 - \beta_1}\right)^2 \|\tilde{w}_{t-1} - w_t\|^2 \end{aligned} \quad (47)$$

516 Plugging (40), (46) and (47) into (36) and taking the expectations on both sides give:

$$\begin{aligned} &\mathbb{E} \left[f(\bar{w}_{t+1}) + \frac{1}{1 - \beta_1} \tilde{M}^2 \left\|\eta_t\hat{v}_t^{-1/2}\right\| - \left(f(\bar{w}_t) + \frac{1}{1 - \beta_1} \tilde{M}^2 \left\|\eta_{t-1}\hat{v}_{t-1}^{-1/2}\right\| \right) \right] \\ &\leq \mathbb{E} \left[-\nabla f(w_t)^\top \eta_{t-1}\hat{v}_{t-1}^{-1/2}\tilde{g}_t - \nabla f(w_t)^\top \eta_t\hat{v}_t^{-1/2}(\beta_1 g_{t-1} + m_{t+1}) \right] \\ &\quad + \mathbb{E} \left[2L \left\|\eta_t\hat{v}_t^{-1/2}\tilde{g}_t\right\|^2 + 4L \left(\frac{\beta_1}{1 - \beta_1}\right)^2 \|\tilde{w}_{t-1} - w_t\|^2 \right] \end{aligned} \quad (48)$$

517 where $\tilde{M}^2 = (a\beta_1^2 - 2a\beta_1 + \beta_1)M^2$. Note that the expectation of \tilde{g}_t conditioned on the filtration \mathcal{F}_t
 518 reads as follows

$$\begin{aligned} \mathbb{E} [\nabla f(w_t)^\top \tilde{g}_t] &= \mathbb{E} [\nabla f(w_t)^\top (g_t - \beta_1 m_t)] \\ &= (1 - a\beta_1) \|\nabla f(w_t)\|^2 \end{aligned} \quad (49)$$

519 Summing from $t = 1$ to $t = T$ leads to

$$\begin{aligned} &\frac{1}{M} \sum_{t=1}^{T_{\max}} ((1 - a\beta_1)\eta_{t-1} + (\beta_1 + a)\eta_t) \|\nabla f(w_t)\|^2 \leq \\ &\mathbb{E} \left[f(\bar{w}_1) + \frac{1}{1 - \beta_1} \tilde{M}^2 \left\|\eta_0\hat{v}_0^{-1/2}\right\| - \left(f(\bar{w}_{T_{\max}+1}) + \frac{1}{1 - \beta_1} \tilde{M}^2 \left\|\eta_{T_{\max}}\hat{v}_{T_{\max}}^{-1/2}\right\| \right) \right] \\ &\quad + 2L \sum_{t=1}^{T_{\max}} \mathbb{E} \left[\left\|\eta_t\hat{v}_t^{-1/2}\tilde{g}_t\right\|^2 \right] + 4L \left(\frac{\beta_1}{1 - \beta_1}\right)^2 \sum_{t=1}^{T_{\max}} \mathbb{E} [\|\tilde{w}_{t-1} - w_t\|^2] \\ &\leq \mathbb{E} \left[\Delta f + \frac{1}{1 - \beta_1} \tilde{M}^2 \left\|\eta_0\hat{v}_0^{-1/2}\right\| \right] + 2L \sum_{t=1}^{T_{\max}} \mathbb{E} \left[\left\|\eta_t\hat{v}_t^{-1/2}\tilde{g}_t\right\|^2 \right] + 4L \left(\frac{\beta_1}{1 - \beta_1}\right)^2 \sum_{t=1}^{T_{\max}} \mathbb{E} [\|\tilde{w}_{t-1} - w_t\|^2] \end{aligned} \quad (50)$$

where $\Delta f = f(\bar{w}_1) - f(\bar{w}_{T_{\max}+1})$. We note that by definition of \hat{v}_t , and a constant learning rate η_t , we have

$$\begin{aligned}\|\tilde{w}_{t-1} - w_t\|^2 &= \left\| \eta_{t-1} \hat{v}_{t-1}^{-1/2} (\theta_{t-1} + h_t) \right\|^2 \\ &= \left\| \eta_{t-1} \hat{v}_{t-1}^{-1/2} (\theta_{t-1} + \beta_1 \theta_{t-2} + (1 - \beta_1) m_t) \right\|^2 \\ &\leq \left\| \eta_{t-1} \hat{v}_{t-1}^{-1/2} \theta_{t-1} \right\|^2 + \left\| \eta_{t-2} \hat{v}_{t-2}^{-1/2} \beta_1 \theta_{t-2} \right\|^2 + (1 - \beta_1)^2 \left\| \eta_{t-1} \hat{v}_{t-1}^{-1/2} m_t \right\|^2\end{aligned}\tag{51}$$

Using Lemma 3 we have

$$\begin{aligned}\sum_{t=1}^{T_{\max}} \mathbb{E} \left[\|\tilde{w}_{t-1} - w_t\|^2 \right] \\ \leq (1 + \beta_1^2) \frac{\eta^2 d T_{\max} (1 - \beta_1)}{(1 - \beta_2)(1 - \gamma)} + (1 - \beta_1)^2 \sum_{t=1}^{T_{\max}} \mathbb{E} \left[\left\| \eta_{t-1} \hat{v}_{t-1}^{-1/2} m_t \right\|^2 \right]\end{aligned}\tag{52}$$

And thus, setting the learning rate to a constant value η and injecting in (50) yields:

$$\begin{aligned}\mathbb{E} [\|\nabla f(w_T)\|^2] &= \frac{1}{\sum_{j=1}^{T_{\max}} \eta_j} \sum_{t=1}^{T_{\max}} \eta_t \|\nabla f(w_t)\|^2 \\ &\leq \frac{M}{(1 - a\beta_1) + (\beta_1 + a)} \frac{1}{\sum_{j=1}^{T_{\max}} \eta_j} \mathbb{E} \left[\Delta f + \frac{1}{1 - \beta_1} \tilde{M}^2 \left\| \eta_0 \hat{v}_0^{-1/2} \right\| \right] \\ &\quad + \frac{4L \left(\frac{\beta_1}{1 - \beta_1} \right)^2 M}{(1 - a\beta_1) + (\beta_1 + a)} \frac{1}{\sum_{j=1}^{T_{\max}} \eta_j} (1 + \beta_1^2) \frac{\eta^2 d T_{\max} (1 - \beta_1)}{(1 - \beta_2)(1 - \gamma)} \\ &\quad + \frac{M}{(1 - a\beta_1) + (\beta_1 + a)} \frac{1}{\sum_{j=1}^{T_{\max}} \eta_j} (1 - \beta_1)^2 \sum_{t=1}^{T_{\max}} \mathbb{E} \left[\left\| \eta_{t-1} \hat{v}_{t-1}^{-1/2} m_t \right\|^2 \right] \\ &\quad + \frac{2LM}{(1 - a\beta_1) + (\beta_1 + a)} \frac{1}{\sum_{j=1}^{T_{\max}} \eta_j} \sum_{t=1}^{T_{\max}} \mathbb{E} \left[\left\| \eta_t \hat{v}_t^{-1/2} \tilde{g}_t \right\|^2 \right]\end{aligned}\tag{53}$$

where T is a random termination number distributed according (6). Setting the stepsize to $\eta = \frac{1}{\sqrt{dT_{\max}}}$ yields :

$$\begin{aligned}\mathbb{E} [\|\nabla f(w_T)\|^2] \\ \leq C_1 \sqrt{\frac{d}{T_{\max}}} + C_2 \frac{1}{T_{\max}} \\ + D_1 \frac{\eta}{T_{\max}} \sum_{t=1}^{T_{\max}} \mathbb{E} \left[\left\| \hat{v}_{t-1}^{-1/2} m_t \right\|^2 \right] + D_2 \frac{\eta}{T_{\max}} \sum_{t=1}^{T_{\max}} \mathbb{E} \left[\left\| \hat{v}_{t-1}^{-1/2} \tilde{g}_t \right\|^2 \right]\end{aligned}\tag{54}$$

where

$$\begin{aligned}C_1 &= \frac{M}{(1 - a\beta_1) + (\beta_1 + a)} \Delta f + \frac{4L \left(\frac{\beta_1}{1 - \beta_1} \right)^2 M}{(1 - a\beta_1) + (\beta_1 + a)} \frac{(1 + \beta_1^2)(1 - \beta_1)}{(1 - \beta_2)(1 - \gamma)} \\ C_2 &= \frac{M}{(1 - \beta_1)((1 - a\beta_1) + (\beta_1 + a))} \tilde{M}^2 \mathbb{E} \left[\left\| \hat{v}_0^{-1/2} \right\| \right]\end{aligned}\tag{55}$$

Simple case as in [43]: if $\beta_1 = 0$ then $\tilde{g}_t = g_t + m_{t+1}$ and $g_t = \theta_t$. Also using Lemma 3 we have that:

$$\sum_{t=1}^{T_{\max}} \eta_t^2 \mathbb{E} \left[\left\| \hat{v}_t^{-1/2} g_t \right\|^2 \right] \leq \frac{\eta^2 d T_{\max}}{(1 - \beta_2)}\tag{56}$$

529 which leads to the final bound:

$$\begin{aligned} & \mathbb{E} [\|\nabla f(w_T)\|^2] \\ & \leq \tilde{C}_1 \sqrt{\frac{d}{T_{\max}}} + \tilde{C}_2 \frac{1}{T_{\max}} \end{aligned} \quad (57)$$

530 where

$$\begin{aligned} \tilde{C}_1 &= C_1 + \frac{M}{(1 - a\beta_1) + (\beta_1 + a)} \left[\frac{a(1 - \beta_1)^2}{1 - \beta_2} + 2L \frac{1}{1 - \beta_2} \right] \\ \tilde{C}_2 &= C_2 = \frac{M}{(1 - \beta_1)((1 - a\beta_1) + (\beta_1 + a))} \tilde{M}^2 \mathbb{E} [\|\hat{v}_0^{-1/2}\|] \end{aligned} \quad (58)$$

531 □

532 E Proof of Lemma 4 (Boundedness of the iterates)

533 **Lemma.** *Given the multilayer model (9), assume the boundedness of the input data and of the loss*
 534 *function, i.e., for any $\xi \in \mathbb{R}^l$ and $y \in \mathbb{R}$ there is a constant $T > 0$ such that:*

$$\|\xi\| \leq 1 \quad \text{a.s.} \quad \text{and} \quad |\mathcal{L}'(\cdot, y)| \leq T \quad (59)$$

where $\mathcal{L}'(\cdot, y)$ denotes its derivative w.r.t. the parameter. Then for each layer $\ell \in [1, L]$, there exist a constant $A_{(\ell)}$ such that:

$$\|w^{(\ell)}\| \leq A_{(\ell)}$$

Proof Recall that for any layer index $\ell \in [1, L]$ we denote the output of layer ℓ by $h^{(\ell)}(w, \xi)$:

$$h^{(\ell)}(w, \xi) = \sigma \left(w^{(\ell)} \sigma \left(w^{(\ell-1)} \dots \sigma \left(w^{(1)} \xi \right) \right) \right)$$

535 Given the sigmoid assumption we have $\|h^{(\ell)}(w, \xi)\| \leq 1$ for any $\ell \in [1, L]$ and any $(w, \xi) \in$
 536 $\mathbb{R}^d \times \mathbb{R}^l$. Observe that at the last layer L :

$$\begin{aligned} \|\nabla_{w^{(L)}} \mathcal{L}(\text{MLN}(w, \xi), y)\| &= \|\mathcal{L}'(\text{MLN}(w, \xi), y) \nabla_{w^{(L)}} \text{MLN}(w, \xi)\| \\ &= \left\| \mathcal{L}'(\text{MLN}(w, \xi), y) \sigma'(w^{(L)} h^{(L-1)}(w, \xi)) h^{(L-1)}(w, \xi) \right\| \\ &\leq \frac{T}{4} \end{aligned} \quad (60)$$

537 where the last equality is due to mild assumptions (59) and to the fact that the norm of the derivative
 538 of the sigmoid function is upperbounded by 1/4.

539 From Algorithm 2, and with $\beta_1 = 0$ for the sake of notation, we have for iteration index $t > 0$:

$$\begin{aligned} \|w_t - \tilde{w}_{t-1}\| &= \left\| -\eta_t \hat{v}_t^{-1/2} (\theta_t + h_{t+1}) \right\| \\ &= \left\| \eta_t \hat{v}_t^{-1/2} (g_t + m_{t+1}) \right\| \\ &\leq \hat{\eta} \left\| \hat{v}_t^{-1/2} g_t \right\| + \hat{\eta} a \left\| \hat{v}_t^{-1/2} g_{t+1} \right\| \end{aligned} \quad (61)$$

where $\hat{\eta} = \max_{t>0} \eta_t$. For any dimension $p \in [1, d]$, using assumption H3, we note that

$$\sqrt{\hat{v}_{t,p}} \geq \sqrt{1 - \beta_2} g_{t,p} \quad \text{and} \quad m_{t+1} \leq a \|g_{t+1}\|$$

540 . Thus:

$$\begin{aligned} \|w_t - \tilde{w}_{t-1}\| &\leq \hat{\eta} \left(\left\| \hat{v}_t^{-1/2} g_t \right\| + a \left\| \hat{v}_t^{-1/2} g_{t+1} \right\| \right) \\ &\leq \hat{\eta} \frac{a + 1}{\sqrt{1 - \beta_2}} \end{aligned} \quad (62)$$

541 In short there exist a constant B such that $\|w_t - \tilde{w}_{t-1}\| \leq B$.

Proof by induction: As in [9], we will prove the containment of the weights by induction. Suppose an iteration index T and a coordinate i of the last layer L such that $w_{T,i}^{(L)} \geq \frac{T}{4\lambda} + B$. Using (60), we have

$$\nabla_i f(w_t^{(L)}) \geq -\frac{T}{4} + \lambda \frac{T}{\lambda 4} \geq 0$$

542 where $f(\cdot)$ is defined by (10) and is the loss of our MLN. This last equation yields $\theta_{T,i}^{(L)} \geq 0$ (given
543 the algorithm and $\beta_1 = 0$) and using the fact that $\|w_t - \tilde{w}_{t-1}\| \leq B$ we have

$$0 \leq w_{T-1,i}^{(L)} - B \leq w_{T,i}^{(L)} \leq w_{T-1,i}^{(L)} \quad (63)$$

which means that $|w_{T,i}^{(L)}| \leq w_{T-1,i}^{(L)}$. So if the first assumption of that induction reasoning holds, i.e., $w_{T-1,i}^{(L)} \geq \frac{T}{4\lambda} + B$, then the next iterates $w_{T,i}^{(L)}$ decreases, see (63) and go below $\frac{T}{4\lambda} + B$. This yields that for any iteration index $t > 0$ we have

$$w_{T,i}^{(L)} \leq \frac{T}{4\lambda} + 2B$$

since B is the biggest jump an iterate can do since $\|w_t - \tilde{w}_{t-1}\| \leq B$. Likewise we can end up showing that

$$|w_{T,i}^{(L)}| \leq \frac{T}{4\lambda} + 2B$$

544 meaning that the weights of the last layer at any iteration is bounded in some matrix norm.

545 Now that we have shown this boundedness property for the last layer L , we will do the same for the
546 previous layers and conclude the verification of assumption H1 by induction.

547 For any layer $\ell \in [1, L - 1]$, we have:

$$\nabla_{w^{(\ell)}} \mathcal{L}(\text{MLN}(w, \xi), y) = \mathcal{L}'(\text{MLN}(w, \xi), y) \left(\prod_{j=1}^{\ell+1} \sigma' \left(w^{(j)} h^{(j-1)}(w, \xi) \right) \right) h^{(\ell-1)}(w, \xi) \quad (64)$$

This last quantity is bounded as long as we can prove that for any layer ℓ the weights $w^{(\ell)}$ are bounded in some matrix norm as $\|w^{(\ell)}\|_F \leq F_\ell$ with the Frobenius norm. Suppose we have shown $\|w^{(r)}\|_F \leq F_r$ for any layer $r > \ell$. Then having this gradient (64) bounded we can use the same lines of proof for the last layer L and show that the norm of the weights at the selected layer ℓ satisfy

$$\|w^{(\ell)}\| \leq \frac{T \prod_{t \geq \ell} F_t}{4^{L-\ell+1}} + 2B$$

548 Showing that the weights of the previous layers $\ell \in [1, L - 1]$ as well as for the last layer L of our
549 fully connected feed forward neural network are bounded at each iteration, leads by induction, to
550 the boundedness (at each iteration) assumption we want to check. \square

F Comparison to some related methods

Comparison to nonconvex optimization works. Recently, [40, 5, 38, 42, 44, 21] provide some theoretical analysis of ADAM-type algorithms when applying them to smooth nonconvex optimization problems. For example, [5] provides a bound, which is $\min_{t \in [T]} \mathbb{E}[\|\nabla f(w_t)\|^2] = O(\log T / \sqrt{T})$. Yet, this data independent bound does not show any advantage over standard stochastic gradient descent. Similar concerns appear in other papers.

To get some adaptive data dependent bound that are in terms of the gradient norms observed along the trajectory) when applying OPTIMISTIC-AMSGRAD to nonconvex optimization, one can follow the approach of [2] or [6]. They provide ways to convert algorithms with adaptive data dependent regret bound for convex loss functions (e.g. ADAGRAD) to the ones that can find an approximate stationary point of non-convex loss functions. Their approaches are modular so that simply using OPTIMISTIC-AMSGRAD as the base algorithm in their methods will immediately lead to a variant of OPTIMISTIC-AMSGRAD that enjoys some guarantee on nonconvex optimization. The variant can outperform the ones instantiated by other ADAM-type algorithms when the gradient prediction m_t is close to g_t . The details are omitted since this is a straightforward application.

Comparison to AO-FTRL [26]. In [26], the authors propose AO-FTRL, which has the update of the form $w_{t+1} = \arg \min_{w \in \Theta} (\sum_{s=1}^t g_s)^\top w + m_{t+1}^\top w + r_{0:t}(w)$, where $r_{0:t}(\cdot)$ is a 1-strongly convex loss function with respect to some norm $\|\cdot\|_{(t)}$ that may be different for different iteration t . Data dependent regret bound was provided in the paper, which is $r_{0:T}(w^*) + \sum_{t=1}^T \|g_t - m_t\|_{(t)}^*$ for any benchmark $w^* \in \Theta$. We see that if one selects $r_{0:t}(w) := \langle w, \text{diag}\{\hat{v}_t\}^{1/2} w \rangle$ and $\|\cdot\|_{(t)} := \sqrt{\langle \cdot, \text{diag}\{\hat{v}_t\}^{1/2} \cdot \rangle}$, then the update might be viewed as an optimistic variant of ADAGRAD. However, no experiments was provided in [26].

Comparison to OPTIMISTIC-ADAM [8]. We are aware that [8] proposed one version of optimistic algorithm for ADAM, which is called OPTIMISTIC-ADAM in their paper. A slightly modified version is summarized in Algorithm 4. Here, OPTIMISTIC-ADAM+ \hat{v}_t is OPTIMISTIC-ADAM in [8] with the additional max operation $\hat{v}_t = \max(\hat{v}_{t-1}, v_t)$ to guarantee that the weighted second moment is monotone increasing.

Algorithm 4 OPTIMISTIC-ADAM [8]+ \hat{v}_t .

- 1: Required: parameter β_1, β_2 , and η_t .
 - 2: Init: $w_1 \in \Theta$ and $\hat{v}_0 = v_0 = \epsilon 1 \in \mathbb{R}^d$.
 - 3: **for** $t = 1$ to T **do**
 - 4: Get mini-batch stochastic gradient vector $g_t \in \mathbb{R}^d$ at w_t .
 - 5: $\theta_t = \beta_1 \theta_{t-1} + (1 - \beta_1) g_t$.
 - 6: $v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$.
 - 7: $\hat{v}_t = \max(\hat{v}_{t-1}, v_t)$.
 - 8: $w_{t+1} = \Pi_k[w_t - 2\eta_t \frac{\theta_t}{\sqrt{\hat{v}_t}} + \eta_t \frac{\theta_{t-1}}{\sqrt{\hat{v}_{t-1}}}]$.
 - 9: **end for**
-

We want to emphasize that the motivations are different. OPTIMISTIC-ADAM in their paper is designed to optimize two-player games (e.g. GANs [14]), while the proposed algorithm in this paper is designed to accelerate optimization (e.g. solving empirical risk minimization quickly). [8] focuses on training GANs [14]. GANs is a two-player zero-sum game. There have been some related works in OPTIMISTIC ONLINE LEARNING like [7, 30, 34]) showing that if both players use some kinds of OPTIMISTIC-update, then accelerating the convergence to the equilibrium of the game is possible. [8] was inspired by these related works and showed that OPTIMISTIC-MIRROR-DESCENT can avoid the cycle behavior in a bilinear zero-sum game, which accelerates the convergence. Furthermore, [8] did not provide theoretical analysis of OPTIMISTIC-ADAM.

G Additional Remarks and Runs on the Gradient Prediction Process

Two illustrative examples. We provide two toy examples to demonstrate how OPTIMISTIC-AMSGRAD works with the chosen extrapolation method. First, consider minimizing a quadratic

function $H(w) := \frac{b}{2}w^2$ with vanilla gradient descent method $w_{t+1} = w_t - \eta_t \nabla H(w_t)$. The gradient $g_t := \nabla H(w_t)$ has a recursive description as $g_{t+1} = bw_{t+1} = b(w_t - \eta_t g_t) = g_t - b\eta_t g_t$. So, the update can be written in the form of (12) with $A = (1 - b\eta)$ and $u_{t-1} = 0$ by setting $\eta_t = \eta$ (constant step size). Therefore, the extrapolation method should predict well.

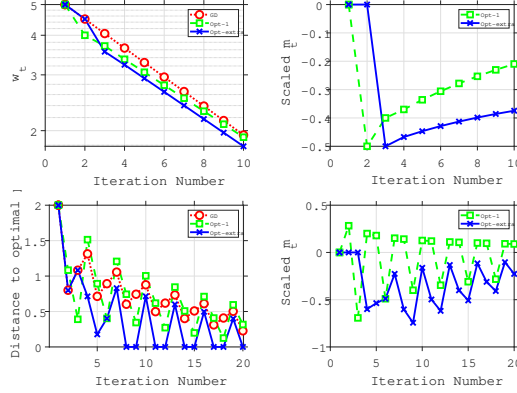


Figure 5: (a): The iterate w_t ; the closer to the optimal point 0 the better. (b): A scaled and clipped version of m_t : $w_t - w_{t-1/2}$, which measures how the prediction of m_t drives the update towards the optimal point. In this scenario, the more negative the better. (c): Distance to the optimal point -1 . The smaller the better. (d): A scaled and clipped version of m_t : $w_t - w_{t-1/2}$, which measures how the prediction of m_t drives the update towards the optimal point. In this scenario, the more negative the better.

Specifically, consider optimizing $H(w) := w^2/2$ by the following three algorithms with the same step size. One is Gradient Descent (GD): $w_{t+1} = w_t - \eta_t g_t$, while the other two are OPTIMISTIC-AMSGRAD with $\beta_1 = 0$ and the second moment term \hat{v}_t being dropped: $w_{t+\frac{1}{2}} = \Pi_{\Theta}[w_{t-\frac{1}{2}} - \eta_t g_t]$, $w_{t+1} = \Pi_{\Theta}[w_{t+\frac{1}{2}} - \eta_{t+1} m_{t+1}]$. We denote the algorithm that sets $m_{t+1} = g_t$ as Opt-1, and denote the algorithm that uses the extrapolation method to get m_{t+1} as Opt-extra. We let $\eta_t = 0.1$ and the initial point $w_0 = 5$ for all the three methods. The simulation results are on Figure 5 (a) and (b). Sub-figure (a) plots update w_t over iteration, where the updates should go towards the optimal point 0. Sub-figure (b) is about a scaled and clipped version of m_t , defined as $w_t - w_{t-1/2}$, which can be viewed as $-\eta_t m_t$ if the projection (if exists) is lifted. Sub-figure (a) shows that Opt-extra converges faster than the other methods. Furthermore, sub-figure (b) shows that the prediction by the extrapolation method is better than the prediction by simply using the previous gradient. The sub-figure shows that $-m_t$ from both methods all point to 0 in all iterations and the magnitude is larger for the one produced by the extrapolation method after iteration 2.²

Now let us consider another problem: an online learning problem proposed in [31]³. Assume the learner's decision space is $\Theta = [-1, 1]$, and the loss function is $\ell_t(w) = 3w$ if $t \bmod 3 = 1$, and $\ell_t(w) = -w$ otherwise. The optimal point to minimize the cumulative loss is $w^* = -1$. We let $\eta_t = 0.1/\sqrt{t}$ and the initial point $w_0 = 1$ for all the three methods. The parameter λ of the extrapolation method is set to $\lambda = 10^{-3} > 0$. The results are on Figure 5 (c) and (d). Sub-figure (c) shows that Opt-extra converges faster than the other methods while Opt-1 is not better than GD. The reason is that the gradient changes from -1 to 3 at $t \bmod 3 = 1$ and it changes from 3 to -1 at $t \bmod 3 = 2$. Consequently, using the current gradient as the guess for the next clearly is not a good choice, since the next gradient is in the opposite direction of the current one. Sub-figure (d) shows that $-m_t$ by the extrapolation method always points to $w^* = -1$, while the one by using the previous negative direction points to the opposite direction in two thirds of rounds. It shows that the extrapolation method is much less affected by the gradient oscillation and always makes the prediction in the right direction, which suggests that the method can capture the aggregate effect.

² The extrapolation method needs at least two gradients for prediction. This is why in the first two iterations, m_t is 0.

³[31] uses this example to show that ADAM [18] fails to converge.