

Pairwise Learning with Adaptive Online Gradient Descent*

Tao Sun

*College of Computer
National University of Defense Technology
Changsha, Hunan, China*

suntao.saltfish@outlook.com

Qingsong Wang

*Department of Mathematics
Scientific Computing and Imaging (SCI) Institute
University of Utah
Salt Lake City, Utah, USA*

wang.8973@osu.edu

Yunwen Lei

*Department of Mathematics
The University of Hong Kong
Pokfulam, Hong Kong*

leiyw@hku.hk

Dongsheng Li

*College of Computer
National University of Defense Technology
Changsha, Hunan, China*

dsli@nudt.edu.cn

Bao Wang

*Department of Mathematics
Scientific Computing and Imaging (SCI) Institute
University of Utah
Salt Lake City, Utah, USA*

wangbaonj@gmail.com

Reviewed on OpenReview: <https://openreview.net/forum?id=rq1SaHQg2k>

Abstract

In this paper, we propose an adaptive online gradient descent method with momentum for pairwise learning, in which the step size is determined by historical information. Due to the structure of pairwise learning, the sample pairs are dependent on the parameters, causing difficulties in the convergence analysis. To this end, we develop novel techniques for the convergence analysis of the proposed algorithm. We show that the proposed algorithm can output the desired solution in strongly convex, convex, and nonconvex cases. Furthermore, we present theoretical explanations for why our proposed algorithm can accelerate previous workhorses for online pairwise learning. All assumptions used in the theoretical analysis are mild and common, making our results applicable to various pairwise learning problems. To demonstrate the efficiency of our algorithm, we compare the proposed adaptive method with the non-adaptive counterpart on the benchmark online AUC maximization problem.

1 Introduction

Let $\mathcal{K} \subseteq \mathbb{R}^d$ be a closed convex set (can be the full space \mathbb{R}^d) representing the parameter space. Given a statistical sample space Ξ with probability distribution \mathcal{P} ; let $F(\cdot; \xi, \xi') : \mathbb{R}^d \rightarrow \mathbb{R}$ be a closed function

*The first and second authors contributed equally to this paper. Dongsheng Li is the corresponding author. Dongsheng Li and Tao Sun are supported in part by the National Science Foundation of China (62025208), Hunan Provincial Natural Science Foundation of China (2022JJ10065), Young Elite Scientists Sponsorship Program by CAST (2022QNRC001), and Continuous Support of PDL (WDZC20235250101).

associated with two samples $\xi, \xi' \in \Xi$. This paper considers the following pairwise learning problem

$$\min_{\mathbf{x} \in \mathcal{K} \subseteq \mathbb{R}^d} \left\{ f(\mathbf{x}) := \mathbb{E}_{(\xi, \xi') \sim \mathcal{P} \oplus \mathcal{P}} F(\mathbf{x}; \xi, \xi') \right\}, \quad (1)$$

where the function $F(\mathbf{x}; \xi, \xi')$ can be either convex or nonconvex in \mathbf{x} . The pairwise learning model (1) describes various classical machine learning tasks, including the metric learning (Weinberger & Saul, 2009; Kulis et al., 2013; Xing et al., 2002; Ying & Li, 2012), ranking (Rejchel, 2012; Agarwal & Niyogi, 2009), two-stage multiple kernel learning (Kumar et al., 2012), neural link prediction (Wang et al., 2021), the minimum error entropy principle (Hu et al., 2013), and can be adapted to the area under ROC curve (AUC) maximization (Zhao et al., 2011; Gao et al., 2013; Ying et al., 2016; Liu et al., 2018) (see Remark 2 in Section 3 for more details).

There are two major kinds of workhorses for the model (1), i.e., *offline* and *online*. The offline one is similar to the empirical risk minimization (ERM): given an i.i.d. sample set $\{\xi^1, \dots, \xi^n\}$, we solve

$$\min_{\mathbf{x} \in \mathcal{K}} \frac{1}{n(n-1)} \sum_{i, j \in [n], i \neq j} F(\mathbf{x}; \xi^i, \xi^j),$$

where $[n] := \{1, 2, \dots, n\}$. The major difference between the offline pairwise learning model and the ERM lies in the efficiency of the samples and whether the objective functions are independent of each other. A n -samples training set outputs a finite-sum minimization with $\mathcal{O}(n)$ sub-functions in ERM, while the same training set results in $\mathcal{O}(n^2)$ in the offline pairwise learning. Furthermore, the objective functions in the ERM are independent of each other, which is broken for the offline pairwise learning¹. One possible modification to circumvent the dependent objective functions is to form the objective function with two new independent samples in each iteration. This method is suggested in (Peel et al., 2010, Section 4.2) and can be regarded as the algorithm proposed by Zhao et al. (2011) with buffer size one. However, as shown in (Zhao et al., 2011), this two-dependent-points version of online learning does not fully utilize the sampling sequence and results in inferior performance compared with algorithms that utilize some historical samples.

The online pairwise learning assumes the i.i.d. samples $(\xi^k)_{k \in [n]}$ are continuously received by the model. In the k th iteration, the online style algorithm proceeds to sample new data ξ^k from \mathcal{P} and reuses the previous samples $(\xi^{j_i})_{1 \leq i \leq s}$ with $\{j_i\}_{1 \leq i \leq s} \subseteq [k-1]$ to get the mini-batch stochastic gradient $\mathbf{g}^k = \frac{1}{s} \sum_{i=1}^s \nabla F(\mathbf{x}^k; \xi^k, \xi^{j_i})$ (Wang et al., 2012; Zhao et al., 2011; Ying & Zhou, 2016). Thus, the online method needs to employ an $\mathcal{O}(s)$ memory to store the previously sampled data and $\mathcal{O}(s)$ computations to calculate the gradient. Nevertheless, the mini-batch version suffers two drawbacks: 1) It has been proved that its excess generalization bound can be as large as $\mathcal{O}\left(\frac{1}{\sqrt{s}} + \frac{1}{\sqrt{n}}\right)$ (Wang et al., 2012), and we need to use a large s to improve generalization of the online method. 2) A large s causes tremendous computational and memory costs that are unacceptable for online settings. To this end, a simple yet efficient online gradient descent (OGD) is proposed (presented as Algorithm 1), in which the stochastic gradient \mathbf{g}^k is set to be $\nabla F(\mathbf{x}^k; \xi^k, \xi^{k-1})$ (Yang et al., 2021b). An interesting finding is that the OGD can achieve $\mathcal{O}\left(\frac{1}{\sqrt{n}}\right)$ excess generalization bound. The favorable memory and computation costs make the OGD applicable to broader online settings.

In this paper, we focus on developing provably convergent online algorithms with adaptive stepsizes for pairwise learning. Considering the efficiency of the sampling method of OGD (Yang et al., 2021b), our algorithm inherits such kind of sampling. Moreover, our algorithm employs momentum.

1.1 The Adaptive Online Gradient Descent for Pairwise Learning

The OGD performs an SGD-style (stochastic gradient descent) iteration but with biased stochastic gradients since ξ^{k-1} is related to \mathbf{x}^k . Motivated by the remarkable success of the adaptive variants of SGD for machine learning (Duchi et al., 2011; McMahan & Streeter, 2010; Tieleman & Hinton, 2012; Kingma & Ba, 2015; Reddi et al., 2018; Ward et al., 2019), we propose the adaptive variant of OGD (AOGD) for pairwise learning,

¹For example, $F(\mathbf{x}; \xi^1, \xi^2)$ and $F(\mathbf{x}; \xi^1, \xi^3)$ are not independent because they have a shared data ξ^1 .

Algorithm 1 Online Gradient Descent (OGD) for Pairwise Learning (Yang et al., 2021b)

Parameters: $\eta > 0$.**Initialization:** $\mathbf{x}^0 = \mathbf{0}$, $\xi^1 \sim \mathcal{P}$ **for** $k = 1, 2, 3, \dots$ **step 1:** receive ξ^k and calculate $\mathbf{g}^k = \nabla F(\mathbf{x}^k; \xi^k, \xi^{k-1})$ **step 2:** $\mathbf{x}^{k+1} = \mathbf{Proj}_{\mathcal{K}}(\mathbf{x}^k - \eta \mathbf{g}^k)$ **End for**

Algorithm 2 Adaptive Online Gradient Descent (AOGD) for Pairwise Learning

Parameters: $\eta > 0$, $0 \leq \theta < 1$.**Initialization:** $\mathbf{x}^0 = \mathbf{m}^0 = \mathbf{0}$, $\xi^1 \sim \mathcal{P}$ **for** $k = 1, 2, 3, \dots$ **step 1:** receive ξ^k and calculate $\mathbf{g}^k = \nabla F(\mathbf{x}^k; \xi^k, \xi^{k-1})$ **step 2:** $\mathbf{m}^k = \theta \mathbf{m}^{k-1} + (1 - \theta) \mathbf{g}^k$ **step 3:** $v^k = v^{k-1} + \|\mathbf{g}^k\|^2$ **step 4:** $\mathbf{x}^{k+1} = \mathbf{Proj}_{\mathcal{K}}(\mathbf{x}^k - \eta \mathbf{m}^k / (v^k)^{\frac{1}{2}})$ **End for**

presented as Algorithm 2. AOGD directly uses the historical sum of the moment rather than the weighted average form used in Adam (Kingma & Ba, 2015), and AOGD can be rewritten in the weighted average form: let $\hat{v}^k := \sum_{i=1}^k \|\mathbf{g}^i\|^2 / k = v^k / k$, then steps 3 and 4 of AOGD can be reformulated as follows:

$$\begin{aligned} \hat{v}^k &= \left(1 - \frac{1}{k}\right) \hat{v}^{k-1} + \frac{1}{k} \|\mathbf{g}^k\|^2, \\ \mathbf{x}^{k+1} &= \mathbf{Proj}_{\mathcal{K}}\left(\mathbf{x}^k - \frac{\eta}{\sqrt{k}} \frac{\mathbf{m}^k}{(\hat{v}^k)^{1/2}}\right). \end{aligned} \quad (2)$$

In reformulation (2), weights of the moment are $\{1/k\}_{k \geq 1}$ and stepsizes are $\{\eta/\sqrt{k}\}_{k \geq 1}$, satisfying the sufficient conditions for the convergence of Adam-type algorithms (Zou et al., 2019; Chen et al., 2019).

Compared with OGD, AOGD employs the momentum and adaptive stepsize generated by the historical information. Thus unlike OGD, AOGD needs a memory of history, but without much computational overhead since AOGD only computes gradient once in each iteration. Another difference between OGD and AOGD lies in the hyper-parameter η : in OGD, η shall be set as small as the desired error ϵ ; while in AOGD, η can be set as a constant that is independent of ϵ .

1.2 Comparison with A Closely Related Work

In (Ding et al., 2015), by adopting the AdaGrad-style adaptive gradient update, Duchi et al. (2011) have proposed an adaptive method for online AUC maximization, which is a kind of pairwise learning. Although both (Ding et al., 2015) and our paper consider the adaptive online method for pairwise learning, there are four major differences between (Ding et al., 2015) and our paper, summarized below.

- 1) We follow the sampling method used by the OGD algorithm in (Yang et al., 2021b).
- 2) We consider more general cases and provide the corresponding theoretical guarantees, including the more general model (pairwise learning rather than only AUC maximization), and convergence for more general settings, including the nonconvex case, and more general schemes e.g., the use of momentum.
- 3) We get rid of using the regret bound because it does not directly tell us whether the algorithm converges to the desired minimizer. Another important reason for not using the regret bound for the analysis is that the regret bound has difficulties in covering the nonconvex cases. Based on the above reasons, a non-regret analysis is necessary.
- 4) We develop new analysis techniques to get the non-regret convergence analysis. Notice that the regret bound is not affected by the stochasticity of the data, and thus the analysis in (Ding et al., 2015) does not need to consider how to deal with the biased stochastic gradients.

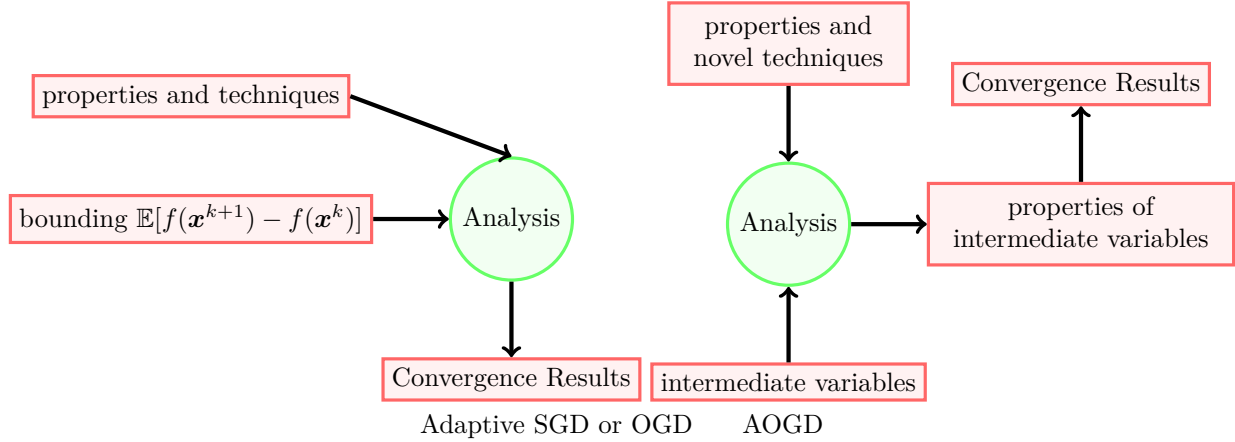


Figure 1: Contrasting the analysis of AOGD against adaptive SGD/OGD. Our analysis of AOGD studies the properties of some intermediate variables. To this end, we develop novel techniques. In the last step, we use the properties of the intermediate variables to derive the convergence guarantee for AOGD.

1.3 Challenges in the Analysis and Difference from Existing Analysis

The primary source of challenges in theoretical analysis comes from the fact that \mathbf{x}^k is dependent on ξ^{k-1} , which immediately breaks the unbiased expectation of the stochastic gradient $\nabla F(\mathbf{x}^k; \xi^k, \xi^{k-1})$, i.e., $\mathbb{E}[\nabla F(\mathbf{x}^k; \xi^k, \xi^{k-1})] \neq \nabla f(\mathbf{x}^k)$. As such, we cannot directly follow the techniques from adaptive SGD (Duchi et al., 2011; Reddi et al., 2018; Chen et al., 2019; Ward et al., 2019; Zou et al., 2019; Li & Orabona, 2019). In paper (Yang et al., 2021b), the authors consider the following decomposition

$$\eta \nabla F(\mathbf{x}^k; \xi^k, \xi^{k-1}) = \eta \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) + \eta \nabla F(\mathbf{x}^k; \xi^k, \xi^{k-1}) - \eta \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}).$$

Notice that \mathbf{x}^{k-1} is independent of ξ^k, ξ^{k-1} , and $\mathbb{E}[\eta \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})] = \eta \nabla f(\mathbf{x}^{k-1})$. The Lipschitz property of the gradient then gives us

$$\left\| \eta \nabla F(\mathbf{x}^k; \xi^k, \xi^{k-1}) - \eta \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) \right\| = \mathcal{O}(\eta \|\mathbf{x}^k - \mathbf{x}^{k-1}\|),$$

and the proof is similar to the “delayed” SGD.

However, our proof cannot directly follow the technique above because AOGD involves two extra recipes, i.e., momentum and adaptive stepsize. When momentum exists, we need to deal with both \mathbf{g}^k and \mathbf{m}^k rather than only \mathbf{g}^k , and we have to modify the techniques from (Yang et al., 2021b) to analyze the effects of momentum². Another difficulty is the use of the adaptive stepsize variable $\eta/(v^k)^{\frac{1}{2}}$. Furthermore, v^k is dependent of the pair (ξ^k, ξ^{k-1}) , making the analysis more challenging.

In contrast to the existing analysis, our approach is not directly establishing the Lyapunov descent starting from bounding $\mathbb{E}[f(\mathbf{x}^{k+1}) - f(\mathbf{x}^k)]$. Instead, we recruit some intermediate variables. Taking the nonconvex case as an example, we first establish a *Lyapunov-like* descent property as follows

$$\mathbb{E} \langle -\nabla f(\mathbf{x}^k), \mathbf{m}^k / (v^k)^{\frac{1}{2}} \rangle \leq \theta \mathbb{E} \langle -\nabla f(\mathbf{x}^{k-1}), \mathbf{m}^{k-1} / (v^{k-1})^{\frac{1}{2}} \rangle + \varphi(\mathbf{x}^k, \mathbf{x}^{k-1}, \xi^k, \xi^{k-1}, \xi^{k-2}),$$

where $\varphi(\mathbf{x}^k, \mathbf{x}^{k-1}, \xi^k, \xi^{k-1}, \xi^{k-2})$ is a function of variables $\mathbf{x}^k, \mathbf{x}^{k-1}, \xi^k, \xi^{k-1}$, and ξ^{k-2} . We stress that the descent is not Lyapunov because $\theta \neq 1$ and $\varphi(\mathbf{x}^k, \mathbf{x}^{k-1}, \xi^k, \xi^{k-1}, \xi^{k-2})$ is not always negative. Then, we build the correspondence between the mathematical convergence measurement and $\mathbb{E}[\langle -\nabla f(\mathbf{x}^k), \mathbf{m}^k / (v^k)^{\frac{1}{2}} \rangle]$. A big picture of the difference in analyzing adaptive SGD/OGD and AOGD is presented in Figure 1.

1.4 Contributions

Our major contributions are threefold, which are summarized below.

²Indeed, in the conclusion part of (Yang et al., 2021b), the authors have listed the momentum variant as future work.

- We propose an adaptive online gradient descent algorithm for pairwise learning with a simple sampling strategy. The proposed algorithm uses adaptive stepsize and momentum, requiring only a small overhead in memory and computational costs.
- We present the convergence results for the proposed algorithm under different settings, including strongly convex, general convex, and nonconvex cases. The use of adaptive stepsize and momentum requires non-trivial techniques for the convergence analysis. We also provide theoretical explanations for why our proposed algorithm can accelerate OGD.
- We verify the efficiency of the proposed AOGD on the benchmark online AUC maximization task, showing that AOGD outperforms OGD.

1.5 Notation

Throughout this paper, we use boldface letters to denote vectors, e.g., $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$. The j th coordinate of the vector \mathbf{x} is denoted by x_j . The L_2 norm of the vector \mathbf{x} is denoted by $\|\mathbf{x}\|$. We denote $\mathbb{E}[\cdot]$ as the expectation with respect to the underlying probability space. We denote the minimum value of the function f over \mathcal{K} as $\min_{\mathcal{K}} f$, and denote $\mathbf{Proj}_{\mathcal{K}}(\mathbf{x})$ as the projection of \mathbf{x} onto the set \mathcal{K} . For two positive sequences $(a_k, b_k)_{k \geq 0}$, $a_k = \mathcal{O}(b_k)$ means that there exists $C > 0$ such that $a_k \leq Cb_k$. The notation $a_k = \Theta(b_k)$ indicates that $a_k = \mathcal{O}(b_k)$ and $b_k = \mathcal{O}(a_k)$. We use $a_k = \tilde{\mathcal{O}}(b_k)$ and $a_k = \tilde{\Theta}(b_k)$ to hide the logarithmic factor but still with the same order. We use $a_k \geq \Theta(b_k)$ to present the relation $a_k \geq Cb_k$ with $C > 0$.

1.6 Organization

We organize this paper as follows: In Section 2, we present assumptions and theoretical convergence results for AOGD under general convex, strongly convex, and nonconvex settings. We numerically verify the efficiency of AOGD and compare it to the benchmark OGD in Section 3. More related works are discussed in Section 4, followed by concluding remarks. The detailed proofs are provided in the supplementary materials.

2 Convergence Analysis

2.1 Assumptions

We first collect several necessary assumptions for the convergence analysis of AOGD:

- **Assumption 1:** *The function $F(\cdot; \xi, \xi')$ is differentiable, and its gradient is L -Lipschitz, i.e.,*

$$\|\nabla F(\mathbf{x}; \xi, \xi') - \nabla F(\mathbf{y}; \xi, \xi')\| \leq L\|\mathbf{x} - \mathbf{y}\|, \forall \mathbf{x}, \mathbf{y} \in \mathcal{K}, \xi, \xi' \in \Xi. \quad (3)$$

- **Assumption 2:** *The gradient of $F(\mathbf{x}; \xi, \xi')$ is uniformly bounded, i.e., $\|\nabla F(\mathbf{x}; \xi, \xi')\| \leq B$ for some constant $B > 0$, $\forall \mathbf{x} \in \mathcal{K}$, and $\xi, \xi' \in \Xi$.*

The Lipschitz smooth gradient assumption is widely used in the (non)convex optimization and pairwise learning communities. While Assumption 2 is frequently used in the adaptive SGD community, see, e.g., (Duchi et al., 2011; Reddi et al., 2018; Chen et al., 2019; Ward et al., 2019; Zou et al., 2019; Li & Orabona, 2019)³. Note that when \mathcal{K} is bounded – Duchi et al. (2011); Reddi et al. (2018) have assumed the boundedness of the constrained set – Assumption 2 directly holds for the continuity of the gradient⁴, but not vice versa. Moreover, Assumption 2 indicates the following estimate of $\nabla f(\mathbf{x})$:

$$\|\nabla f(\mathbf{x})\| = \|\mathbb{E}_{\xi, \xi' \sim \mathcal{P} \oplus \mathcal{P}} \nabla F(\mathbf{x}; \xi, \xi')\| \leq \mathbb{E}_{\xi, \xi' \sim \mathcal{P} \oplus \mathcal{P}} \|\nabla F(\mathbf{x}; \xi, \xi')\| \leq B.$$

Using mathematical induction, we can see that the momentum \mathbf{m}^k also enjoys the uniform bound according to Assumption 2. Furthermore, we stress that we do not need to assume the variance is bounded, which is indicated by Assumption 2 since we have

$$\mathbb{E}\|\nabla F(\mathbf{x}; \xi, \xi') - \nabla f(\mathbf{x})\|^2 \leq \mathbb{E}\|\nabla F(\mathbf{x}; \xi, \xi')\|^2 \leq B^2.$$

³The uniform assumption presented in (Li & Orabona, 2019) enjoys another presentation, i.e., (Assumption (H2)) presented as $|f(\mathbf{x}) - f(\mathbf{y})| \leq G\|\mathbf{x} - \mathbf{y}\|$. Note that when f is differentiable, it is equivalent to $\sup_{\mathbf{x}} \|\nabla f(\mathbf{x})\| \leq G$. With bounded assumption for $F(\mathbf{x}; \xi) - \nabla f(\mathbf{x})$ in (Li & Orabona, 2019), which is indeed the uniform bound assumption.

⁴Any continuous function is uniformly bounded over a closed bounded subset of \mathbb{R}^d .

Assumptions 1 and 2 will be used in the analysis of AOGD for all different scenarios in the subsequent analysis.

2.2 General Convex Cases

In this subsection, we present the convergence result of AOGD for the general convex case, i.e., $F(\mathbf{x}; \xi, \xi')$ is convex with respect to \mathbf{x} and any fixed ξ, ξ' . Note that the convexity of $F(\mathbf{x}; \xi, \xi')$ indicates the convexity of $f(\mathbf{x})$ in (1), but not vice versa.

Theorem 1 (General Convexity) *Let Assumption 1 hold, and $\|\mathbf{g}^0\|^2 \geq \delta > 0$ for some constant δ , and $F(\mathbf{x}; \xi, \xi')$ be convex. Assume $\{\mathbf{x}^k\}_{k \geq 1}$ is generated by AOGD for pairwise learning, $\mathbf{x}^* := \arg \min_{\mathbf{x} \in \mathcal{K}} f(\mathbf{x})$, and \mathcal{K} is additionally bounded, i.e., $\max_{\mathbf{x}, \mathbf{y} \in \mathcal{K}} \|\mathbf{x} - \mathbf{y}\| \leq D$. Then we have*

$$\mathbb{E} \left[f \left(\frac{\sum_{k=1}^K \mathbf{x}^k}{K} \right) - \min_{\mathcal{K}} f \right] \leq \frac{c_1 + c_2(K)}{K}, \quad (4)$$

where $c_1 := \frac{\mathbb{E} \sqrt{v^1} \|\mathbf{x}^* - \mathbf{x}^1\|^2}{2(1-\theta)\eta} + \frac{3B^2}{\sqrt{\delta}}$, and $c_2(K) := \left[\frac{\eta + 2D^2}{2(1-\theta)} + (\eta\theta + 1) \right] \cdot \mathbb{E} \sqrt{v^{K+1}} + \left[4L + \frac{2L^2}{\sqrt{\delta}} \right] \ln \frac{\mathbb{E} \sqrt{v^{K+1}}}{\sqrt{\delta}}$.

Assumption 2 is implicitly used in Theorem 1 because we have assumed the boundedness of the constrained set \mathcal{K} . The bounded constrained set is indeed stronger than the uniform bounded gradient assumption. We leave how to relax the bounded set assumption as future work. From Theorem 1, we can see that the convergence rate is dependent on $\mathbb{E}[\sqrt{v^K}]$. The boundedness of the stochastic gradients indicates that $\mathbb{E}[\sqrt{v^K}] = \mathcal{O}(\sqrt{K})$, which means the worst convergence rate of AOGD is $\tilde{\mathcal{O}}\left(\frac{1}{\sqrt{K}}\right)$. We notice that the convergence rate $\tilde{\mathcal{O}}\left(\frac{1}{\sqrt{K}}\right)$ coincides with the rate of OGD for pairwise learning in the general convex case (Yang et al., 2021b). However, in some cases $\mathbb{E}[\sqrt{v^K}]$ can decay faster than $\mathcal{O}(\sqrt{K})$, based on which we can establish an accelerated rate of AOGD.

Proposition 1 *Assume the conditions of Theorem 1 hold, and assume $\mathbb{E}[\sqrt{v^K}] = \mathcal{O}(K^\alpha)$ with $0 < \alpha \leq \frac{1}{2}$. Then we have*

$$\mathbb{E} \left[f \left(\frac{\sum_{k=1}^K \mathbf{x}^k}{K} \right) - \min_{\mathcal{K}} f \right] = \tilde{\mathcal{O}} \left(\frac{1}{K^{1-\alpha}} \right). \quad (5)$$

If $\alpha < \frac{1}{2}$, the convergence rate of AOGD is thus better than $\tilde{\mathcal{O}}\left(\frac{1}{\sqrt{K}}\right)$. Proposition 1 then provides a theoretical interpretation of why AOGD is possible to be faster than OGD (Yang et al., 2021b).

Remark 1 *Note that the condition $\alpha = 1/2$ directly holds due to the boundedness of the stochastic gradients. The fast decaying condition $\alpha < 1/2$ has been commonly used as a standard explanation for why adaptive stochastic optimization algorithms often outperform non-adaptive schemes, as shown in several studies (Reddi et al., 2018; Chen et al., 2018; 2019; Liu et al., 2019). As far as our current knowledge goes, the superiority of adaptive SGD over non-adaptive approaches is not well-explained, apart from the assumption that $\alpha < 1/2$. This assumption is commonly used in previous works because many training tasks involve sparse stochastic gradients. However, it is important to note that this assumption is not a logical consequence of sparse gradients resulting in $\alpha < 1/2$. Instead, the assumption is more of a hypothetical analysis. In summary, while the assumption $\alpha < 1/2$ is widely employed in the literature, we do not have a definitive explanation for the effectiveness of adaptive SGD beyond this assumption. Further research is required to gain a comprehensive understanding of the benefits of adaptive SGD compared to non-adaptive methods.*

2.3 Strongly Convex Cases

In this subsection, we consider the case that the function $F(\mathbf{x}; \xi, \xi')$ is ν -strongly convex for some constant $\nu > 0$, i.e., $F(\mathbf{x}; \xi, \xi') - F(\mathbf{y}; \xi, \xi') - \langle \nabla F(\mathbf{y}; \xi, \xi'), \mathbf{x} - \mathbf{y} \rangle \geq \frac{\nu}{2} \|\mathbf{x} - \mathbf{y}\|^2$. In particular, if $\nu = 0$, $F(\mathbf{x}; \xi, \xi')$ then reduces to general convex.

Stepsize rule for the strongly convex case. Before developing a convergence guarantee of AOGD for strongly convex cases, we need to select the appropriate stepsize rule for AOGD. We first need to explain the stepsize rule of AOGD for the general convex case, i.e., $\eta/\sqrt{v^k}$ for some constant $\eta > 0$, is inappropriate for solving strongly convex problems. The boundedness of gradient directly gives us $\frac{1}{\sqrt{v^k}} \geq \Theta\left(\frac{1}{\sqrt{k}}\right)$. However, such a stepsize rule causes the accumulation of stochastic noise; it will not improve the convergence rate of AOGD under strong convexity compared to the rate of AOGD under general convexity. A proper stepsize choice for the strongly convex case is $\Theta\left(\frac{1}{k}\right)$ (Bach & Moulines, 2013). Indeed, in papers (Duchi et al., 2011; Sun et al., 2020), the authors use the $\frac{1}{\sqrt{k}\sqrt{v^k}}$ stepsize rule for AdaGrad when the problem is strongly convex, and we follow this stepsize rule for the strongly convex online pairwise learning. Note that $\frac{1}{\sqrt{k}\sqrt{v^k}} \geq \Theta\left(\frac{1}{k}\right)$, which coincides with the stepsize used for SGD when the underlying problem is strongly convex. In summary, in the strongly convex case, we replace **step 4** of Algorithm 2 with **step 4'**, which is given below

$$\text{step 4'} : \mathbf{x}^{k+1} = \text{Proj}_{\mathcal{K}}\left(\mathbf{x}^k - \frac{\eta}{\sqrt{k}} \mathbf{m}^k / \sqrt{v^k}\right). \quad (6)$$

Next, we present the convergence rate of AOGD for pairwise learning with the strong convexity assumption.

Theorem 2 *Let Assumption 1 hold, and $\|\mathbf{g}^0\|^2 \geq \delta > 0$, and $F(\mathbf{x}; \xi, \xi')$ be strongly convex. Assume $\{\mathbf{x}^k\}_{k \geq 1}$ is generated by the AOGD for pairwise learning with $\theta = 0$ using stepsize rule **step 4'**, and $\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathcal{K}} f(\mathbf{x})$. By setting $\eta = \frac{B}{2\nu}$, then we have*

$$\mathbb{E}[\|\mathbf{x}^K - \mathbf{x}^*\|^2] = \mathcal{O}\left(\frac{\ln K}{K}\right). \quad (7)$$

In the strongly convex case, Assumption 2 is equivalent to the boundedness assumption of the constrained set \mathcal{K} . This is because the function $f(\mathbf{x})$ is also ν -strongly convex⁵, yielding $\|\nabla f(\mathbf{x})\| = \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{x}^\dagger)\| \geq \nu\|\mathbf{x} - \mathbf{x}^\dagger\|$ with \mathbf{x}^\dagger being the global minimizer of f . When Assumption 2 holds, in Subsection 2.1 we have shown that $\|\nabla f(\mathbf{x})\| \leq B$ over \mathcal{K} . Thus, we get $\|\mathbf{x}\| \leq \|\mathbf{x} - \mathbf{x}^\dagger\| + \|\mathbf{x}^\dagger\| \leq \frac{B}{\nu} + \|\mathbf{x}^\dagger\|$ when $\mathbf{x} \in \mathcal{K}$. Notice that \mathbf{x}^\dagger is fixed, and the set \mathcal{K} is then uniformly bounded.

Theorem 2 above shows that AOGD can achieve a faster convergence rate under the strong convexity assumption than that for general convex cases. Theorem 2 also shows that in the strongly convex case, AOGD achieves an almost optimal convergence rate of SGD under strong convexity, specifically $\tilde{\mathcal{O}}\left(\frac{1}{K}\right)$ (the optimal convergence rate is $\mathcal{O}\left(\frac{1}{K}\right)$ according to (Rakhlin et al., 2012)). The result in Theorem 2 does not generalize to the general convex case since we set $\eta = \frac{B}{2\nu}$, which is infinity when $\nu = 0$. For technical reasons, we set $\theta = 0$ in Theorem 2, i.e., we only consider the momentum-free case. We will consider how to build the convergence rate $\tilde{\mathcal{O}}\left(\frac{1}{K}\right)$ for AOGD with momentum in our future work.

2.4 Nonconvex Cases

In this part, we consider the case when $F(\mathbf{x}; \xi, \xi')$ is nonconvex. The assumptions for the nonconvex case are much milder than the convex and strongly convex cases. We can even get rid of using the projection operator $\text{Proj}_{\mathcal{K}}(\cdot)$ for AOGD. The convergence result of AOGD for the nonconvex case is presented as follows.

Theorem 3 *Let Assumptions 1, 2 hold, and let $\{\mathbf{x}^k\}_{k \geq 1}$ be generated by AOGD, and $\|\mathbf{g}^0\|^2 \geq \delta$ for some constant $\delta > 0$. Suppose $\sqrt{v^k} \leq C \cdot k^\alpha$ for two constants $C > 0$ and $0 < \alpha \leq \frac{1}{2}$, and \mathcal{K} is the full space. Then, we have*

$$\min_{1 \leq k \leq K} \left\{ \mathbb{E} \|\nabla f(\mathbf{x}^k)\|^2 \right\} \leq \frac{c_3 + c_4(K)}{K^{1-\alpha}}, \quad (8)$$

⁵Taking expectation of both sides of the inequality $F(\mathbf{x}; \xi, \xi') - F(\mathbf{y}; \xi, \xi') - \langle \nabla F(\mathbf{y}; \xi, \xi'), \mathbf{x} - \mathbf{y} \rangle \geq \frac{\nu}{2} \|\mathbf{x} - \mathbf{y}\|^2$ gives us the strong convexity of f .

where $c_3 := \frac{2(7-6\theta)B^2C}{(1-\theta)\sqrt{\delta}} + \frac{4Cf(\mathbf{x}^1)}{\eta}$, and $c_4(K) := \left(\frac{(1+\theta)\eta}{1-\theta} + 2L^2/\sqrt{\delta} + L^2/\delta + 5/2 \right) \ln \left(\frac{CK}{\delta} \right)$.

From Theorem 3, we can see that the convergence rate of AOGD is $\tilde{O}\left(\frac{1}{K^{1-\alpha}}\right)$, and $\alpha = \frac{1}{2}$ is the worst case due to the boundedness of the stochastic gradient. When $\alpha < \frac{1}{2}$, we get a faster convergence of AOGD compared with OGD or SGD in the general nonconvex case.

The conditions of Theorem 3 can also be satisfied by the general convex case without projection. Therefore, (8) also holds when $F(\mathbf{x}; \xi, \xi')$ is convex. However, (8) is weaker than (4) since the convergence rate in (8) is not established with respect to the function values.

3 Numerical Results

Table 1: Statistics of the dataset used for contrasting the performance of AOGD and OGD, where n is the number of samples in each dataset, and d is the number of features of each instance in a given dataset. All datasets come from the LIBSVM website (Chang & Lin, 2011), and they are used in (Yang et al., 2021b).

	diabetes	german	ijcnn1	letter	mnist	usps
n	768	1,000	49,990	15,000	60,000	7,291
d	8	24	22	161	780	256

In this section, we numerically validate our theoretical findings for both convex and nonconvex cases. To this end, we compare our proposed AOGD against the baseline OGD proposed in (Yang et al., 2021b)⁶ for pairwise learning in terms of generalization and rate of convergence with respect to the number of iteration. We also include the results of AdaOAM in (Ding et al., 2015), an adaptive online algorithm for AUC maximization. However, we note that the algorithm in (Ding et al., 2015) is not designed for general pairwise learning. Following (Yang et al., 2021b), we consider six benchmark datasets, summarized in Table 1. Also, following the data split strategy used in (Yang et al., 2021b), for the dataset with multiple classes, we convert the first half of classes to the positive class and the second half of classes to the negative class.

For each dataset, we use 80% of the data for training and the remaining 20% for testing. All the reported results are based on 25 runs with random shuffling. The generalization performance is reported using the average AUC score and standard deviation on the test data. To determine proper hyperparameters for OGD, AOGD, and AdaOAM, we conduct 5-fold cross-validation on the training sets: 1) for OGD, we select stepsizes $\eta_t = \eta \in 10^{[-5:5]}$ ⁷ and the parameter space \mathcal{K} is set to be the L^2 -ball centered at the origin with radius $R \in 10^{[-3,3]}$; 2) for AOGD, we let $\theta = 0.9$ and we select stepsizes $\eta_t = \eta \in 10^{[-5:5]}$ and the parameter space \mathcal{K} is also the L^2 -ball centered at the origin with radius $R \in 10^{[-3,3]}$; 3) for AdaOAM, we select stepsizes $\eta_t = \eta \in 10^{[-5:5]}$ and the parameter space \mathcal{K} is set to be the L^2 -ball centered at the origin with radius $R \in 10^{[-3,3]}$.

Convex case: We run experiments on AUC maximization using the following convex loss function

$$f(\mathbf{w}; (\mathbf{x}, y), (\mathbf{x}', y')) = \ell(\mathbf{w}^\top(\mathbf{x} - \mathbf{x}'))\mathbb{I}_{[y=1 \wedge y'=-1]} + \ell(\mathbf{w}^\top(\mathbf{x}' - \mathbf{x}))\mathbb{I}_{[y=-1 \wedge y'=1]},$$

where

$$\mathbb{I}_{[y=1 \wedge y'=-1]} = \begin{cases} 1 & \text{if } y = 1 \text{ and } y' = -1, \\ 0 & \text{otherwise.} \end{cases}$$

and similarly for $\mathbb{I}_{[y=-1 \wedge y'=1]}$. Also, ℓ is a surrogate loss function, e.g., the hinge loss

$$\ell(t) = (1-t)_+ = \begin{cases} 0 & \text{if } 1-t < 0, \\ 1-t & \text{otherwise.} \end{cases}$$

⁶In (Yang et al., 2021b), the authors have shown that OGD can remarkably outperform existing algorithms, including OLP (Kar et al., 2013), OAM_{gra} (Zhao et al., 2011), SGD_{pair} (Lei et al., 2020), and SPAUC (Lei & Ying, 2021).

⁷ $[-5 : 5]$ stands for integers in the interval $[-5, 5]$.

Remark 2 The loss function above follows the setup in (Yang et al., 2021b) which is designed for the AUC maximization problem with data stream (\mathbf{x}, y) sampled from a distribution that contains both positive and negative samples by utilizing the indicator function $\mathbb{I}_{[y=1 \wedge y'=-1]}$ and $\mathbb{I}_{[y=-1 \wedge y'=1]}$. We want to mention that the second term is not presented in the main text of (Yang et al., 2021b) but implicitly utilized their code ⁸ on which our implementation is based. In particular, in both OGD and AOGD, the loss is activated or non-zero whenever the consecutive sample pairs (\mathbf{x}, y) and (\mathbf{x}', y') are of opposite labels, i.e., $y = 1$ and $y' = -1$ or $y = -1$ and $y' = 1$.

We now provide further clarification on the pairwise model and AUC maximization: In model (1), the distribution for ξ and ξ' is assumed to be the same. This is in contrast to the early AUC maximization formulation, as seen in (Zhao et al., 2011), which does not have this property. In the early formulation, the positive and negative classes are treated differently, with potentially different distributions. However, we can introduce indicator functions and show that the AUC maximization can be represented as a special form of (1) as above.

A natural question that arises is why we would choose to reformulate the AUC maximization problem into the model (1). At first glance, there does not seem to be any apparent advantage to such a reformulation, as the algorithm only works when a negative sample is encountered along with a positive sample. In the worst-case scenario, if all the negative labels are encountered before the positive labels, the algorithm would almost output no predictions.

We provide a necessary explanation here:

- a) The worst-case scenario mentioned above is highly unlikely to occur because we assume that the data is independently and identically distributed from the underlying distribution.
- b) In the online setting, the labels are unknown in advance. While the work of Zhao et al. (2011) discusses the online scenario, they employ an “update buffer” data pre-processing step to divide the data and subsequently run the algorithm offline. Therefore, we cannot directly apply the formulation described in (Zhao et al., 2011). To address this, it is more flexible to utilize the formulation presented in (Yang et al., 2021b). This formulation is better suited for online learning, as it does not require an offline processing step like the one used in (Zhao et al., 2011).
- c) Existing results from (Yang et al., 2021b) demonstrate that the formulation (1) outperforms previous methods, even when a negative sample is encountered along with a positive sample. This is because previous methods often have high gradient complexity at each iteration, while the online complexity is very small in the proposed formulation.

In summary, the reformulation (1) allows for more flexibility in the optimization process and has the potential to yield better results in practical scenarios.

Figure 2 plots the AUC scores of AOGD, OGD, and AdaOAM against the number of iterations (in log scale)

9

Table 2: Average AUC scores \pm standard deviation with convex loss function across the six benchmark datasets listed in Table 1. The best results are highlighted in boldface.

Algorithm	diabetes	german	ijcnn1	letter	mnist	usps
AOGD	.831 \pm .027	.795 \pm .026	.934 \pm .002	.814 \pm .006	.931 \pm .002	.925 \pm .004
OGD	.831 \pm .030	.793 \pm .021	.934 \pm .002	.810 \pm .007	.932 \pm .001	.926 \pm .006
AdaOAM	.829 \pm .027	.792 \pm .028	.934 \pm .003	.814 \pm .009	.932 \pm .002	.924 \pm .005

Table 2 summarizes the generalization performance between AOGD and OGD. The results for AOGD and AdaOAM are obtained from our above experiment and the results for OGD are adapted from (Yang et al.,

⁸<https://github.com/zhenhuan-yang/simple-pairwise>

⁹Number of iterations N is given by $2^{N_{log}/2+4}$ where N_{log} is the log-scaled number of iterations. on the six benchmark datasets listed in Table 1. The numerical results on the six benchmark datasets show that AOGD converges faster than OGD and AdaOAM in general, confirming our theoretical results.

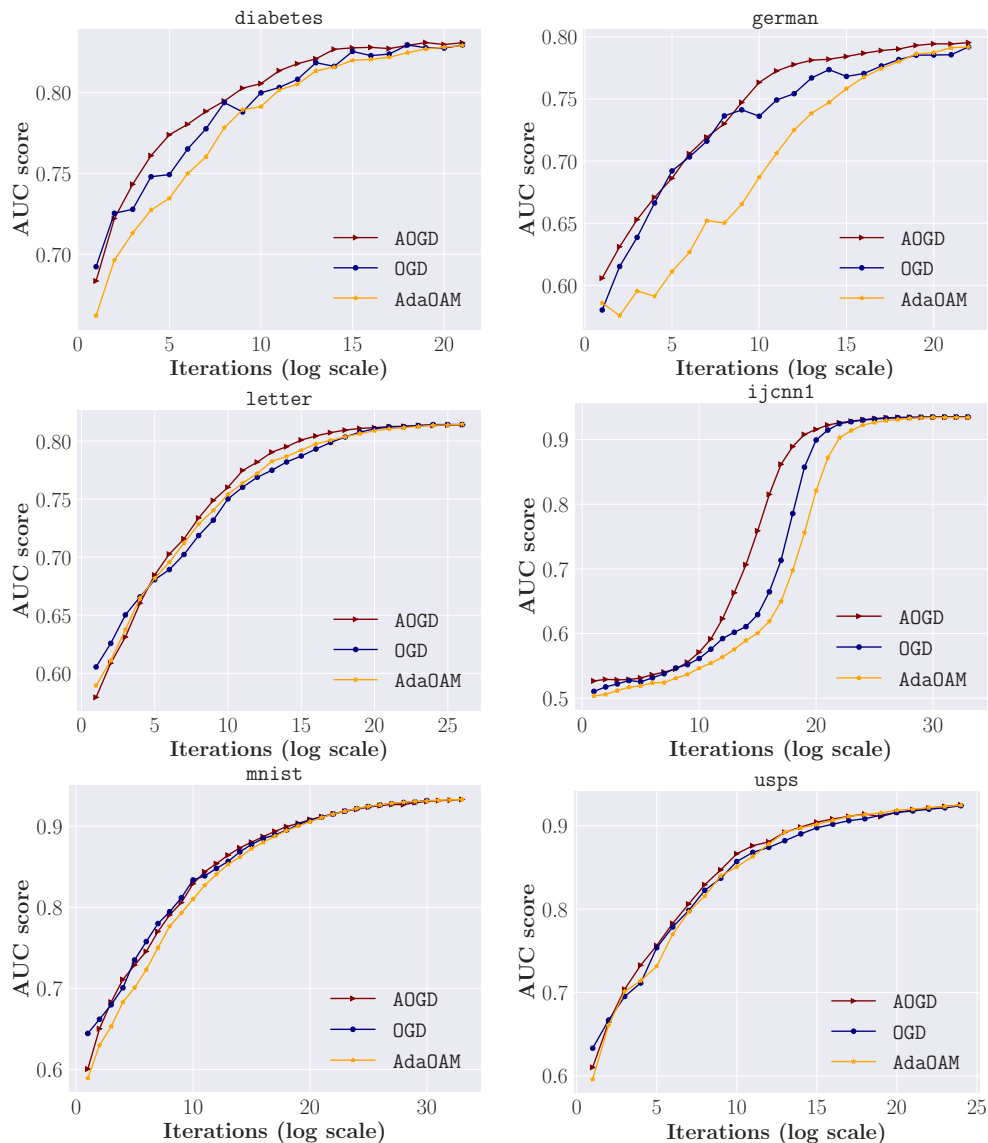


Figure 2: The AUC score of AOGD and OGD against the number of iterations (in log scale) for AUC maximization with **convex loss** (hinge loss). It is evident that AOGD converges faster than OGD and AdaOAM in general, confirming our established theoretical results for AOGD.

2021b). Overall, AOGD generalizes as well as OGD and AdaOAM as the peak performance of these three methods are comparable. Establishing the generalization of OAGD is an interesting future direction.

Non-convex case: We follow the setting in Appendix G of (Yang et al., 2021b) on AUC maximization using the logistic link function $\text{logit}(t) = (1 + \exp(-t))^{-1}$ and then the square loss function $\ell(t) = (1 - t)^2$. Hence, the loss function for the AUC maximization problem is given by the non-convex function

$$f(\mathbf{w}; (\mathbf{x}, y), (\mathbf{x}', y')) = (1 - \text{logit}(\mathbf{w}^\top (\mathbf{x} - \mathbf{x}')))^2 \mathbb{I}_{[y=1 \wedge y'=-1]} + (1 - \text{logit}(\mathbf{w}^\top (\mathbf{x}' - \mathbf{x})))^2 \mathbb{I}_{[y=-1 \wedge y'=1]}.$$

We plot in Figure 3 the AUC scores of AOGD, OGD, and AdaOAM against the number of iterations on the six benchmark datasets listed in Table 1. The numerical results on the six benchmark datasets show that AOGD converges faster than OGD and AdaOAM in general, confirming our theoretical results.

Table 3: Average AUC scores \pm standard deviation with nonconvex loss function across the six benchmark datasets listed in Table 1. The best results are highlighted in boldface.

Algorithm	diabetes	german	ijcnn1	letter	mnist	usps
AOGD	.834 \pm .022	.797 \pm .023	.935 \pm .002	.815 \pm .005	.932 \pm .002	.928 \pm .004
OGD	.829 \pm .033	.794 \pm .022	.934 \pm .002	.815 \pm .003	.931 \pm .002	.926 \pm .005
AdaOAM	.831 \pm .025	.789 \pm .021	.931 \pm .002	.815 \pm .006	.932 \pm .003	.926 \pm .007

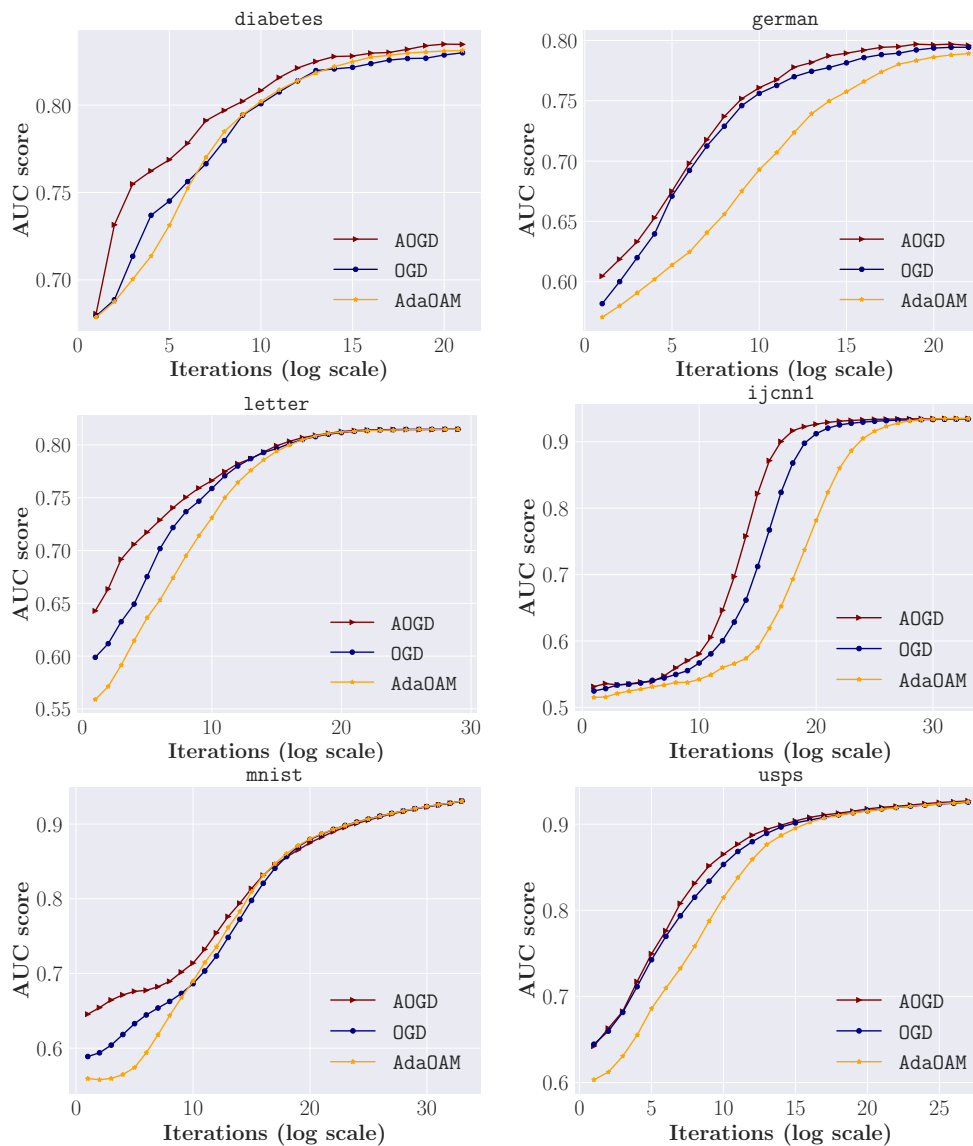


Figure 3: The AUC score of AOGD, OGD, and AdaOAM against the number of iterations (in log scale) with non-convex loss function. In this case, AOGD converges faster than OGD and AdaOAM in general, confirming our established theoretical results for AOGD. In particular, the slower convergence of AdaOAM is more prominent in the non-convex case.

4 More Related Works

Because offline methods employ the ERM-like policy, the core problem of most offline methods is establishing the generalization bound of the finite-sum model with statistical learning theory or algorithmic stability

(Agarwal & Niyogi, 2009; Jin et al., 2009; Wang et al., 2019; Gao & Zhou, 2013; Lei et al., 2020). It is worth mentioning that the difficulty in the generalization analysis for pairwise learning lies in that the objective functions fail to be i.i.d. with each other, which breaks the fundamental assumption in statistical learning theory and the algorithmic stability communities.

The online methods for pairwise learning assume the model accesses a data stream of i.i.d. samples, including online AUC maximization algorithms (Zhao et al., 2011; Ying et al., 2016; Liu et al., 2018; Natole et al., 2018; Lei & Ying, 2021; Guo et al., 2020), online metric learning algorithms (Shalev-Shwartz et al., 2004; Davis et al., 2007; Jain et al., 2008; Jin et al., 2009), online learning to rank (Rejchel, 2012; Schuth et al., 2013; Zoghi et al., 2017; Li et al., 2019), neural link prediction (Wang et al., 2021), etc. The online AUC maximization is proposed by Zhao et al. (2011) with theoretical guarantees. In the paper (Ying et al., 2016), a stochastic online AUC maximization algorithm is proposed from the perspective of a saddle representation. The main advantage of the algorithm in (Ying et al., 2016) is to avoid storing all previous examples and second-order covariance matrices. Leveraging saddle representation, Liu et al. (2018) propose a faster online AUC maximization algorithm with provably improved statistical convergence rates. The stochastic proximal algorithms for AUC maximization with non-differentiable regularization are proposed and studied in (Natole et al., 2018). To make the algorithm scalable to large-scale streaming data, Lei & Ying (2021) propose a new stochastic proximal algorithm. In the paper (Guo et al., 2020), the authors consider the distributed setting and propose a communication-efficient stochastic AUC maximization with deep neural networks. Shalev-Shwartz et al. (2004) propose an online algorithm for supervised learning of pseudo-metrics. In (Davis et al., 2007), the authors present an information-theoretic approach for online metric learning. In (Jain et al., 2008), leveraging the LogDet regularization, the authors propose a fast online metric learning for the similarity search. The generalization bound of regularized distance metric learning is established in (Jin et al., 2009). Rejchel (2012) consider ranking estimators that minimize the convex empirical risks and prove their generalization bounds. A framework of online learning to rank is proposed by Schuth et al. (2013). In the paper (Zoghi et al., 2017), the authors investigate online learning to rank in stochastic click models. Paper (Li et al., 2019) introduces a new model for online ranking with features. The differentially private pairwise learning has been recently studied, and representative works include (Huai et al., 2020; Yang et al., 2021a; Xue et al., 2021; Yang et al., 2021b).

5 Conclusions

In this paper, we propose adaptive online gradient descent algorithms to solve pairwise learning problems and establish their theoretical performance bounds in strongly convex, convex, and nonconvex settings. Our theoretical results explain why the convergence speed of adaptive online gradient descent can outperform the one without adaptive stepsize for pairwise learning. We also provide numerical experiments to demonstrate the efficiency of the proposed algorithm.

Limitation and future work. There are two major limitations in our analysis: 1) we assume the set \mathcal{K} is bounded in establishing Theorem 1, and 2) the convergence rate in Theorem 2 is analyzed for the adaptive online gradient descent without momentum. We leave how to overcome the above two limitations as future work. There are numerous other avenues for future work, including 1) Can we establish the lower bound of the convergence rates for the adaptive online gradient descent applied to pairwise learning? 2) Can we extend the online adaptive gradient descent to the proximal settings to solve nonsmooth pairwise learning problems?

Broader Impact Statement

N/A

References

- Shivani Agarwal and Partha Niyogi. Generalization bounds for ranking algorithms via algorithmic stability. *Journal of Machine Learning Research*, 10(2), 2009.
- Francis Bach and Eric Moulines. Non-strongly-convex smooth stochastic approximation with convergence rate $O(1/n)$. *Advances in neural information processing systems*, 26, 2013.
- Chih-Chung Chang and Chih-Jen Lin. LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2:27:1–27:27, 2011. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- Xiangyi Chen, Sijia Liu, Ruoyu Sun, and Mingyi Hong. On the convergence of a class of Adam-type algorithms for non-convex optimization. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=H1x-x309tm>.
- Zaiyi Chen, Zhuoning Yuan, Jinfeng Yi, Bowen Zhou, Enhong Chen, and Tianbao Yang. Universal stagewise learning for non-convex problems with convergence on averaged solutions. In *International Conference on Learning Representations*, 2018.
- Jason V Davis, Brian Kulis, Prateek Jain, Suvrit Sra, and Inderjit S Dhillon. Information-theoretic metric learning. In *Proceedings of the 24th international conference on Machine learning*, pp. 209–216, 2007.
- Yi Ding, Peilin Zhao, Steven Hoi, and Yew-Soon Ong. An adaptive gradient method for online auc maximization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 29, 2015.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(Jul):2121–2159, 2011.
- Wei Gao and Zhi-Hua Zhou. Uniform convergence, stability and learnability for ranking problems. In *Twenty-Third International Joint Conference on Artificial Intelligence*. Citeseer, 2013.
- Wei Gao, Rong Jin, Shenghuo Zhu, and Zhi-Hua Zhou. One-pass AUC optimization. In *International conference on machine learning*, pp. 906–914. PMLR, 2013.
- Zhishuai Guo, Mingrui Liu, Zhuoning Yuan, Li Shen, Wei Liu, and Tianbao Yang. Communication-efficient distributed stochastic AUC maximization with deep neural networks. In *International Conference on Machine Learning*, pp. 3864–3874. PMLR, 2020.
- Ting Hu, Jun Fan, Qiang Wu, and Ding-Xuan Zhou. Learning theory approach to minimum error entropy criterion. *Journal of Machine Learning Research*, 14(2), 2013.
- Mengdi Huai, Di Wang, Chenglin Miao, Jinhui Xu, and Aidong Zhang. Pairwise learning with differential privacy guarantees. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 694–701, 2020.
- Prateek Jain, Brian Kulis, Inderjit Dhillon, and Kristen Grauman. Online metric learning and fast similarity search. *Advances in neural information processing systems*, 21, 2008.
- Rong Jin, Shijun Wang, and Yang Zhou. Regularized distance metric learning: Theory and algorithm. *Advances in neural information processing systems*, 22, 2009.
- Purushottam Kar, Bharath Sriperumbudur, Prateek Jain, and Harish Karnick. On the generalization ability of online learning algorithms for pairwise loss functions. In Sanjoy Dasgupta and David McAllester (eds.), *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pp. 441–449, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR. URL <https://proceedings.mlr.press/v28/kar13.html>.

- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.
- Brian Kulis et al. Metric learning: A survey. *Foundations and Trends® in Machine Learning*, 5(4):287–364, 2013.
- Abhishek Kumar, Alexandru Niculescu-Mizil, Koray Kavukcoglu, and Hal Daumé. A binary classification framework for two-stage multiple kernel learning. In *Proceedings of the 29th International Conference on Machine Learning*, pp. 1331–1338, 2012.
- Yunwen Lei and Yiming Ying. Stochastic proximal AUC maximization. *Journal of Machine Learning Research*, 22(61):1–45, 2021.
- Yunwen Lei, Antoine Ledent, and Marius Kloft. Sharper generalization bounds for pairwise learning. *Advances in Neural Information Processing Systems*, 33:21236–21246, 2020.
- Shuai Li, Tor Lattimore, and Csaba Szepesvári. Online learning to rank with features. In *International Conference on Machine Learning*, pp. 3856–3865. PMLR, 2019.
- Xiaoyu Li and Francesco Orabona. On the convergence of stochastic gradient descent with adaptive stepsizes. *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 983–992, 2019.
- Mingrui Liu, Xiaoxuan Zhang, Zaiyi Chen, Xiaoyu Wang, and Tianbao Yang. Fast stochastic AUC maximization with $O(1/n)$ -convergence rate. In *International Conference on Machine Learning*, pp. 3189–3197. PMLR, 2018.
- Mingrui Liu, Youssef Mroueh, Jerret Ross, Wei Zhang, Xiaodong Cui, Payel Das, and Tianbao Yang. Towards better understanding of adaptive gradient algorithms in generative adversarial nets. In *International Conference on Learning Representations*, 2019.
- H Brendan McMahan and Matthew Streeter. Adaptive bound optimization for online convex optimization. *COLT 2010*, pp. 244, 2010.
- Michael Natole, Yiming Ying, and Siwei Lyu. Stochastic proximal algorithms for AUC maximization. In *International Conference on Machine Learning*, pp. 3710–3719. PMLR, 2018.
- Thomas Peel, Sandrine Anthoine, and Liva Ralaivola. Empirical bernstein inequalities for u-statistics. *Advances in Neural Information Processing Systems*, 23, 2010.
- Alexander Rakhlin, Ohad Shamir, and Karthik Sridharan. Making gradient descent optimal for strongly convex stochastic optimization. In *Proceedings of the 29th International Conference on Machine Learning*, pp. 1571–1578, 2012.
- Sashank J. Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of Adam and beyond. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=ryQu7f-RZ>.
- Wojciech Rejchel. On ranking and generalization bounds. *Journal of Machine Learning Research*, 13(5), 2012.
- Anne Schuth, Katja Hofmann, Shimon Whiteson, and Maarten De Rijke. Lerot: An online learning to rank framework. In *Proceedings of the 2013 workshop on Living labs for information retrieval evaluation*, pp. 23–26, 2013.
- Shai Shalev-Shwartz, Yoram Singer, and Andrew Y Ng. Online and batch learning of pseudo-metrics. In *Proceedings of the twenty-first international conference on Machine learning*, pp. 94, 2004.
- Tao Sun, Linbo Qiao, Qing Liao, and Dongsheng Li. Novel convergence results of adaptive stochastic gradient descents. *IEEE Transactions on Image Processing*, 30:1044–1056, 2020.
- Tijmen Tieleman and Geoffrey Hinton. Lecture 6.5-RMSProp, coursera: Neural networks for machine learning. *University of Toronto, Technical Report*, 2012.

- Boyu Wang, Hejia Zhang, Peng Liu, Zebang Shen, and Joelle Pineau. Multitask metric learning: Theory and algorithm. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 3362–3371. PMLR, 2019.
- Yuyang Wang, Roni Khardon, Dmitry Pechyony, and Rosie Jones. Generalization bounds for online learning algorithms with pairwise loss functions. In *Conference on Learning Theory*, pp. 13–1. JMLR Workshop and Conference Proceedings, 2012.
- Zhitao Wang, Yong Zhou, Litao Hong, Yuanhang Zou, and Hanjing Su. Pairwise learning for neural link prediction. *arXiv preprint arXiv:2112.02936*, 2021.
- Rachel Ward, Xiaoxia Wu, and Leon Bottou. Adagrad stepsizes: Sharp convergence over nonconvex landscapes. In *International Conference on Machine Learning*, pp. 6677–6686. PMLR, 2019.
- Kilian Q Weinberger and Lawrence K Saul. Distance metric learning for large margin nearest neighbor classification. *Journal of machine learning research*, 10(2), 2009.
- Eric Xing, Michael Jordan, Stuart J Russell, and Andrew Ng. Distance metric learning with application to clustering with side-information. *Advances in neural information processing systems*, 15, 2002.
- Zhiyu Xue, Shaoyang Yang, Mengdi Huai, and Di Wang. Differentially private pairwise learning revisited. *IJCAI*, 2021.
- Zhenhuan Yang, Yunwen Lei, Siwei Lyu, and Yiming Ying. Stability and differential privacy of stochastic gradient descent for pairwise learning with non-smooth loss. In *International Conference on Artificial Intelligence and Statistics*, pp. 2026–2034. PMLR, 2021a.
- Zhenhuan Yang, Yunwen Lei, Puyu Wang, Tianbao Yang, and Yiming Ying. Simple stochastic and online gradient descent algorithms for pairwise learning. *Advances in Neural Information Processing Systems*, 34, 2021b.
- Yiming Ying and Peng Li. Distance metric learning with eigenvalue optimization. *The Journal of Machine Learning Research*, 13(1):1–26, 2012.
- Yiming Ying and Ding-Xuan Zhou. Online pairwise learning algorithms. *Neural computation*, 28(4):743–777, 2016.
- Yiming Ying, Longyin Wen, and Siwei Lyu. Stochastic online AUC maximization. *Advances in neural information processing systems*, 29, 2016.
- Peilin Zhao, Steven CH Hoi, Rong Jin, and Tianbo Yang. Online AUC maximization. In *Proceedings of the 28th International Conference on Machine Learning ICML*, pp. 233–240, 2011.
- Masrour Zoghi, Tomas Tunys, Mohammad Ghavamzadeh, Branislav Kveton, Csaba Szepesvari, and Zheng Wen. Online learning to rank in stochastic click models. In *International Conference on Machine Learning*, pp. 4199–4208. PMLR, 2017.
- Fangyu Zou, Li Shen, Zequn Jie, Weizhong Zhang, and Wei Liu. A sufficient condition for convergences of Adam and RMSProp. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 11127–11135, 2019.

A Proofs of Results in the Convex Scenario

A.1 Technical Lemmas

Given any $\mathbf{x}^* \in \arg_{\mathcal{K}} \min f$, in the (strongly) convex setting, we introduce the following notation

$$\left\{ \begin{array}{l} \phi_k := \mathbb{E} \left[\frac{\|\mathbf{g}^k\|^2}{v^k} \right] + 2B^2 \mathbb{E} \left[\frac{1}{\sqrt{v^{k-2}}} - \frac{1}{\sqrt{v^{k-1}}} \right] + \frac{L^2 \mathbb{E} \|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2}{\sqrt{\delta}} + B^2 \mathbb{E} \left[\frac{1}{\sqrt{v^{k-2}}} - \frac{1}{\sqrt{v^k}} \right], \\ A_k := \mathbb{E} (\|\mathbf{m}^k\|^2 / (v^k)^{\frac{1}{2}}), \\ B_k := \mathbb{E} (\langle \mathbf{x}^* - \mathbf{x}^k, \mathbf{m}^k \rangle), \\ C_k := (\eta\theta + \frac{1-\theta}{2})A_{k-1} + (1-\theta)\phi_k \end{array} \right. \quad (9)$$

We have the following lemmas.

Lemma 1 [Lemma 9 in the appendix, (Li & Orabona, 2019)] Let a_1, a_2, \dots, a_K be non-negative, and h be a non-increasing function. Then we have

$$\sum_{k=1}^K a_k h(a_0 + \sum_{i=1}^k a_i) \leq \int_{a_0}^{\sum_{k=0}^K a_k} h(t) dt.$$

Lemma 2 Assume $\{\mathbf{x}^k\}_{k \geq 1}$ is generated by Algorithm 2, then we have

$$\sum_{k=1}^K A_k \leq \sum_{k=1}^K \mathbb{E} [\|\mathbf{g}^k\|^2 / (v^k)^{\frac{1}{2}}] \leq \mathbb{E} (v^k)^{\frac{1}{2}}.$$

Lemma 3 Assume $\{\mathbf{x}^k\}_{k \geq 1}$ is generated by Algorithm 2 and $\|\mathbf{g}^0\| \geq \sqrt{\delta} > 0$, then we have

$$\sum_{k=1}^K \mathbb{E} \|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2 \leq \sum_{k=1}^K \mathbb{E} [\|\mathbf{g}^k\|^2 / v^k] \leq \mathbb{E} \ln \frac{v^k}{\delta}.$$

Lemma 4 Assume $\{\mathbf{x}^k\}_{k \geq 1}$ is generated by Algorithm 2, then we have

$$\mathbb{E} \|\nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})\|^2 \leq 2L^2 \mathbb{E} \|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2 + 2\mathbb{E} \|\mathbf{g}^k\|^2.$$

Lemma 5 Assume $\{\mathbf{x}^k\}_{k \geq 1}$ is generated by Algorithm 2, and Assumptions 1 and 2 hold, then the following result holds

$$B_k \leq \theta B_{k-1} + C_k + (1-\theta) \mathbb{E} [f(\mathbf{x}^*) - f(\mathbf{x}^k)].$$

A.2 Proof of Theorem 1

Given $k \in \mathbb{Z}^+$, with Lemma 5 and mathematical induction, we have the following inequality

$$B_k \leq \theta^k B_1 + \sum_{i=1}^{k-1} \theta^{k-1-i} C_i + \sum_{i=1}^{k-1} (1-\theta) \theta^{k-1-i} \mathbb{E} (f(\mathbf{x}^*) - f(\mathbf{x}^i)).$$

Notice that $\mathbf{m}^1 = \mathbf{0}$ and $B_1 = 0$, we get

$$B_k \leq \sum_{i=1}^{k-1} \theta^{k-1-i} C_i + \sum_{i=1}^{k-1} (1-\theta) \theta^{k-1-i} \mathbb{E} (f(\mathbf{x}^*) - f(\mathbf{x}^i)). \quad (10)$$

Summing the inequality (10) from $k = 1$ to K gives us

$$\begin{aligned} \sum_{k=1}^K B_k &\leq \sum_{k=1}^K \sum_{i=1}^{k-1} \theta^{k-1-i} C_i + \sum_{k=1}^K \sum_{i=1}^{k-1} (1-\theta) \theta^{k-1-i} \mathbb{E}(f(\mathbf{x}^*) - f(\mathbf{x}^i)) \\ &\leq \frac{\sum_{k=1}^K C_k}{1-\theta} + (1-\theta) \sum_{k=1}^K \mathbb{E}(f(\mathbf{x}^*) - f(\mathbf{x}^k)). \end{aligned}$$

Therefore, we have

$$\sum_{k=1}^K \mathbb{E}(f(\mathbf{x}^k) - f(\mathbf{x}^*)) \leq - \sum_{k=1}^K \frac{B_k}{(1-\theta)} + \frac{\sum_{k=1}^K C_k}{(1-\theta)^2}. \quad (11)$$

The scheme of the algorithm indicates that

$$\mathbf{x}^{k+1} = \mathbf{Proj}_{\mathcal{K}}(\mathbf{x}^k - \eta \mathbf{m}^k / (v^k)^{\frac{1}{2}}).$$

We then get

$$\begin{aligned} \|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2 &= \|\mathbf{Proj}_{\mathcal{K}}(\mathbf{x}^k - \eta \mathbf{m}^k / (v^k)^{\frac{1}{2}}) - \mathbf{x}^*\|^2 \\ &= \|\mathbf{Proj}_{\mathcal{K}}(\mathbf{x}^k - \eta \mathbf{m}^k / (v^k)^{\frac{1}{2}}) - \mathbf{Proj}_{\mathcal{K}}(\mathbf{x}^*)\|^2 \\ &\leq \|\mathbf{x}^k - \eta \mathbf{m}^k / (v^k)^{\frac{1}{2}} - \mathbf{x}^*\|^2. \end{aligned}$$

Multiplying both sides with $(v^k)^{\frac{1}{2}}$, we are then led to

$$(v^k)^{\frac{1}{2}} \|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2 \leq (v^k)^{\frac{1}{2}} \|\mathbf{x}^* - \mathbf{x}^k\|^2 + 2\eta \langle \mathbf{m}^k, \mathbf{x}^* - \mathbf{x}^k \rangle + \eta^2 \|\mathbf{m}^k\|^2 / (v^k)^{\frac{1}{2}},$$

which is equivalent to

$$\begin{aligned} (v^{k+1})^{\frac{1}{2}} \|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2 &\leq (v^k)^{\frac{1}{2}} \|\mathbf{x}^* - \mathbf{x}^k\|^2 + 2\eta \langle \mathbf{m}^k, \mathbf{x}^* - \mathbf{x}^k \rangle \\ &\quad + \eta^2 \|\mathbf{m}^k\|^2 / (v^k)^{\frac{1}{2}} + ((v^{k+1})^{\frac{1}{2}} - (v^k)^{\frac{1}{2}}) \|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2. \end{aligned}$$

Taking the total expectation of both sides of the above equation gives us

$$\begin{aligned} -2\eta B_k &\leq \mathbb{E}(v^k)^{\frac{1}{2}} \|\mathbf{x}^k - \mathbf{x}^*\|^2 - \mathbb{E}(v^{k+1})^{\frac{1}{2}} \|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2 \\ &\quad + \eta^2 A_k + (\mathbb{E}(v^{k+1})^{\frac{1}{2}} - \mathbb{E}(v^k)^{\frac{1}{2}}) D^2. \end{aligned} \quad (12)$$

Summing (12) from $k = 1$ to K gives us

$$\sum_{k=1}^K (-B_k) / (1-\theta) \leq \frac{\mathbb{E}\sqrt{v^1} \|\mathbf{x}^* - \mathbf{x}^1\|^2}{2\eta(1-\theta)} + \frac{\eta}{2(1-\theta)} \sum_{k=1}^K A_k + \frac{D^2}{(1-\theta)} \mathbb{E}\sqrt{v^{K+1}}.$$

Together with (11), we then have

$$\begin{aligned} \sum_{k=1}^K \mathbb{E}(f(\mathbf{x}^k) - f(\mathbf{x}^*)) &\leq \frac{\mathbb{E}\sqrt{v^1} \|\mathbf{x}^* - \mathbf{x}^1\|^2}{2\eta(1-\theta)} \\ &\quad + \frac{\eta}{2(1-\theta)} \sum_{k=1}^K A_k + \frac{D^2}{(1-\theta)} \mathbb{E}\sqrt{v^{K+1}} + \frac{\sum_{k=1}^K C_k}{(1-\theta)^2}. \end{aligned} \quad (13)$$

We turn to bound the right-hand side of (13) and get the following bound

$$\frac{\eta}{2(1-\theta)} \sum_{k=1}^K A_k \leq \frac{\eta}{2(1-\theta)} \mathbb{E}(v^{K+1})^{\frac{1}{2}}. \quad (14)$$

On the other hand, we can get

$$\begin{aligned} \frac{\sum_{k=1}^K C_k}{(1-\theta)^2} &\leq \frac{\eta\theta}{(1-\theta)^2} \mathbb{E}(v^{K+1})^{\frac{1}{2}} + \frac{2L}{(1-\theta)^2} \mathbb{E} \ln \frac{v^{K+1}}{\delta} + \frac{1}{(1-\theta)^2} \sum_{k=1}^K \mathbb{E} \phi_k \\ &\leq \frac{\eta\theta}{(1-\theta)^2} \mathbb{E}(v^{K+1})^{\frac{1}{2}} + \frac{2L}{(1-\theta)^2} \mathbb{E} \ln \frac{v^{K+1}}{\delta} + \frac{\frac{3B^2}{\sqrt{\delta}} + \frac{L^2}{\sqrt{\delta}} \mathbb{E} \ln \frac{v^{K+1}}{\delta} + \mathbb{E}(v^{K+1})^{\frac{1}{2}}}{(1-\theta)^2} \end{aligned} \quad (15)$$

Substituting the bounds (14) and (15) into (13), we are then led to

$$\begin{aligned} &\left[\sum_{k=1}^K \mathbb{E}(f(\mathbf{x}^k) - f(\mathbf{x}^*)) \right] \\ &\leq \frac{\mathbb{E}\sqrt{v^1} \|\mathbf{x}^* - \mathbf{x}^1\|^2}{2\eta(1-\theta)} + \left[\frac{\eta + 2D^2}{2(1-\theta)} + \frac{\eta\theta + 1}{(1-\theta)^2} \right] \cdot \mathbb{E}\sqrt{v^{K+1}} \\ &\quad + \left[\frac{4L}{(1-\theta)^2} + \frac{2L^2}{(1-\theta)^2\sqrt{\delta}} \right] \mathbb{E} \ln \frac{\sqrt{v^{K+1}}}{\sqrt{\delta}} + \frac{3B^2}{(1-\theta)^2\sqrt{\delta}}. \end{aligned} \quad (16)$$

Notice that $-\ln(\cdot)$ and $-\frac{1}{\cdot}$ are both convex when \cdot is positive, using Jensen's inequality gives us

$$\begin{aligned} -\ln \mathbb{E}\sqrt{v^{K+1}} &\leq -\mathbb{E} \ln \sqrt{v^{K+1}}, \quad -\frac{1}{\mathbb{E}(\sqrt{v^k})} \leq -\mathbb{E} \frac{1}{\sqrt{v^k}} \\ \Rightarrow \mathbb{E} \ln \sqrt{v^{K+1}} &\leq \ln \mathbb{E}\sqrt{v^{K+1}}, \quad \mathbb{E} \frac{1}{\sqrt{v^k}} \leq \frac{1}{\mathbb{E}\sqrt{v^k}}. \end{aligned}$$

Therefore, (16) can be presented as

$$\mathbb{E} \left[f \left(\frac{\sum_{k=1}^K \mathbf{x}^k}{K} \right) - \min f \right] \leq \frac{c_1 + c_2(K)}{K},$$

where

$$c_1 := \frac{\mathbb{E}\sqrt{v^1} \|\mathbf{x}^* - \mathbf{x}^1\|^2}{2(1-\theta)\eta} + \frac{3B^2}{(1-\theta)^2\sqrt{\delta}},$$

and

$$c_2(K) := \left[\frac{\eta + 2D^2}{2(1-\theta)} + \frac{\eta\theta + 1}{(1-\theta)^2} \right] \cdot \mathbb{E}\sqrt{v^{K+1}} + \left[\frac{4L}{(1-\theta)^2} + \frac{2L^2}{(1-\theta)^2\sqrt{\delta}} \right] \ln \frac{\mathbb{E}\sqrt{v^{K+1}}}{\sqrt{\delta}}.$$

B Proofs of Results in the Strongly Convex Scenario

B.1 Proof of Proposition 2

The boundedness of the gradient and the strong convexity indicate \mathcal{K} is bounded, whose radius is assumed to be $D > 0$. Notice that the operator $\mathbf{Proj}_{\mathcal{K}}(\cdot)$ is contractive, as $\theta = 0$, we get

$$\begin{aligned} \mathbb{E} \|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2 &= \mathbb{E} \|\mathbf{Proj}_{\mathcal{K}}(\mathbf{x}^k - \frac{\eta}{\sqrt{k}} \mathbf{m}^k / \sqrt{v^k}) - \mathbf{x}^*\|^2 \\ &= \mathbb{E} \|\mathbf{Proj}_{\mathcal{K}}(\mathbf{x}^k - \frac{\eta}{\sqrt{k}} \mathbf{m}^k / \sqrt{v^k}) - \mathbf{Proj}_{\mathcal{K}}(\mathbf{x}^*)\|^2 \\ &\leq \mathbb{E} \|\mathbf{x}^k - \frac{\eta}{\sqrt{k}} \mathbf{m}^k / \sqrt{v^k} - \mathbf{x}^*\|^2 \\ &= \mathbb{E} \|\mathbf{x}^k - \mathbf{x}^*\|^2 - 2 \frac{\eta}{\sqrt{k}} \mathbb{E} \left(\frac{\langle \nabla F(\mathbf{x}^k; \xi^k, \xi^{k-1}), \mathbf{x}^k - \mathbf{x}^* \rangle}{\sqrt{v^k}} \right) + \frac{\eta^2}{k} \mathbb{E} \|\mathbf{m}^k\|^2 / v^k \\ &\leq \mathbb{E} \left[(1 - 2\eta\nu / \sqrt{kv^k}) \mathbb{E} \|\mathbf{x}^k - \mathbf{x}^*\|^2 \right] + \frac{\eta^2}{k} \mathbb{E} \frac{\|\mathbf{m}^k\|^2}{v^k} \\ &\quad + 2 \frac{\eta}{\sqrt{k}} \mathbb{E} ([F(\mathbf{x}^*; \xi^k, \xi^{k-1}) - F(\mathbf{x}^k; \xi^k, \xi^{k-1})] / \sqrt{v^k}), \end{aligned}$$

where we used the strong convexity of $F(\mathbf{x}; \xi^k, \xi^{k-1})$. Now, we turn to bound

$$\mathbb{E}([F(\mathbf{x}^*; \xi^k, \xi^{k-1}) - F(\mathbf{x}^k; \xi^k, \xi^{k-1})]/\sqrt{v^k}).$$

With direct computations, we have the decomposition

$$\begin{aligned} & \mathbb{E}([F(\mathbf{x}^*; \xi^k, \xi^{k-1}) - F(\mathbf{x}^k; \xi^k, \xi^{k-1})]/\sqrt{v^k}) \\ &= \mathbb{E}([F(\mathbf{x}^*; \xi^k, \xi^{k-1}) - F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})]/\sqrt{v^{k-2}}) \\ &+ \mathbb{E}([F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) - F(\mathbf{x}^k; \xi^k, \xi^{k-1})]/\sqrt{v^{k-2}}) \\ &+ \frac{\eta}{\sqrt{k}} \mathbb{E}\left([F(\mathbf{x}^*; \xi^k, \xi^{k-1}) - F(\mathbf{x}^k; \xi^k, \xi^{k-1})] \times (1/\sqrt{v^k} - 1/\sqrt{v^{k-2}})\right), \end{aligned}$$

consequently,

$$\begin{aligned} & \mathbb{E}([F(\mathbf{x}^*; \xi^k, \xi^{k-1}) - F(\mathbf{x}^k; \xi^k, \xi^{k-1})]/\sqrt{v^k}) \\ & \leq \mathbb{E}\left(\frac{\langle \mathbf{m}^{k-1}, \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) \rangle}{\sqrt{k-1}\sqrt{v^{k-1}}\sqrt{v^{k-2}}}\right) + \frac{DB\eta}{\sqrt{k}} \mathbb{E}(1/\sqrt{v^{k-2}} - 1/\sqrt{v^k}) \\ & \leq \frac{1}{2\sqrt{k-1}} \mathbb{E}\frac{\|\mathbf{m}^{k-1}\|}{v^{k-1}} + \frac{1}{2\sqrt{k-1}} \mathbb{E}\frac{\|\nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})\|^2}{v^{k-2}} \\ & \quad + \frac{DB\eta}{\sqrt{k}} \mathbb{E}(1/\sqrt{v^{k-2}} - 1/\sqrt{v^k}) + \frac{B}{\sqrt{k-1}\sqrt{\delta}} \mathbb{E}(1/\sqrt{v^{k-2}} - 1/\sqrt{v^{k-1}}) \\ & \leq \frac{1}{\sqrt{\delta}} \mathbb{E}\left(\frac{\|\mathbf{m}^{k-1}\|^2}{\sqrt{k-1}v^{k-1}}\right) + \frac{DB\eta}{\sqrt{k}} \mathbb{E}(1/\sqrt{v^{k-2}} - 1/\sqrt{v^k}) \\ & \quad + \frac{B}{\sqrt{k-1}\sqrt{\delta}} \mathbb{E}(1/\sqrt{v^{k-2}} - 1/\sqrt{v^{k-1}}). \end{aligned}$$

where we used $\mathbb{E}([F(\mathbf{x}^*; \xi^k, \xi^{k-1}) - F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})]/\sqrt{v^{k-2}}) = [f(\mathbf{x}^*) - f(\mathbf{x}^{k-1})]/\sqrt{v^{k-2}} \leq 0$. Letting $\eta = \frac{B}{2\nu}$ and $\beta_k := \frac{\eta}{\sqrt{k}\sqrt{k-1}} \mathbb{E}\frac{\|\mathbf{m}^{k-1}\|}{v^{k-1}} + \frac{\eta}{\sqrt{k}\sqrt{k-1}} \mathbb{E}\frac{\|\nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})\|^2}{v^{k-2}} + \frac{2DB\eta^2}{k} \mathbb{E}(1/\sqrt{v^{k-2}} - 1/\sqrt{v^k}) + \frac{2B\eta}{\sqrt{k-1}\sqrt{k}\sqrt{\delta}} \mathbb{E}(1/\sqrt{v^{k-2}} - 1/\sqrt{v^{k-1}}) + \frac{\eta^2}{k} \mathbb{E}\frac{\|\mathbf{m}^k\|^2}{v^k}$,

$$a_{k+1} \leq \left(1 - \frac{1}{k}\right)a_k + \beta_k.$$

With direct computations, we have

$$\begin{aligned} a_3 &\leq \frac{1}{2}a_2 + \beta_2, \\ a_4 &\leq \frac{2}{3}a_3 + \beta_3 \leq \frac{1}{3}a_2 + \frac{2}{3}\beta_2 + \beta_3, \\ a_5 &\leq \frac{3}{4}a_4 + \beta_4 \leq \frac{1}{4}a_2 + \frac{2}{4}\beta_2 + \frac{3}{4}\beta_3 + \beta_4, \\ &\vdots \\ a_{K+1} &\leq \frac{a_2}{K} + \sum_{k=2}^K \frac{k\beta_k}{K}. \end{aligned}$$

Notice that

$$\begin{aligned} k\beta_k &= \mathcal{O}\left(\mathbb{E}\frac{\|\mathbf{m}^{k-1}\|}{v^{k-1}} + \mathbb{E}\frac{\|\nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})\|^2}{v^{k-2}} + \mathbb{E}(1/\sqrt{v^{k-2}} - 1/\sqrt{v^k})\right. \\ &\quad \left.+ \mathbb{E}(1/\sqrt{v^{k-2}} - 1/\sqrt{v^{k-1}}) + \mathbb{E}\frac{\|\mathbf{m}^k\|^2}{v^k}\right), \end{aligned}$$

we just need to compute $\sum_{k=1}^{\infty} \mathbb{E} \frac{\|\nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})\|^2}{v^{k-2}}$. By using Lemma 4, we have

$$\mathbb{E} \|\nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})\|^2 / v^{k-2} \leq 2L^2 \mathbb{E} \|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2 / \delta + 2\mathbb{E} \|\mathbf{g}^k\|^2 / v^{k-2}.$$

With the fact

$$\mathbb{E} \frac{\|\mathbf{g}^k\|^2}{v^{k-2}} \leq \mathbb{E} \frac{\|\mathbf{g}^k\|^2}{v^k} + B^2 (\mathbb{E} 1/v^{k-2} - \mathbb{E} 1/v^k),$$

we then get

$$\sum_{k=1}^{\infty} \mathbb{E} \frac{\|\nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})\|^2}{v^{k-2}} = \mathcal{O}(\ln v^K).$$

With Lemma 3, we can get

$$\sum_{k=2}^K k\beta_k = \mathcal{O}(\ln v^K) = \mathcal{O}(\ln K) \Rightarrow a_{K+1} = \frac{\ln K}{K}.$$

C Proofs of Results in Nonconvex Scenario

C.1 Additional Technical Lemmas

In the nonconvex case, we denote the following items to simplify the presentations of the technical lemmas

$$\left\{ \begin{array}{l} \hat{A}_k := \mathbb{E} \left[\|\mathbf{m}^k\|^2 / v^k \right], \\ \hat{B}_k := \mathbb{E} \left(\langle -\nabla f(\mathbf{x}^k), \mathbf{m}^k / (v^k)^{\frac{1}{2}} \rangle \right), \\ \hat{C}_k := \theta \eta \hat{A}_k + 2(1-\theta) \mathbb{E} \left[\|\mathbf{g}^k\|^2 / v^k \right] + 6(1-\theta) B^2 \mathbb{E} [1/(v^{k-2})^{\frac{1}{2}} - 1/(v^k)^{\frac{1}{2}}] \\ \quad + (1-\theta)(2L^2/\sqrt{\delta} + L^2/\delta + 1/2) \mathbb{E} \|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2. \end{array} \right. \quad (17)$$

Lemma 6 Assume $\{\mathbf{x}^k\}_{k \geq 1}$ is generated by Algorithm 2, then we have

$$\sum_{k=1}^K \hat{A}_k \leq \sum_{k=1}^K \mathbb{E} \|\mathbf{g}^k\|^2 / v^k \leq \mathbb{E} \ln \frac{v^k}{\delta}.$$

Lemma 7 Assume $\{\mathbf{x}^k\}_{k \geq 1}$ is generated by Algorithm 2 and the functions are nonconvex. Let \mathcal{K} be the full space and Assumptions 1 and 2 hold, then the following result holds

$$\hat{B}_k + \frac{(1-\theta)}{2} \mathbb{E} \left(\|\nabla f(\mathbf{x}^k)\|^2 / (v^k)^{\frac{1}{2}} \right) \leq \theta \hat{B}_{k-1} + \hat{C}_k.$$

C.2 Proof of Theorem 3

According to Lemma 7, we have

$$\begin{aligned} \frac{(1-\theta)}{2} \sum_{k=1}^K \mathbb{E} \left(\|\nabla f(\mathbf{x}^k)\|^2 / (v^k)^{\frac{1}{2}} \right) &\leq -\hat{B}_K + (\theta-1) \sum_{k=1}^{K-1} \hat{B}_k + \sum_{k=1}^K \hat{C}_k \\ &\leq (\theta-1) \sum_{k=1}^{K-1} \hat{B}_k + \sum_{k=1}^K \hat{C}_k + \frac{B^2}{\sqrt{\delta}}. \end{aligned} \quad (18)$$

The Lipschitz property of the gradients gives

$$\mathbb{E} f(\mathbf{x}^{k+1}) - \mathbb{E} f(\mathbf{x}^k) \leq \mathbb{E} \langle \nabla f(\mathbf{x}^k), \mathbf{x}^{k+1} - \mathbf{x}^k \rangle + \frac{L \mathbb{E} \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2}{2} = \eta \hat{B}_k + \frac{L\eta^2}{2} \hat{A}_k. \quad (19)$$

Combining with (19), we get the following estimate

$$\sum_{k=1}^{K-1} -\hat{B}_k \leq L\eta \sum_{k=1}^{K-1} \hat{A}_k + \frac{2f(\mathbf{x}^1)}{\eta}. \quad (20)$$

On the other hand, with Lemma 3, we have the following bound

$$\frac{2}{1-\theta} \sum_{k=1}^K \hat{C}_k \leq \frac{2\theta\eta}{1-\theta} \sum_{k=1}^K \hat{A}_k + \frac{6B^2}{\sqrt{\delta}} + (2L^2/\sqrt{\delta} + L^2/\delta + 5/2)\mathbb{E} \ln\left(\frac{v^K}{\delta}\right). \quad (21)$$

Using [Lemma 2, (Li & Orabona, 2019)] and Lemma 6, we have

$$\sum_{k=1}^K \hat{A}_k \leq \mathbb{E} \ln\left(\frac{v^K}{\delta}\right). \quad (22)$$

Substituting (22), (21) and (20) into (18), then we get

$$\begin{aligned} & \sum_{k=1}^K \mathbb{E} \left(\|\nabla f(\mathbf{x}^k)\|^2 / (v^k)^{\frac{1}{2}} \right) \\ & \leq \left(\frac{(1+\theta)\eta}{1-\theta} + 2L^2/\sqrt{\delta} + L^2/\delta + 5/2 \right) \mathbb{E} \ln\left(\frac{v^K}{\delta}\right) + \frac{(7-6\theta)B^2}{(1-\theta)\sqrt{\delta}} + \frac{2f(\mathbf{x}^1)}{\eta}. \end{aligned} \quad (23)$$

Notice that $\frac{1}{(v^k)^{\frac{1}{2}}} \geq \frac{1}{Ck^\alpha}$, then we have

$$\sum_{k=1}^K \mathbb{E} \left(\|\nabla f(\mathbf{x}^k)\|^2 / (v^k)^{\frac{1}{2}} \right) \geq \left(\sum_{k=1}^K \frac{1}{k^\alpha C} \right) \cdot \min_{1 \leq k \leq K} \{ \mathbb{E} \|\nabla f(\mathbf{x}^k)\|^2 \}.$$

We then complete the proof by using the fact that $0 < \frac{1}{1-\alpha} \leq 2$ when $0 < \alpha \leq 1/2$. Thus, we can get the following estimate

$$\min_{1 \leq k \leq K} \{ \mathbb{E} \|\nabla f(\mathbf{x}^k)\|^2 \} \leq \frac{c_3 + c_4(K)}{K^{1-\alpha}},$$

where

$$c_3 := \frac{2(7-6\theta)B^2C}{(1-\theta)\sqrt{\delta}} + \frac{4Cf(\mathbf{x}^1)}{\eta},$$

and

$$c_4(K) := \left(\frac{(1+\theta)\eta}{1-\theta} + 2L^2/\sqrt{\delta} + L^2/\delta + 5/2 \right) \ln\left(\frac{CK}{\delta}\right).$$

D Proofs of the Technical Lemmas

D.1 Proof of Lemma 2

With the fact that $\mathbf{m}^k = (1-\theta) \sum_{j=1}^k \theta^{k-j} \mathbf{g}^j$ when $k \geq 1$, we have

$$\begin{aligned} \|\mathbf{m}^k\|^2 / (v^k)^{\frac{1}{2}} &= \sum_{i=1}^d |\mathbf{m}_i^k / (v^k)^{\frac{1}{4}}|^2 \leq \sum_{i=1}^d (1-\theta)^2 \left| \sum_{j=1}^k \theta^{k-j} \mathbf{g}_i^j / (v^k)^{\frac{1}{4}} \right|^2 \\ &\stackrel{a)}{\leq} \sum_{i=1}^d (1-\theta)^2 \left(\sum_{j=1}^k \theta^{k-j} (v^k)^{\frac{1}{2}} \right) \times \sum_{j=1}^k \theta^{k-j} (\mathbf{g}_i^j)^2 / (v^k) \\ &\leq \sum_{i=1}^d (1-\theta)^2 \cdot \frac{(v^k)^{\frac{1}{2}}}{1-\theta} \cdot \sum_{j=1}^k \theta^{k-j} (\mathbf{g}_i^j)^2 / (v^k) \\ &= (1-\theta) \cdot \sum_{j=1}^k \theta^{k-j} \|\mathbf{g}^j\|^2 / (v^k)^{\frac{1}{2}} \stackrel{b)}{\leq} (1-\theta) \cdot \sum_{j=1}^k \theta^{k-j} \|\mathbf{g}^j\|^2 / (v^j)^{\frac{1}{2}} \end{aligned}$$

where a) uses the Cauchy's inequality $(\sum_{j=1}^k a_j b_j)^2 \leq (\sum_{j=1}^k a_j^2) \cdot (\sum_{j=1}^k b_j^2)$ with $a_j = \theta^{\frac{k-j}{2}} (v^k)^{\frac{1}{4}}$ and $b_j = \theta^{\frac{k-j}{2}} \mathbf{g}_i^j / (v^k)^{\frac{1}{2}}$; b) is due to $v^j \leq v^k$ when $j \leq k$. Thus, we are led to

$$\begin{aligned} \sum_{k=1}^K \sum_{j=1}^k \theta^{k-j} \|\mathbf{g}^j\|^2 / (v^j)^{\frac{1}{2}} &= \sum_{j=1}^K \sum_{k=j}^K \theta^{k-j} \|\mathbf{g}^j\|^2 / (v^j)^{\frac{1}{2}} \\ &= \sum_{j=1}^K \sum_{k=j}^K \theta^{k-j} \|\mathbf{g}^j\|^2 / (v^j)^{\frac{1}{2}} \leq \frac{1}{1-\theta} \sum_{j=1}^K \|\mathbf{g}^j\|^2 / (v^j)^{\frac{1}{2}}. \end{aligned}$$

Thus, we can get

$$\sum_{k=1}^K A_k \leq \sum_{k=1}^K \mathbb{E}[\|\mathbf{g}^k\|^2 / (v^k)^{\frac{1}{2}}].$$

Using Lemma 1 with $a_k = \|\mathbf{g}^k\|^2$ and $h(\cdot) = \frac{1}{\sqrt{\cdot}}$, we get

$$\sum_{k=1}^K \|\mathbf{g}^k\|^2 / (v^k)^{\frac{1}{2}} \leq \sqrt{v^K} - \sqrt{v^0} \leq \sqrt{v^K},$$

where we used the fact that $v^0 \geq 0$. The proof is then completed.

D.2 Proof of Lemma 3

Similar to the proof of Lemma 2, we have

$$\begin{aligned} \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2 &= \|\mathbf{Proj}_{\mathcal{K}}(\mathbf{x}^k - \eta \mathbf{m}^k) - \mathbf{x}^k\|^2 = \|\mathbf{Proj}_{\mathcal{K}}(\mathbf{x}^k - \eta \mathbf{m}^k) - \mathbf{Proj}_{\mathcal{K}}(\mathbf{x}^k)\|^2 \\ &\leq \|\mathbf{m}^k / (v^k)^{\frac{1}{2}}\|^2 = \sum_{i=1}^d |\mathbf{m}_i^k / (v^k)^{\frac{1}{2}}|^2 \leq \sum_{i=1}^d (1-\theta)^2 \left| \sum_{j=1}^k \theta^{k-j} \mathbf{g}_i^j / (v^k)^{\frac{1}{2}} \right|^2 \\ &\stackrel{a)}{\leq} \sum_{i=1}^d (1-\theta)^2 \left(\sum_{j=1}^k \theta^{k-j} \right) \cdot \sum_{j=1}^k \theta^{k-j} \frac{(\mathbf{g}_i^j)^2}{(v^k)} \leq \sum_{i=1}^d (1-\theta)^2 \cdot \frac{1}{1-\theta} \cdot \sum_{j=1}^{k-1} \theta^{k-j} (\mathbf{g}_i^j)^2 / (v^k) \\ &= (1-\theta) \cdot \sum_{j=1}^k \theta^{k-j} \|\mathbf{g}^j\|^2 / (v^k) \stackrel{b)}{=} (1-\theta) \cdot \sum_{j=1}^k \theta^{k-j} \|\mathbf{g}^j\|^2 / v^j \end{aligned}$$

where a) uses the fact that $(\sum_{j=1}^k a_j b_j)^2 \leq \sum_{j=1}^k a_j^2 \sum_{j=1}^k b_j^2$ with $a_j = \theta^{\frac{k-j}{2}}$ and $b_j = \theta^{\frac{k-j}{2}} \mathbf{g}_i^j / (v^k)^{\frac{1}{2}}$, and b) is due to $v^j \leq v^k$ when $j \leq k$. Thus, we can get

$$\sum_{k=1}^K \|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2 \leq \sum_{k=1}^K \|\mathbf{g}^k\|^2 / v^k.$$

Using Lemma 1 with $a_k = \|\mathbf{g}^k\|^2$ and $h(\cdot) = \frac{1}{\cdot}$, we are then led to

$$\sum_{k=1}^K \|\mathbf{g}^k\|^2 / v^k \leq \ln \frac{v^K}{v^0} \leq \ln \frac{v^K}{\delta}.$$

The proof is then completed.

D.3 Proof of Lemma 4

Direct computations give us

$$\begin{aligned}
\mathbb{E}\|\mathbf{g}^k\|^2 &= \mathbb{E}\|\nabla F(\mathbf{x}^k; \xi^k, \xi^{k-1})\|^2 \\
&= \mathbb{E}\|\nabla F(\mathbf{x}^k; \xi^k, \xi^{k-1}) - \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) + \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})\|^2 \\
&= \mathbb{E}\|\nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})\|^2 + \mathbb{E}\|\nabla F(\mathbf{x}^k; \xi^k, \xi^{k-1}) - \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})\|^2 \\
&\quad + 2\mathbb{E}\langle \nabla F(\mathbf{x}^k; \xi^k, \xi^{k-1}) - \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}), \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) \rangle \\
&\stackrel{a)}{\geq} \frac{1}{2}\mathbb{E}\|\nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})\|^2 - \mathbb{E}\|\nabla F(\mathbf{x}^k; \xi^k, \xi^{k-1}) - \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})\|^2 \\
&\geq \frac{1}{2}\mathbb{E}\|\nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})\|^2 - L^2\mathbb{E}\|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2,
\end{aligned}$$

where $a)$ uses the inequality $2\mathbb{E}\langle \nabla F(\mathbf{x}^k; \xi^k, \xi^{k-1}) - \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}), \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) \rangle \geq -\frac{1}{2}\mathbb{E}\|\nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})\|^2 - 2\mathbb{E}\|\nabla F(\mathbf{x}^k; \xi^k, \xi^{k-1}) - \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})\|^2$. Thus we can get

$$\mathbb{E}\|\nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})\|^2 \leq 2L^2\mathbb{E}\|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2 + 2\mathbb{E}\|\mathbf{g}^k\|^2.$$

D.4 Proof of Lemma 5

The convexity of $f_{i_k}(\mathbf{x})$ with respect to \mathbf{x} and $\mathbf{g}^k = \nabla F(\mathbf{x}^k; \xi^k, \xi^{k-1})$ gives us

$$\begin{aligned}
\mathbb{E}\langle \mathbf{x}^* - \mathbf{x}^k, \mathbf{g}^k \rangle &\leq \mathbb{E}[F(\mathbf{x}^*; \xi^k, \xi^{k-1}) - F(\mathbf{x}^k; \xi^k, \xi^{k-1})] \\
&= \mathbb{E}[F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) - F(\mathbf{x}^k; \xi^k, \xi^{k-1}) + F(\mathbf{x}^*; \xi^k, \xi^{k-1}) - F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})].
\end{aligned} \tag{24}$$

The Lipschitz gradient continuity of $F(\mathbf{x}; \xi^k, \xi^{k-1})$ with respect to \mathbf{x} indicates

$$\begin{aligned}
&\mathbb{E}[F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) - F(\mathbf{x}^k; \xi^k, \xi^{k-1})] \\
&\leq \frac{L}{2}\mathbb{E}\|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2 + \mathbb{E}\langle \mathbf{x}^{k-1} - \mathbf{x}^k, \nabla F(\mathbf{x}^k; \xi^k, \xi^{k-1}) \rangle \\
&\leq \frac{L}{2}\mathbb{E}\|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2 + \mathbb{E}\langle \mathbf{x}^{k-1} - \mathbf{x}^k, \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) \rangle \\
&\quad + \mathbb{E}\langle \mathbf{x}^{k-1} - \mathbf{x}^k, \nabla F(\mathbf{x}^k; \xi^k, \xi^{k-1}) - \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) \rangle \\
&\leq \frac{3L}{2}\mathbb{E}\|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2 + \mathbb{E}\langle \mathbf{x}^{k-1} - \mathbf{x}^k, \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) \rangle,
\end{aligned} \tag{25}$$

where we used the Cauchy's inequality $\mathbb{E}\langle \mathbf{x}^{k-1} - \mathbf{x}^k, \nabla F(\mathbf{x}^k; \xi^k, \xi^{k-1}) - \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) \rangle \leq L\|\mathbf{x}^{k-1} - \mathbf{x}^k\|^2$. Because ξ^k, ξ^{k-1} are independent of \mathbf{x}^{k-1} ,

$$\mathbb{E}(F(\mathbf{x}^*; \xi^k, \xi^{k-1}) - F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})) = f(\mathbf{x}^*) - f(\mathbf{x}^{k-1}). \tag{26}$$

Turning back to (24), we get

$$\mathbb{E}\langle \mathbf{x}^* - \mathbf{x}^k, \mathbf{g}^k \rangle \leq \mathbb{E}[F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) - F(\mathbf{x}^k; \xi^k, \xi^{k-1})] + f(\mathbf{x}^*) - f(\mathbf{x}^{k-1}). \tag{27}$$

Notice that \mathbf{x}^{k-1} is independent of (ξ^k, ξ^{k-1}) ,

$$\begin{aligned}
&\mathbb{E}\langle \mathbf{x}^{k-1} - \mathbf{x}^k, \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) \rangle \\
&= \mathbb{E}\langle \mathbf{x}^{k-1} - \mathbf{x}^k, \nabla f(\mathbf{x}^{k-1}) \rangle \\
&\quad + \mathbb{E}\langle \mathbf{x}^{k-1} - \mathbf{x}^k, \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) - \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) \rangle \\
&\leq \mathbb{E}\langle \mathbf{x}^{k-1} - \mathbf{x}^k, \nabla f(\mathbf{x}^{k-1}) \rangle + \frac{1}{2}\mathbb{E}\frac{\|\mathbf{m}^{k-1}\|^2}{\sqrt{v^{k-1}}} \\
&\quad + \frac{1}{2}\mathbb{E}\frac{\|\nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) - \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})\|^2}{\sqrt{v^{k-1}}},
\end{aligned} \tag{28}$$

where we used $|\mathbb{E}\langle X, Y \rangle| \leq \mathbb{E}\|X\|^2 + \mathbb{E}\|Y\|^2$ with $X = \sqrt[4]{v^{k-1}}(\mathbf{x}^{k-1} - \mathbf{x}^k)$, and $Y = [\nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) - \mathbb{E}\nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})] / \sqrt[4]{v^{k-1}}$. Now, we turn to the upper bound of $\frac{1}{2}\mathbb{E}\frac{\|\nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) - \mathbb{E}\nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})\|^2}{\sqrt{v^{k-1}}}$:

$$\begin{aligned}
& \frac{1}{2}\mathbb{E}\frac{\|\nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) - \mathbb{E}\nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})\|^2}{\sqrt{v^{k-1}}} \\
&= \frac{1}{2}\mathbb{E}\frac{\|\nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) - \mathbb{E}\nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})\|^2}{\sqrt{v^{k-2}}} \\
&\quad + \frac{1}{2}\mathbb{E}\left[\left(\|\nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) - \mathbb{E}\nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})\|^2\right) \times \left(\frac{1}{\sqrt{v^{k-1}}} - \frac{1}{\sqrt{v^{k-2}}}\right)\right] \\
&\leq \frac{1}{2}\frac{\mathbb{E}\|\nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})\|^2}{\sqrt{v^{k-2}}} + 2B^2\mathbb{E}\left(\frac{1}{\sqrt{v^{k-2}}} - \frac{1}{\sqrt{v^{k-1}}}\right) \\
&\stackrel{a)}{\leq} \frac{L^2\mathbb{E}\|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2 + \mathbb{E}\|\mathbf{g}^k\|^2}{\sqrt{v^{k-2}}} + 2B^2\mathbb{E}\left(\frac{1}{\sqrt{v^{k-2}}} - \frac{1}{\sqrt{v^{k-1}}}\right) \\
&\leq \phi_k,
\end{aligned} \tag{29}$$

where $\phi_k := \mathbb{E}\frac{\|\mathbf{g}^k\|^2}{\sqrt{v^k}} + 2B^2\mathbb{E}\left(\frac{1}{\sqrt{v^{k-2}}} - \frac{1}{\sqrt{v^{k-1}}}\right) + \frac{L^2\mathbb{E}\|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2}{\sqrt{\delta}} + B^2\mathbb{E}\left(\frac{1}{\sqrt{v^{k-2}}} - \frac{1}{\sqrt{v^k}}\right)$, and $a)$ depends on Lemma 4. Thus, we have

$$\begin{aligned}
& \mathbb{E}\langle \mathbf{x}^{k-1} - \mathbf{x}^k, \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) \rangle \\
&\leq \mathbb{E}\langle \mathbf{x}^{k-1} - \mathbf{x}^k, \nabla f(\mathbf{x}^{k-1}) \rangle + \frac{1}{2}\mathbb{E}\frac{\|\mathbf{m}^{k-1}\|^2}{\sqrt{v^{k-1}}} + \phi_k.
\end{aligned} \tag{30}$$

Once with the Lipchitz property,

$$\langle \nabla f(\mathbf{x}^{k-1}), \mathbf{x}^{k-1} - \mathbf{x}^k \rangle \leq f(\mathbf{x}^{k-1}) - f(\mathbf{x}^k) + \frac{L}{2}\|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2. \tag{31}$$

Combing (31) and (30), we then get

$$\begin{aligned}
& \mathbb{E}\langle \mathbf{x}^{k-1} - \mathbf{x}^k, \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) \rangle \leq f(\mathbf{x}^{k-1}) - f(\mathbf{x}^k) \\
&\quad + \frac{L}{2}\|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2 + \frac{1}{2}\mathbb{E}\frac{\|\mathbf{m}^{k-1}\|^2}{\sqrt{v^{k-1}}} + \phi_k.
\end{aligned} \tag{32}$$

Substituting (32) into (25),

$$\begin{aligned}
& \mathbb{E}[F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) - F(\mathbf{x}^k; \xi^k, \xi^{k-1})] \leq 2L\mathbb{E}\|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2 \\
&\quad + f(\mathbf{x}^{k-1}) - f(\mathbf{x}^k) + \frac{1}{2}\mathbb{E}\frac{\|\mathbf{m}^{k-1}\|^2}{\sqrt{v^{k-1}}} + \phi_k.
\end{aligned} \tag{33}$$

Substituting (33) into (27), we then get

$$\mathbb{E}\langle \mathbf{x}^* - \mathbf{x}^k, \mathbf{g}^k \rangle \leq 2L\mathbb{E}\|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2 + f(\mathbf{x}^*) - f(\mathbf{x}^k) + \frac{1}{2}\mathbb{E}\frac{\|\mathbf{m}^{k-1}\|^2}{\sqrt{v^{k-1}}} + \phi_k. \tag{34}$$

According to our algorithm and we denote $\Lambda := \mathbb{E}\langle \mathbf{x}^* - \mathbf{x}^k, \mathbf{g}^k \rangle$, then we have

$$\begin{aligned}
& \mathbb{E}\langle \mathbf{x}^* - \mathbf{x}^k, \mathbf{m}^k \rangle = \mathbb{E}\langle \mathbf{x}^* - \mathbf{x}^k, \theta\mathbf{m}^{k-1} + (1-\theta)\mathbf{g}^k \rangle \\
&= (1-\theta) \cdot \Lambda + \theta\mathbb{E}\langle \mathbf{x}^* - \mathbf{x}^k, \mathbf{m}^{k-1} \rangle \\
&= (1-\theta) \cdot \Lambda + \theta\mathbb{E}\langle \mathbf{x}^* - \mathbf{x}^{k-1}, \mathbf{m}^{k-1} \rangle + \theta\mathbb{E}\langle \mathbf{x}^k - \mathbf{x}^{k-1}, \mathbf{m}^{k-1} \rangle \\
&\stackrel{b)}{\leq} (1-\theta) \cdot \Lambda + \theta\mathbb{E}\langle \mathbf{x}^* - \mathbf{x}^{k-1}, \mathbf{m}^{k-1} \rangle + \eta\theta\mathbb{E}\|\mathbf{m}^{k-1}\|^2 / (v^{k-1})^{\frac{1}{2}},
\end{aligned}$$

where $b)$ depends on that $\langle \mathbf{x}^k - \mathbf{x}^{k-1}, \mathbf{m}^{k-1} \rangle \leq \|\mathbf{x}^k - \mathbf{x}^{k-1}\| \cdot \|\mathbf{m}^{k-1}\| = \|\mathbf{Proj}_{\mathcal{K}}(\mathbf{x}^{k-1} - \eta\mathbf{m}^k) - \mathbf{Proj}_{\mathcal{K}}(\mathbf{x}^{k-1})\| \cdot \|\mathbf{m}^{k-1}\| \leq \|\mathbf{m}^{k-1}\|^2 / (v^{k-1})^{\frac{1}{2}}$. Then, we get

$$B_k \leq (1-\theta)\mathbb{E}\langle \mathbf{x}^* - \mathbf{x}^k, \mathbf{g}^k \rangle + \theta B_{k-1} + \eta\theta A_{k-1}. \tag{35}$$

Substituting (34) into (35), we then proved the desired result.

D.5 Proof of Lemma 6

This proof is identical to the proof of Lemma 3 and will not be reproduced.

D.6 Proof of Lemma 7

Notice that $\mathbb{E}\nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) = \nabla f(\mathbf{x}^{k-1})$, we then get

$$\begin{aligned}
& \mathbb{E}\langle -\nabla f(\mathbf{x}^k)/(v^k)^{1/2}, \mathbf{g}^k \rangle \\
&= \mathbb{E}\langle -\nabla f(\mathbf{x}^k)/(v^k)^{1/2}, \nabla F(\mathbf{x}^k; \xi^k, \xi^{k-1}) \rangle \\
&= -\mathbb{E}\|\nabla f(\mathbf{x}^k)\|^2/(v^k)^{1/2} + \mathbb{E}\langle \nabla f(\mathbf{x}^k)/(v^k)^{\frac{1}{2}}, \nabla f(\mathbf{x}^k) - \nabla f(\mathbf{x}^{k-1}) \rangle \\
&+ \underbrace{\mathbb{E}\langle \nabla f(\mathbf{x}^k)/(v^k)^{\frac{1}{2}}, \nabla f(\mathbf{x}^{k-1}) - \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) \rangle}_{:= (\dagger)} \\
&+ \mathbb{E}\left\langle \frac{\nabla f(\mathbf{x}^k)}{(v^k)^{\frac{1}{2}}}, \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) - \nabla F(\mathbf{x}^k; \xi^k, \xi^{k-1}) \right\rangle.
\end{aligned} \tag{36}$$

The Cauchy's inequality together with the smooth assumption gives us

$$\begin{aligned}
& \mathbb{E}\langle \nabla f(\mathbf{x}^k)/(v^k)^{\frac{1}{2}}, \nabla f(\mathbf{x}^k) - \nabla f(\mathbf{x}^{k-1}) \rangle \\
&\leq |\mathbb{E}\langle \nabla f(\mathbf{x}^k)/(v^k)^{\frac{1}{4}}, [\nabla f(\mathbf{x}^{k-1}) - \nabla f(\mathbf{x}^k)]/(v^k)^{\frac{1}{4}} \rangle| \\
&\leq \frac{1}{4}\mathbb{E}\|\nabla f(\mathbf{x}^k)\|^2/(v^k)^{\frac{1}{2}} + \mathbb{E}\|\nabla f(\mathbf{x}^{k-1}) - \nabla f(\mathbf{x}^k)\|^2/(v^k)^{\frac{1}{2}} \\
&\leq \frac{1}{4}\mathbb{E}\|\nabla f(\mathbf{x}^k)\|^2/(v^k)^{\frac{1}{2}} + L^2\mathbb{E}\|\mathbf{x}^{k-1} - \mathbf{x}^k\|^2/(v^k)^{\frac{1}{2}}.
\end{aligned} \tag{37}$$

Similarly, the Lipschitz property of the stochastic gradient yields

$$\begin{aligned}
& \mathbb{E}\langle \nabla f(\mathbf{x}^k)/(v^k)^{\frac{1}{2}}, \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) - \nabla F(\mathbf{x}^k; \xi^k, \xi^{k-1}) \rangle \\
&\leq \mathbb{E}\left| \left\langle \frac{\nabla f(\mathbf{x}^k)}{(v^k)^{\frac{1}{4}}}, \frac{[\nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) - \nabla F(\mathbf{x}^k; \xi^k, \xi^{k-1})]}{(v^k)^{\frac{1}{4}}} \right\rangle \right| \\
&\leq \frac{1}{4}\mathbb{E}\|\nabla f(\mathbf{x}^k)\|^2/(v^k)^{\frac{1}{2}} + L^2\mathbb{E}\|\mathbf{x}^{k-1} - \mathbf{x}^k\|^2/(v^k)^{\frac{1}{2}}.
\end{aligned} \tag{38}$$

Substituting (38) and (37) into (36), we are then led to

$$\begin{aligned}
& \mathbb{E}\langle -\nabla f(\mathbf{x}^k)/(v^k)^{\frac{1}{2}}, \mathbf{g}^k \rangle \leq -\mathbb{E}\left(\|\nabla f(\mathbf{x}^k)\|^2/(v^k)^{\frac{1}{2}}\right) + (\dagger) \\
&+ \frac{1}{2}\mathbb{E}(\|\nabla f(\mathbf{x}^k)\|^2/(v^k)^{\frac{1}{2}}) + 2L^2\mathbb{E}(\|\mathbf{x}^{k-1} - \mathbf{x}^k\|^2/(v^k)^{\frac{1}{2}}).
\end{aligned} \tag{39}$$

Now, we turn to bound (\dagger) :

$$\begin{aligned}
(\dagger) &= \underbrace{\mathbb{E}\langle \nabla f(\mathbf{x}^{k-1})/(v^{k-1})^{\frac{1}{2}}, \nabla f(\mathbf{x}^{k-1}) - \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) \rangle}_{=0} \\
&+ \mathbb{E}\left\langle \nabla f(\mathbf{x}^k)/(v^k)^{\frac{1}{2}} - \nabla f(\mathbf{x}^{k-1})/(v^k)^{\frac{1}{2}} \right. \\
&+ \left. \nabla f(\mathbf{x}^{k-1})/(v^k)^{\frac{1}{2}} - \nabla f(\mathbf{x}^{k-1})/(v^{k-1})^{\frac{1}{2}}, \nabla f(\mathbf{x}^{k-1}) - \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) \right\rangle \\
&\leq \mathbb{E}\langle [\nabla f(\mathbf{x}^k) - \nabla f(\mathbf{x}^{k-1})]/(v^k)^{\frac{1}{2}}, \nabla f(\mathbf{x}^{k-1}) - \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1}) \rangle \\
&+ 2B^2\mathbb{E}[1/(v^{k-2})^{\frac{1}{2}} - 1/(v^k)^{\frac{1}{2}}] \\
&\leq \frac{\mathbb{E}\|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2}{2} + 2B^2\mathbb{E}[1/(v^{k-2})^{\frac{1}{2}} - 1/(v^k)^{\frac{1}{2}}] \\
&+ \frac{1}{2}\mathbb{E}\|\nabla f(\mathbf{x}^{k-1}) - \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})\|^2/v^k.
\end{aligned}$$

Furthermore, we have

$$\begin{aligned}
& \frac{1}{2} \mathbb{E} \|\nabla f(\mathbf{x}^{k-1}) - \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})\|^2 / v^k \\
&= \frac{1}{2} \mathbb{E} \|\nabla f(\mathbf{x}^{k-1}) - \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})\|^2 / v^{k-2} \\
&+ \frac{1}{2} \mathbb{E} \|\nabla f(\mathbf{x}^{k-1}) - \nabla F(\mathbf{x}^{k-1}; \xi^k, \xi^{k-1})\|^2 (1/v^k - 1/v^{k-2}) \\
&\stackrel{\text{Lemma 4}}{\leq} L^2 \mathbb{E} \|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2 / v^{k-2} + 2 \mathbb{E} \|\mathbf{g}^k\|^2 / v^{k-2} + 2B^2 \mathbb{E} (1/v^k - 1/v^{k-2}) \\
&\leq \frac{L^2}{\delta} \mathbb{E} \|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2 + 2 \mathbb{E} \|\mathbf{g}^k\|^2 / v^k + 4B^2 \mathbb{E} (1/v^k - 1/v^{k-2}).
\end{aligned}$$

Thus, (39) can also be bounded as

$$\begin{aligned}
\mathbb{E} \langle -\nabla f(\mathbf{x}^k) / (v^k)^{\frac{1}{2}}, \mathbf{g}^k \rangle &\leq -\frac{1}{2} \mathbb{E} \left(\|\nabla f(\mathbf{x}^k)\|^2 / (v^k)^{\frac{1}{2}} \right) \\
&+ 6B^2 \mathbb{E} [1/(v^{k-2})^{\frac{1}{2}} - 1/(v^k)^{\frac{1}{2}}] + 2 \mathbb{E} \|\mathbf{g}^k\|^2 / v^k \\
&+ (2L^2/\sqrt{\delta} + L^2/\delta + 1/2) \mathbb{E} \|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2.
\end{aligned} \tag{40}$$

We also use a shorthand notation $\Lambda := \mathbb{E} \langle -\nabla f(\mathbf{x}^k) / (v^k)^{\frac{1}{2}}, \mathbf{g}^k \rangle$ and then

$$\begin{aligned}
& \mathbb{E} \left(\langle -\nabla f(\mathbf{x}^k), \mathbf{m}^k / (v^k)^{\frac{1}{2}} \rangle \right) \\
&= \mathbb{E} \left(\langle -\nabla f(\mathbf{x}^k) / (v^k)^{\frac{1}{2}}, \theta \mathbf{m}^{k-1} + (1 - \theta) \mathbf{g}^k \rangle \right) \\
&= (1 - \theta) \cdot \Lambda + \theta \mathbb{E} \langle -\nabla f(\mathbf{x}^k) / (v^k)^{\frac{1}{2}}, \mathbf{m}^{k-1} \rangle \\
&\stackrel{a)}{\leq} (1 - \theta) \cdot \Lambda + \theta \mathbb{E} \langle -\nabla f(\mathbf{x}^{k-1}) / (v^{k-1})^{\frac{1}{2}}, \mathbf{m}^{k-1} \rangle + \theta \eta \mathbb{E} \|\mathbf{m}^{k-1}\|^2 / v^{k-1}.
\end{aligned} \tag{41}$$

where $a)$ uses the Cauchy-Schwarz inequality and the Lipschitz property of f and $v^{k-1} \leq v^k$. Substituting (41) into (40), we then proved the result.

E Ablation study of the parameter θ

In this section, we investigate the influence of the parameter θ on the convergence rate of AOGD. The experimental setup adheres to the one detailed in Section 3, focusing on the hinge loss applied to the `ijcnn1` dataset.

We fix the parameters η and R at values 1.0 and 10.0, respectively, while varying θ over the set $\{0.1, 0.3, 0.5, 0.7, 0.9\}$. The rationale behind these chosen values for η and R stems from optimization results observed at $\theta = 0.9$, aligning with the configurations presented in Table 2 and Fig. 2.

As depicted in Figure 4, the convergence rate of AOGD across different θ values is displayed. Notably, the convergence rate appears robust to the specific selection of θ with slightly better performance observed at an intermediate θ value, specifically $\theta = 0.7$.

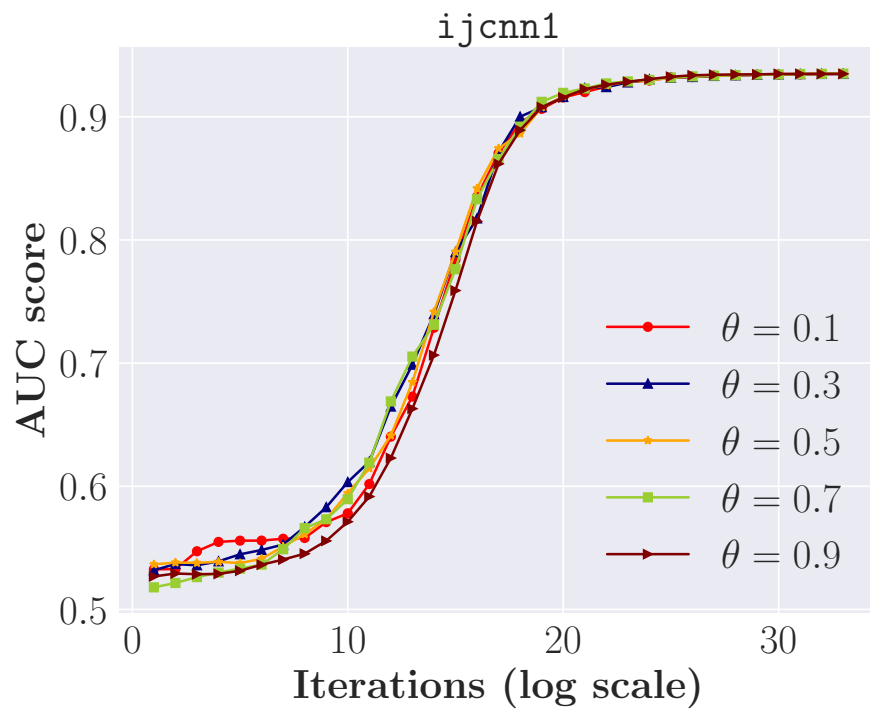


Figure 4: The effect of θ on the convergence rate of AOGD with hinge loss on `ijcnn1` dataset. The result suggests AOGD achieves robust performance across different θ values.