
Divergence Results and Convergence of a Variance Reduced Version of ADAM

Ruiqi Wang
Northwestern University

Diego Klabjan
Northwestern University

Abstract

Stochastic optimization algorithms using exponential moving averages of the past gradients, such as ADAM, RMSProp and AdaGrad, have been having great successes in many applications, especially in training deep neural networks. ADAM in particular stands out as efficient and robust. Despite of its outstanding performance, ADAM has been proved to be divergent for some specific problems. We revisit the divergent question and provide divergent examples under stronger conditions such as in expectation or high probability. Under a variance reduction assumption, we show that an ADAM-type algorithm converges, which means that it is the variance of gradients that causes the divergence of original ADAM. To this end, we propose a variance reduced version of ADAM and provide a convergent analysis of the algorithm. Numerical experiments show that the proposed algorithm has as good performance as ADAM. Our work suggests a new direction for fixing the convergence issues.

1 Introduction

Stochastic optimization based on mini-batch is a common training procedure in machine learning. Suppose we have finitely many differentiable objectives $\{f_n(w)\}_{n=1}^N$ defined on \mathbb{R}^d with N being the size of the training set. In each iteration, a random index set \mathcal{B}_t is selected from $\{1, \dots, N\}$ and the update is made based on the mini-batch loss $F^{\mathcal{B}_t}(w) = \frac{1}{b} \sum_{n \in \mathcal{B}_t} f_n(w)$, where $b = |\mathcal{B}_t|$ is the batch size. The goal is to minimize the empirical risk $\min_{w \in \mathbb{R}^d} F(w) := \frac{1}{N} \sum_{n=1}^N f_n(w)$.

First order methods, which make updates based on the information of the gradient of mini-batch loss functions, prevail in practice, [Goodfellow et al., 2016]. A simple method is stochastic gradient descent (SGD), where the model parameters are updated at the negative direction of the mini-batch loss gradient in each iteration. Although SGD is straightforward and is proved to be convergent, the steps of SGD near the minima are very noisy and take longer to

converge. Several adaptive variants of SGD, such as AdaGrad [Duchi et al., 2011], RMSProp [Hinton et al., 2012] and ADAM [Kingma and Ba, 2015], are proved to converge faster than SGD in practice. These methods take the historical gradients into account. Specifically, instead of using a predefined learning rate schema, they adjust the step size automatically based on the information from the past mini-batch losses. AdaGrad is the earliest algorithm in the adaptive method family and performs better than SGD when gradients are sparse. Although AdaGrad has great theoretical properties for convex loss, it does not work well practically in training. RMSProp replaces the sum of square scaling in AdaGrad with exponential moving average and fixes the rapid decay of the learning rate in AdaGrad. ADAM-type algorithms combine the exponential moving average of both first and second order moments. The original ADAM enjoys the advantages of AdaGrad in sparse problems and RMSProp in non-stationary problems and became one of the most popular optimization methods in practice.

Yet, ADAM may fail to solve some problems. Reddi et al. [Reddi et al., 2018] found a flaw in the proof of convergence in [Kingma and Ba, 2015] and proposed a divergent example for online ADAM. Based on the divergent example, they pointed out that when some large, informative but rare gradients occur, the exponential moving average would make them decay quickly and hence would lead to the failure of convergence. To this end, Reddi et al. proposed two variants of ADAM to fix this problem. The first proposal, known as AMSGrad, suggests taking the historical maximum of the ADAM state v_t in order to obtain ‘long-term memories’ and prevent the large and informative gradients from being forgotten. Although this helps keeping the information of large gradients, it hurts the adaptability of ADAM. If the algorithm is exposed to a large gradient at early iterations, the v_t parameter will stay constant, hence the algorithm will not automatically adapt the step size, and it will degenerate to a momentum method. Another intuitive criticism is that keeping v_t increasing is against what one expects, since if the algorithm converges, the norm of gradients should decrease and $v_{t+1} - v_t = (1 - \beta_2)(g_t^2 - v_t)$ is more likely to be negative, where g_t is the stochastic gradient in step t and β_2 is a hyper parameter.

Several other proposals tried to fix the divergent problem of

ADAM. The second variant proposed in [Reddi et al., 2018], called ADAMNC, requires the second order moment hyper-parameter β_2 to increase and to satisfy several conditions. However the conditions are hard to check. Although they claim that $\beta_{2,t} = 1 - 1/t$ satisfies the conditions, this case is actually AdaGrad, which is already well-known for its convergence. Zhou et al. [Zhou et al., 2019] analyzed the divergent example in [Reddi et al., 2018], and pointed out that the correlation of v_t and g_t causes divergence of ADAM, and proposed a decorrelated variant of ADAM. The theoretical analysis in [Zhou et al., 2019] is based on complex assumptions and they do not provide a convergence analysis of their algorithm. Several other works, such as [Guo et al., 2021, Shi et al., 2020, Wang et al., 2019, Zou et al., 2019] suggested properly tuning the hyper-parameters of ADAM-type algorithms had helped with convergence in practice.

It is empirically well-known that larger batch size reduces the variance of the loss of a stochastic optimization algorithm. [Qian and Klabjan, 2020] gave a theoretical proof that the variance of the stochastic gradient is proportional to $1/b$. Although several works connected the convergence of ADAM with the mini-batch size, the direct connection between convergence and variance is wanted. For the full-batch case (i.e., where there is no variance), [De et al., 2018] showed that ADAM converges under some specific scheduling of learning rates. [Shi et al., 2020] showed the convergence of full gradient ADAM and RMSProp with the learning rate schedule $\alpha_t = \alpha/\sqrt{t}$ and constants β_1 and β_2 satisfying $\beta_1 < \sqrt{\beta_2}$. For the stochastic setting with a fixed batch size, Zaheer et al. [Zaheer et al., 2018] proved that the expected norm of the gradient can be bounded into a neighborhood of 0, whose size is proportional to $1/b$. They suggested to increase the batch size with the number of iterations in order to establish convergence. One question is that whether there exists a threshold of batch size $b^* < N$, such that any batch size larger than b^* guarantees convergence. We show that even when $b = N - 1$, there still exist divergent examples of ADAM. This means that although large batch size helps tighten the optimality gap, the convergence issue is not solved as long as the variance exists. Another possible convergent result is to analyze the convergence in expectation or high probability under a stochastic starting point. However our divergent result holds for any initial point, which rules out this possibility.

Without relying on the mini-batch size, we make a direct analysis of variance and the convergence of ADAM. We first show a motivating result which points out that the convergence of an ADAM-type algorithm can be implied by reducing the variance. Motivated by this, we propose a variance reduced version of ADAM, called VRADAM, and show that VRADAM converges. We provide two options regarding to resetting of ADAM states during the full gradient steps, and recommend the resetting option based on a

theoretical analysis herein and computational experiments. Finally, we conduct several computational experiments, and show that our algorithm performs as well as the original version of ADAM.

In Section 3, we show a divergent example. Using contradiction by assuming the algorithm converges, we show that the expected update of iterates is larger than a positive constant, which means that it is impossible for the algorithm to converge to an optimal solution, which contradicts with the assumption. In Section 5, we prove the convergence of VRADAM. The main proof technique applied is to properly bound the difference between the estimated gradients and the true value of gradients. By bounding the update of the objective function in each iterate, we can further employ the strong convexity assumption and conclude convergence.

Our contributions are as follows.

1. We provide an unconstrained and strongly convex stochastic optimization problem on which the original ADAM diverges. We show that the divergence holds for any initial point, which rules out all of the possible weaker convergent results under stochastic starting point.
2. We construct a divergent mini-batch problem with $b = N - 1$, and conclude that there does not exist a convergent threshold for the mini-batch size.
3. We propose a variance reduced version of ADAM. We provide convergence results of the variance reduced version for strongly convex objectives to optimality or non-convex objectives. We show by experiments that the variance reduction does not harm the numerical performance of ADAM.

In Section 2, we review the literature on the topics of the convergence/divergence issue of ADAM and variance reduction optimization methods. In Section 3 we provide divergent examples for stochastic ADAM. We show that the example is divergent for large batch sizes, which disproves the existence of a convergence threshold of mini-batch size. In Section 4 we start from a reducing variance condition and prove the convergence of an ADAM-type algorithm under this condition. In Section 5 we propose a variance reduced version of ADAM. We show that resetting the states in the algorithm helps with the performance. We also provide a convergence result of our variance reduced ADAM. In Section 6 we conduct several numerical experiments and show the convergence and sensitivity of the proposed algorithm.

2 Literature Review

Convergence of ADAM: Reddi et al. [Reddi et al., 2018] firstly pointed out the convergence issue of ADAM and proposed two convergent variants: (a) AMSGrad takes the

historical maximum value of v_t to keep the step size decreasing and (b) ADAMNC requires the hyper-parameters to satisfy specific conditions. Both of the approaches require that β_1 varies with time, which is inconsistent with practice. Fang and Klabjan [Fang and Klabjan, 2019] gave a convergence proof for AMSGrad with constant β_1 and [Alacaoglu et al., 2020] provided a tighter bound. Enlarging the mini-batch size is another direction. [De et al., 2018] and [Shi et al., 2020] proved the convergence of ADAM for full batch gradients and [Zaheer et al., 2018] showed the convergence of ADAM as long as the batch size is of the same order as the maximum number of iterations, but one criticism is that such a setting for the batch size is very inefficient in practice since the calculation of a large batch gradient is expensive. Several works, such as [Guo et al., 2021, Zou et al., 2019, Wang et al., 2019] proposed guidelines on setting hyper-parameters in order to obtain convergent results. [Guo et al., 2021] showed that as long as β_1 is close enough to 1, in particular, $1 - \beta_{1,t} \propto 1/\sqrt{t}$, ADAM establishes a convergent rate of $\mathcal{O}(1/\sqrt{T})$. However, since [Reddi et al., 2018] proposed the divergent example for any fixed β_1 and β_2 such that $\beta_1 < \sqrt{\beta_2}$, there is no hope to extend the results of [Guo et al., 2021] to constant momentum parameters. [Zou et al., 2019] also provided a series of conditions under which ADAM could converge. Specifically, they require the quantity $\alpha_t/\sqrt{1 - \beta_{2,t}}$ to be ‘almost’ non-increasing. [Wang et al., 2019] proposed to set the denominator hyper-parameter ϵ to be $1/t$, and showed the convergence of ADAM for strongly convex objectives. The aforementioned works focus on setting the hyper-parameters in ADAM. On contrary, our work proposes a new algorithm that only requires basic and common conditions. We show a $\mathcal{O}(T^{-p})$ convergence rate for $0 < p < 1$ where p is dependent on hyper-parameters.

Variance reduction: The computational efficiency issues of full gradient descent methods get more severe with a large data size, but employing stochastic gradient descent may cause divergence because of the issue of variance. One classic method for variance reduction is to use mini-batch losses with a larger batch size, which however does not guarantee the variance to converge to zero. As an estimation of the full gradient, the stochastic average gradient (SAG) method [Le Roux et al., 2012] uses an average of $\nabla f_i(x_{k_i})$, where k_i is the most recent step index when sample i is picked. Although the convergence analysis of SAG provided in [Schmidt et al., 2017] showed its remarkable linear convergence, the estimator of the descent direction is biased and the analysis of SAG is complicated. SAGA [Defazio et al., 2014], an unbiased variant to SAG introduced a concept called ‘covariates’ and guarantees linear convergence as well. Both SAG and SAGA require the memory of $\mathcal{O}(Nd)$, which is expensive when the data set is large. SVRG [Johnson and Zhang, 2013] constructs two layers of iterations and calculates the full gradient as an auxiliary vector for variance reduction before starting each inner loop.

It only requires a memory of $\mathcal{O}(d)$. Most of the literature on variance reduction focus on the convergence rate and memory requirement on the plain SGD algorithm. Recently, [Dubois-Taine et al., 2021] combined AdaGrad with SVRG for robustness in the learning rate. Our work introduces the idea of variance reduction to the convergence analysis of ADAM. It initiates the idea of the dynamic learning rate to SVRG.

3 Divergent examples for stochastic ADAM with large batch size

Several recent works [Shi et al., 2020, Zaheer et al., 2018, De et al., 2018] have suggested increasing the mini-batch size may help with convergence of ADAM. In particular, vanilla ADAM is convergent if the mini-batch size b is equal to the size of the training set, or it increases in the same order as training iterates. An interesting question is whether there exists a threshold of the batch size $b^* = b(N)$, which is smaller than N , such that $b > b^*$ implies convergence of ADAM. If such a threshold exists, the convergence can be guaranteed by a sufficiently large, but neither increasing nor as large as the training set size, batch size. Unfortunately, such a threshold does not exist. In fact, we show in this section that as long as the algorithm is not full batch, one can find a divergent example of ADAM.

Another aspect of interest is if ADAM converges on average or with high probability. Our example establishes non-convergence for any initial data point (even starting with an optimal one). We conclude that a probabilistic statement is impossible if stochasticity comes from either sampling or the initial point.

Reddi et al. firstly proposed a divergent example for ADAM in [Reddi et al., 2018]. The example, which is under the population loss minimization framework, consists of two linear functions defined on a finite interval. One drawback of this example is that the optimization problem is constrained, yet training in machine learning is usually an unconstrained problem. Under the unconstrained framework, the example proposed in [Reddi et al., 2018] does not have a minimum solution, hence it does not satisfy the basic requirements. We firstly propose an unconstrained problem under the population loss minimization framework.

Let a random variable ξ take discrete value from the set $\{1, 2\}$, and set $\mathbb{P}(\xi = 1) = \frac{1+\delta}{1+\delta^4}$ for some $\delta > 1$. Furthermore, we define the estimation of gradients by

$$\mathcal{G}(w; 1) = \frac{w}{\delta} + \delta^4 \quad \text{and} \quad \mathcal{G}(w; 2) = \frac{w}{\delta} - 1,$$

which implies the stochastic optimization problem with the loss functions

$$f_1(w) = \frac{w^2}{2\delta} + \delta^4 w \quad \text{and} \quad f_2(w) = \frac{w^2}{2\delta} - w$$

with the corresponded probability distribution with respect to ξ . The population loss is given as $F(w) = \mathbb{E}_\xi [f_\xi(w)]$. We call this stochastic optimization problem the **Original Problem**(δ), or **OP**(δ) for short. We should note that OP(δ) is defined on \mathbb{R} , thus it is unconstrained. In addition, it is a strongly convex problem. As a divergent property of OP(δ), we show the following result.

Theorem 1 *There exists a $\delta^* > 2$ such that for any $\delta > \delta^*$ and any initial point w_1 , ADAM diverges in expectation on OP(δ), i.e., $\mathbb{E}[F(w_t)] \not\rightarrow F^*$ where F^* is the optimal value of $F(w)$.*

The proof of Theorem 1 is given in the appendix, where we show that for large enough δ , the expectation of the ADAM update between two consequential iterates is always positive. As a consequence, the iterates keep drifting from the optimal solution. The divergent example also tells us that strong convexity and the relaxation of constraints cannot help with the convergence of ADAM.

Based on the construction of OP(δ), we can give the divergent examples for any fixed mini-batch size.

Theorem 2 *For any fixed b , there exists an N_b^* , such that for any $N > N_b^*$, there exists a mini-batch problem with sample size N and batch size b where ADAM diverges for any initial point.*

Even if the batch size is unreasonably large, say $b = N - 1$, we can still construct the divergent example based on OP(δ) as stated next.

Theorem 3 *There exists an N^* such that for any $N > N^*$, there exists a mini-batch problem with sample size N and batch size $b = N - 1$ where ADAM diverges for any initial point.*

In conclusion, Theorem 2 and Theorem 3 extinguish the hope of finding a large enough batch size for stochastic ADAM to converge. Among the related works regarding the convergence of ADAM and batch size, larger batch size is always suggested, but the results in this section have enlightened the limitations of such approaches.

4 Motivation

In this section, we stick with the general ADAM algorithm described in Algorithm 1. To analyze, we make several assumptions on the gradient estimator and objective.

Assumption 1 *The gradient estimator $\mathcal{G} : \mathbb{R}^d \times \Omega \rightarrow \mathbb{R}^d$ and objective $F : \mathbb{R}^d \rightarrow \mathbb{R}$ satisfy the following:*

1. \mathcal{G} is unbiased, i.e., for any $w \in \mathbb{R}^d$, $\mathbb{E}_\xi [\mathcal{G}(w; \xi)] = \nabla F(w)$.

Algorithm 1 General ADAM

Require: Gradient estimation $\mathcal{G}(\cdot; \cdot)$, seed generation rule \mathbb{P}_ξ , initial point w_1 , mini-batch size b , learning rate α_t , exponential decay rates $\beta_1, \beta_2 \in [0, 1)$, denominator hyperparameter $\epsilon > 0$.

$m_0 \leftarrow 0, v_0 \leftarrow 0$

for $t \in 1, \dots, T$ **do**

 Sample $\xi_t \sim \mathbb{P}_\xi$

$g_t \leftarrow \mathcal{G}(w_t; \xi_t)$

$m_t \leftarrow \beta_1 m_{t-1} + (1 - \beta_1) g_t$

$v_t \leftarrow \beta_2 v_{t-1} + (1 - \beta_2) g_t \odot g_t$

$V_t \leftarrow \text{diag}(v_t) + \epsilon I_d$

$w_{t+1} \leftarrow w_t - \alpha_t V_t^{-1/2} m_t$

end for

2. *There exists a constant $0 < L < +\infty$, such that for any $\xi \in \Omega$ and $w, \bar{w} \in \mathbb{R}^d$, we have $\|\mathcal{G}(w; \xi) - \mathcal{G}(\bar{w}; \xi)\|_2 \leq L \|w - \bar{w}\|_2$ and $\|\nabla F(w) - \nabla F(\bar{w})\|_2 \leq L \|w - \bar{w}\|_2$.*

3. *There exists a constant $0 < G < +\infty$, such that for any $\xi \in \Omega$ and $w \in \mathbb{R}^d$, we have $\|\mathcal{G}(w; \xi)\|_2 \leq G$ and $\|F(w)\|_2 \leq G$.*

As this point convexity is not needed. We mainly focus on the variance of the gradient estimator. The common assumptions in the literature are that the variance is bounded by a constant, [Zaheer et al., 2018], or a linear function of the square of the norm of the objective $\text{Var}(\mathcal{G}_i(w; \xi)) \leq C_1 + C_2 \|\nabla F(w)\|_2^2$, [Bottou et al., 2018]. Another assumption made in [Shi et al., 2020, Vaswani et al., 2019] is called the ‘strongly growth condition’ which is $\sum_{n=1}^N \|\nabla f_n(w)\|_2^2 \leq C \|\nabla F(w)\|_2^2$ for some $C > 0$. Note that for vanilla ADAM where $\mathcal{G}(w; \xi) = \nabla F^{\mathcal{B}}(w)$ the strongly growth condition implies that $\nabla F^{\mathcal{B}}(w^*) = 0$ if and only if $\nabla F(w^*) = 0$. As a result the strongly growth condition implies that $\text{Var}(\mathcal{G}(w; \xi)) \leq 2LE \|\|w - w^*\|_2^2$, given Lipschitz smooth gradients for full-batch and mini-batch losses. For those iterates close to a saddle point, the variance is automatically reduced, because $\|w - w^*\|_2^2$ is small. However, the strongly growth condition is so strong that the majority of practical problems do not satisfy it. In fact, one observation of OP(δ) is that the variance is a constant, which also breaks the strongly growth condition.

In this section, as a motivative result, let us assume the variance of the gradient estimator is reduced a priori. Let us denote a series of positive constants $\{\lambda_t\}_{t=1}^T$ such that for any $t = 1, \dots, T$, we have $\text{Var}(\mathcal{G}(w_t; \xi_t)) \leq \lambda_t$. For the objective with a finite lower bound, we have the following result.

Theorem 4 *Let Assumption 1 be satisfied, and assume that $F(w)$ is lower bounded by $F_{\text{inf}} > -\infty$. Then for any initial*

point w_1 , ADAM satisfies

$$\min_{1 \leq t \leq T} \mathbb{E} \left[\|\nabla F(w_t)\|_2^2 \right] \leq \mathcal{O} \left(\frac{\sum_{t=1}^T \alpha_t^2}{\sum_{t=1}^T \alpha_t} + \frac{\sum_{t=1}^T \alpha_t \lambda_t}{\sum_{t=1}^T \alpha_t} \right).$$

The proof is in the appendix. Let us assume that the two common conditions $\sum_{t=1}^{\infty} \alpha_t = \infty$ and $\sum_{t=1}^{\infty} \alpha_t^2 < \infty$ are satisfied. Theorem 4 shows that ADAM converges if $\sum_{t=1}^{\infty} \alpha_t \lambda_t < +\infty$. In fact, $\lambda_t \rightarrow 0$ as $t \rightarrow \infty$ implies that $\sum_{t=1}^T \alpha_t \lambda_t / \sum_{t=1}^T \alpha_t \rightarrow 0$, and hence it leads to convergence of the algorithm.

We emphasize that since the assumption on variance is made on the algorithmic iterates $\{w_t\}_{t=1}^T$, it is very difficult to be checked for a specific problem in advance. However, we showed that if the variance is convergent, an ADAM-type algorithm converges. We show next that the algorithm we propose has convergent variance and furthermore is convergent.

5 Variance Reduced ADAM

Algorithm 2 Variance Reduced ADAM

Require: Loss functions $\{f_n(w)\}_{n=1}^N$, initial point \tilde{w}_1 , learning rate α_t , exponential decay rates $\beta_1, \beta_2 \in [0, 1)$, denominator hyper-parameter $\epsilon > 0$, inner iteration size m . Initialize $m_m^{(0)} \leftarrow 0, v_m^{(0)} \leftarrow 0$.

for $t = 1, \dots, T$ **do**

 Compute full-batch gradient $\nabla F(\tilde{w}_t)$

$w_1^{(t)} \leftarrow \tilde{w}_t$

Option A: (Resetting) $m_0^{(t)} \leftarrow 0, v_0^{(t)} \leftarrow 0$

Option B: (No Resetting) $m_0^{(t)} \leftarrow m_m^{(t-1)}$ and

$v_0^{(t)} \leftarrow v_m^{(t-1)}$

for $k = 1, \dots, m$ **do**

 Sample $\mathcal{B}_k^{(t)}$ from $\{1, \dots, N\}$ with $|\mathcal{B}_k^{(t)}| = b$.

$g_k^{(t)} \leftarrow \nabla F^{\mathcal{B}_k^{(t)}}(w_k^{(t)}) - \nabla F^{\mathcal{B}_k^{(t)}}(\tilde{w}_t) + \nabla F(\tilde{w}_t)$

$m_k^{(t)} \leftarrow \beta_1 m_{k-1}^{(t)} + (1 - \beta_1) g_k^{(t)}$

$v_k^{(t)} \leftarrow \beta_2 v_{k-1}^{(t)} + (1 - \beta_2) g_k^{(t)} \odot g_k^{(t)}$

Option A: $\tilde{m}_k^{(t)} \leftarrow \frac{m_k^{(t)}}{1 - \beta_1^k}, \tilde{v}_k^{(t)} \leftarrow \frac{v_k^{(t)}}{1 - \beta_2^k}$

Option B: $\tilde{m}_k^{(t)} \leftarrow \frac{m_k^{(t)}}{1 - \beta_1^{k+(t-1)m}}, \tilde{v}_k^{(t)} \leftarrow$

$\frac{v_k^{(t)}}{1 - \beta_2^{k+(t-1)m}}$

$V_k^{(t)} \leftarrow \text{diag}(\tilde{v}_k^{(t)} + \epsilon)$

$w_{k+1}^{(t)} \leftarrow w_k^{(t)} - \alpha_t (V_k^{(t)})^{-1/2} \tilde{m}_k^{(t)}$

end for

$\tilde{w}_{t+1} \leftarrow w_{m+1}^{(t)}$

end for

Variance reduction for random variables is a common topic in many fields. In general, an unbiased variance reduction

of a random variable X is $\tilde{X} = X - Y + \mathbb{E}Y$, which establishes the variance $\text{Var}(\tilde{X}) = \text{Var}(X) + \text{Var}(Y) - 2\text{Cov}(X, Y) < \text{Var}(X)$ given $\text{Cov}(X, Y) > \text{Var}(Y)/2$, i.e., X and Y are positively correlated at a sufficient level. In the context of stochastic gradient descent, the random variable for variance reduction is $\mathcal{G}(w_t; \xi_t)$, the gradient of mini-batch loss $\nabla F^{\mathcal{B}_t}(w_t)$. Johnson and Zhang [Johnson and Zhang, 2013] proposed a solution for SGD. They suggested the associate random variable to be the gradient of the same mini-batch loss at a previous iterate \tilde{w} . Since the expectation of a mini-batch gradient is the full-batch gradient, the descent direction becomes $g_t = \nabla F^{\mathcal{B}_t}(w_t) - \nabla F^{\mathcal{B}_t}(\tilde{w}) + \nabla F(\tilde{w})$. Vector \tilde{w} is known as the snapshot model. Since calculation of the full batch gradient at \tilde{w} is required, [Johnson and Zhang, 2013] proposed to save the snapshot model every m iterations, which is known as the SVRG algorithm. Inspired by SVRG and motivated by the result in Section 4, we propose the combination of the variance reduce method and ADAM, called VRADAM (Algorithm 2).

An intuitive analysis of variance of the update direction

$$\begin{aligned} \text{Var}(g_{k,i}^{(t)}) &= \text{Var}(\nabla_i F^{\mathcal{B}_k^{(t)}}(w_k^{(t)}) - \nabla_i F^{\mathcal{B}_k^{(t)}}(\tilde{w}_t)) \\ &\leq \mathbb{E} \left[\left(\nabla_i F^{\mathcal{B}_k^{(t)}}(w_k^{(t)}) - \nabla_i F^{\mathcal{B}_k^{(t)}}(\tilde{w}_t) \right)^2 \right] \\ &\leq L^2 \mathbb{E} \left[\left\| w_k^{(t)} - \tilde{w}_t \right\|_2^2 \right]. \end{aligned}$$

is that as the iterates become close to the optimal point, the variance is reduced simultaneously, which guarantees a similar condition of variance as the strongly growth condition.

5.1 Resetting/No resetting options

We provide two options with regard to the update of ADAM states. In one option, we reinitialize the ADAM states at the beginning of each outer iteration, while the other option keeps the state through the whole training process. Although for the original ADAM, resetting the states harms the performance of the algorithm, we computationally found that the resetting option works better in VRADAM. Intuitively, this is because in each inner loop, the first step $g_1^{(t)}$ is always the full gradient direction, which makes a more efficient update than the direction adapted by previous ADAM states. In order to support our argument, we provide a theoretical analysis of an example. If we fix the initial point $w_1 \in \mathbb{R}$ and the mini-batch losses $F^{\mathcal{B}_1^{(1)}}, F^{\mathcal{B}_2^{(1)}}, \dots, F^{\mathcal{B}_m^{(1)}}$, the iterates are identical between the two options through the $t = 1$ iteration. We consider the objective values after the first update in the second outer iteration, i.e. $F(w_2^{(2)})$. At the end of $t = 1$ iteration, we obtain $w_{m+1}^{(1)} = w_1^{(2)} = \tilde{w}_2$ and the ADAM states $m_{m+1}^{(1)}$ and $v_{m+1}^{(1)}$. Then while $t = 2$, the

algorithm makes the first update as follows.

Option A:	Option B:
$m_1^{(2)} = (1 - \beta_1)g_1^{(2)}$	$\hat{m}_1^{(2)} = \beta_1 m_{m+1}^{(1)} + (1 - \beta_1)g_1^{(2)}$
$\tilde{m}_1^{(2)} = g_1^{(2)}$	$\hat{\tilde{m}}_1^{(2)} = \hat{m}_1^{(2)} / (1 - \beta_1^{m+1})$
$v_1^{(2)} = (1 - \beta_2) \left(g_1^{(2)}\right)^2$	$\hat{v}_1^{(2)} = \beta_2 v_{m+1}^{(1)} + (1 - \beta_2) \left(g_1^{(2)}\right)^2$
$\tilde{v}_1^{(2)} = \left(g_1^{(2)}\right)^2$	$\hat{\tilde{v}}_1^{(2)} = \hat{v}_1^{(2)} / (1 - \beta_2^{m+1})$
$w_2^{(2)} = w_1^{(2)} - \alpha_2 \frac{\tilde{m}_1^{(2)}}{\sqrt{\tilde{v}_1^{(2)} + \epsilon}}$	$\hat{w}_2^{(2)} = w_1^{(2)} - \alpha_2 \frac{\hat{\tilde{m}}_1^{(2)}}{\sqrt{\hat{\tilde{v}}_1^{(2)} + \epsilon}}$

We make the following assumptions.

Assumption 2 *The framework described in this section satisfies:*

1. $F(w)$ is c -strongly convex, and its gradient is L -smooth. Each one of $|g_1^{(1)}|, \dots, |g_m^{(1)}|$ and $|g_1^{(2)}|$ is bounded above by $G > 0$.
2. The algorithm makes progress in the $t = 1$ iteration, specifically, $|m_{m+1}^{(1)}| \geq |F'(w_1^{(2)})|$.
3. The hyper-parameters satisfy

$$L\alpha_2 \geq 2\sqrt{G^2 + \epsilon} \quad \text{and} \quad \frac{L}{c} \leq \frac{2\beta_1 - 1}{1 - \beta_1^{m+1}} \sqrt{\frac{\epsilon}{G^2 + \epsilon}}.$$

Notice that the second assumption assumes that the exponential moving average of the steps in the first loop is larger than the full gradient at the beginning of the second loop, which reflects that the algorithm makes progress in the first iteration. The third assumption can be satisfied by β_1 that is close enough to 1 and appropriately selected α_2 . Our concern is to compare the objective values $F(w_2^{(2)})$ and $F(\hat{w}_2^{(2)})$ of the two options. The following theorem shows that option A, i.e., the option where the states of ADAM are reset at the beginning of each outer iteration, makes more efficient descent, hence works better.

Theorem 5 *Given Assumption 2, we have $F(\hat{w}_2^{(2)}) \geq F(w_2^{(2)})$.*

While Theorem 5 allows the preference of option A within the two outer iterations, the computational experiments confirm this choice in general.

5.2 Convergence results for VRADAM with resetting

In the previous section we show the advantage of resetting of the ADAM states every outer iteration over not doing

so by comparing the values of the objectives in the two options. In this section, we provide a convergence proof of the resetting option of VRADAM. Similar to Assumption 1 we make under the general ADAM framework, we make the following specific assumptions for the convergence proof of VRADAM.

Assumption 3 *The loss functions $f_1(w), \dots, f_N(w)$ satisfy the following conditions.*

1. There exists a constant $0 < L < +\infty$, such that for any $n = 1, \dots, N$ and $w, \bar{w} \in \mathbb{R}^d$, we have $\|\nabla f_n(w) - \nabla f_n(\bar{w})\|_2 \leq L\|w - \bar{w}\|_2$.
2. There exists a constant $0 < G < +\infty$, such that for any $n = 1, \dots, N$ and $w \in \mathbb{R}^d$, we have $\|\nabla f_n(w)\|_2 \leq G$.

We show the convergence result of VRADAM for strongly convex functions first.

Theorem 6 *Let Assumption 3 be satisfied and assume that $F(w)$ is strongly convex with parameter c , and let F^* be the unique minimum of F . Let $\alpha_t = \alpha/t$ and we require $C_2 = 2c(1 - \beta_1)/\sqrt{9G^2 + \epsilon} < 1/\alpha m$. Then for any initial point w_1 , Algorithm 2 with Option A satisfies $F(\tilde{w}_T) - F^* \leq \mathcal{O}(T^{-C_2 m \alpha})$ almost surely.*

We remark that the requirement $C_2 m \alpha < 1$ can be satisfied by properly selected α and β_1 . Specifically, when β_1 is close to 1 and α is small, the assumption is more likely to be satisfied. The proof starts by bounding the update of the objective function in one iterate and then applies strong convexity.

As for objectives that are not necessarily convex, we exhibit the following result.

Theorem 7 *Let Assumption 3 be satisfied and assume that $F(w)$ is lower bounded by $F_{\inf} > -\infty$. We require $\sum_{t=1}^{\infty} \alpha_t^2 < +\infty$ and $\sum_{t=1}^{\infty} \alpha_t = +\infty$. Then for any initial point w_1 , Algorithm 2 with Option A satisfies $\liminf_{t \rightarrow \infty} \|\nabla F(\tilde{w}_t)\|_2 = 0$ almost surely.*

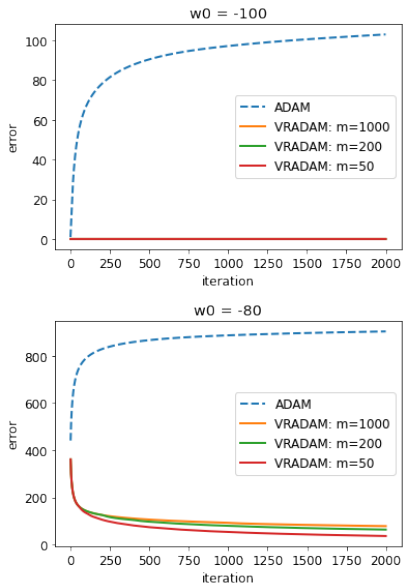
Under the general objective setting, as a corollary of Theorem 7, we can bound the variance of the norm of the gradient.

Corollary 1 *Given the same conditions as in Theorem 7, we have $\liminf_{t \rightarrow \infty} \text{Var}(\|\nabla F(\tilde{w}_t)\|_2) = 0$.*

The proof of the corollary simply comes from

$$\liminf_{t \rightarrow \infty} \text{Var}(\|\nabla F(\tilde{w}_t)\|_2) \leq \liminf_{t \rightarrow \infty} \mathbb{E}(\|\nabla F(\tilde{w}_t)\|_2^2)$$

where the right hand side is 0 because almost sure convergence implies L2 convergence.

Figure 1: The numerical experiments of $OP(\delta)$

6 Numerical Experiments

6.1 Divergent example

In Section 3, we provide an unconstrained stochastic optimization problem where ADAM diverges for any initial point. The goal of this section is to numerically show that ADAM diverges on this problem. Since the construction of the mini-batch problem can be equivalently transformed into a stochastic optimization problem with population loss, we stick to the experiments of $OP(\delta)$ defined in Section 3. Letting $\delta = 10$, the optimal solution of $OP(\delta)$ is $w^* = -\delta^2 = -100$. We consider two cases: when the algorithms start from the optimal solution, i.e., $w_0 = -100$ and when they start from somewhere far from the optimal solution, in the experiments we set $w_0 = -80$. We simulate 1,000 trials for each case and plot the expected L2 error $\mathbb{E}[(w_t - w^*)^2]$ of w . We notice from Figure 1 that even starting from the optimal point, ADAM still diverges, while VRADAM converges well. When $w_0 = -80$, VRADAM eventually converges to the optimal solution while ADAM diverges. This result rules out all of the possible temptations of solving the divergence problem of ADAM by using a stochastic initial point.

6.2 Experiments on machine learning problems

6.2.1 Datasets and implementation

In this section, we compare the numerical performances of VRADAM and ADAM on several real-world classification tasks. The experiments are conducted on the following data sets.

- **Coverage Type** [Blackard et al., 1998]: A dataset predicting forest cover type from 54 cartographic variables. The dataset contains 581,012 data points and assigns them into 7 different categories.
- **MNIST** [Deng, 2012]: A handwritten digit dataset containing 60,000 grey level images with size 28×28 pixels.
- **NSL-KDD** [Tavallaee et al., 2009]: A selected subset of the KDD CUP 99 dataset, which is a public dataset used to train a network intrusion detection system. The subset eliminates repeated data samples and avoids several shortcomings in the original dataset.
- **Embedded CIFAR-10**: CIFAR-10 [Krizhevsky and Hinton, 2009] consists of 60,000 color images in 10 classes. We feed each sample to a pretrained ResNet model [He et al., 2016] and obtain a 1,000 dimensional embedding vector for each image.

We use logistic regression on the four previously mentioned data sets and non-convex deep neural networks on MNIST and Coverage Type datasets. The structures of the deep neural networks used in this section are explained in the appendix. Cross entropy is the underlying loss. While training these models, we fix the batch size $|\mathcal{B}_t| = 64$ and the ADAM hyper-parameters $\beta_1 = 0.9, \beta_2 = 0.999$, which are commonly recommended values in the literature and fine tune the learning rate schedules among $\alpha_t = \alpha_0, \alpha_t = \alpha_0/t$ and $\alpha_t = \alpha_0\gamma^t$, i.e., the constant learning rate, inverse-proportional learning rate and the exponentially decaying learning rate. We perform a grid search among $\alpha_0 \in \{0.0005, 0.001, 0.005, 0.01, 0.05\}$ and $\gamma \in \{0.6, 0.8, 0.95\}$.

In order to eliminate luck from randomness, we ran each experiment setting with 3 different random seeds. We report the average loss of each experiment. We train each setting for 15 epochs for VRADAM and 50 epochs for ADAM. While comparing the performances, we display the loss functions up to convergence. The experiments were conducted in PyTorch 1.12.1 on the Google Colab cloud service.

6.2.2 Main results

As for VRADAM, the number of inner loop iterations m in Algorithm 2 is also a hyper parameter to decide. If m is too large, the variance reduce procedure does not make a real difference, while an improperly small m would lead to a very frequent computation of the full gradients, which is computationally expensive. We recommend that the length of the inner loop should be about the size of an epoch. In other word, $m \approx N/b$, where N is the number of samples and b is the mini-batch size. In our experiments, we test the performance of VRADAM with $m \in \{N/2b, N/b, 2N/b, 4N/b\}$.

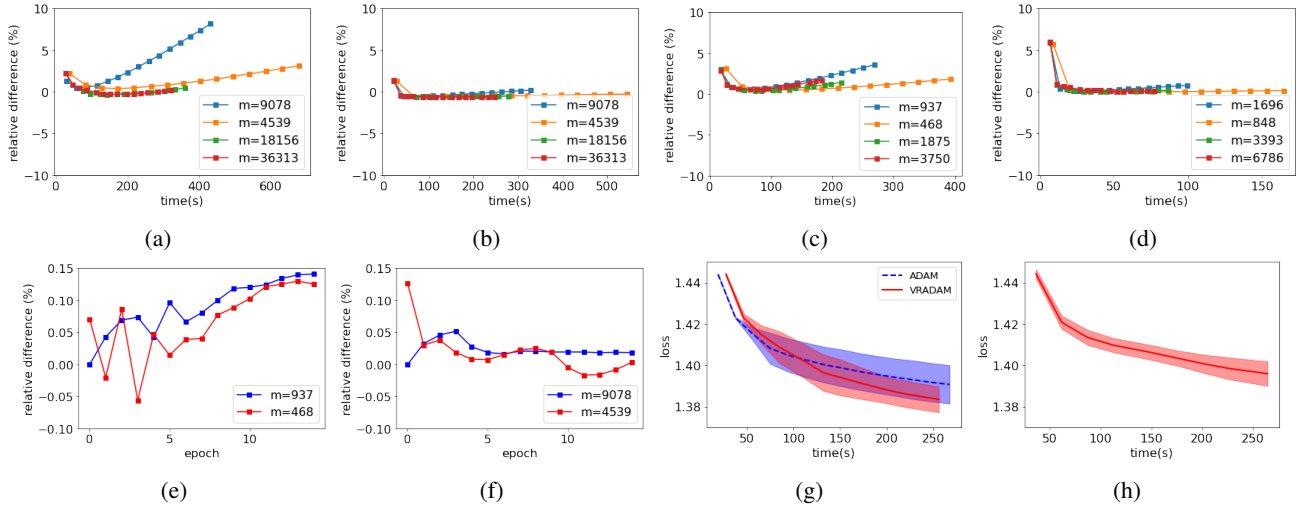


Figure 2: Top: Relative differences of loss function of VRADAM over ADAM on classification tasks of (a) CovType dataset with a feedforward network and (b) CovType with logistic regression (c) MNIST with logistic regression (d) NSL-KDD with logistic regression. Bottom left: Relative differences of loss function of VRADAM without reset over VRADAM with reset on classification tasks of CovType dataset with (e) a feedforward network and (f) logistic regression. Bottom right: Deviation of loss. The shaded areas mark the maximal and minimal losses among the three seeds.

It is unfair to compare the convergence of the training loss with regard to the number of epochs, since the time VRADAM spends on training one epoch is longer than ADAM. Instead, we compare the convergence rate with regard to the computational time.

Figures 2a, 2b, 2c and 2d show the relative difference of the loss of VRADAM with regard to ADAM. It shows that our approach of variance reduction has as good convergence rate as ADAM. Although VRADAM converges slower than ADAM in a few starting iterations, due to the first full gradient computation, the variance reduction approach can catch up and reach a similar convergence rate as ADAM. Specifically, we observe that VRADAM works better on large datasets, such as CovType and NSL-KDD. We also find that VRADAM works better on convex problems by comparing Figures 2a and 2b.

As we mentioned in Section 5, we recommend resetting the optimizer states every inner loop. Figures 2e and 2f display the performance of the resetting option in a few experiments and show that the resetting option helps with the convergence of VRADAM. Additional experiments are shown in the appendix.

In conclusion, we observe that VRADAM outperforms ADAM on large datasets, such as CovType. Such datasets may contain extreme values, which may harm the convergence of ADAM. We find that the computational cost when calculating the full gradients can be compensated by the benefits of quick convergence of VRADAM. We recommend VRADAM over ADAM for tasks with large datasets, in particular if loss is convex.

6.2.3 Sensitivity

We consider the CovType classification task with the FFN model as an example to analyze the sensitivity of our algorithm. As we stated previously, our algorithm fixes the convergence issue of ADAM by reducing the variance. Our experiments take three different seeds and the deviation of the results reflects the variance of the algorithm. Figure 2g shows that VRADAM reduces the noise comparing with ADAM. We also studied the sensitivity of VRADAM over the different initial points, which is shown in Figure 2h.

7 Conclusions

We started from an analysis of the divergence of the original ADAM algorithm and concluded that even strong convexity can not help with the convergence of ADAM. We then gave a convergence analysis under a high-level variance reduction assumption, and concluded that an ADAM-type algorithm converges if its variance is reduced. Inspired by this motivating result and the idea of SVRG, we proposed a variance reduced approach which fixes the original convergence issue of ADAM. We finally showed with numerical experiments that even though our approach requires more gradient computation than ADAM, VRADAM converges as quickly as ADAM after several initial iterations. The motivating results we provided can also lead to other variance reduction approaches, which are possible future research directions on the convergence issue of ADAM.

References

- [Alacaoglu et al., 2020] Alacaoglu, A., Malitsky, Y., Mertikopoulos, P., and Cevher, V. (2020). A new regret analysis for ADAM-type algorithms. In *International Conference on Machine Learning*.
- [Blackard et al., 1998] Blackard, J. A., Dean, D. J., and Anderson, C. W. (1998). UCI machine learning repository: Covertypes data set. <https://archive.ics.uci.edu/ml/datasets/covertypes>.
- [Bottou et al., 2018] Bottou, L., Curtis, F. E., and Nocedal, J. (2018). Optimization methods for large-scale machine learning. *Siam Review*, 60(2):223–311.
- [De et al., 2018] De, S., Mukherjee, A., and Ullah, E. (2018). Convergence guarantees for RMSprop and ADAM in non-convex optimization and an empirical comparison to Nesterov acceleration. *arXiv preprint arXiv:1807.06766*.
- [Defazio et al., 2014] Defazio, A., Bach, F., and Lacoste-Julien, S. (2014). SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives. *Advances in Neural Information Processing Systems*.
- [Deng, 2012] Deng, L. (2012). The MNIST database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine*, 29(6):141–142.
- [Dubois-Taine et al., 2021] Dubois-Taine, B., Vaswani, S., Babanezhad, R., Schmidt, M., and Lacoste-Julien, S. (2021). SVRG meets AdaGrad: Painless variance reduction. *arXiv preprint arXiv:2102.09645*.
- [Duchi et al., 2011] Duchi, J., Hazan, E., and Singer, Y. (2011). Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(7).
- [Fang and Klabjan, 2019] Fang, B. and Klabjan, D. (2019). Convergence analyses of online adam algorithm in convex setting and two-layer relu neural network. *arXiv preprint arXiv:1905.09356*.
- [Goodfellow et al., 2016] Goodfellow, I., Bengio, Y., and Courville, A. (2016). *Deep Learning*. MIT Press. <http://www.deeplearningbook.org>.
- [Guo et al., 2021] Guo, Z., Xu, Y., Yin, W., Jin, R., and Yang, T. (2021). A novel convergence analysis for algorithms of the ADAM family. *arXiv preprint arXiv:2112.03459*.
- [He et al., 2016] He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*.
- [Hinton et al., 2012] Hinton, G., Srivastava, N., and Swersky, K. (2012). Neural networks for machine learning: lecture 6a overview of mini-batch gradient descent. https://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf.
- [Johnson and Zhang, 2013] Johnson, R. and Zhang, T. (2013). Accelerating stochastic gradient descent using predictive variance reduction. *Advances in Neural Information Processing Systems*.
- [Kingma and Ba, 2015] Kingma, D. P. and Ba, J. (2015). ADAM: A method for stochastic optimization. *International Conference on Learning Representations*.
- [Krizhevsky and Hinton, 2009] Krizhevsky, A. and Hinton, G. (2009). Learning multiple layers of features from tiny images. <https://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.222.9220&rep=rep1&type=pdf>.
- [Le Roux et al., 2012] Le Roux, N., Schmidt, M., and Bach, F. (2012). A stochastic gradient method with an exponential convergence rate for finite training sets. *Advances in Neural Information Processing Systems*.
- [Qian and Klabjan, 2020] Qian, X. and Klabjan, D. (2020). The impact of the mini-batch size on the variance of gradients in stochastic gradient descent. *arXiv preprint arXiv:2004.13146*.
- [Reddi et al., 2018] Reddi, S. J., Kale, S., and Kumar, S. (2018). On the convergence of ADAM and beyond. *International Conference on Learning Representations*.
- [Schmidt et al., 2017] Schmidt, M., Le Roux, N., and Bach, F. (2017). Minimizing finite sums with the stochastic average gradient. *Mathematical Programming*, 162(1):83–112.
- [Shi et al., 2020] Shi, N., Li, D., Hong, M., and Sun, R. (2020). RMSprop converges with proper hyperparameter. In *International Conference on Learning Representations*.
- [Tavallae et al., 2009] Tavallae, M., Bagheri, E., Lu, W., and Ghorbani, A. A. (2009). A detailed analysis of the KDD cup 99 data set. In *2009 IEEE symposium on Computational Intelligence for Security and Defense Applications*.
- [Vaswani et al., 2019] Vaswani, S., Bach, F., and Schmidt, M. (2019). Fast and faster convergence of SGD for overparameterized models and an accelerated perceptron. In *The 22nd International Conference on Artificial Intelligence and Statistics*.

- [Wang et al., 2019] Wang, G., Lu, S., Tu, W., and Zhang, L. (2019). SADAM: A variant of ADAM for strongly convex functions. *arXiv preprint arXiv:1905.02957*.
- [Zaheer et al., 2018] Zaheer, M., Reddi, S., Sachan, D., Kale, S., and Kumar, S. (2018). Adaptive methods for nonconvex optimization. *Advances in Neural Information Processing Systems*.
- [Zhou et al., 2019] Zhou, Z., Zhang, Q., Lu, G., Wang, H., Zhang, W., and Yu, Y. (2019). ADAshift: Decorrelation and convergence of adaptive learning rate methods. *International Conference on Learning Representations*.
- [Zou et al., 2019] Zou, F., Shen, L., Jie, Z., Zhang, W., and Liu, W. (2019). A sufficient condition for convergences of ADAM and RMSprop. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.

Divergence Results and Convergence of a Variance Reduced Version of ADAM: Supplementary Materials

A Proofs

A.1 Technical Lemmas

Lemma 1 *In a probability space $(\Omega, \mathcal{F}, \mathbb{P})$, let there be an \mathcal{F} -measurable random variable $X(w)$ and an event $A \in \mathcal{F}$ such that $\mathbb{P}\{A\} > 0$. For a convex function $\phi(x)$, we have*

$$\frac{\mathbb{E}[\phi(X)\mathbf{1}\{A\}]}{\mathbb{P}\{A\}} \geq \phi\left(\frac{\mathbb{E}[X\mathbf{1}\{A\}]}{\mathbb{P}\{A\}}\right).$$

Proof: Let

$$x_0 = \frac{\mathbb{E}[X\mathbf{1}\{A\}]}{\mathbb{P}\{A\}} = \frac{1}{\mathbb{P}\{A\}} \int_A X(w) d\mathbb{P}(w).$$

Since ϕ is convex, there exists a sub-gradient of ϕ at x_0 , i.e., there exists an a such that

$$\phi(x) \geq \phi(x_0) + a(x - x_0)$$

for any $x \in \mathbb{R}$. Then we have

$$\begin{aligned} \frac{\mathbb{E}[\phi(X)\mathbf{1}\{A\}]}{\mathbb{P}\{A\}} &= \frac{1}{\mathbb{P}\{A\}} \int_A \phi(X(w)) d\mathbb{P}(w) \\ &\geq \frac{1}{\mathbb{P}\{A\}} \int_A a(X(w) - x_0) + \phi(x_0) d\mathbb{P}(w) \\ &= \frac{a}{\mathbb{P}\{A\}} \int_A X(w) d\mathbb{P}(w) + \frac{\phi(x_0) - ax_0}{\mathbb{P}\{A\}} \int_A 1 d\mathbb{P}(w) \\ &= ax_0 + \phi(x_0) - ax_0 \\ &= \phi(x_0), \end{aligned}$$

which finishes the proof.

Lemma 2 *For any $x, y > 0$,*

$$(x + y)^3 \leq 4(x^3 + y^3).$$

Proof: For $t \geq 0$, let

$$h(t) := \frac{(1+t)^3}{1+t^3} = 1 + 3\frac{t+t^2}{1+t^3}.$$

The derivative of $h(t)$ reads

$$h'(t) = 3\frac{(1+2t)(1+t^3) - 3t^2(t+t^2)}{(1+t^3)^2} = -3\frac{(t-1)(t+1)^3}{(t^3+1)^2}.$$

Apparently $h(t)$ achieves the maximum at $t = 1$, thus $h(x/y) \leq h(1) = 4$, and we have

$$\frac{(x+y)^3}{x^3+y^3} \leq 4.$$

Lemma 3 Given $\alpha_t = \alpha/t$ and $\beta_1 \in [0, 1)$, there exists a constant $\bar{C} > 0$, such that for any $t \geq 2$ we have

$$\sum_{j=1}^{t-1} \alpha_j \beta_1^{t-j} \leq \bar{C} \alpha_t.$$

Proof: Letting $t^* = \lfloor \frac{t-1}{2} \rfloor$, we have

$$\begin{aligned} \sum_{j=1}^{t-1} \alpha_j \beta_1^{t-j} &= \sum_{j=1}^{t^*} \alpha_j \beta_1^{t-j} + \sum_{j=t^*+1}^{t-1} \alpha_j \beta_1^{t-j} \\ &\leq \alpha \sum_{j=1}^{t^*} \beta_1^{t-j} + \frac{\alpha}{t^*+1} \sum_{j=t^*+1}^{t-1} \beta_1^{t-j} \\ &\leq \frac{\alpha \beta_1^{t-t^*}}{1-\beta_1} + \frac{\alpha}{t^*+1} \frac{\beta_1}{1-\beta_1} \\ &\leq \frac{\alpha \beta_1^{(t+1)/2}}{1-\beta_1} + \frac{2\alpha \beta_1}{(1-\beta_1)} \frac{1}{t-1} = \mathcal{O}(t^{-1}). \end{aligned}$$

Thus there exists a positive constant \bar{C} such that for any $t \geq 2$,

$$\sum_{j=1}^{t-1} \alpha_j \beta_1^{t-j} \leq \bar{C} \alpha_t.$$

Lemma 4 Consider $0 < A < 1$ and $T \geq 2$, and let

$$\begin{aligned} \lambda_{T-1} &= \prod_{t=1}^{T-1} \left(1 - \frac{A}{t}\right), \\ \nu_{T-1} &= \sum_{t=1}^{T-1} \frac{1}{t^2} \prod_{j=t+1}^{T-1} \left(1 - \frac{A}{t}\right). \end{aligned}$$

Then we have

$$\lambda_{T-1} \leq \mathcal{O}(T^{-A})$$

and

$$\nu_{T-1} \leq \mathcal{O}(T^{-A}).$$

Proof: Notice that

$$\log \lambda_{T-1} = \sum_{t=1}^{T-1} \log \left(1 - \frac{A}{t}\right) \leq -A \sum_{t=1}^{T-1} \frac{1}{t} \leq -A \log T,$$

where the first inequality comes from $\log(1-x) \leq -x$ for $x \geq 0$ and the second inequality uses the integral approximation

$$\sum_{t=1}^{T-1} \frac{1}{t} \geq \sum_{t=1}^{T-1} \int_t^{t+1} \frac{1}{s} ds = \int_1^T \frac{1}{s} ds = \log T.$$

Thus $\lambda_{T-1} \leq T^{-A}$. Similarly, we have

$$\log \frac{\lambda_{T-1}}{\lambda_t} = \sum_{k=t+1}^{T-1} \log \left(1 - \frac{A}{k}\right) \leq -A \log \frac{T}{t+1},$$

and then,

$$\nu_{T-1} \leq \sum_{t=1}^{T-1} \frac{1}{t^2} \left(\frac{T}{t+1} \right)^{-A} \leq T^{-A} \sum_{t=1}^{T-1} \frac{(t+1)^A}{t^2} \leq 2^A T^{-A} \sum_{t=1}^{T-1} t^{-2+A}.$$

Again, applying the integral approximation yields

$$\sum_{t=1}^{T-1} t^{-2+A} \leq 1 + \sum_{t=2}^{T-1} \int_{t-1}^t s^{-2+A} ds = 1 + \int_1^{T-1} s^{-2+A} ds \leq 1 + \frac{1}{1-A} = \frac{2-A}{1-A} < \infty.$$

Then we have $\nu_{T-1} \leq \mathcal{O}(T^{-A})$.

A.2 Proof of Theorem 1

Let $p = \mathbb{P}(\xi = 1)$. In each step, the update value is

$$\begin{aligned} \Delta_t &= -\frac{\alpha g_t}{\sqrt{\beta_2 v_{t-1} + (1-\beta_2)g_t^2}} \\ &= \begin{cases} -\frac{\alpha(w_t/\delta + \delta^4)}{\sqrt{\beta_2 v_{t-1} + (1-\beta_2)(w_t/\delta + \delta^4)^2}} & \text{with probability } p \\ \frac{\alpha(1-w_t/\delta)}{\sqrt{\beta_2 v_{t-1} + (1-\beta_2)(1-w_t/\delta)^2}} & \text{with probability } 1-p. \end{cases} \end{aligned}$$

Apparently,

$$F(w) = \frac{w^2}{2\delta} + \delta w,$$

and

$$w^* = -\delta^2.$$

We use contradiction to prove the theorem. Assume that $\mathbb{E}[F(w_t) - F(w^*)] \rightarrow 0$. Notice that

$$F(w_t) - F(w^*) = \frac{1}{2\delta}(w_t - w^*)^2,$$

which means that $\mathbb{E}[F(w_t) - F(w^*)] \rightarrow 0$ is equivalent to $\mathbb{E}[(w_t - w^*)^2] \rightarrow 0$. Let us select $0 < \epsilon < 1/2$, and we choose T_ϵ such that $t > T_\epsilon$ implies $\mathbb{E}[(w_t - w^*)^2] < \epsilon$. The following discussion is based on w_t such that $t > T_\epsilon$.

We have

$$|\Delta_t| = \frac{\alpha|g_t|}{\sqrt{\beta_2 v_{t-1} + (1-\beta_2)g_t^2}} \leq \frac{\alpha|g_t|}{\sqrt{(1-\beta_2)g_t^2}} = \frac{\alpha}{\sqrt{1-\beta_2}}, \quad (1)$$

where the inequality is due to v_k being non-negative for any k .

Let \mathcal{F}_t be the filtration including all the information obtained until the update of w_t , including w_t . We define the following event

$$E := \{|w_t - w^*| < \delta^2\}$$

which is known given \mathcal{F}_t . We have

$$\mathbb{P}\{E^c\} = \mathbb{P}\{|w_t - w^*| \geq \delta^2\} \leq \frac{\mathbb{E}[(w_t - w^*)^2]}{\delta^4} < \frac{\epsilon}{\delta^4}.$$

Given E^c , we simply bound the step size with the lower bound

$$\mathbb{E}[\Delta_t \mathbf{1}\{E^c\}] \geq -\frac{\alpha}{\sqrt{1-\beta_2}} \mathbb{E}[\mathbf{1}\{E^c\}] \geq -\frac{\epsilon}{\delta^4} \frac{\alpha}{\sqrt{1-\beta_2}}.$$

For the samples in E , we have

$$\begin{aligned}
 \mathbb{E}[\Delta_t \mathbf{1}\{E\}] &= \mathbb{E}[\mathbb{E}[\Delta_t \mathbf{1}\{E\} | \mathcal{F}_t]] = \mathbb{E}[\mathbb{E}[\Delta_t | \mathcal{F}_t] \mathbf{1}\{E\}] \\
 &= \mathbb{E} \left[\mathbf{1}\{E\} \left\{ (1-p) \frac{\alpha(1-w_t/\delta)}{\sqrt{\beta_2 v_{t-1} + (1-\beta_2)(1-w_t/\delta)^2}} \right\} \right] \\
 &\quad - \mathbb{E} \left[\left\{ p \frac{\alpha(w_t/\delta + \delta^4)}{\sqrt{\beta_2 v_{t-1} + (1-\beta_2)(w_t/\delta + \delta^4)^2}} \right\} \right] \\
 &\geq \mathbb{E} \left[\mathbf{1}\{E\} (1-p) \frac{\alpha}{\sqrt{\beta_2 v_{t-1} + (1-\beta_2)(1+2\delta)^2}} \right] \\
 &\quad - p \frac{\alpha}{\sqrt{1-\beta_2}} \mathbb{P}\{E\}. \tag{2}
 \end{aligned}$$

In the inequality, the first term is bounded because $-2\delta < w_t/\delta < 0$ by the definition of E and the second term is bounded by the bound of the step length in (1).

By applying Lemma 1 to (2), we have

$$\mathbb{E}[\Delta_t \mathbf{1}\{E\}] \geq (1-p)\mathbb{P}(E) \frac{\alpha}{\sqrt{\beta_2 \mathbb{E}[v_{t-1} \mathbf{1}\{E\}] / \mathbb{P}\{E\} + (1-\beta_2)(1+2\delta)^2}} - p\mathbb{P}(E) \frac{\alpha}{\sqrt{1-\beta_2}}.$$

We next focus on the conditional expectation

$$\mathbb{E}[v_{t-1} \mathbf{1}\{E\}] = (1-\beta_2) \sum_{k=1}^{t-1} \beta_2^{t-1-k} \mathbb{E}[\mathbf{1}\{E\} g_k^2].$$

We claim that for any trajectory in E and for any $k < t$, we have

$$|w_t - w_k| = \left| \sum_{j=k}^{t-1} \Delta_j \right| \leq \sum_{j=k}^{t-1} |\Delta_j| \leq \frac{\alpha(t-k)}{\sqrt{1-\beta_2}}.$$

The last inequality comes from the bound of the step length in (1). Then we have

$$w_t - \frac{\alpha(t-k)}{\sqrt{1-\beta_2}} \leq w_k \leq w_t + \frac{\alpha(t-k)}{\sqrt{1-\beta_2}}.$$

Let us recall that given E , we have $-2\delta^2 < w_k < 0$, and hence

$$-2\delta^2 - \frac{\alpha(t-k)}{\sqrt{1-\beta_2}} \leq w_k \leq \frac{\alpha(t-k)}{\sqrt{1-\beta_2}}. \tag{3}$$

Then for each $k = 1, \dots, t-1$, we obtain

$$\begin{aligned}
 \mathbb{E}[g_k^2 \mathbf{1}\{E\}] &= \mathbb{E}[\mathbb{E}[g_k^2 \mathbf{1}\{E\} | \mathcal{F}_k]] \\
 &\leq \mathbb{E} \left[\left(\mathbb{E}[g_k^{2(1+\mu)} | \mathcal{F}_k] \right)^{1/(1+\mu)} (\mathbb{E}[\mathbf{1}\{E\}])^{\mu/(1+\mu)} \right] \\
 &\leq \mathbb{E} \left[\left(\mathbb{E}[g_k^{2(1+\mu)} | \mathcal{F}_k] \right)^{1/(1+\mu)} \right]
 \end{aligned}$$

where the inequality holds for any μ according to the Holder inequality. Let $0 < \mu < 1/2$. According to the bound given previously in (3) and $\delta \geq 2$, we have

$$\begin{aligned}
 \left(\frac{w_k}{\delta} + \delta^4 \right)^2 &\leq \left(\frac{\alpha(t-k)}{\delta\sqrt{1-\beta_2}} + 2\delta + \delta^4 \right)^2 \\
 \left(1 - \frac{w_k}{\delta} \right)^2 &\leq \left(\frac{\alpha(t-k)}{\delta\sqrt{1-\beta_2}} + 2\delta + 1 \right)^2.
 \end{aligned}$$

Then we derive,

$$\begin{aligned}
 \mathbb{E} \left[g_k^{2(1+\mu)} | \mathcal{F}_k \right] &= p \left(\frac{w_k}{\delta} + \delta^4 \right)^{2(1+\mu)} + (1-p) \left(\frac{w_k}{\delta} - 1 \right)^{2(1+\mu)} \\
 &\leq \frac{1+\delta}{\delta^4} \left(\frac{\alpha(t-k)}{\delta\sqrt{1-\beta_2}} + 2\delta + \delta^4 \right)^{2(1+\mu)} + \left(\frac{\alpha(t-k)}{\delta\sqrt{1-\beta_2}} + 2\delta + 1 \right)^{2(1+\mu)} \\
 &= (1+\delta) \left(\frac{\alpha(t-k)}{\delta^5\sqrt{1-\beta_2}} + \frac{2}{\delta^3} + 1 \right)^{2(1+\mu)} \delta^{4+8\mu} + \left(\frac{\alpha(t-k)}{\delta\sqrt{1-\beta_2}} + 2\delta + 1 \right)^{2(1+\mu)} \\
 &\leq 2\delta \cdot \left(\frac{\alpha(t-k)}{\delta^5\sqrt{1-\beta_2}} + \frac{2}{\delta^3} + 1 \right)^3 \cdot \delta^{4+8\mu} + \left(\frac{\alpha(t-k)}{\delta\sqrt{1-\beta_2}} + 2\delta + 1 \right)^3 \\
 &\leq 2 \left(\frac{\alpha(t-k)}{\sqrt{1-\beta_2}} + 2 \right)^3 \delta^{5+8\mu} + \left(\frac{\alpha(t-k)}{\sqrt{1-\beta_2}} + 3\delta \right)^3 \\
 &\leq \left(\frac{8\alpha^3(t-k)^3}{(1-\beta_2)^{3/2}} + 64 \right) \delta^{5+8\mu} + \frac{4\alpha^3(t-k)^3}{(1-\beta_2)^{3/2}} + 108\delta^3.
 \end{aligned}$$

We have used $\delta > 2$ and $0 < \mu < 1/2$. The last inequality holds because of Lemma 2. Then we obtain

$$\begin{aligned}
 \left(\mathbb{E} \left[g_k^{2(1+\mu)} | \mathcal{F}_k \right] \right)^{1/(1+\mu)} &\leq \left[\left(\frac{8\alpha^3(t-k)^3}{(1-\beta_2)^{3/2}} + 64 \right) \delta^{5+8\mu} + \frac{4\alpha^3(t-k)^3}{(1-\beta_2)^{3/2}} + 108\delta^3 \right]^{1/(1+\mu)} \\
 &= \left(\frac{8\alpha^3(t-k)^3}{(1-\beta_2)^{3/2}} + 64 + \frac{4\alpha^3(t-k)^3}{(1-\beta_2)^{3/2}} \delta^{-5-8\mu} + 108\delta^{-2-8\mu} \right)^{1/(1+\mu)} \delta^{(5+8\mu)/(1+\mu)} \\
 &\leq \left(\frac{12\alpha^3(t-k)^3}{(1-\beta_2)^{3/2}} + 172 \right)^{1/(1+\mu)} \delta^{(5+8\mu)/(1+\mu)} \\
 &\leq \left(\frac{12\alpha^3(t-k)^3}{(1-\beta_2)^{3/2}} + 172 \right) \delta^{(5+8\mu)/(1+\mu)}.
 \end{aligned}$$

The third inequality uses $\delta > 1$ and thus

$$\begin{aligned}
 \mathbb{E}[v_{t-1} \mathbf{1}\{E\}] &\leq (1-\beta_2) \sum_{k=1}^{t-1} \beta_2^{t-1-k} \left(\frac{12\alpha^3(t-k)^3}{(1-\beta_2)^{3/2}} + 172 \right) \delta^{(5+8\mu)/(1+\mu)} \\
 &\leq \delta^{(5+8\mu)/(1+\mu)} \left\{ \frac{12\alpha^3}{\sqrt{1-\beta_2}} \sum_{k=1}^{\infty} \beta_2^{k-1} k^3 + 172 \right\} \\
 &\leq \delta^{(5+8\mu)/(1+\mu)} \left\{ \frac{72\alpha^3}{(1-\beta_2)^{9/2}} + 172 \right\} := M_1 \delta^{(5+8\mu)/(1+\mu)}
 \end{aligned}$$

where the last inequality is because $\sum_{k=1}^{\infty} \beta_2^{k-1} k^3 = (1+4\beta_2+\beta_2^2)/(1-\beta_2)^4 < 6/(1-\beta_2)^4$. Thus, we have

$$\begin{aligned}
 \mathbb{E}[\Delta_t \mathbf{1}\{E\}] &\geq \mathbb{P}\{E\} \left\{ \left(1 - \frac{1+\delta}{1+\delta^4} \right) \frac{\alpha}{\sqrt{\beta_2 M_1 \delta^{(5+8\mu)/(1+\mu)} / \mathbb{P}\{E\} + (1-\beta_2)(1+2\delta)^2}} \right. \\
 &\quad \left. - \frac{1+\delta}{1+\delta^4} \frac{\alpha}{\sqrt{1-\beta_2}} \right\} \\
 &\geq \frac{1}{2} \left\{ \left(1 - \frac{1+\delta}{1+\delta^4} \right) \frac{\alpha}{\sqrt{2\beta_2 M_1 \delta^{(5+8\mu)/(1+\mu)} + (1-\beta_2)(1+2\delta)^2}} \right. \\
 &\quad \left. - \frac{1+\delta}{1+\delta^4} \frac{\alpha}{\sqrt{1-\beta_2}} \right\},
 \end{aligned}$$

where the second inequality follows from

$$\mathbb{P}\{E\} > 1 - \frac{\epsilon}{\delta^4} > 1 - \epsilon > \frac{1}{2}.$$

Then the full expectation of Δ_t is

$$\mathbb{E}[\Delta_t] \geq \frac{1}{2} \left\{ \left(1 - \frac{1+\delta}{1+\delta^4} \right) \frac{\alpha}{\underbrace{\sqrt{2\beta_2 M_1 \delta^{(5+8\mu)/(1+\mu)} + (1-\beta_2)(1+2\delta)^2}}_{T_1}} - \underbrace{\frac{1+\delta}{1+\delta^4} \frac{\alpha}{\sqrt{1-\beta_2}}}_{T_2} \right\} - \underbrace{\frac{1}{2\delta^4} \frac{\alpha}{\sqrt{1-\beta_2}}}_{T_3}.$$

Notice that $T_1 = \Omega(\delta^{-(5+8\mu)/(2+2\mu)})$, $T_2 = \mathcal{O}(\delta^{-3})$ and $T_3 = \mathcal{O}(\delta^{-4})$. As long as we set $\mu < 1/2$, we have $(5+8\mu)/(2+2\mu) < 3 < 4$, thus the right hand side can be positive for sufficiently large δ , only dependent on α and β_2 . We conclude that we can assume $\mathbb{E}[\Delta_t] > c_0 > 0$. This means w_t keeps drifting in the positive direction. Then for any $k \geq 1$ and $t > T_\epsilon$, we have

$$\begin{aligned} \mathbb{E}[(w_{t+k} - w^*)^2] &= \mathbb{E}[(w_{t+k} - w_t)^2] + 2\mathbb{E}[(w_{t+k} - w_t)(w_t - w^*)] + \mathbb{E}[(w_t - w^*)^2] \\ &\geq \mathbb{E}[(w_{t+k} - w_t)^2] - 2\sqrt{\mathbb{E}[(w_{t+k} - w_t)^2] \mathbb{E}[(w_t - w^*)^2]} + \mathbb{E}[(w_t - w^*)^2] \\ &= \left(\sqrt{\mathbb{E}[(w_{t+k} - w_t)^2]} - \sqrt{\mathbb{E}[(w_t - w^*)^2]} \right)^2, \end{aligned} \quad (4)$$

where the inequality is the Cauchy-Schwartz inequality for random variables. If we select k large enough such that $k > 3\sqrt{\epsilon}/c_0$, which implies $kc_0 - \sqrt{\epsilon} > 2\sqrt{\epsilon}$, then $\mathbb{E}[(w_{t+k} - w_t)^2] \geq (\mathbb{E}[w_{t+k} - w_t])^2 \geq k^2 c_0^2$, and thus from (4) we have

$$\mathbb{E}[(w_{t+k} - w^*)^2] \geq (kc_0 - \sqrt{\epsilon})^2 \geq 4\epsilon > \epsilon,$$

which contradicts the convergence assumption. Thus, ADAM diverges for this unconstrained stochastic optimization problem. This completes the proof of Theorem 1.

A.3 Proof of Theorem 2

Consider function $\pi(\delta) = (1+\delta)/(1+\delta^4)$. We notice that $\pi(1) = 1$, $\pi(\delta) \leq 1$ for $\delta \geq 1$, and $\pi(+\infty) = 0$. Since π has only a finite number of stationary points, there exists a $\bar{\delta}$ such that π is decreasing on $[\bar{\delta}, \infty)$. Thus for any b , there exists an N_b^* such that for any $N \geq N_b^*$, $N > b$ there exists a $\delta_{N,b} > \max(\delta^*, \bar{\delta}) > \delta^*$ with

$$\frac{b}{N} = \pi(\delta_{N,b}).$$

Let us consider the following mini-batch problem with sample size $N > N_b^*$.

$$\begin{aligned} f_n(w) &= \frac{w^2}{2\delta_{N,b}} - w \quad \text{for } n = 1, \dots, N-1, \\ f_N(w) &= \frac{w^2}{2\delta_{N,b}} + (b\delta_{N,b}^4 + b - 1)w. \end{aligned}$$

Apparently, the selection of mini-batch \mathcal{B}_t satisfies

$$\mathbb{P}\{N \in \mathcal{B}_t\} = \frac{\binom{N-1}{b-1}}{\binom{N}{b}} = \frac{b}{N} = \frac{1 + \delta_{N,b}}{1 + \delta_{N,b}^4} := p.$$

If $M \in \mathcal{B}_t$, we have

$$F^{\mathcal{B}_t}(w) = \frac{1}{b} (f_N(w) + (b-1)f_1(w)) = \frac{w^2}{2\delta_{N,b}} + \delta_{N,b}^4 w.$$

Otherwise, it is clear that

$$F^{\mathcal{B}_t}(w) = f_1(w) = \frac{w^2}{2\delta_{N,b}} - w.$$

To summarize, the mini-batch loss reads

$$F^{\mathcal{B}_t}(w) = \begin{cases} \frac{w^2}{2\delta_{N,b}^2} + \delta_{N,b}^4 w & \text{with probability } p \\ \frac{w^2}{2\delta_{N,b}^2} - w & \text{with probability } 1 - p \end{cases}$$

which is an OP($\delta_{N,b}$), since $\delta_{N,b} > \delta^*$. By Theorem 1, ADAM diverges on this problem.

A.4 Proof of Theorem 3

Similarly to the proof of Theorem 2, there exists N^* such that for each $N > N^*$, there exists a $\delta_N > \delta^*$ such that

$$\frac{1}{N} = \frac{1 + \delta_N}{1 + \delta_N^4}.$$

We let

$$\begin{aligned} f_n(w) &= \frac{w^2}{2\delta_N} + \delta_N^4 w \text{ for } n = 1, \dots, N-1 \\ f_N(w) &= \frac{w^2}{2\delta_N} - ((N-1) + (N-2)\delta_N^4)w. \end{aligned}$$

The selection of mini-batch \mathcal{B}_t satisfies

$$\mathbb{P}\{N \notin \mathcal{B}_t\} = \frac{1}{N}.$$

If $N \notin \mathcal{B}_t$, we have

$$F^{\mathcal{B}_t}(w) = f_1(w) = \frac{w^2}{2\delta_N} + \delta_N^4 w,$$

and otherwise

$$F^{\mathcal{B}_t}(w) = \frac{N-2}{N-1}f_1(w) + \frac{1}{N-1}f_N(w) = \frac{w^2}{2\delta_N} - w.$$

This is an OP(δ_N), which is divergent according to Theorem 1.

A.5 Proof of Theorem 4

We first introduce the following lemma.

Lemma 5 *Given Assumption 1 is satisfied, there exist positive constants Q_1 and Q_2 such that for any t ,*

$$\mathbb{E}[F(w_{t+1})] - \mathbb{E}[F(w_t)] \leq -\frac{\alpha_t}{4\sqrt{G^2 + \epsilon}} \mathbb{E}[\|\nabla F(w_t)\|_2^2] + Q_1 \alpha_t \lambda_t + Q_2 \alpha_t \sum_{k=1}^{t-1} \beta_1^{t-k} \lambda_k + Q_3 \alpha_t^2.$$

Proof: Let us start from the application of L -smoothness of gradient of $F(w)$ as follows.

$$\begin{aligned}
 F(w_{t+1}) &\leq F(w_t) + \nabla F(w_t)^\top (w_{t+1} - w_t) + \frac{L}{2} \|w_{t+1} - w_t\|_2^2 \\
 &= F(w_t) - \frac{\alpha_t}{1 - \beta_1^t} \nabla F(w_t)^\top V_t^{-1/2} m_t + \frac{\alpha_t^2 L}{2(1 - \beta_1^t)^2} \|V_t^{-1/2} m_t\|_2^2 \\
 &= F(w_t) - \frac{\alpha_t(1 - \beta_1)}{1 - \beta_1^t} \nabla F(w_t)^\top V_t^{-1/2} \sum_{k=1}^t \beta_1^{t-k} g_k + \frac{\alpha_t^2 L}{2(1 - \beta_1^t)^2} \|V_t^{-1/2} m_t\|_2^2 \\
 &= F(w_t) - \alpha_t \nabla F(w_t)^\top V_t^{-1/2} g_t - \frac{\alpha_t(1 - \beta_1)}{1 - \beta_1^t} \sum_{k=1}^{t-1} \beta_1^{t-k} \nabla F(w_t)^\top V_t^{-1/2} (g_k - g_t) \\
 &\quad + \frac{\alpha_t^2 L}{2(1 - \beta_1^t)^2} \|V_t^{-1/2} m_t\|_2^2 \\
 &= F(w_t) - \alpha_t \underbrace{\nabla F(w_t)^\top V_t^{-1/2} \mathcal{G}(w_t; \xi_t)}_{T_1} \\
 &\quad - \frac{\alpha_t(1 - \beta_1)}{1 - \beta_1^t} \underbrace{\sum_{k=1}^{t-1} \beta_1^{t-k} \nabla F(w_t)^\top V_t^{-1/2} (\mathcal{G}(w_k; \xi_k) - \mathcal{G}(w_k; \xi_t))}_{T_2} \\
 &\quad - \frac{\alpha_t(1 - \beta_1)}{1 - \beta_1^t} \underbrace{\sum_{k=1}^{t-1} \beta_1^{t-k} \nabla F(w_t)^\top V_t^{-1/2} (\mathcal{G}(w_k; \xi_t) - \mathcal{G}(w_t; \xi_t))}_{T_3} \\
 &\quad + \frac{\alpha_t^2 L}{2(1 - \beta_1^t)^2} \underbrace{\|V_t^{-1/2} m_t\|_2^2}_{T_4}.
 \end{aligned}$$

Bounding T_1 : We start from

$$\begin{aligned}
 \mathbb{E}[T_1] &= \mathbb{E} \left[\left\| V_t^{-1/4} \nabla F(w_t) \right\|_2^2 \right] + \mathbb{E} \left[\nabla F(w_t)^\top V_t^{-1/2} (\mathcal{G}(w_t; \xi_t) - \nabla F(w_t)) \right] \\
 &\geq \frac{1}{2} \mathbb{E} \left[\left\| V_t^{-1/4} \nabla F(w_t) \right\|_2^2 \right] - \frac{1}{2} \mathbb{E} \left[\left\| V_t^{-1/4} (\mathcal{G}(w_t; \xi_t) - \nabla F(w_t)) \right\|_2^2 \right] \\
 &\geq \frac{1}{2\sqrt{G^2 + \epsilon}} \mathbb{E} \left[\left\| \nabla F(w_t) \right\|_2^2 \right] - \frac{1}{2\sqrt{\epsilon}} \mathbb{E} \left[\left\| \mathcal{G}(w_t; \xi_t) - \nabla F(w_t) \right\|_2^2 \right],
 \end{aligned}$$

where the first inequality applies the Cauchy-Schwartz inequality and the second inequality is due to

$$\left\| V_t^{-1/4} \nabla F(w_t) \right\|_2^2 = \sum_{i=1}^d \frac{(\nabla_i F(w_t))^2}{\sqrt{\tilde{v}_{t,i} + \epsilon}} \geq \frac{1}{\sqrt{G^2 + \epsilon}} \sum_{i=1}^d (\nabla_i F(w_t))^2 = \frac{1}{\sqrt{G^2 + \epsilon}} \left\| \nabla F(w_t) \right\|_2^2$$

and

$$\begin{aligned}
 \left\| V_t^{-1/4} (\mathcal{G}(w_t; \xi_t) - \nabla F(w_t)) \right\|_2^2 &= \sum_{i=1}^d \frac{(\nabla_i F(w_t) - \mathcal{G}_i(w_t; \xi_t))^2}{\sqrt{\tilde{v}_{t,i} + \epsilon}} \\
 &\leq \frac{1}{\sqrt{\epsilon}} \sum_{i=1}^d (\nabla_i F(w_t) - \mathcal{G}_i(w_t; \xi_t))^2 \\
 &= \frac{1}{\sqrt{\epsilon}} \left\| \mathcal{G}(w_t; \xi_t) - \nabla F(w_t) \right\|_2^2.
 \end{aligned}$$

According to the unbiased assumption, we have

$$\mathbb{E} \left[\left\| \mathcal{G}(w_t; \xi_t) - \nabla F(w_t) \right\|_2^2 \right] = \sum_{i=1}^d \mathbb{E} \left[(\mathcal{G}_i(w_t; \xi_t) - \nabla_i F(w_t))^2 \right] = \sum_{i=1}^d \text{Var}(\mathcal{G}_i(w_t; \xi_t)) \leq d\lambda_t. \quad (5)$$

Then we can lower bound the expectation of T_1 as

$$\mathbb{E}[T_1] \geq \frac{1}{2\sqrt{G^2 + \epsilon}} \mathbb{E} \left[\|\nabla F(w_t)\|_2^2 \right] - \frac{d}{2\sqrt{\epsilon}} \lambda_t. \quad (6)$$

Bounding T_2 : Notice that for two random vectors X and Y , and a constant a we have

$$\left\| aX + \frac{1}{a}Y \right\|_2^2 = a^2 \|X\|_2^2 + \frac{1}{a^2} \|Y\|_2^2 + 2Y^\top X,$$

and thus

$$\mathbb{E} [Y^\top X] \geq -\frac{a^2}{2} \mathbb{E}[\|X\|_2^2] - \frac{1}{2a^2} \mathbb{E}[\|Y\|_2^2].$$

If $\mathcal{F}_k = \{\xi_1, \dots, \xi_{k-1}\}$, then w_k is known given \mathcal{F}_k . We then have

$$\begin{aligned} \mathbb{E}[T_2] &= \sum_{k=1}^{t-1} \beta_1^{t-k} \mathbb{E} \left[\nabla F(w_t)^\top V_t^{-1/2} (\mathcal{G}(w_k; \xi_k) - \mathcal{G}(w_k; \xi_t)) \right] \\ &= \sum_{k=1}^{t-1} \beta_1^{t-k} \mathbb{E} \left[\mathbb{E} \left[\nabla F(w_t)^\top V_t^{-1/2} (\mathcal{G}(w_k; \xi_k) - \mathcal{G}(w_k; \xi_t)) \mid \mathcal{F}_k \right] \right] \\ &\geq -\frac{1}{2} \sum_{k=1}^{t-1} \beta_1^{t-k} \mathbb{E} \left[a^2 \mathbb{E} \left[\|V_t^{-1/2} \nabla F(w_t)\|_2^2 \mid \mathcal{F}_k \right] + \frac{1}{a^2} \mathbb{E} \left[\|\mathcal{G}(w_k; \xi_k) - \mathcal{G}(w_k; \xi_t)\|_2^2 \mid \mathcal{F}_k \right] \right] \\ &\geq -\frac{1}{2} \sum_{k=1}^{t-1} \beta_1^{t-k} \left\{ \frac{a^2}{\epsilon} \mathbb{E} \left[\|\nabla F(w_t)\|_2^2 \right] + \frac{1}{a^2} \mathbb{E} \left[\mathbb{E} \left[\|\mathcal{G}(w_k; \xi_k) - \mathcal{G}(w_k; \xi_t)\|_2^2 \mid \mathcal{F}_k \right] \right] \right\} \\ &= -\frac{1}{2} \sum_{k=1}^{t-1} \beta_1^{t-k} \left\{ \frac{a^2}{\epsilon} \mathbb{E} \left[\|\nabla F(w_t)\|_2^2 \right] + \frac{2}{a^2} \mathbb{E} \left[\mathbb{E} \left[\|\nabla \mathcal{G}(w_k; \xi_k) - \nabla F(w_k)\|_2^2 \mid \mathcal{F}_k \right] \right] \right\} \\ &\geq -\frac{a^2}{2\epsilon} \frac{1}{1 - \beta_1} \mathbb{E} \left[\|\nabla F(w_t)\|_2^2 \right] - \frac{1}{a^2} \sum_{k=1}^{t-1} \beta_1^{t-k} \mathbb{E} \left[\|\nabla \mathcal{G}(w_k; \xi_k) - \nabla F(w_k)\|_2^2 \right] \\ &\geq -\frac{a^2}{2\epsilon} \frac{1}{1 - \beta_1} \mathbb{E} \left[\|\nabla F(w_t)\|_2^2 \right] - \frac{d}{a^2} \sum_{k=1}^{t-1} \beta_1^{t-k} \lambda_k \end{aligned}$$

for any positive constant a , where the third equality holds because $\mathcal{G}(w_k; \xi_k)$ and $\mathcal{G}(w_k; \xi_t)$ are i.i.d. given \mathcal{F}_k , and thus

$$\begin{aligned} \mathbb{E} \left[\|\mathcal{G}(w_k; \xi_k) - \mathcal{G}(w_k; \xi_t)\|_2^2 \mid \mathcal{F}_k \right] &= \mathbb{E} \left[\|\mathcal{G}(w_k; \xi_k) - \nabla F(w_k)\|_2^2 \mid \mathcal{F}_k \right] + \mathbb{E} \left[\|\mathcal{G}(w_k; \xi_t) - \nabla F(w_k)\|_2^2 \mid \mathcal{F}_k \right] \\ &\quad - 2\mathbb{E} \left[(\mathcal{G}(w_k; \xi_k) - \nabla F(w_k)) \mid \mathcal{F}_k \right]^\top \mathbb{E} \left[(\mathcal{G}(w_k; \xi_t) - \nabla F(w_k)) \mid \mathcal{F}_k \right] \\ &= 2\mathbb{E} \left[\|\mathcal{G}(w_k; \xi_k) - \nabla F(w_k)\|_2^2 \mid \mathcal{F}_k \right]. \end{aligned}$$

The last inequality applies (5).

If $a = \sqrt{\epsilon(1 - \beta_1)/2\sqrt{G^2 + \epsilon}}$, then we have

$$\mathbb{E}[T_2] \geq -\frac{1}{4\sqrt{G^2 + \epsilon}} \mathbb{E} \left[\|\nabla F(w_t)\|_2^2 \right] - \frac{2d\sqrt{G^2 + \epsilon}}{\epsilon(1 - \beta_1)} \sum_{k=1}^{t-1} \beta_1^{t-k} \lambda_k. \quad (7)$$

Bounding T_3 : We derive

$$\begin{aligned}
 T_3 &= \sum_{k=1}^{t-1} \beta_1^{t-k} \nabla F(w_t) V_t^{-1/2} (\mathcal{G}(w_k; \xi_t) - \mathcal{G}(w_t; \xi_t)) \\
 &\geq - \sum_{k=1}^{t-1} \beta_1^{t-k} \left\| \nabla F(w_t) V_t^{-1/2} \right\|_2 \|\mathcal{G}(w_k; \xi_t) - \mathcal{G}(w_t; \xi_t)\|_2 \\
 &\geq - \frac{LG}{\sqrt{\epsilon}} \sum_{k=1}^{t-1} \beta_1^{t-k} \|w_t - w_k\| \\
 &= - \frac{LG}{\sqrt{\epsilon}} \sum_{k=1}^{t-1} \beta_1^{t-k} \left\| \sum_{j=k}^{t-1} \alpha_j V_j^{-1/2} \tilde{m}_j \right\|_2 \\
 &\geq - \frac{LG}{\sqrt{\epsilon}} \sum_{k=1}^{t-1} \beta_1^{t-k} \sum_{j=k}^{t-1} \alpha_j \left\| V_j^{-1/2} \tilde{m}_j \right\|_2 \\
 &\geq - \frac{LG^2}{\epsilon \sqrt{1-\beta_1}} \sum_{k=1}^{t-1} \sum_{j=k}^{t-1} \beta_1^{t-k} \alpha_j \\
 &= - \frac{LG^2}{\epsilon \sqrt{1-\beta_1}} \sum_{j=1}^{t-1} \alpha_j \sum_{k=1}^j \beta_1^{t-k} \\
 &\geq - \frac{LG^2}{\epsilon (1-\beta_1)^{3/2}} \sum_{j=1}^{t-1} \alpha_j \beta_1^{t-j} \\
 &\geq - \frac{LG^2 \bar{C}}{\epsilon (1-\beta_1)^{3/2}} \alpha_t, \tag{8}
 \end{aligned}$$

where the first inequality is the Cauchy-Schwartz inequality, the second inequality applies L smoothness of $G(\cdot; \xi)$ for any ξ , the forth inequality holds because

$$\left\| V_j^{-1/2} \tilde{m}_j \right\|_2 = \sqrt{\frac{1}{1-\beta_1^j} \sum_{i=1}^d \frac{m_{j,i}^2}{v_{j,i} + \epsilon}} \leq \frac{G}{\sqrt{\epsilon(1-\beta_1)}}$$

and the last inequality comes from Lemma 3.

Bounding T_4 : It is easy to show that

$$T_4 = \sum_{i=1}^d \frac{m_{t,i}^2}{v_{t,i} + \epsilon} \leq \frac{G^2}{\epsilon}. \tag{9}$$

According to the bounds in (6), (7), (8) and (9), we get

$$\begin{aligned}
 \mathbb{E}[F(w_{t+1})] - \mathbb{E}[F(w_t)] &\leq -\alpha_t \left\{ \frac{1}{2\sqrt{G^2 + \epsilon}} \mathbb{E} \left[\|\nabla F(w_t)\|_2^2 \right] - \frac{d}{2\sqrt{\epsilon}} \lambda_t \right\} \\
 &\quad - \frac{\alpha_t (1-\beta_1)}{1-\beta_1^t} \left\{ -\frac{1}{4\sqrt{G^2 + \epsilon}} \mathbb{E} \left[\|\nabla F(w_t)\|_2^2 \right] - \frac{2d\sqrt{G^2 + \epsilon}}{\epsilon(1-\beta_1)} \sum_{k=1}^{t-1} \beta_1^{t-k} \lambda_k \right\} \\
 &\quad + \frac{LG^2 \bar{C}}{\epsilon \sqrt{1-\beta_1} (1-\beta_1^t)} \alpha_t^2 + \frac{LG^2}{2\epsilon(1-\beta_1^t)^2} \alpha_t^2 \\
 &\leq -\alpha_t \frac{1}{4\sqrt{G^2 + \epsilon}} \mathbb{E} \left[\|\nabla F(w_t)\|_2^2 \right] + \frac{d}{2\sqrt{\epsilon}} \alpha_t \lambda_t + \frac{2d\sqrt{G^2 + \epsilon}}{\epsilon(1-\beta_1)} \alpha_t \sum_{k=1}^{t-1} \beta_1^{t-k} \lambda_k \\
 &\quad + \left\{ \frac{LG^2 \bar{C}}{\epsilon(1-\beta_1)^{3/2}} + \frac{LG^2}{2\epsilon(1-\beta_1)^2} \right\} \alpha_t^2.
 \end{aligned}$$

Letting

$$\begin{aligned} Q_1 &= \frac{d}{2\sqrt{\epsilon}} \\ Q_2 &= \frac{2d\sqrt{G^2 + \epsilon}}{\epsilon(1 - \beta_1)} \\ Q_3 &= \frac{LG^2\bar{C}}{\epsilon(1 - \beta_1)^{3/2}} + \frac{LG^2}{2\epsilon(1 - \beta_1)^2} \end{aligned}$$

completes the proof.

Proof of Theorem 4: According to Lemma 5, we have

$$\begin{aligned} F_{\text{inf}} - F(w_1) &\leq \mathbb{E}[F(w_{T+1})] - F(w_1) \\ &= \sum_{t=1}^T \mathbb{E}[F(w_{t+1})] - \mathbb{E}[F(w_t)] \\ &\leq -\frac{1}{4\sqrt{G^2 + \epsilon}} \sum_{i=1}^T \alpha_t \mathbb{E} [\|\nabla F(w_t)\|_2^2] + Q_1 \sum_{i=1}^T \alpha_t \lambda_t + Q_2 \sum_{i=1}^T \alpha_t \sum_{k=1}^{t-1} \beta_1^{t-k} \lambda_k \\ &\quad + Q_3 \sum_{i=1}^T \alpha_t^2. \end{aligned}$$

Then we obtain

$$\begin{aligned} \frac{1}{4\sqrt{G^2 + \epsilon}} \sum_{t=1}^T \alpha_t \mathbb{E} [\|\nabla F(w_t)\|_2^2] &\leq F(w_1) - F_{\text{inf}} + Q_1 \sum_{t=1}^T \alpha_t \lambda_t \\ &\quad + Q_2 \sum_{t=1}^T \alpha_t \sum_{k=1}^{t-1} \beta_1^{t-k} \lambda_k + Q_3 \sum_{t=1}^T \alpha_t^2 \\ &\leq F(w_1) - F_{\text{inf}} + Q_1 \sum_{t=1}^T \alpha_t \lambda_t \\ &\quad + Q_2 \sum_{k=1}^T \lambda_k \sum_{t=k}^T \beta_1^{t-k} \alpha_t + Q_3 \sum_{t=1}^T \alpha_t^2 \\ &\leq F(w_1) - F_{\text{inf}} + Q_1 \sum_{t=1}^T \alpha_t \lambda_t \\ &\quad + \frac{Q_2}{1 - \beta_1} \sum_{k=1}^T \lambda_k \alpha_k + Q_3 \sum_{t=1}^T \alpha_t^2. \end{aligned}$$

Noticing that the left-hand side can be bounded as

$$\sum_{t=1}^T \alpha_t \mathbb{E} [\|\nabla F(w_t)\|_2^2] \geq \sum_{t=1}^T \alpha_t \min_{1 \leq t \leq T} \mathbb{E} [\|\nabla F(w_t)\|_2^2],$$

we obtain

$$\begin{aligned} \min_{1 \leq t \leq T} \mathbb{E} [\|\nabla F(w_t)\|_2^2] &\leq \frac{4\sqrt{G^2 + \epsilon}}{\sum_{t=1}^T \alpha_t} + 4\sqrt{G^2 + \epsilon} \left(Q_1 + \frac{Q_2}{1 - \beta_2} \right) \frac{\sum_{t=1}^T \lambda_t \alpha_t}{\sum_{t=1}^T \alpha_t} \\ &\quad + 4\sqrt{G^2 + \epsilon} Q_3 \frac{\sum_{t=1}^T \alpha_t^2}{\sum_{t=1}^T \alpha_t}. \end{aligned}$$

A.6 Proof of Theorem 5

Applying L-smoothness of gradients of F and strong convexity of F , we have

$$\begin{aligned} F(\hat{w}_2^{(2)}) &\geq F(w_1^{(2)}) + F'(w_1^{(2)})(\hat{w}_2^{(2)} - w_1^{(2)}) + \frac{c}{2}(\hat{w}_2^{(2)} - w_1^{(2)})^2 \\ F(w_2^{(2)}) &\leq F(w_1^{(2)}) + F'(w_1^{(2)})(w_2^{(2)} - w_1^{(2)}) + \frac{L}{2}(w_2^{(2)} - w_1^{(2)})^2. \end{aligned}$$

By definition, we have

$$\begin{aligned} F(\hat{w}_2^{(2)}) - F(w_2^{(2)}) &\geq F'(w_1^{(2)})(\hat{w}_2^{(2)} - w_2^{(2)}) + \frac{c}{2}(\hat{w}_2^{(2)} - w_1^{(2)})^2 - \frac{L}{2}(w_2^{(2)} - w_1^{(2)})^2 \\ &= F'(w_1^{(2)}) \left(\alpha_2 \frac{\tilde{m}_1^{(2)}}{\sqrt{\tilde{v}_1^{(2)} + \epsilon}} - \alpha_2 \frac{\hat{m}_1^{(2)}}{\sqrt{\hat{v}_1^{(2)} + \epsilon}} \right) + \frac{c\alpha_2^2}{2} \frac{\tilde{m}_1^{(2)}}{\tilde{v}_1^{(2)} + \epsilon} - \frac{L\alpha_2^2}{2} \frac{\hat{m}_1^{(2)}}{\hat{v}_1^{(2)} + \epsilon} \\ &= \alpha_2 F'(w_1^{(2)}) \left(\frac{F'(w_1^{(2)})}{\sqrt{Q_3}} - \gamma \frac{(1 - \beta_1)F'(w_1^{(2)}) + \beta_1 m_{m+1}^{(1)}}{\sqrt{Q_4}} \right) \\ &\quad + \frac{c\alpha_2^2 \gamma^2}{2} \frac{((1 - \beta_1)F'(w_1^{(2)}) + \beta_1 m_{m+1}^{(1)})^2}{Q_4} - \frac{L\alpha_2^2}{2} \frac{(F'(w_1^{(2)}))^2}{Q_3} \end{aligned}$$

where

$$\begin{aligned} Q_3 &= \tilde{v}_1^{(2)} + \epsilon \\ Q_4 &= \hat{v}_1^{(2)} + \epsilon \\ \gamma &= \frac{1}{1 - \beta_1^{m+1}}. \end{aligned}$$

Thus we have

$$F(\hat{w}_2^{(2)}) - F(w_2^{(2)}) \geq (F'(w_1^{(2)}))^2 q \left(\frac{m_{m+1}^{(1)}}{F'(w_1^{(2)})} \right)$$

where $q(x) = Q_5 x^2 + Q_6 x + Q_7$ is a function with parameters

$$\begin{aligned} Q_5 &= \frac{c\alpha_2^2 \gamma^2 \beta_1^2}{2Q_4} \\ Q_6 &= \frac{c\alpha_2^2 \gamma^2 \beta_1 (1 - \beta_1)}{Q_4} - \frac{\alpha_2 \gamma \beta_1}{\sqrt{Q_2}} \\ Q_7 &= \frac{\alpha_2}{\sqrt{Q_3}} - \frac{\alpha_2 \gamma (1 - \beta_1)}{\sqrt{Q_4}} + \frac{c\alpha_2^2 \gamma^2 (1 - \beta_1)^2}{2Q_4} - \frac{L\alpha_2^2}{2Q_3}. \end{aligned}$$

Apparently, from

$$\begin{aligned} \tilde{v}_1^{(2)} &= (g_1^{(2)})^2 \leq G^2 \\ \hat{v}_1^{(2)} &= \frac{1 - \beta_1}{1 - \beta_2^{m+1}} \left(\sum_{k=1}^m \beta_2^{m+1-k} (g_k^{(1)})^2 + (g_1^{(2)})^2 \right) \leq \frac{1 - \beta_1}{1 - \beta_2^{m+1}} \left(\sum_{k=1}^m \beta_2^{m+1-k} + 1 \right) G^2 = G^2, \end{aligned}$$

we have

$$\begin{aligned} \epsilon &\leq Q_3 \leq \epsilon + G^2 \\ \epsilon &\leq Q_4 \leq \epsilon + G^2. \end{aligned}$$

Noticing that

$$\begin{aligned}
 \Delta &= Q_6^2 - 4Q_5Q_7 \\
 &= \frac{\alpha_2^2\gamma^2\beta_1^2}{Q_4} \left(1 - \frac{2c\alpha_2}{\sqrt{Q_3}} + \frac{cL\alpha_2^2}{Q_3} \right) \\
 &> \frac{\alpha_2^2\gamma^2\beta_1^2}{Q_4} \left(1 - \frac{2c\alpha_2}{\sqrt{Q_3}} + \frac{c^2\alpha_2^2}{Q_3} \right) \\
 &= \frac{\alpha_2^2\gamma^2\beta_1^2}{Q_4} \left(1 - \frac{c\alpha_2}{\sqrt{Q_3}} \right)^2 \geq 0,
 \end{aligned}$$

where the first inequality uses the property of the strong convexity parameter and the L -smoothness gradient parameter $c < L$, we have that there exists

$$\begin{aligned}
 x_1 &= \frac{-Q_6 + \sqrt{\Delta}}{2Q_5} \\
 x_2 &= \frac{-Q_6 - \sqrt{\Delta}}{2Q_5}
 \end{aligned}$$

such that $q(x_1) = q(x_2) = 0$. We claim that $|x_1| \leq 1$ and $|x_2| \leq 1$, which is implied by

$$\sqrt{\Delta} \leq \min\{2Q_5 + Q_6, 2Q_5 - Q_6\}. \quad (10)$$

We notice that

$$\begin{aligned}
 2Q_5 + Q_6 &= \frac{\alpha_2\gamma\beta_1}{\sqrt{Q_4}} \left(\frac{c\alpha_2\gamma}{\sqrt{Q_4}} - 1 \right) \\
 2Q_5 - Q_6 &= \frac{\alpha_2\gamma\beta_1}{\sqrt{Q_4}} \left(\frac{c\alpha_2\gamma(2\beta_1 - 1)}{\sqrt{Q_4}} + 1 \right)
 \end{aligned}$$

and

$$\begin{aligned}
 \Delta &\leq \frac{\alpha_2^2\gamma^2\beta_1^2}{Q_4} \left(1 - \frac{2L\alpha_2}{\sqrt{Q_3}} + \frac{L^2\alpha_2^2}{Q_3} \right) \\
 &= \frac{\alpha_2^2\gamma^2\beta_1^2}{Q_4} \left(1 - \frac{L\alpha_2}{\sqrt{Q_3}} \right)^2
 \end{aligned}$$

where the inequality holds according to Assumption 2 $L\alpha_2 \geq 2\sqrt{G^2 + \epsilon} \geq 2\sqrt{Q_3}$. Thus we have

$$\sqrt{\Delta} \leq \frac{\alpha_2\gamma\beta_1}{\sqrt{Q_4}} \left(\frac{L\alpha_2}{\sqrt{Q_3}} - 1 \right) \leq \frac{\alpha_2\gamma\beta_1}{\sqrt{Q_4}} \left(\frac{c\alpha_2\gamma(2\beta_1 - 1)}{\sqrt{Q_4}} - 1 \right) \leq \min\{2Q_5 + Q_6, 2Q_5 - Q_6\},$$

where the second inequality holds according to Assumption 2, $\frac{L}{c} \leq \frac{2\beta_1 - 1}{1 - \beta_1^{m+1}} \sqrt{\frac{\epsilon}{G + \epsilon}} \leq \frac{2\beta_1 - 1}{1 - \beta_1^{m+1}} \sqrt{\frac{Q_3}{Q_4}}$.

Hence we obtain (10), which implies that $q(x) \geq 0$ where $|x| \geq 1$. As we assume

$$\left| m_{m+1}^{(1)} \right| \geq \left| F' \left(w_1^{(2)} \right) \right|,$$

we have

$$F \left(\hat{w}_2^{(2)} \right) - F \left(w_2^{(2)} \right) \geq 0,$$

which finishes the proof.

A.7 Proof of Theorem 6 and Theorem 7

We start proving the following lemmas.

Lemma 6 *Given Assumption 3, we have that for any $1 \leq k \leq m$ and $1 \leq t \leq T$, the ADAM states in Algorithm 2 with option A satisfy*

$$\begin{aligned}\|m_k^{(t)}\|_2 &\leq 3G, \\ \|v_k^{(t)}\|_2 &\leq 9G^2.\end{aligned}$$

Proof: By definition, we have

$$\begin{aligned}m_k^{(t)} &= (1 - \beta_1) \sum_{j=1}^k \beta_1^{k-j} g_j^{(t)} \\ v_k^{(t)} &= (1 - \beta_2) \sum_{j=1}^k \beta_2^{k-j} g_j^{(t)} \odot g_j^{(t)}.\end{aligned}$$

Applying the Cauchy-Schwartz inequality, we obtain

$$\begin{aligned}\|m_k^{(t)}\|_2 &\leq (1 - \beta_1) \sum_{j=1}^k \beta_1^{k-j} \|g_j^{(t)}\|_2 \\ &\leq (1 - \beta_1) \sum_{j=1}^k \beta_1^{k-j} \left(\|\nabla F^{\mathcal{B}_j^{(t)}}(w_j^{(t)})\|_2 + \|\nabla F^{\mathcal{B}_j^{(t)}}(\tilde{w}_t)\|_2 + \|\nabla F(\tilde{w}_t)\|_2 \right) \\ &\leq (1 - \beta_1) \sum_{j=1}^k \beta_1^{k-j} 3G \leq 3G\end{aligned}$$

and

$$\begin{aligned}\|v_k^{(t)}\|_2 &\leq (1 - \beta_2) \sum_{j=1}^k \beta_2^{k-j} \|g_j^{(t)} \odot g_j^{(t)}\|_2 \\ &= (1 - \beta_2) \sum_{j=1}^k \beta_2^{k-j} \|g_j^{(t)}\|_2^2 \\ &\leq (1 - \beta_2) \sum_{j=1}^k \beta_2^{k-j} \left(\|\nabla F^{\mathcal{B}_j^{(t)}}(w_j^{(t)})\|_2 + \|\nabla F^{\mathcal{B}_j^{(t)}}(\tilde{w}_t)\|_2 + \|\nabla F(\tilde{w}_t)\|_2 \right)^2 \\ &\leq (1 - \beta_2) \sum_{j=1}^k \beta_2^{k-j} 9G^2 \leq 9G^2.\end{aligned}$$

Lemma 7 *Given Assumption 3, there exist positive constants Q_8 and Q_9 such that Algorithm 2 with option A satisfies that for any t ,*

$$F(\tilde{w}_{t+1}) - F(\tilde{w}_t) \leq -Q_8 \alpha_t m \|\nabla F(\tilde{w}_t)\|_2^2 + Q_9 \alpha_t^2$$

holds almost surely.

Proof: We start from the application of L -smoothness of gradient of $F(w)$ as follows

$$\begin{aligned}F(\tilde{w}_{t+1}) &\leq F(\tilde{w}_t) + \nabla F(\tilde{w}_t)^\top (\tilde{w}_{t+1} - \tilde{w}_t) + \frac{L}{2} \|\tilde{w}_{t+1} - \tilde{w}_t\|_2^2 \\ &= F(\tilde{w}_t) + \nabla F(\tilde{w}_t)^\top \sum_{k=1}^m (w_{k+1}^{(t)} - w_k^{(t)}) + \frac{L}{2} \left\| \sum_{k=1}^m (w_{k+1}^{(t)} - w_k^{(t)}) \right\|_2^2 \\ &= F(\tilde{w}_t) - \alpha_t \nabla F(\tilde{w}_t)^\top \sum_{k=1}^m (V_k^{(t)})^{-1/2} \tilde{m}_k^{(t)} + \frac{\alpha_t^2 L}{2} \left\| \sum_{k=1}^m (V_k^{(t)})^{-1/2} \tilde{m}_k^{(t)} \right\|_2^2.\end{aligned}$$

By definition and the resetting option, we have

$$\tilde{m}_k^{(t)} = \frac{1 - \beta_1}{1 - \beta_1^k} \left(V_k^{(t)} \right)^{-1/2} \sum_{j=1}^k \beta_1^{k-j} g_j^{(t)},$$

and thus

$$\begin{aligned} F(\tilde{w}_{t+1}) &\leq F(\tilde{w}_t) - \alpha_t(1 - \beta_1) \nabla F(\tilde{w}_t)^\top \sum_{k=1}^m \frac{1}{1 - \beta_1^k} \left(V_k^{(t)} \right)^{-1/2} \sum_{j=1}^k \beta_1^{k-j} g_j^{(t)} \\ &\quad + \frac{\alpha_t^2 L}{2} \left\| \sum_{k=1}^m \left(V_k^{(t)} \right)^{-1/2} \tilde{m}_k^{(t)} \right\|_2^2 \\ &= F(\tilde{w}_t) - \alpha_t(1 - \beta_1) \nabla F(\tilde{w}_t)^\top \sum_{j=1}^m \left(\sum_{k=j}^m \frac{\beta_1^{k-j}}{1 - \beta_1^k} \left(V_k^{(t)} \right)^{-1/2} \right) g_j^{(t)} \\ &\quad + \frac{\alpha_t^2 L}{2} \left\| \sum_{k=1}^m \left(V_k^{(t)} \right)^{-1/2} \tilde{m}_k^{(t)} \right\|_2^2 \\ &= F(\tilde{w}_t) - \alpha_t(1 - \beta_1) \\ &\quad \underbrace{\times \nabla F(\tilde{w}_t)^\top \sum_{j=1}^m \left(\sum_{k=j}^m \frac{\beta_1^{k-j}}{1 - \beta_1^k} \left(V_k^{(t)} \right)^{-1/2} \right) \left(\nabla F^{\mathcal{B}_j^{(t)}}(w_j^{(t)}) - \nabla F^{\mathcal{B}_j^{(t)}}(\tilde{w}_t) \right)}_{T_1} \\ &\quad - \alpha_t(1 - \beta_1) \nabla F(\tilde{w}_t)^\top \underbrace{\left(\sum_{j=1}^m \sum_{k=j}^m \frac{\beta_1^{k-j}}{1 - \beta_1^k} \left(V_k^{(t)} \right)^{-1/2} \right)}_{T_2} \nabla F(\tilde{w}_t) \\ &\quad + \frac{\alpha_t^2 L}{2} \underbrace{\left\| \sum_{k=1}^m \left(V_k^{(t)} \right)^{-1/2} \tilde{m}_k^{(t)} \right\|_2^2}_{T_3}. \end{aligned}$$

Bounding T_1 : We have

$$T_1 \geq - \sum_{j=1}^m \left\| \sum_{k=j}^m \frac{\beta_1^{k-j}}{1 - \beta_1^k} \left(V_k^{(t)} \right)^{-1/2} \nabla F(\tilde{w}_t) \right\|_2 \left\| \nabla F^{\mathcal{B}_j^{(t)}}(w_j^{(t)}) - \nabla F^{\mathcal{B}_j^{(t)}}(\tilde{w}_t) \right\|_2. \quad (11)$$

The first thing to notice is that

$$\begin{aligned} \left\| \sum_{k=j}^m \frac{\beta_1^{k-j}}{1 - \beta_1^k} \left(V_k^{(t)} \right)^{-1/2} \nabla F(\tilde{w}_t) \right\|_2 &\leq \sum_{k=j}^m \frac{\beta_1^{k-j}}{1 - \beta_1^k} \left\| \left(V_k^{(t)} \right)^{-1/2} \nabla F(\tilde{w}_t) \right\|_2 \\ &= \sum_{k=j}^m \frac{\beta_1^{k-j}}{1 - \beta_1^k} \sqrt{\sum_{i=1}^d \frac{\nabla_i F(\tilde{w}_t)^2}{v_{k,i}^{(t)} + \epsilon}} \\ &\leq \sum_{k=j}^m \frac{\beta_1^{k-j}}{1 - \beta_1^k} \frac{\|\nabla F(\tilde{w}_t)\|_2}{\sqrt{\epsilon}} \\ &\leq \frac{G}{\sqrt{\epsilon}} \sum_{k=j}^m \frac{\beta_1^{k-j}}{1 - \beta_1^k} \\ &\leq \frac{G}{\sqrt{\epsilon}} \frac{1}{1 - \beta_1} \sum_{k=j}^m \beta_1^{k-j} \leq \frac{G}{(1 - \beta_1)^2 \sqrt{\epsilon}}, \end{aligned} \quad (12)$$

where the second inequality employs the assumption that the gradients of F are bounded. Secondly, according to L -smoothness of gradients of every loss function, we derive

$$\begin{aligned}
 \left\| \nabla F^{\mathcal{B}_j^{(t)}} \left(w_j^{(t)} \right) - \nabla F^{\mathcal{B}_j^{(t)}} \left(\tilde{w}_t \right) \right\|_2 &\leq L \left\| w_j^{(t)} - w_1^{(t)} \right\|_2 \\
 &= L \left\| \sum_{l=1}^{j-1} \alpha_t \left(V_l^{(t)} \right)^{-1/2} \tilde{m}_l^{(t)} \right\|_2 \\
 &\leq \alpha_t L \sum_{l=1}^{j-1} \left\| \left(V_l^{(t)} \right)^{-1/2} \tilde{m}_l^{(t)} \right\|_2 \\
 &= \alpha_t L \sum_{l=1}^{j-1} \frac{1}{1 - \beta_1^l} \sqrt{\sum_{i=1}^d \frac{\left(m_{l,i}^{(t)} \right)^2}{v_{l,i}^{(t)} + \epsilon}} \\
 &\leq \frac{\alpha_t L}{1 - \beta_1} \sum_{l=1}^{j-1} \frac{\left\| m_l^{(t)} \right\|_2}{\sqrt{\epsilon}} \leq \frac{3GL}{(1 - \beta_1)\sqrt{\epsilon}} (j - 1)\alpha_t, \tag{13}
 \end{aligned}$$

where the first inequality applies the Cauchy-Schwartz inequality and the last one applies Lemma 6. By plugging equations (12) and (13) into equation (11), we obtain

$$T_1 \geq - \sum_{j=1}^m \frac{3G^2L}{(1 - \beta_1)^3\epsilon} (j - 1)\alpha_t = - \frac{3G^2L}{2(1 - \beta_1)^3\epsilon} m(m - 1)\alpha_t. \tag{14}$$

Bounding T_2 : We have

$$\begin{aligned}
 T_2 &= \sum_{j=1}^m \sum_{k=j}^m \frac{\beta_1^{k-j}}{1 - \beta_1^k} \nabla F(\tilde{w}_t)^\top \left(V_k^{(t)} \right)^{-1/2} \nabla F(\tilde{w}_t) \\
 &= \sum_{j=1}^m \sum_{k=j}^m \frac{\beta_1^{k-j}}{1 - \beta_1^k} \sum_{i=1}^d \frac{\nabla_i F(\tilde{w}_t)^2}{\sqrt{v_{k,i}^{(t)} + \epsilon}} \\
 &\geq \sum_{j=1}^m \sum_{k=j}^m \beta_1^{k-j} \sum_{i=1}^d \frac{\nabla_i F(\tilde{w}_t)^2}{\sqrt{9G^2 + \epsilon}} \\
 &= \frac{\|\nabla F(\tilde{w}_t)\|_2^2}{\sqrt{9G^2 + \epsilon}} \sum_{j=1}^m \sum_{k=j}^m \beta_1^{k-j} \\
 &\geq \frac{\|\nabla F(\tilde{w}_t)\|_2^2}{\sqrt{9G^2 + \epsilon}} \sum_{j=1}^m 1 \\
 &= \frac{1}{\sqrt{9G^2 + \epsilon}} m \|\nabla F(\tilde{w}_t)\|_2^2. \tag{15}
 \end{aligned}$$

Bounding T_3 : We obtain

$$\begin{aligned}
 T_3 &\leq \left(\sum_{k=1}^m \left\| \left(V_k^{(t)} \right)^{-1/2} \tilde{m}_k^{(t)} \right\|_2 \right)^2 \\
 &= \left(\sum_{k=1}^m \frac{1}{1 - \beta_1^k} \sqrt{\sum_{i=1}^d \frac{\left(m_{k,i}^{(t)} \right)^2}{v_{k,i}^{(t)} + \epsilon}} \right)^2 \\
 &\leq \left(\sum_{k=1}^m \frac{1}{\sqrt{\epsilon}} \left\| m_k^{(t)} \right\|_2 \right)^2 \leq \left(\sum_{k=1}^m \frac{3G}{\sqrt{\epsilon}} \right)^2 \leq \frac{9G^2 m^2}{\epsilon}. \tag{16}
 \end{aligned}$$

In summary, we get

$$\begin{aligned} F(\tilde{w}_{t+1}) &\leq F(\tilde{w}_t) - \alpha_t \frac{m(1-\beta_1)}{\sqrt{9G^2 + \epsilon}} \|\nabla F(\tilde{w}_t)\|_2^2 + \frac{3G^2 Lm(m-1)/(1-\beta_1)^2 + 9G^2 Lm^2}{2\epsilon} \alpha_t^2 \\ &= F(\tilde{w}_t) - Q_8 \alpha_t m \|\nabla F(\tilde{w}_t)\|_2^2 + Q_9 \alpha_t^2, \end{aligned}$$

where

$$\begin{aligned} Q_8 &= \frac{1-\beta_1}{\sqrt{9G^2 + \epsilon}} \\ Q_9 &= \frac{3G^2 Lm(m-1)/(1-\beta_1)^2 + 9G^2 Lm^2}{2\epsilon}. \end{aligned}$$

Proof of Theorem 6: If $F(w)$ is c -strongly convex, we have

$$\|\nabla F(w)\|_2^2 \geq 2c(F(w) - F^*),$$

and thus according to Lemma 7, we have

$$F(\tilde{w}_{t+1}) \leq F(\tilde{w}_t) - 2cQ_8\alpha_t m (F(\tilde{w}_t) - F^*) + Q_9\alpha_t^2,$$

which is equivalent to

$$F(\tilde{w}_{t+1}) - F^* \leq \left(1 - \frac{C_2 m \alpha}{t}\right) (F(\tilde{w}_t) - F^*) + Q_9 \alpha_t^2.$$

We obtain recursively

$$F(\tilde{w}_T) - F^* \leq \prod_{t=1}^{T-1} \left(1 - \frac{C_2 m \alpha}{t}\right) (F(\tilde{w}_1) - F^*) + \sum_{t=1}^{T-1} \alpha_t \prod_{j=t+1}^{T-1} \left(1 - \frac{C_2 m \alpha}{j}\right).$$

By definition, we have $C_2 m \alpha < 1$, and thus we can use Lemma 4 to obtain

$$F(\tilde{w}_T) - F^* \leq \mathcal{O}(T^{-C_2 m \alpha}).$$

Proof of Theorem 7: Let us consider the set of indices $A = \{t \in \mathbb{N} : \|\nabla F(\tilde{w}_t)\| = 0\}$. If the set is infinite, there exists a sequence $\{t_k\}_{k=1}^{+\infty}$ such that $\|\nabla F(\tilde{w}_{t_k})\| = 0$ for all k . Then we have

$$\liminf_{t \rightarrow \infty} \|\nabla F(\tilde{w}_t)\|_2 = 0.$$

Otherwise, A is finite, and thus its maximum exists. For all $t > \tau := \max A$, we have $\|\nabla F(\tilde{w}_t)\|_2 > 0$. Applying Lemma 7, we have

$$F(\tilde{w}_{t+1}) - F(\tilde{w}_t) \leq -\alpha_t Q_8 m \|\nabla F(\tilde{w}_t)\|_2^2 + Q_9 \alpha_t^2.$$

Then it follows

$$\begin{aligned} F_{\inf} - F(\tilde{w}_{\tau+1}) &\leq F(\tilde{w}_{T+1}) - F(\tilde{w}_{\tau+1}) \\ &\leq \sum_{t=\tau+1}^T -Q_8 m \alpha_t \|\nabla F(\tilde{w}_t)\|_2^2 + Q_9 \sum_{t=\tau+1}^T \alpha_t^2, \end{aligned}$$

and thus

$$\min_{\tau+1 \leq t \leq T} \|\nabla F(\tilde{w}_t)\|_2^2 \sum_{t=\tau+1}^T \alpha_t \leq \sum_{t=\tau+1}^T \alpha_t \|\nabla F(\tilde{w}_t)\|_2^2 \leq \frac{F(\tilde{w}_{\tau+1}) - F_{\inf}}{Q_8 m} + \frac{Q_9}{Q_8 m} \sum_{t=\tau+1}^T \alpha_t^2.$$

Then, we have

$$\min_{\tau+1 \leq t \leq T} \|\nabla F(\tilde{w}_t)\|_2^2 \leq \frac{1}{\sum_{t=\tau+1}^T \alpha_t} \left\{ \frac{F(\tilde{w}_{\tau+1}) - F_{\inf}}{Q_8 m} + \frac{Q_9}{Q_8 m} \sum_{t=\tau+1}^T \alpha_t^2 \right\},$$

which yields

$$\lim_{T \rightarrow +\infty} \min_{\tau+1 \leq t \leq T} \|\nabla F(\tilde{w}_t)\|_2^2 = 0. \quad (17)$$

For all $r > \tau$, there must exist an $s > r$ such that

$$\|\nabla F(\tilde{w}_s)\|_2 < \|\nabla F(\tilde{w}_r)\|_2.$$

Otherwise, if there exists an $r_0 > \tau$, such that for all $s > r_0$ we have

$$\|\nabla F(\tilde{w}_s)\|_2 \geq \|\nabla F(\tilde{w}_{r_0})\|_2,$$

then for all $T \geq r_0$, we have

$$\min_{\tau+1 \leq t \leq T} \|\nabla F(\tilde{w}_t)\|_2^2 = \min_{\tau+1 \leq t \leq r_0} \|\nabla F(\tilde{w}_t)\|_2^2 = A > 0,$$

which contradicts (17), since A is a positive constant.

Let $t_1 = \tau + 1$ and for all $k \in \mathbb{N}$, let $t_{k+1} = \inf \{s > t_k : \|\nabla F(\tilde{w}_s)\|_2 < \|\nabla F(\tilde{w}_{t_k})\|_2\}$. This implies a sub-sequence $\{\|\nabla F(\tilde{w}_{t_k})\|_2\}_k$ of sequence $\{\|\nabla F(\tilde{w}_t)\|_2\}_t$. Since

$$\|\nabla F(\tilde{w}_{t_k})\|_2 = \min_{\tau+1 \leq t \leq t_k} \|\nabla F(\tilde{w}_t)\|_2,$$

by employing (17), we have

$$\lim_{k \rightarrow \infty} \|\nabla F(\tilde{w}_{t_k})\|_2 = 0,$$

which implies that

$$\liminf_{t \rightarrow \infty} \|\nabla F(\tilde{w}_t)\|_2 = 0.$$

This completes the proof.

B Experiments

B.1 Network Structure

Dataset	Input dimension	Hidden dimension	Output Dimension
CovType	98	100	7
MNIST	784	100	10

Table 1: Feedforward network structure

The feedforward networks used in the experiments have two fully connected layers with the dimensions described in Table 1. The structure of the CNN used in the experiments is described as follows. The CNN is mainly composed of two convolution layers, two max pooling layers and one fully connected layer. The kernel size of the convolution layers is 4 and the kernel size of the pooling layers is 2. The numbers of channels of the two convolution layers are 16 and 32, respectively, and the dimensions of the fully connected layer are 32 for input and 10 for output.

B.2 Additional Results

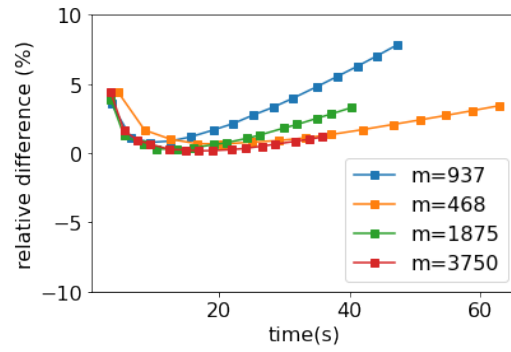


Figure 3: Relative difference of VRADAM in classifying Embedded CIFAR10 with Logistic regression

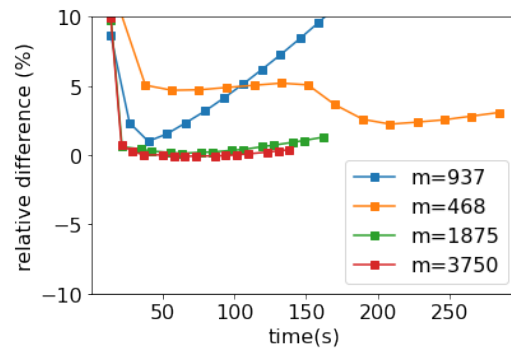


Figure 4: Relative difference of VRADAM in classifying MNIST with CNN

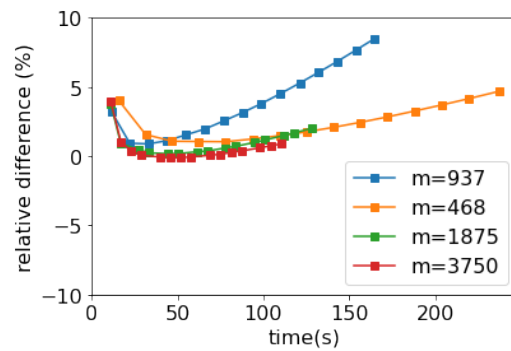


Figure 5: Relative difference of VRADAM in classifying MNIST with FFN

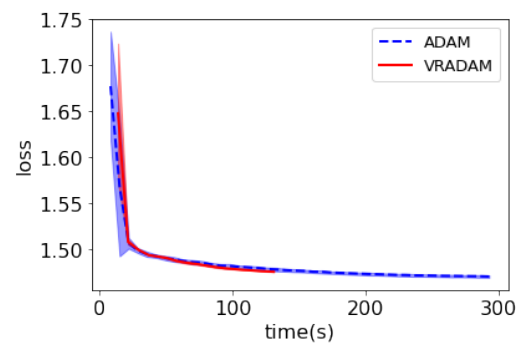


Figure 6: Deviation of VRADAM and ADAM for MNIST with CNN

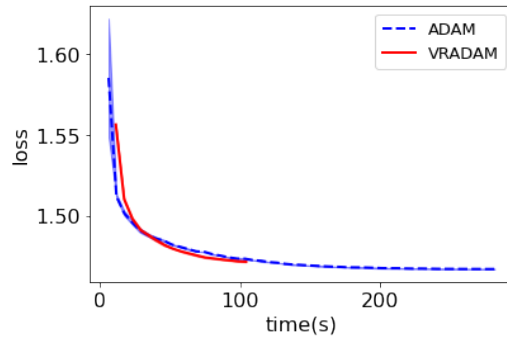


Figure 7: Deviation of VRADAM and ADAM for MNIST with FFN

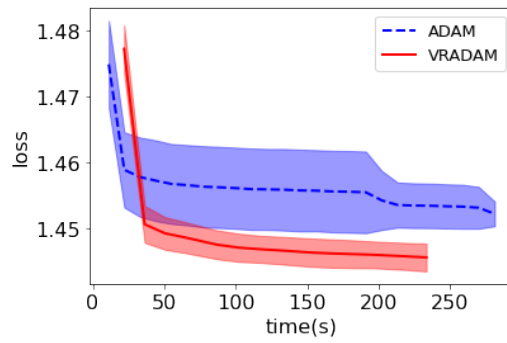


Figure 8: Deviation of VRADAM and ADAM for CovType with logistic regression

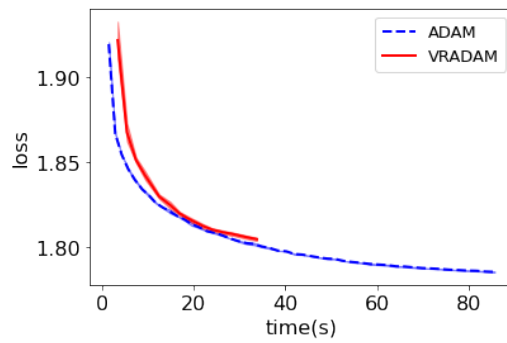


Figure 9: Deviation of VRADAM and ADAM for Embedded CIFAR-10 with logistic regression

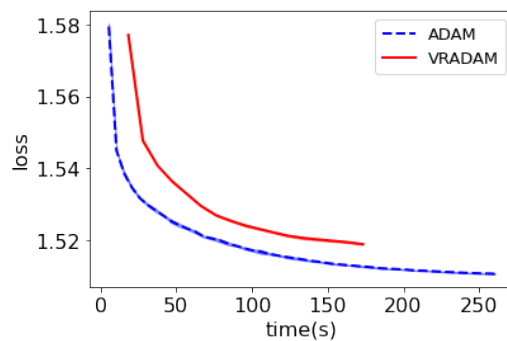


Figure 10: Deviation of VRADAM and ADAM for MNIST with logistic regression

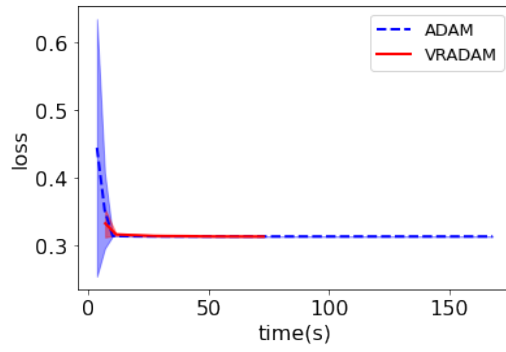


Figure 11: Deviation of VRADAM and ADAM for NSL-KDD with logistic regression

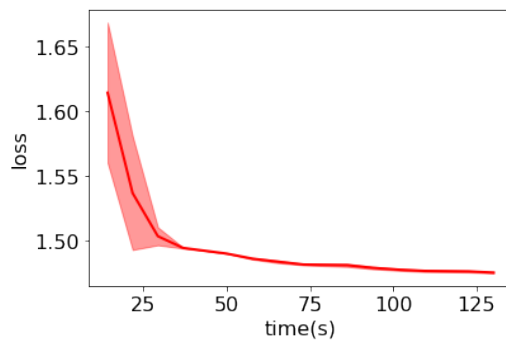


Figure 12: Sensitivity on initial point for MNIST with CNN

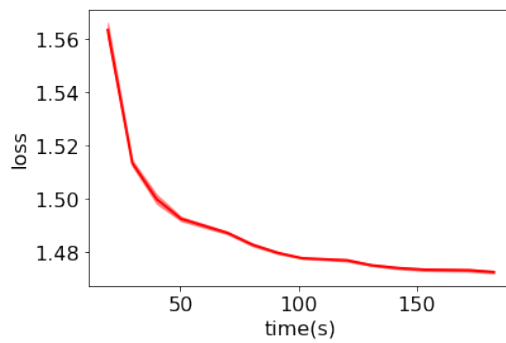


Figure 13: Sensitivity on initial point for MNIST with FFN

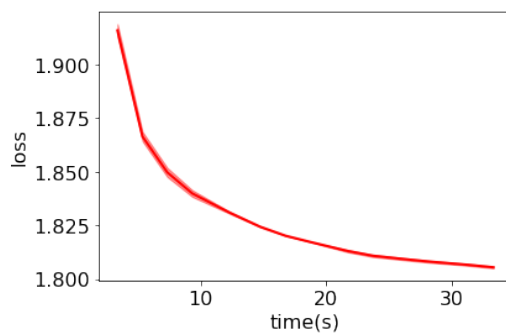


Figure 14: Sensitivity on initial point for Embedded CIFAR-10 with logistic regression

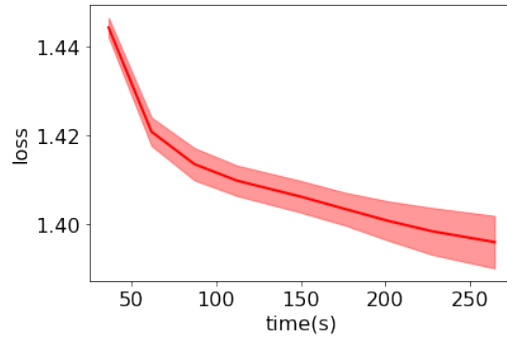


Figure 15: Sensitivity on initial point for CovType with logistic regression

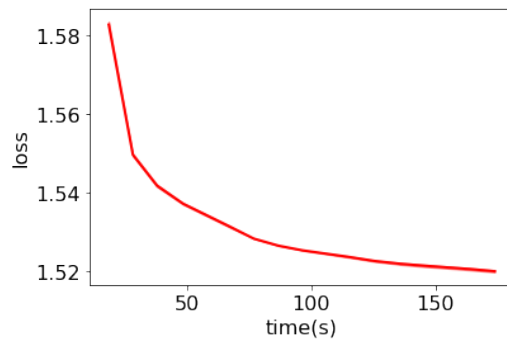


Figure 16: Sensitivity on initial point for MNIST with logistic regression

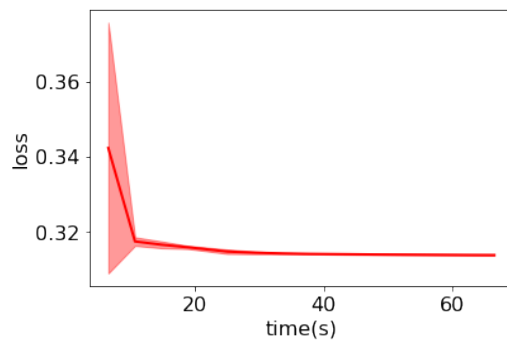


Figure 17: Sensitivity on initial point for NSL-KDD with logistic regression