

# DOES MOMENTUM CHANGE THE IMPLICIT BIAS ON SEPARABLE DATA?

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

The momentum acceleration technique is widely adopted in many optimization algorithms. However, the theoretical understanding on how the momentum affects the generalization performance of the optimization algorithms is still unknown. In this paper, we answer this question through analyzing the implicit bias of momentum-based optimization. We prove that **on the linear classification problem with separable data**, both SGD with momentum and Adam (**without stochasticity**) converge to the  $L_2$  max-margin solution for exponential-tailed loss, which is the same as vanilla gradient descent. That means, these optimizers with momentum acceleration still converge to a model with low complexity, which provides guarantees on their generalization. Technically, to overcome the difficulty brought by the error accumulation in analyzing the momentum, we construct new Lyapunov functions as a tool to analyze the gap between the model parameter and the max-margin solution.

## 1 INTRODUCTION

It is widely believed that the optimizers have implicit bias in terms of selecting output parameters among all the local minima on the landscape (Neyshabur et al., 2015; Keskar et al., 2017; Wilson et al., 2017). It is shown in the analysis of Adaboost that the *coordinate descent* would converge to the  $L^1$  max-margin solution for the linear classification task with exponential-tailed loss ((Schapire & Freund, 2013; Telgarsky, 2013)). Latter, Soudry et al. (2018) shows that *gradient descent* would converge to the  $L^2$  max-margin solution under the same setting, which mirrors its good generalization property in practice. Since then, many efforts have been taken on analyzing the implicit bias of various local-search optimizers, including stochastic gradient descent (Nacson et al., 2019), steepest descent (Gunasekar et al., 2018a), AdaGrad (Qian & Qian, 2019) and optimizers for homogeneous neural networks (Lyu & Li, 2019; Ji & Telgarsky, 2020; Wang et al., 2021).

However, though the momentum acceleration technique is widely adopted in the optimization algorithms in both convex and non-convex learning tasks (Sutskever et al., 2013; Vaswani et al., 2017; Tan & Le, 2019), the understanding on how the momentum would affect generalization performance of the optimization algorithms is still unclear. A natural question is:

*Can we theoretically analyze the implicit bias of momentum-based optimizers?*

In this paper, we take the first step to analyze the convergence of momentum based optimizers and unveil their implicit bias. Specifically, we study the classification problem with linear model and exponential-tailed loss using Stochastic Gradient Descent with Momentum (SGDM) and Adam optimizers. We consider the optimizers with constant learning rate and constant momentum hyper-parameters, which are widely adopted in practice, e.g., the default setting in popular machine learning frameworks (Paszke et al., 2019) and in experiments (Xie et al., 2017). We note that Gradient Descent with Momentum (GDM) can be viewed as a special case of SGDM, and naturally share the properties for SGDM. Our main results are summarized in Theorem 1.

**Theorem 1** (informal). *With linear separable dataset  $\mathcal{S}$ , for SGDM and Adam (**without stochasticity, abbreviated as w/s latter**), the loss converges to 0 with rate  $\mathcal{O}(\frac{1}{t})$  where  $t$  is the number of iterations, the parameter norm diverges to infinity, and direction of parameters converges to the direction of the  $L^2$  max-margin solution.*

Theorem 1 states SGDM converges to the  $L^2$  max-margin solution, which is the same as SGD, indicating that momentum does not affect the convergent direction. The good generalization behavior of the output parameters of SGDM is well validated as the margin of a classifier is positively correlated with its generalization error (Jiang et al., 2019). This is supported by existing experimental observations (c.f. Figure 1, (Soudry et al., 2018) and Figure 2, (Nacson et al., 2019)). Similar claims hold for Adam (w/s), which is also well supported by empirical results in Wang et al. (2021).

Our contributions are significant in terms of the following aspects:

- We establish the implicit bias of the momentum based optimizers, an open problem since the initial work Soudry et al. (2018). The momentum based optimizers are widely used in practice and our theoretical characterization deepens the understanding on their generalization property, which is important by its own.
- Technically, we propose a *new Lyapunov function* to analyze the convergence of SGDM, which helps to bound the sum of squared gradients along the training trajectory. Compared to the usual one, the new Lyapunov function depends on a middle variable of an alternative update rule of SGDM, which helps to capture the historical dependence in the momentum update. To our knowledge, such a technique has not been exploited ever before and can be of independent interest for convergence analysis of momentum-based optimizers. We then construct a new Lyapunov function to bound the difference of learned parameters and the scaled max-margin solution, which finally leads to the direction convergence. This Lyapunov function provides a direct way to establish the convergence to the desired direction.

**Organization of This Paper.** Section 2 collects further related works on the implicit bias of first order optimizers and convergence of momentum-based optimizers. Section 3 shows basic settings and assumptions which will be used throughout this paper. Section 4 studies the implicit bias of GDM as a warm up, while Section 5 and Section 6 explore respectively the implicit bias of SGDM and Adam (w/s). Discussions of these results are put in Section 7.

## 2 FURTHER RELATED WORKS

**Implicit Bias of First-order Optimization Methods.** Soudry et al. (2018) prove that gradient descent on linear classification problem with exponential-tailed loss converges to the direction of the max  $L^2$  margin solution of the corresponding hard-margin Support Vector Machine. Nacson et al. (2019) extend the results in (Soudry et al., 2018) to the stochastic case, proving that the convergent direction of SGD is the same as GD almost surely. Qian & Qian (2019) go beyond the vanilla gradient descent methods and consider the AdaGrad optimizer instead. They prove that the convergent direction of AdaGrad has a dependency on the optimizing trajectory, which varies according to the initialization. Ji & Telgarsky (2021) propose a primal-dual analysis framework for the linear classification models, and prove a faster convergent rate of the margin by increasing the learning rate according to the loss. Based on (Ji & Telgarsky, 2021), (Ji et al., 2021) design another algorithm with an even faster convergent rate of margin by applying the Nesterov’s Acceleration Method on the dual space. However, the corresponding form of the algorithm on the primal space is no longer a Nesterov’s Acceleration Method nor GDM, which is significantly different from our settings.

On the other hand, there is another line of work trying to extend the result in the linear case to deep neural networks. Ji & Telgarsky (2018); Gunasekar et al. (2018b) study the deep linear network and Soudry et al. (2018) study the two-layer neural network with ReLU activation. Lyu & Li (2019) propose a framework to analyze the asymptotic direction of GD on homogeneous neural networks, proving that given there exists a time the network achieves 100% training accuracy, GD will converge to some KKT point of the  $L^2$  max-margin problem. Wang et al. (2021) extend the framework of Lyu & Li (2019) to adaptive optimizers, and prove RMSProp and Adam without momentum have the same convergent direction as GD, while AdaGrad doesn’t. The results (Lyu & Li, 2019; Wang et al., 2021) indicate that results in the linear model can be extended to deep homogeneous neural networks, and suggest that the linear model is a proper start point to study the implicit bias. **There are also works on the implicit bias of regression problems with bounded optimal points and interesting readers can refer to (Rosasco & Villa, 2015; Lin & Rosasco, 2017; Ali et al., 2020) etc. for details.**

**Convergence of Momentum-Based Optimization Methods.** For convex optimization problems, the convergence rate of Nesterov’s Acceleration Method has been proved in Nesterov (1983). In

contrast, although GDM (Polyak’s Heavy-Ball Method) was proposed in (Polyak, 1964) prior to the Nesterov’s Acceleration Method, the convergence of GDM on convex loss with Lipschitz gradient was not solved until Ghadimi et al. (2015) provides an ergodic convergent result for GDM, i.e., the convergent result for the running average of the iterates. However, the ergodic result is undesired under many learning scenarios, e.g., in classification tasks, the optimization algorithms usually output the parameters of the last step. To the best of our knowledge, there is only two works on the non-ergodic analysis of (S)GDM: Sun et al. (2019) and Tao et al. (2021). Sun et al. (2019) prove that if the training loss is coercive (the training loss goes to infinity whenever parameter norm goes to infinity), convex, and globally smooth, then, with constant momentum hyper-parameter, the training loss converges to minima with rate  $\mathcal{O}(t^{-1})$ . Tao et al. (2021) analyze a case when momentum coefficient increases to 1 and the gradient is bounded all over the parameter space, showing SGDM can improve the convergence rate of SGD by a factor  $\log(t)$ .

There are also works on the gradient norm convergence of SGDM under various settings (Yan et al., 2018; Yu et al., 2019; Liu et al., 2020) and on the investigation of momentum-based method from the view point of dynamics (Sarao Mannelli & Urbani, 2021). However, there is no existing work on the implicit bias of momentum-based optimizers for the classification problem, which is first analyzed by this paper.

### 3 PRELIMINARIES

In this paper, we focus on the linear model with exponential-tailed loss. We first derive the results for binary classification, then we show that the methodology can be easily extended to the multi-class classification problem.

**Problem setting.** The dataset used for training is defined as  $\mathcal{S} = (\mathbf{x}_i, \mathbf{y}_i)_{i=1}^N$ , where  $\mathbf{x}_i \in \mathbb{R}^d$  is the  $i$ -th feature, and  $\mathbf{y}_i \in \mathbb{R}$  is the  $i$ -th label ( $i = 1, 2, \dots, N$ ). We will use the linear model to fit the label: for any feature  $\mathbf{x} \in \mathbb{R}^d$  and parameter  $\mathbf{w} \in \mathbb{R}^d$ , the prediction is given by  $\langle \mathbf{w}, \mathbf{x} \rangle$ .

For binary classification, given any data  $\mathbf{z}_i = (\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{S}$ , the individual loss for parameter  $\mathbf{w}$  is given as  $\ell(\mathbf{y}_i \langle \mathbf{w}, \mathbf{x}_i \rangle)$ . A common setting is to ensemble the feature and label together as  $\tilde{\mathbf{x}}_i = \mathbf{y}_i \mathbf{x}_i$  (there is a mapping  $\mathcal{T} : (\mathbf{x}, \mathbf{y}) \rightarrow \mathbf{y}\mathbf{x}$ , and  $\tilde{\mathbf{x}}_i$  is  $\mathcal{T}((\mathbf{x}_i, \mathbf{y}_i))$ ). The individual loss can be rewritten as

$$\tilde{\ell}(\mathbf{w}, (\mathbf{x}_i, \mathbf{y}_i)) = \ell(\mathbf{y}_i \langle \mathbf{w}, \mathbf{x}_i \rangle) = \ell(\langle \mathbf{w}, \tilde{\mathbf{x}}_i \rangle).$$

The optimization target is defined as the averaged loss:

$$\mathcal{L}(\mathbf{w}) = \frac{\sum_{i=1}^N \tilde{\ell}(\mathbf{w}, (\mathbf{x}_i, \mathbf{y}_i))}{N}.$$

Without loss of generality, we consider the case with normalized data<sup>1</sup>, that is,  $\|\tilde{\mathbf{x}}_i\| \leq 1, \forall i \in [N]$ .

**Optimizer.** Here we will introduce the update rules of SGDM and Adam (w/s). SGDM can be viewed as a stochastic version of GDM by randomly choosing a subset of the dataset to update. Specifically, the update rule of SGDM is given as

$$(SGDM) : \mathbf{w}(t+1) - \mathbf{w}(t) = -\eta \nabla \mathcal{L}_{\mathcal{B}(t)}(\mathbf{w}(t)) + \beta(\mathbf{w}(t) - \mathbf{w}(t-1)), \forall t \geq 1, \quad (1)$$

where  $\mathcal{B}(t)$  is a subset of  $\mathcal{S}$  with size  $b$  which is sampled independently and uniformly with replacement, and  $\mathcal{L}_{\mathcal{B}(t)}$  is defined as  $\mathcal{L}_{\mathcal{B}(t)}(\mathbf{w}) = \frac{\sum_{\mathbf{z} \in \mathcal{B}(t)} \tilde{\ell}(\mathbf{w}, \mathbf{z})}{b}$ . We also define  $\mathcal{F}_t$  as the sub-sigma field such that  $\{\mathbf{w}(t)\}_{t=1}^{\infty}$  is adapted with respect to Filtration  $\{\mathcal{F}_t\}_{t=1}^{\infty}$ .

<sup>1</sup>The proof can be naturally applied to unnormalized data by letting  $\ell(x) = \ell(\max_{\tilde{\mathbf{x}} \in \mathcal{T}(\mathcal{S})} \|\tilde{\mathbf{x}}\| \cdot x)$  and  $\tilde{\mathbf{x}}_i = \frac{\tilde{\mathbf{x}}_i}{\max_{\tilde{\mathbf{x}} \in \mathcal{T}(\mathcal{S})} \|\tilde{\mathbf{x}}\|}$ .

The Adam (**w/s**) can be viewed as a variant of GDM in which the preconditioner is adopted, whose form is characterized as follows:

$$\begin{aligned} \mathbf{m}(0) &= \mathbf{0}, \mathbf{m}(t) = \beta_1 \mathbf{m}(t-1) + (1 - \beta_1) \nabla \mathcal{L}(\mathbf{w}(t)), \hat{\mathbf{m}}(t) = \frac{1}{1 - \beta_1^t} \mathbf{m}(t), \forall t \geq 1, \\ \nu(0) &= 0, \nu(t) = \beta_2 \nu(t-1) + (1 - \beta_2) (\nabla \mathcal{L}(\mathbf{w}(t)))^2, \hat{\nu}(t) = \frac{1}{1 - \beta_2^t} \nu(t), \forall t \geq 1, \\ (\text{Adam } (\mathbf{w}/s)) : \mathbf{w}(t) &= \mathbf{w}(t-1) - \eta \frac{\hat{\mathbf{m}}(t-1)}{\sqrt{\hat{\nu}(t-1) + \varepsilon \mathbb{1}_d}}, \forall t \geq 1, \end{aligned} \quad (2)$$

where  $\frac{1}{\sqrt{\hat{\nu}(t-1) + \varepsilon \mathbb{1}_d}}$  is called the preconditioner at step  $t$ .

**Assumptions:** The analysis of this paper are based on **three common assumptions in existing literature (first proposed by (Soudry et al., 2018))**, respectively on the separability of the dataset, the individual loss behaviour at the tail, and the smoothness of the individual loss. We list them as follows:

**Assumption 1** (Linearly Separable Dataset). *There exists one parameter  $\mathbf{w} \in \mathbb{R}^d$ , such that*

$$\langle \mathbf{w}, \tilde{\mathbf{x}}_i \rangle > 0, \forall i \in [N].$$

**Assumption 2** (Exponential-tailed Loss). *The individual loss  $\ell$  is exponential-tailed, i.e.,*

- *Differentiable and monotonically decreasing to zero, with its derivative also converging to zero, i.e.,  $\lim_{x \rightarrow \infty} \ell(x) = \lim_{x \rightarrow \infty} \ell'(x) = 0$ , and  $\ell'(x) < 0 \forall x$ ;*
- *Close to exponential loss when  $x$  is large enough, i.e., there exist positive constants  $c, a, K, \mu_+, \mu_-, x_+, x_-$  and  $x_0$ , such that,*

$$\forall x > x_+ : -\ell'(x) \leq c(1 + e^{-\mu_+ x})e^{-ax}, \quad (3)$$

$$\forall x > x_- : -\ell'(x) \geq c(1 - e^{-\mu_- x})e^{-ax}. \quad (4)$$

**Assumption 3** (Smooth Loss). *Either of the following assumptions holds regarding the case:*

**(D): (Deterministic Case)** *The individual loss  $\ell$  is locally smooth, i.e., for any  $s_0 \in \mathbb{R}$ , there exists a positive real  $H_{s_0}$ , such that  $\forall x, y \geq s_0, |\ell'(x) - \ell'(y)| \leq H_{s_0} |x - y|$ .*

**(S): (Stochastic Case)** *The individual loss  $\ell$  is globally smooth, i.e., there exists a positive real  $H$ , such that  $\forall x, y \in \mathbb{R}, |\ell'(x) - \ell'(y)| \leq H |x - y|$ .*

We provide explanations of these three assumptions respectively. Based on Assumption 1, we can formally define the margin and the maximum margin solution of an optimization problem, which is introduced in Definition 1

**Definition 1.** *Let the margin  $\gamma(\mathbf{w})$  of parameter  $\mathbf{w}$  defined as the lowest score of the prediction of  $\mathbf{w}$  over the dataset  $\mathcal{S}$ , i.e.,  $\gamma(\mathbf{w}) = \min_{\tilde{\mathbf{x}}_i \in \mathcal{T}(\mathcal{S})} \langle \mathbf{w}, \tilde{\mathbf{x}}_i \rangle$ . By Assumption 1 and the positive homogeneous of  $\gamma$ ,  $\gamma(\frac{\hat{\mathbf{w}}}{\gamma(\hat{\mathbf{w}})}) = 1$ , and thus we define the maximum margin solution  $\hat{\mathbf{w}}$  and the  **$L^2$  max margin  $\gamma$**  of the dataset  $\mathcal{S}$  as follows:*

$$\hat{\mathbf{w}} \triangleq \arg \min_{\gamma(\mathbf{w}) \geq 1} \|\mathbf{w}\|^2, \quad \gamma \triangleq \frac{1}{\|\hat{\mathbf{w}}\|}$$

Since  $\|\cdot\|^2$  is strongly convex and set  $\{\mathbf{w} : \gamma(\mathbf{w}) \geq 1\}$  is convex,  $\hat{\mathbf{w}}$  is uniquely defined.

Assumption 2 constraints the loss to be exponential-tailed, which is satisfied by many popular choices of  $\ell$ , including the ( $\ell_{exp}(x) = e^{-x}$ ) and the logistic loss ( $\ell_{log}(x) = \log(1 + e^{-x})$ ). Also, as  $c$  and  $a$  can be respectively absorbed by resetting the learning rate and data as  $\eta = c\eta$  and  $\mathbf{x}_i = a\mathbf{x}_i$ , without loss of generality, in this paper we only analyze the case that  $c = a = 1$ .

The globally smooth assumption (Assumption 3. (S)) is strictly stronger than the locally smooth assumption (Assumption 3. (D)). One can easily verify that both the exponential loss and the logistic loss meet Assumption 3. (D), and the logistic loss also meets Assumption 3. (S).

#### 4 WARM UP: IMPLICIT BIAS OF GDM

In this section, we study the implicit bias of GDM as a warm up. The update rule of GDM is

$$\mathbf{w}(t+1) - \mathbf{w}(t) = -\eta \nabla \mathcal{L}(\mathbf{w}(t)) + \beta(\mathbf{w}(t) - \mathbf{w}(t-1)), \forall t \geq 1 \quad (5)$$

which is a determined case of SGDM, with batch size  $b = N$  and thus without randomness. Although the analysis of GDM is essentially easier than that of SGDM, it helps us gain a better understanding of this problem, and helps us to demonstrate several techniques which will be applied latter. The formal theorem of the implicit bias of GDM is as follows:

**Theorem 2.** *Let Assumptions 1, 2, and 3. (D) hold. Let  $\beta \in [0, 1)$  and  $\eta < 2 \frac{1-\beta}{H_{\ell^{-1}(N\mathcal{L}(\mathbf{w}_1))}}$ . Then, for almost every data set  $S^2$ , with arbitrary initialization point  $\mathbf{w}(1)$ , GDM (Eq. (5)) satisfies that  $\mathbf{w}(t) - \log(t)\hat{\mathbf{w}}$  is bounded as  $t \rightarrow \infty$ , and  $\lim_{t \rightarrow \infty} \frac{\mathbf{w}(t)}{\|\mathbf{w}(t)\|} = \frac{\hat{\mathbf{w}}}{\|\hat{\mathbf{w}}\|}$ .*

Before we move on to the proof, we would like to provide some explanations of the results in Theorem 2. To begin with, Theorem 2 adopts an constant momentum hyper-parameter, which agrees with the practice use (e.g.,  $\beta$  is set to be 0.9 (Xie et al., 2017)). Also, Theorem 2 puts no restriction on the range of  $\beta$ , which allows wider choices of hyper-parameter tuning. Furthermore, Theorem 2 shows the implicit bias of GDM agrees with GD in linear classification with exponential-tailed loss (c.f. Soudry et al. (2018) for results on GD), and this consistency can be verified by existing results (c.f. Section 7 for detailed discussions).

We then present a proof sketch of Theorem 2, which is divided into two parts: we first prove that the sum of squared gradients is bounded, which indicates both the loss and the norm of gradient converge to 0; we then show the difference between  $\mathbf{w}(t)$  and  $\log(t)\hat{\mathbf{w}}$  is bounded, and therefore, the direction of  $\hat{\mathbf{w}}$  dominates as  $t \rightarrow \infty$ .

**Stage I: Bound the sum of squared gradients.** The core of Stage I is to select a proper Lyapunov function  $\xi(t)$ , which is required to correlated with the training loss  $\mathcal{L}$  and be non-increasing along the optimization trajectory. For GD, since  $\mathcal{L}$  itself is non-increasing with properly chosen learning rate, we can just pick  $\xi(t) = \mathcal{L}(t)$ . However, as the update of GDM doesn't align with the direction of negative gradient, training loss  $\mathcal{L}(t)$  in GDM is no longer monotonously decreasing, and the Lyapunov function requires special construction. A choice of  $\xi(t)$  is proposed as the follows:

**Lemma 1.** *Let all conditions in Theorem 2 hold. Define  $\xi(t) \triangleq \mathcal{L}(\mathbf{w}(t)) + \frac{\beta}{2\eta} \|\mathbf{w}(t) - \mathbf{w}(t-1)\|^2$ . Let  $C_1$  be a positive real, s.t.,  $\eta = 2 \frac{1-\beta}{H_{\ell^{-1}(N\mathcal{L}(\mathbf{w}_1))}} C_1$ . We then have*

$$\xi(t) \geq \xi(t+1) + \frac{(1-\beta)(1-C_1)}{\eta} \|\mathbf{w}(t+1) - \mathbf{w}(t)\|^2. \quad (6)$$

**Remark 1.** *Although this Lyapunov function is obtained by (Sun et al., 2019) by directly examining the Taylor's expansion at  $\mathbf{w}(t)$ , the proof here is non-trivial as we only requires the loss to be locally smooth instead of globally smooth in (Sun et al., 2019), and the Taylor's expansion can only be applied to  $\mathbf{w}(t+1)$  if it's ensured that all parameters on the line  $\alpha\mathbf{w}(t) + (1-\alpha)\mathbf{w}(t+1)$  ( $\alpha \in (0, 1)$ ) have training loss no larger than  $\mathcal{L}(\mathbf{w}(0))$ .*

To prove Lemma 1, we define  $\mathbf{w}(t+\alpha) = \alpha\mathbf{w}(t) + (1-\alpha)\mathbf{w}(t+1)$  for any  $t \in \mathbb{Z}^+$  and  $\alpha \in (0, 1)$ , and prove a generalized version of Eq. (6):

$$\xi(t) \geq \xi(t+\alpha) + \frac{(1-\beta)(1-c)\alpha^2}{\eta} \|\mathbf{w}(t+1) - \mathbf{w}(t)\|^2. \quad (7)$$

The proof idea is that as long as Eq. (7) holds for all time in  $[1, t+\alpha)$ , the training loss across  $[1, t+\alpha]$  will be smaller than  $\mathcal{L}(\mathbf{w}(1))$ , and the strict inequality in Eq. (7) holds. Consequently, there exists some small enough positive real  $\varepsilon$ , such that Eq. (7) holds for all time in  $[1, t+\alpha+\varepsilon)$ , and we are able to extend the feasible set where Eq. (7) holds. The formal description of the generalized lemma and its corresponding proof is deferred to Appendix B.1.1.

<sup>2</sup>Here "almost everywhere" means the conclusion holds except a zero-measure set in  $\mathbb{R}^{d \times N}$ .

By Lemma 1, we have that  $\xi(t)$  is monotonously decreasing by gap  $\frac{(1-\beta)(1-c)\alpha^2}{\eta} \|\mathbf{w}(t+1) - \mathbf{w}(t)\|^2$ . As  $\xi(1) = \mathcal{L}(\mathbf{w}(1))$  is a finite number, we have  $\sum_{t=1}^{\infty} \|\mathbf{w}(t+1) - \mathbf{w}(t)\|^2 < \infty$ , which further leads to the following corollary by the negative derivative of the individual loss and separable data:

**Corollary 1.** *Let all conditions in Theorem 2 hold. We have,  $\sum_{t=1}^{\infty} \|\nabla \mathcal{L}(\mathbf{w}(t))\|^2 < \infty$ .*

The proof of Corollary 1 can be found in Appendix B.1.1. Resulted from Corollary 1, we have  $\|\nabla \mathcal{L}(\mathbf{w}(t))\|^2 \rightarrow 0$  as  $t \rightarrow \infty$ , which by the exponential-tailed loss assumption further leads to  $\lim_{t \rightarrow \infty} \mathcal{L}(\mathbf{w}(t)) = 0$ .

**Stage II. Bound the difference between  $\mathbf{w}(t)$  and  $\log(t)\hat{\mathbf{w}}$ .** It seems natural to directly bound  $\|\mathbf{w}(t) - \log(t)\hat{\mathbf{w}}\|$  across all iterations. Soudry et al. (2018) and Nacson et al. (2019) indeed follow this routine by showing the norm of  $\mathbf{r}(t) \triangleq \mathbf{w}(t) - \log(t)\hat{\mathbf{w}} - \tilde{\mathbf{w}}$  is bounded, where  $\tilde{\mathbf{w}}$  is some constant vector satisfying  $e^{\langle \tilde{\mathbf{w}}, \hat{\mathbf{x}} \rangle}$  recovers the coefficient of support vector  $\hat{\mathbf{x}}$  in  $\hat{\mathbf{w}}$ . However, their analyses rely on the fact that  $\mathbf{w}(t+1) - \mathbf{w}(t)$  equals  $-\eta \nabla \mathcal{L}(\mathbf{w}(t))$ , which has a simple form, i.e.,

$$\text{For GD: } \|\mathbf{r}(t+1)\|^2 - \|\mathbf{r}(t)\|^2 = \|\mathbf{r}(t+1) - \mathbf{r}(t)\|^2 - 2\langle \log \frac{t+1}{t} \hat{\mathbf{w}} + \eta \nabla \mathcal{L}(\mathbf{w}(t)), \mathbf{r}(t) \rangle.$$

However, it no longer holds for GDM, as  $\mathbf{w}(t+1) - \mathbf{w}(t)$  is a exponentially-decayed sum of gradients at  $\mathbf{w}(s)$ ,  $s \leq t$ . To handle this problem, we propose a novel Lyapunov's function for the direction convergence, concluded as the following Lemma.

**Lemma 2.** *Let all conditions in Theorem 2 hold. Then,  $\|\mathbf{r}(t)\|$  is bounded if and only if the function  $g(t)$  is upper bounded, where  $g: \mathbb{Z}^+ \rightarrow \mathbb{R}$  is defined as*

$$g(t) \triangleq \frac{1}{2} \|\mathbf{r}(t)\|^2 + \frac{\beta}{1-\beta} \langle \mathbf{r}(t), \mathbf{w}(t) - \mathbf{w}(t-1) \rangle - \frac{\beta}{1-\beta} \sum_{\tau=2}^t \langle \mathbf{r}(\tau) - \mathbf{r}(\tau-1), \mathbf{w}(\tau) - \mathbf{w}(\tau-1) \rangle. \quad (8)$$

Furthermore, we have  $\sum_{t=1}^{\infty} (g(t+1) - g(t))$  is upper bounded.

The first claim in Lemma 2 provides an alternative approach to verify that  $\|\mathbf{r}(t)\|$  is bounded, while the second claim shows this approach can be fulfilled. The motivation of this lemma is to construct a easy-to-verify criterion (i.e.,  $g(t)$  is upper bounded) of  $\|\mathbf{r}(t)\|$  is upper bounded. To demonstrate how  $g(t)$  is constructed intuitively, we consider the continuous dynamics approximation of GDM, i.e.,

$$\frac{\beta}{1-\beta} \frac{d^2 \mathbf{w}(t)}{dt^2} + \frac{d\mathbf{w}(t)}{dt} + \frac{\eta}{1-\beta} \nabla \mathcal{L}(\mathbf{w}(t)) = 0.$$

We then calculate the  $\frac{1}{2} \|\mathbf{r}(t)\|^2 - \frac{1}{2} \|\mathbf{r}(1)\|^2$  by applying the integral of the dynamics as

$$\begin{aligned} & \int_1^t \frac{1}{2} \frac{d\|\mathbf{r}(s)\|^2}{ds} ds = \int_1^t \left\langle \mathbf{r}(s), -\frac{\eta}{1-\beta} \nabla \mathcal{L}(\mathbf{w}(s)) - \frac{1}{s} \hat{\mathbf{w}} \right\rangle ds + \int_1^t \frac{\beta}{1-\beta} \left\langle \mathbf{r}(s), -\frac{d^2 \mathbf{w}(s)}{ds^2} \right\rangle ds \\ & = \int_1^t \left\langle \mathbf{r}(s), -\frac{\eta}{1-\beta} \nabla \mathcal{L}(\mathbf{w}(s)) - \frac{1}{s} \hat{\mathbf{w}} \right\rangle ds + \frac{\beta}{1-\beta} \left( \left\langle \mathbf{r}(t), -\frac{d\mathbf{w}(t)}{dt} \right\rangle - \left\langle \mathbf{r}(1), -\frac{d\mathbf{w}(t)}{dt} \Big|_{t=1} \right\rangle + \int_1^t \left\langle \frac{d\mathbf{r}(s)}{ds}, \frac{d\mathbf{w}(s)}{ds} \right\rangle ds \right), \end{aligned}$$

where the last equation is due to Integration by part. We denote  $\tilde{g}(t) = \frac{\beta}{1-\beta} (\langle \mathbf{r}(t), \frac{d\mathbf{w}(t)}{dt} \rangle - \int_1^t \langle \frac{d\mathbf{r}(s)}{ds}, \frac{d\mathbf{w}(s)}{ds} \rangle ds) + \frac{1}{2} \|\mathbf{r}(t)\|^2$ , and it can be verified the derivative of  $\tilde{g}(t)$  is  $\langle \mathbf{r}(t), -\frac{\eta}{1-\beta} \nabla \mathcal{L}(\mathbf{w}(t)) - \frac{1}{t} \hat{\mathbf{w}} \rangle$ . By replacing  $\frac{d\mathbf{w}(t)}{dt}$  as  $\mathbf{w}(t) - \mathbf{w}(t-1)$ , and  $\frac{d\mathbf{r}(t)}{dt}$  as  $\mathbf{r}(t) - \mathbf{r}(t-1)$  in  $\tilde{g}(t)$ , we obtain  $g(t)$ , and  $g(t+1) - g(t)$  has the form  $\langle \mathbf{r}(t), -\frac{\eta}{1-\beta} \nabla \mathcal{L}(\mathbf{w}(t)) - \log \frac{t+1}{t} \hat{\mathbf{w}} \rangle$ , which can be analyzed following the similar routine as (Soudry et al., 2018). The proof of the second claim is completed.

On the other hand, the core of the proof of Lemma 2 is that  $\frac{1}{2} \|\mathbf{r}(t)\|^2 + \frac{\beta}{1-\beta} \langle \mathbf{r}(t), \mathbf{w}(t) - \mathbf{w}(t-1) \rangle$  is bounded if and only if  $\frac{1}{2} \|\mathbf{r}(t)\|^2$  is bounded, while  $\langle \mathbf{r}(\tau) - \mathbf{r}(\tau-1), \mathbf{w}(\tau) - \mathbf{w}(\tau-1) \rangle$  is an absolute convergent based on  $\sum_{t=1}^{\infty} \|\mathbf{w}(t+1) - \mathbf{w}(t)\|^2 < \infty$ .

## 5 TACKLE THE DIFFICULTY BROUGHT BY RANDOM SAMPLING

In this section, we analyze the implicit bias of SGDM. The randomness introduced by random sampling makes the analysis of GDM no longer work for SGDM. As an example, if we want to

follow the same routine of the GDM case to show  $\mathcal{L}(\mathbf{w}(t)) + \frac{\beta}{2\eta}\mathbb{E}\|\mathbf{w}(t) - \mathbf{w}(t-1)\|^2$  is a Lyapunov function of SGDM, we can only have

$$\begin{aligned} \mathbb{E}\mathcal{L}(\mathbf{w}(t)) + \frac{\beta}{2\eta}\mathbb{E}\|\mathbf{w}(t) - \mathbf{w}(t-1)\|^2 &\geq \mathbb{E}\mathcal{L}(\mathbf{w}(t+1)) + \frac{\beta}{2\eta}\mathbb{E}\|\mathbb{E}[\mathbf{w}(t+1) - \mathbf{w}(t)|\mathcal{F}_t]\|^2 \\ &\quad + \frac{(1-\beta)(1-c)}{\eta}\mathbb{E}\|\mathbb{E}[\mathbf{w}(t+1) - \mathbf{w}(t)|\mathcal{F}_t]\|^2. \end{aligned}$$

As the squared first momentum is smaller than the second momentum,  $\mathbb{E}\mathcal{L}(\mathbf{w}(t)) + \frac{\beta}{2\eta}\mathbb{E}\|\mathbf{w}(t) - \mathbf{w}(t-1)\|^2$  may no longer be monotonously decreasing for every  $\beta \in [0, 1)$ . Furthermore, if the margin  $\gamma$  is small,  $\beta$  will be required to be upper bounded by a small number, which contradicts to the relatively large choice of  $\beta$  in practice, e.g. 0.9 (We defer a detailed discussion to Appendix B.2.3). To tackle this problem, we propose a new Lyapunov function, with which we derive the following theorem for the implicit bias of SGDM:

**Theorem 3.** *Let Assumption 1, 2, and 3. (S) hold. Let  $\beta \in [0, 1)$  and  $\eta < \min\{\frac{1-\beta}{1+\frac{H}{b\gamma^2}}, \frac{(1-\beta)^3\gamma^4 b}{2H^2N^3\beta^2}\}$ .*

*Then, for almost every data set  $S$ , with arbitrary initialization point  $\mathbf{w}(1)$ , SGDM (Eq. (1)) satisfies  $\mathbf{w}(t) - \log(t)\hat{\mathbf{w}}$  is bounded as  $t \rightarrow \infty$  and  $\lim_{t \rightarrow \infty} \frac{\mathbf{w}(t)}{\|\mathbf{w}(t)\|} = \frac{\hat{\mathbf{w}}}{\|\hat{\mathbf{w}}\|}$ , almost surely (a.s.).*

To the best of our knowledge, this is the first convergence analysis of SGDM with no restriction on the gradient norm and with constant learning rates. Also, as both  $\frac{1-\beta}{1+\frac{H}{b\gamma^2}}$  and  $\frac{(1-\beta)^3\gamma^4 b}{2H^2N^3\beta^2}$  are monotonously increasing with respect to batch size  $b$ , Theorem 3 also sheds light on the learning rate tuning, i.e., the larger the batch size is, the larger the learning rate is. Furthermore, similar to the GDM case, Theorem 3 shows the implicit bias of SGDM under this setting is consistent with SGD (c.f. (Nacson et al., 2019) for the implicit bias of SGD). This matches the observations in practice (c.f. Section 7 for details).

**Remark 2.** *In practice, mini-batch SGDM are more widely adopted, whose update rule differs from SGDM only by the way obtaining  $\mathbf{B}(t)$ . Specifically, within the  $T$ -th epoch  $E(T) = \{KT + 1, \dots, KT + T\}$  ( $K$  is the length of one epoch),  $\{\mathbf{B}(t)\}_{t \in E(T)}$  is a randomly and uniformly partition of  $S$ . One may wonder whether the same result hold for mini-batch SGDM. The answer is "Yes", with the a.s. condition removed in this case. We defer the detailed description of the corresponding theorem together with the proof to Appendix D.1.*

The proof sketch follows the same framework as that for GDM by dividing the proof into two stages. However, the implementations in both stages differ, among which the proof of Stage I is significantly distinctive, as a new Lyapunov function is proposed.

**Stage I: Bound the sum of squared gradients.** To obtain the new Lyapunov function, we first present a lemma to provide an alternative form of update rule of SGDM:

**Lemma 3.** *Define  $\tilde{\eta} \triangleq \frac{\eta}{1-\beta}$ , and  $\mathbf{u}(1) \triangleq \mathbf{w}(1)$ . The update rule of SGDM (Eq. (1)) is equivalent to*

$$\begin{aligned} \mathbf{u}(t+1) &= -\tilde{\eta}\nabla\mathcal{L}_{\mathbf{B}(t)}(\mathbf{w}(t)) + \mathbf{u}(t), \\ \mathbf{w}(t+1) &= \beta\mathbf{w}(t) + (1-\beta)\mathbf{u}(t+1). \end{aligned} \tag{9}$$

By a simple rearrangement of the second equation, we have  $\mathbf{u}(t+1) = \mathbf{w}(t) + \frac{1}{1-\beta}(\mathbf{w}(t+1) - \mathbf{w}(t))$ , which differs from  $\mathbf{w}(t+1)$  only by a larger step size updated from  $\mathbf{w}(t)$ . The following lemma then indicates that  $\mathbb{E}\mathcal{L}(\mathbf{u}(t))$  is a proper selection of Lyapunov function.

**Lemma 4.** *Let all conditions in Theorem 3 hold. Then, we have*

$$\mathbb{E}[\mathcal{L}(\mathbf{u}(t+1))] \leq \mathbb{E}\mathcal{L}(\mathbf{u}(t)) - \frac{\tilde{\eta}}{2}\mathbb{E}\|\nabla\mathcal{L}(\mathbf{w}(t))\|^2 + \frac{1-\beta}{4}\tilde{\eta} \left( \sum_{s=1}^{t-1} \beta^{t-1-s} \mathbb{E}\|\nabla\mathcal{L}(\mathbf{w}(s))\|^2 \right),$$

and

$$\mathbb{E}[\mathcal{L}(\mathbf{u}(t+1))] \leq \mathcal{L}(\mathbf{u}(1)) - \sum_{s=1}^t \frac{\tilde{\eta}}{4}\mathbb{E}\|\nabla\mathcal{L}(\mathbf{w}(s))\|^2.$$

$\mathbb{E}\mathcal{L}(\mathbf{u}(t))$  may be not monotonously decreasing. However, the upper bound of  $\mathbb{E}\mathcal{L}(\mathbf{u}(t))$ , i.e.,  $\mathcal{L}(\mathbf{u}(1)) - \sum_{s=1}^t \frac{\eta}{4} \mathbb{E}\|\nabla\mathcal{L}(\mathbf{w}(s))\|^2$ , is monotonously decreasing by  $\frac{\eta}{4}\mathbb{E}\|\nabla\mathcal{L}(\mathbf{w}(t))\|^2$  at step  $t$ , and leads to the sum of expected squared gradients being finite along the trajectory (which further indicates the sum of squared gradients is finite, a.s.).

The proof idea of Lemma 4 is the expectation of the first order Taylor's expansion of  $\mathcal{L}$  at  $\mathbf{w}(t)$  has the form

$$\mathbb{E}[\mathcal{L}(\mathbf{u}(t+1))|\mathcal{F}_t] \leq \mathcal{L}(\mathbf{u}(t)) - \tilde{\eta}\langle\nabla\mathcal{L}(\mathbf{u}(t)), \nabla\mathcal{L}(\mathbf{w}(t))\rangle + \frac{L\tilde{\eta}^2}{2}\mathbb{E}[\|\nabla\mathcal{L}_{\mathcal{B}(t)}(\mathbf{w}(t))\|^2|\mathcal{F}_t].$$

As  $\mathcal{L}$  is  $H$  smooth, we have  $\nabla\mathcal{L}(\mathbf{u}(t)) = \nabla\mathcal{L}(\mathbf{w}(t)) + \mathcal{O}(\|\mathbf{u}(t) - \mathbf{w}(t)\|)$ , where  $\mathbf{u}(t) - \mathbf{w}(t) = \frac{\beta}{1-\beta}(\mathbf{w}(t) - \mathbf{w}(t-1)) = -\eta\frac{\beta}{1-\beta}\sum_{s=1}^{t-1}\beta^{t-1-s}\nabla\mathcal{L}_{\mathcal{B}(s)}(\mathbf{w}(s))$  is proportional to  $\eta$ . Therefore, when  $\eta$  is small,  $-\langle\nabla\mathcal{L}(\mathbf{u}(t)), \nabla\mathcal{L}(\mathbf{w}(t))\rangle$  is close to  $-\|\nabla\mathcal{L}(\mathbf{w}(t))\|^2$ , and the coefficient of  $-\|\nabla\mathcal{L}(\mathbf{w}(t))\|^2$  becomes dominant.

**Stage II. Bound the difference between  $\mathbf{w}(t)$  and  $\log(t)\hat{\mathbf{w}}$ .** The Lyapunov function used by Stage II of SGDM is the same as the  $g(t)$  in the GDM case. As we have proved the sum of squared gradient along the trajectory is bounded, one can easily obtain that  $\sum_{t=1}^{\infty}\|\mathbf{w}(t+1) - \mathbf{w}(t)\|^2$  is bounded, and  $\|\mathbf{r}(t)\|$  being bounded is still equivalent to that  $g(t)$  is upper bounded. However, when it comes to the detailed calculation of analyzing  $g(t+1) - g(t)$ , it differs from the GDM case due to the random subset. We defer the proof to Appendix B.2.2.

## 6 ANALYZE THE EFFECT OF PRECONDITIONERS

When it comes to the analysis of Adam (**w/s**), the effect of the preconditioner should be taken into consideration when designing the corresponding Lyapunov functions. To tackle this problem, we incorporate the preconditioner into the Lyapunov function, and obtain the implicit bias of Adam (**w/s**) as follows:

**Theorem 4.** *Let Assumption 1, 2, and 3. (D) hold. Let  $1 > \beta_2 > \beta_1^4 \geq 0$ , and the learning rate  $\eta$  is a small enough constant (The upper bound of learning rate is complex, and we defer it to Appendix C.1). Then, for almost every data set  $S$ , with arbitrary initialization point  $\mathbf{w}(1)$ , Adam (**w/s**) (Eq. (2)) satisfies that  $\mathbf{w}(t) - \log(t)\hat{\mathbf{w}}$  is bounded as  $t \rightarrow \infty$ , and  $\lim_{t \rightarrow \infty} \frac{\mathbf{w}(t)}{\|\mathbf{w}(t)\|} = \frac{\hat{\mathbf{w}}}{\|\hat{\mathbf{w}}\|}$ .*

**Remark 3.** *Existing literature usually assume a time-decaying hyperparameter choice of  $\beta_1$  or  $\beta_2$  (c.f., (Kingma & Ba, 2014; Chen et al., 2018)). To the best of our knowledge, the only work analyzing Adam with constant  $\beta_1$  and  $\beta_2$  is (Reddi et al., 2019), which assumes  $\beta_2 > \beta_1^2$  (stronger than our assumption). Our result indicates that the ranges of  $\beta_1$  and  $\beta_2$  can be broader in hyper-parameter selection.*

We still start the proof sketch by proving the sum of squared gradients is finite in Stage I. Compared to GDM, the difference is that in Stage I, we will also prove the loss converges to 0 with rate  $\Theta(t^{-1})$ . This property will be used in Stage II to analyze the effect of preconditioner on implicit bias.

**Stage I: Bound the sum of squared gradients.** The next lemma can be viewed as an extension of Lemma 1 by taking the preconditioner into consideration, and characterizes the one-step loss update in Adam (**w/s**).

**Lemma 5.** *Let all conditions in Theorem 4 hold. Then, for any  $t \geq 1$ ,*

$$\begin{aligned} & \mathcal{L}(t+1) + \frac{1}{2} \frac{1 - \beta_1^t}{\eta(1 - \beta_1)} \left\| \sqrt[4]{\varepsilon \mathbb{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) \right\|^2 \\ & \leq \mathcal{L}(\mathbf{w}(t)) + \frac{1 - \beta_1^{t-1}}{2c\eta(1 - \beta_1)} \frac{1 - (c\beta_1)^t}{1 - (c\beta_1)^{t-1}} \left\| \sqrt[4]{\varepsilon \mathbb{1}_d + \hat{\nu}(t-1)} \odot (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\|^2. \end{aligned} \quad (10)$$

The difference between the proof of Lemma 5 and Lemma 1 is that we need to handle the gap between the preconditioners at step  $t$  and step  $t+1$ , this leads to a amplifying factor  $\sqrt[4]{\frac{1 - \beta_2^t}{\beta_2(1 - \beta_2^{t-1})}}$  of term  $\left\| \sqrt[4]{\varepsilon \mathbb{1}_d + \hat{\nu}(t-1)} \odot (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\|^2$ , which, however, is smaller than the shrinking factor  $\frac{1 - \beta_1^t}{\beta_2(1 - \beta_2^{t-1})}$  introduced by the gap between the coefficients of momentum



$\hat{\mathbf{m}}(t)$  and  $\hat{\mathbf{m}}(t-1)$ . As  $t \rightarrow \infty$ , both  $\beta_1^t \rightarrow 0$  and  $(c\beta_1)^t \rightarrow 0$ , and we can obtain that  $\xi(t) \triangleq \mathcal{L}(\mathbf{w}(t)) + \frac{1}{2\sqrt{c\eta}(1-\beta_1)} \left\| \sqrt[4]{\varepsilon \mathbb{1}_d + \hat{\nu}(t-1)} \odot (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\|^2$  is a Lyapunov function. Based on Lemma 5, we can prove the following asymptotic rate of  $\mathcal{L}$ .

**Lemma 6.** *Let all conditions in Theorem 4 hold. Then,*

$$\mathcal{L}(\mathbf{w}(t)) = \Theta(t^{-1}), \|\mathbf{w}(t)\| = \Theta(\log(t)), \text{ and } \|\mathbf{w}(t) - \mathbf{w}(t-1)\| = \Theta(t^{-1}).$$

The proof idea of Lemma 6 is that by the exponential-tailed property, when time is large enough, the training loss can be bounded by the gradient, and further bounded by  $\left\| \sqrt[4]{\varepsilon \mathbb{1}_d + \hat{\nu}(t-1)} \odot (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\|$ . Consequently,  $\xi(t+1) - \xi(t) \leq -\mathcal{O}(\xi(t))$ , and  $\xi(t) = \mathcal{O}(t^{-1})$ .

**Stage II. Bound the difference between  $\mathbf{w}(t)$  and  $\log(t)\hat{\mathbf{w}}$ .** The Lyapunov function for the direction difference is also modified according to the preconditioner, as introduced in Lemma 7:

**Lemma 7.** *Let all conditions in Theorem 4 hold. Then,  $\|\mathbf{r}(t)\|$  is bounded if and only if  $g(t)$  is upper bounded, where  $g(t)$  is defined as follows.*

$$\begin{aligned} g(t) \triangleq & \frac{\sqrt{\varepsilon}}{2} \|\mathbf{r}(t)\|^2 + \frac{\beta_1}{1-\beta_1} \left\langle \mathbf{r}(t), (1-\beta_1^{t-1}) \sqrt{\varepsilon \mathbb{1}_d + \hat{\nu}(t-1)} \odot (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\rangle \\ & - \frac{\beta_1}{1-\beta_1} \sum_{\tau=2}^t \left\langle \mathbf{r}(\tau) - \mathbf{r}(\tau-1), (1-\beta_1^{\tau-1}) \sqrt{\varepsilon \mathbb{1}_d + \hat{\nu}(\tau-1)} \odot (\mathbf{w}(\tau) - \mathbf{w}(\tau-1)) \right\rangle. \end{aligned}$$

Furthermore, we have  $\sum_{s=1}^{\infty} (g(t+1) - g(t))$  is upper bounded.

The first claim of Lemma 7 is similar to that in GDM and SGDM case. However, when we come to the second claim, simple calculation of  $g(t+1) - g(t)$  leads to

$$\left\langle \mathbf{r}(t), (\sqrt{\varepsilon} - (1-\beta_1^t) \sqrt{\varepsilon \mathbb{1}_d + \hat{\nu}(t)}) \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) \right\rangle + \left\langle \mathbf{r}(t), -\sqrt{\varepsilon} \log\left(\frac{t+1}{t}\right) \hat{\mathbf{w}} - \frac{\eta}{1-\beta} \nabla \mathcal{L}(\mathbf{w}(t)) \right\rangle.$$

This is where we need Lemma 6, which bounds the gap between  $\sqrt{\varepsilon \mathbb{1}_d + \hat{\nu}(t)}$  and  $\sqrt{\varepsilon}$  by  $\mathcal{O}(t^{-2})$  and makes  $\sum_{t=1}^{\infty} \left\langle \mathbf{r}(t), (\sqrt{\varepsilon} - (1-\beta_1^t) \sqrt{\varepsilon \mathbb{1}_d + \hat{\nu}(t)}) \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) \right\rangle$  an absolute continuous sequence. The second term  $\left\langle \mathbf{r}(t), -\sqrt{\varepsilon} \log\left(\frac{t+1}{t}\right) \hat{\mathbf{w}} - \frac{\eta}{1-\beta} \nabla \mathcal{L}(\mathbf{w}(t)) \right\rangle$  can be tackled by following the same routine as the GDM/SGDM case. Consequently, the proof of Lemma 7 is completed.

Till now, we have obtained the implicit bias for GDM, SGDM and Adam (w/s), one may wonder whether the implicit bias of stochastic Adam can be obtained. Here, we make some discussions here and put its further investigation for future works. First, our analysis can be extended to prove that the stochastic adaptive heavy-ball algorithm converges to the  $L^2$  max margin solution, which is proposed by (Tao et al., 2021) and has the following form<sup>3</sup>:

$$\mathbf{w}(t+1) = \mathbf{w}(t) - \eta \frac{\nabla \mathcal{L}_{\mathbf{B}(t)}(\mathbf{w}(t))}{\sqrt{\varepsilon \mathbb{1}_d + \hat{\nu}(t)}} + \beta(\mathbf{w}(t) - \mathbf{w}(t-1)), t \geq 1,$$

where  $\hat{\nu}(t)$  is defined as Eq. (2). The proof is a simple combination of the proof techniques in SGDM and Adam by observing that the stochastic adaptive heavy-ball algorithm has the following equivalent update rule ( $\mathbf{u}(1) = \mathbf{w}(1)$ )

$$\begin{aligned} \mathbf{u}(t+1) &= -\frac{\eta}{1-\beta} \frac{\nabla \mathcal{L}_{\mathbf{B}(t)}(\mathbf{w}(t))}{\sqrt{\varepsilon \mathbb{1}_d + \hat{\nu}(t)}} + \mathbf{u}(t), \\ \mathbf{w}(t+1) &= \beta \mathbf{w}(t) + (1-\beta) \mathbf{u}(t+1). \end{aligned}$$

However, Adam is different from the stochastic adaptive heavy-ball algorithm as it combine the conditioner and momentum as a whole, which makes our constructed Lyapunov function can not be applied. We will put further investigation on stochastic Adam in future work.

<sup>3</sup>Compared to (Tao et al., 2021), we put the  $\varepsilon \mathbb{1}_d$  inside the square-root and use  $\hat{\nu}(t)$  instead of  $\nu(t)$  to demonstrate the difference between Adaptive Heavy Ball and Adam, while these changes will not influence the convergent direction.

## 7 DISCUSSIONS

**Extension to the Multi-Class Classification Problem.** As mentioned in the introduction, despite all the previous analyses are aimed at the binary classification problem, they can be naturally extended to the analyses multi-class classification problem. Specifically, in the linear multi-class classification problem, for any  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{d_x} \times \{1, \dots, C\}$  in the sample set  $\mathcal{S}$ , the (individual) logistic loss with parameter  $\mathbf{W} \in \mathbb{R}^{C \times d_x}$  is denoted as

$$\ell(\mathbf{y}, \mathbf{W}\mathbf{x}) = \log \frac{e^{\mathbf{W}_{\mathbf{y},\mathbf{x}}}}{\sum_{i=1}^C e^{\mathbf{W}_{i,\mathbf{x}}}}.$$

Correspondingly, dataset  $\mathcal{S}$  is separable if there exists a parameter  $\mathbf{W}$ , such that  $\forall (\mathbf{x}, \mathbf{y}) \in \mathcal{S}$ , we have  $\mathbf{W}_{\mathbf{y},\mathbf{x}} > \mathbf{W}_{i,\mathbf{x}}, \forall i \neq \mathbf{y}$ . The multi-class  $L^2$  max-margin problem is then defined as

$$\min \|\mathbf{W}\|_F, \text{ subject to : } \mathbf{W}_{\mathbf{y},\mathbf{x}} \geq \mathbf{W}_{i,\mathbf{x}} + 1, \forall (\mathbf{x}, \mathbf{y}) \in \mathcal{S}, i \neq \mathbf{y},$$

where  $\|\cdot\|_F$  denotes the Frobenius norm. Denote  $\hat{\mathbf{W}}$  as the  $L^2$  max-margin solution, we have (mini-batch) SGDM and Adam (w/s) still converges to the direction of  $\hat{\mathbf{W}}$ , the proof of which is deferred to the Appendix D.2.

**Theorem 5.** *For linear multi-class classification problem using logistic loss and almost every separable data, with a small enough learning rate, and  $1 > \beta_2 > \beta_1^4 \geq 0$  (for Adam (w/s)), (mini-batch) SGDM and Adam (w/s) converge to the multi-class  $L^2$  max-margin solution (a.s. for SGDM).*

**Consistency with the Existing Experimental Results.** Our results stand with existing experimental observations. Specifically, Soudry et al. (2018) conduct experiments using GD and GDM on the same linear separable data (c.f., Figure 1 in (Soudry et al., 2018)), and it is observed that the training behaviors of GD and GDM are quite similar in terms of the direction  $w(t)$ , the training loss, and the margin, which supports our Theorem 2. Nacson et al. (2019) extend the experiment to the stochastic setting (c.f. Figure 2 in (Nacson et al., 2019)), and observe the same similarity between SGD and SGDM, which agrees with our Theorem 3. Wang et al. (2021) conduct the experiments of SGD, SGDM, Adam (without momentum) and Adam on MNIST using homogeneous neural networks (c.f. Appendix F.1.2 in (Wang et al., 2021)), and observe such similarity for deep neural networks. Our theorems apply to linear models, which is a special case of homogeneous neural network, and meet their observation.

**Gap Between The Linear Model and Deep Neural Networks.** While our results only hold for the linear classification problem, extending the results to the deep neural networks is possible. Specifically, Lyu & Li (2019) construct a Lyapunov function bounding the rate between loss and parameter norm of GD on the homogeneous deep neural networks, and prove the parameter direction converges to some KKT point of the  $L^2$  max-margin problem, which extends the result of GD on linear model. Wang et al. (2021) further show the proof techniques in (Lyu & Li, 2019) can be generalized, and successfully extend the Lyapunov function to AdaGrad, RMSProp and Adam (without momentum) on deep homogeneous neural networks. It will be interesting to see whether such a Lyapunov function can be constructed for GDM and Adam in homogeneous neural networks.

## 8 CONCLUSION

In this paper, we study the implicit bias of momentum-based optimizers in linear classification with exponential-tailed loss. Our results indicates that for SGD and Adam, adding momentum will not influence the implicit bias, and the direction of parameter converges to the  $L^2$  max-margin solution. Our theoretical results stands with existing experimental observations.

## REFERENCES

Alnur Ali, Edgar Dobriban, and Ryan Tibshirani. The implicit regularization of stochastic gradient flow for least squares. In *International Conference on Machine Learning*, pp. 233–244. PMLR, 2020.

- Xiangyi Chen, Sijia Liu, Ruoyu Sun, and Mingyi Hong. On the convergence of a class of adam-type algorithms for non-convex optimization. In *International Conference on Learning Representations*, 2018.
- Euhanna Ghadimi, Hamid Reza Feyzmahdavian, and Mikael Johansson. Global convergence of the heavy-ball method for convex optimization. In *2015 European control conference (ECC)*, pp. 310–315. IEEE, 2015.
- Suriya Gunasekar, Jason Lee, Daniel Soudry, and Nathan Srebro. Characterizing implicit bias in terms of optimization geometry. In *International Conference on Machine Learning*, pp. 1832–1841. PMLR, 2018a.
- Suriya Gunasekar, Jason D Lee, Daniel Soudry, and Nati Srebro. Implicit bias of gradient descent on linear convolutional networks. In *Advances in Neural Information Processing Systems*, pp. 9461–9471, 2018b.
- Ziwei Ji and Matus Telgarsky. Gradient descent aligns the layers of deep linear networks. *arXiv preprint arXiv:1810.02032*, 2018.
- Ziwei Ji and Matus Telgarsky. Directional convergence and alignment in deep learning. *arXiv preprint arXiv:2006.06657*, 2020.
- Ziwei Ji and Matus Telgarsky. Characterizing the implicit bias via a primal-dual analysis. In *Algorithmic Learning Theory*, pp. 772–804. PMLR, 2021.
- Ziwei Ji, Nathan Srebro, and Matus Telgarsky. Fast margin maximization via dual acceleration. In *International Conference on Machine Learning*, pp. 4860–4869. PMLR, 2021.
- Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic generalization measures and where to find them. *arXiv preprint arXiv:1912.02178*, 2019.
- Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. In *International Conference on Learning Representations*, 2017.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Junhong Lin and Lorenzo Rosasco. Optimal rates for multi-pass stochastic gradient methods. *The Journal of Machine Learning Research*, 18(1):3375–3421, 2017.
- Yanli Liu, Yuan Gao, and Wotao Yin. An improved analysis of stochastic gradient descent with momentum. *Advances in Neural Information Processing Systems*, 33, 2020.
- Kaifeng Lyu and Jian Li. Gradient descent maximizes the margin of homogeneous neural networks. *arXiv preprint arXiv:1906.05890*, 2019.
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*, 2018.
- Mor Shpigel Nacson, Nathan Srebro, and Daniel Soudry. Stochastic gradient descent on separable data: Exact convergence with a fixed learning rate. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 3051–3059. PMLR, 2019.
- Y Nesterov. A method for solving a convex programming problem with convergence rate  $o(1/k^2)$ . In *Soviet Mathematics. Doklady*, volume 27, pp. 367–372, 1983.
- Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. In search of the real inductive bias: On the role of implicit regularization in deep learning. In *ICLR*, 2015.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32: 8026–8037, 2019.

- Boris T Polyak. Some methods of speeding up the convergence of iteration methods. *Ussr computational mathematics and mathematical physics*, 4(5):1–17, 1964.
- Qian Qian and Xiaoyuan Qian. The implicit bias of adagrad on separable data. In *Advances in Neural Information Processing Systems*, pp. 7761–7769, 2019.
- Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. *arXiv preprint arXiv:1904.09237*, 2019.
- Lorenzo Rosasco and Silvia Villa. Learning with incremental iterative regularization. In *Proceedings of the 28th International Conference on Neural Information Processing Systems-Volume 1*, pp. 1630–1638, 2015.
- Stefano Sarao Mannelli and Pierfrancesco Urbani. Analytical study of momentum-based acceleration methods in paradigmatic high-dimensional non-convex problems. *Advances in Neural Information Processing Systems*, 34, 2021.
- Robert E Schapire and Yoav Freund. Boosting: Foundations and algorithms. *Kybernetes*, 2013.
- Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit bias of gradient descent on separable data. *The Journal of Machine Learning Research*, 19(1): 2822–2878, 2018.
- Tao Sun, Penghang Yin, Dongsheng Li, Chun Huang, Lei Guan, and Hao Jiang. Non-ergodic convergence analysis of heavy-ball algorithms. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 5033–5040, 2019.
- Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initialization and momentum in deep learning. In *International conference on machine learning*, pp. 1139–1147. PMLR, 2013.
- Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International Conference on Machine Learning*, pp. 6105–6114. PMLR, 2019.
- Wei Tao, Sheng Long, Gaowei Wu, and Qing Tao. The role of momentum parameters in the optimal convergence of adaptive polyak’s heavy-ball methods. *arXiv preprint arXiv:2102.07314*, 2021.
- Matus Telgarsky. Margins, shrinkage, and boosting. In *International Conference on Machine Learning*, pp. 307–315. PMLR, 2013.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pp. 5998–6008, 2017.
- Bohan Wang, Qi Meng, Wei Chen, and Tie-Yan Liu. The implicit bias for adaptive optimization algorithms on homogeneous neural networks. In *International Conference on Machine Learning*, pp. 10849–10858. PMLR, 2021.
- Ashia C Wilson, Rebecca Roelofs, Mitchell Stern, Nati Srebro, and Benjamin Recht. The marginal value of adaptive gradient methods in machine learning. In *Advances in neural information processing systems*, pp. 4148–4158, 2017.
- Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1492–1500, 2017.
- Y Yan, T Yang, Z Li, Q Lin, and Y Yang. A unified analysis of stochastic momentum methods for deep learning. In *IJCAI International Joint Conference on Artificial Intelligence*, 2018.
- Hao Yu, Rong Jin, and Sen Yang. On the linear speedup analysis of communication efficient momentum sgd for distributed non-convex optimization. In *International Conference on Machine Learning*, pp. 7184–7193. PMLR, 2019.

## Supplementary materials for “Does momentum Change the Implicit Bias on separable data?”

### A PREPARATIONS

This section collect definitions and lemmas which will be used throughout the proofs.

#### A.1 CHARACTERIZATION OF THE MAX-MARGIN SOLUTION

This section collects several commonly-used characterization of the max-margin solution from Nacson et al. (2019) and Soudry et al. (2018).

To start with, we define support vectors and support set, which are two common terms in margin analysis.

**Definition 2** (Support vectors and support set). *For any  $i \in [N]$ ,  $\tilde{\mathbf{x}}_i$  is called a support vector of the dataset  $\mathcal{S}$ , if*

$$\langle \tilde{\mathbf{x}}_i, \hat{\mathbf{w}} \rangle = 1.$$

*Correspondingly,  $\tilde{\mathbf{x}}_i$  is called a non-support vector if  $\langle \tilde{\mathbf{x}}_i, \hat{\mathbf{w}} \rangle > 1$ . The support set of  $\mathcal{S}$  is then defined as*

$$\mathcal{S}_s = \{(\mathbf{x}_i, \mathbf{y}_i) : \langle \mathbf{y}_i \mathbf{x}_i, \hat{\mathbf{w}} \rangle = 1\}.$$

The following lemma delivers  $\hat{\mathbf{w}}$  as an linear combination of support vectors.

**Lemma 8** (Lemma 12, Soudry et al. (2018)). *For almost every datasets  $\mathcal{S}$ , there exists a unique vector  $\mathbf{v} = (v_1, \dots, v_N)$ , such that  $\hat{\mathbf{w}}$  can be represented as*

$$\hat{\mathbf{w}} = \sum_{i=1}^N v_i \tilde{\mathbf{x}}_i, \quad (11)$$

where  $\mathbf{v}$  satisfies  $v_i = 0$  if  $\tilde{\mathbf{x}}_i \notin \mathcal{T}(\mathcal{S}_s)$ , and  $v_i > 0$  if  $\tilde{\mathbf{x}}_i \in \mathcal{T}(\mathcal{S}_s)$ . Furthermore, the size of  $\tilde{\mathcal{S}}_s$  is at most  $d$ .

By Lemma 8, we further have the following corollary:

**Corollary 2.** *For almost every datasets  $\mathcal{S}$ , the unique  $\mathbf{v}$  given by Lemma 8 further satisfies that for any positive constant  $C_2$ , there exists a non-zero vector  $\tilde{\mathbf{w}}$ , such that,  $\tilde{\mathbf{x}}_i \in \mathcal{T}(\mathcal{S}_s)$ , we have*

$$C_2 e^{-\langle \tilde{\mathbf{x}}_i, \tilde{\mathbf{w}} \rangle} = v_i. \quad (12)$$

*Proof.* For almost every datasets  $\mathcal{S}$ , any subsets with size  $d$  of  $\mathcal{S}$  is linearly independent. Since  $\tilde{\mathcal{S}}_s$  has size no larger than  $d$  (by Lemma 8), and Eq. (12) is equivalent to linear equations, the proof is completed.  $\square$

For the stochastic case, we will also need the following lemma when we calculate the form of parameter at time  $t$ .

**Lemma 9** (Lemma 5, Nacson et al. (2019)). *Let  $\mathcal{B}(s)$  be the random subset used in SGD (i.e., the one used in SGDM). Almost surely, there exists a vector  $\tilde{\mathbf{w}}$*

$$\frac{N}{b} \sum_{s=1}^{t-1} \frac{1}{s} \sum_{i: \tilde{\mathbf{x}}_i \in \mathcal{T}(\mathcal{B}(s) \cap \mathcal{S}_s)} v_i \tilde{\mathbf{x}}_i = \log \left( \frac{bt}{N} \right) \hat{\mathbf{w}} + \mathbf{n}(t),$$

where  $\mathbf{n}(t)$  satisfies  $\|\mathbf{n}(t)\| = \mathcal{O}(t^{-0.5+\varepsilon})$  for any  $\varepsilon > 0$ , and  $\|\mathbf{n}(t+1) - \mathbf{n}(t)\| = \mathcal{O}(t^{-1})$ . If the  $\mathcal{B}(s)$  is the random subset used in SGDM (i.e., the one used in mini-batch SGDM), then the a.s. condition can be removed.

## A.2 PREPARATIONS OF THE OPTIMIZATION ANALYSIS

This section collects technical lemmas which will be used in latter proofs. We begin with a lemma bounding the smooth constants if the loss is bounded.

**Lemma 10.** *If loss  $\ell$  satisfies (D) in Assumption 3, then for any  $\mathbf{w}_0$ , if  $\mathcal{L}(\mathbf{w}) \leq \mathcal{L}(\mathbf{w}_0)$ , then we have  $\mathcal{L}$  is  $H_{s_0}$  smooth at point  $\mathbf{w}$ , where  $s_0 = \ell^{-1}(N\mathcal{L}(\mathbf{w}_0))$ . Furthermore,  $\mathcal{L}$  is globally  $H_{s_0}$  smooth over the set  $\{\mathbf{w} : \mathcal{L}(\mathbf{w}) \leq \mathcal{L}(\mathbf{w}_0)\}$ .*

*Proof.* Since  $\ell$  is positive, we have  $\forall i \in [N]$ ,

$$\frac{\tilde{\ell}(\mathbf{w}, \mathbf{z}_i)}{N} < \frac{\sum_{j=1}^N \tilde{\ell}(\mathbf{w}, \mathbf{z}_j)}{N} = \mathcal{L}(\mathbf{w}) \leq \mathcal{L}(\mathbf{w}_0),$$

which leads to  $\tilde{\ell}(\mathbf{w}, \mathbf{z}_i) < N\mathcal{L}(\mathbf{w}_0)$ , and  $\ell$  is  $H_{s_0}$  smooth at  $\langle \mathbf{w}, \tilde{\mathbf{x}}_i \rangle$ .

Furthermore, since  $\nabla_{\mathbf{w}} \tilde{\ell}(\mathbf{w}, \mathbf{z}_i) = \nabla_{\mathbf{w}} \ell(\langle \mathbf{w}, \tilde{\mathbf{x}}_i \rangle) = \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}}_i \rangle) \tilde{\mathbf{x}}_i$ , for any two parameters  $\mathbf{w}_1$  and  $\mathbf{w}_2$  close enough to  $\mathbf{w}$ ,

$$\begin{aligned} \|\nabla_{\mathbf{w}} \tilde{\ell}(\mathbf{w}_1, \mathbf{z}_i) - \nabla_{\mathbf{w}} \tilde{\ell}(\mathbf{w}_2, \mathbf{z}_i)\| &= \|(\ell'(\langle \mathbf{w}_1, \tilde{\mathbf{x}}_i \rangle) - \ell'(\langle \mathbf{w}_2, \tilde{\mathbf{x}}_i \rangle)) \tilde{\mathbf{x}}_i\| \\ &\leq |\ell'(\langle \mathbf{w}_1, \tilde{\mathbf{x}}_i \rangle) - \ell'(\langle \mathbf{w}_2, \tilde{\mathbf{x}}_i \rangle)| \leq H_{s_0} |\langle \mathbf{w}_1 - \mathbf{w}_2, \tilde{\mathbf{x}}_i \rangle| \leq H_{s_0} \|\mathbf{w}_1 - \mathbf{w}_2\|, \end{aligned}$$

and thus,

$$\|\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}_1) - \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}_2)\| \leq \frac{1}{N} \sum_{i=1}^N \|\nabla_{\mathbf{w}} \tilde{\ell}(\mathbf{w}_1, \mathbf{z}_i) - \nabla_{\mathbf{w}} \tilde{\ell}(\mathbf{w}_2, \mathbf{z}_i)\| \leq H_{s_0} \|\mathbf{w}_1 - \mathbf{w}_2\|,$$

which completes the proof that  $\mathcal{L}$  is locally  $H_{s_0}$  smooth at  $\mathbf{w}$ .

Now if  $\mathbf{w}_1$  and  $\mathbf{w}_2$  both belong to  $\{\mathbf{w} : \mathcal{L}(\mathbf{w}) \leq \mathcal{L}(\mathbf{w}_0)\}$ , we have for any  $\tilde{\mathbf{x}}_i \in \mathcal{T}(\mathcal{S})$ ,  $\langle \mathbf{w}_1, \tilde{\mathbf{x}}_i \rangle > \ell^{-1}(N\mathcal{L}(\mathbf{w}_0))$ , and  $\langle \mathbf{w}_2, \tilde{\mathbf{x}}_i \rangle > \ell^{-1}(N\mathcal{L}(\mathbf{w}_0))$ . Following the same routine as the locally smooth proof, we complete the second argument.

The proof is completed.  $\square$

Based on Assumption 2, we also have the following lemma characterizing the relationship between loss  $\ell$  and its derivative  $\ell'$  when  $x$  is large enough.

**Lemma 11.** *Let loss  $\ell$  satisfy Assumption 2. Then, there exists an large enough  $x_0$  and a positive real  $K$ , such that,  $\forall x > x_0$ , we have*

$$-\frac{1}{K} \ell'(x) \leq \ell(x) \leq -K \ell'(x).$$

*Proof.* By Assumption 2, there exists a large enough  $x_0$ , such that  $\forall x > x_0$ , we have

$$\frac{1}{2} e^{-x} \leq -\ell'(x) \leq 2e^{-x}. \quad (13)$$

On the other hand, as  $\lim_{t \rightarrow \infty} \ell(x) = 0$ , we have

$$\ell(x) = \int_{s=x}^{\infty} -\ell'(s) ds,$$

which by Eq. (13) leads to

$$\frac{1}{2} e^{-x} = \frac{1}{2} \int_x^{\infty} e^{-s} ds \leq \ell(x) \leq 2 \int_x^{\infty} e^{-s} ds = 2e^{-x}.$$

Therefore, setting  $K = 4$  completes the proof.  $\square$

The following lemma bridges the second moment of  $\nabla \mathcal{L}_{\mathcal{B}(t)}$  with its squared first moment.

**Lemma 12.** *Let the dataset  $\mathcal{S}$  satisfies the separable assumption 1. Let  $\mathcal{B}$  be a random subset of  $\mathcal{S}$  with size  $b$  sampled independently and uniformly without replacement. Then, at any point  $\mathbf{w}$ , we have*

$$\|\nabla\mathcal{L}(\mathbf{w})\|^2 \leq \mathbb{E}_{\mathcal{B}} [\|\nabla\mathcal{L}_{\mathcal{B}}(\mathbf{w})\|^2] \leq \frac{N}{\gamma^2 b} \|\nabla\mathcal{L}(\mathbf{w})\|^2.$$

*Proof.* To start with, notice that

$$\|\nabla\mathcal{L}(\mathbf{w})\| = \|\mathbb{E}_{\mathcal{B}}\mathcal{L}_{\mathcal{B}}(\mathbf{w})\| \leq \mathbb{E}_{\mathcal{B}}\|\mathcal{L}_{\mathcal{B}}(\mathbf{w})\|.$$

Therefore, the first inequality can be directly obtained by Cauchy-Schwartz's inequality. To prove the second inequality, we first calculate the explicit form of  $\nabla\mathcal{L}_{\mathcal{B}}(\mathbf{w})$ .

$$\begin{aligned} \|\nabla\mathcal{L}_{\mathcal{B}}(\mathbf{w})\|^2 &= \frac{1}{b^2} \left\| \nabla \sum_{\mathbf{z} \in \mathcal{B}} \tilde{\ell}(\mathbf{w}, \mathbf{z}) \right\|^2 \\ &= \frac{1}{b^2} \left\| \sum_{\tilde{\mathbf{x}} \in \mathcal{T}(\mathcal{B})} \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}} \rangle) \tilde{\mathbf{x}} \right\|^2 = \frac{1}{b^2} \sum_{\tilde{\mathbf{x}}, \tilde{\mathbf{x}}' \in \mathcal{T}(\mathcal{B})} \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}} \rangle) \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}}' \rangle) \langle \tilde{\mathbf{x}}, \tilde{\mathbf{x}}' \rangle \\ &\leq \frac{1}{b^2} \sum_{\tilde{\mathbf{x}}, \tilde{\mathbf{x}}' \in \mathcal{T}(\mathcal{B})} \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}} \rangle) \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}}' \rangle). \end{aligned}$$

Therefore,

$$\begin{aligned} &\mathbb{E}_{\mathcal{B}} \|\nabla\mathcal{L}_{\mathcal{B}}(\mathbf{w})\|^2 \\ &\leq \mathbb{E}_{\mathcal{B}} \frac{1}{b^2} \sum_{\tilde{\mathbf{x}}, \tilde{\mathbf{x}}' \in \mathcal{T}(\mathcal{B})} \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}} \rangle) \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}}' \rangle) \end{aligned} \quad (14)$$

$$\begin{aligned} &= \mathbb{E}_{\mathcal{B}} \frac{1}{b^2} \sum_{\tilde{\mathbf{x}}, \tilde{\mathbf{x}}' \in \mathcal{T}(\mathcal{S})} \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}} \rangle) \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}}' \rangle) \mathbb{1}_{\tilde{\mathbf{x}}, \tilde{\mathbf{x}}' \in \mathcal{T}(\mathcal{B})} \\ &= \frac{1}{b^2} \sum_{\tilde{\mathbf{x}}, \tilde{\mathbf{x}}' \in \mathcal{T}(\mathcal{S})} \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}} \rangle) \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}}' \rangle) \mathbb{E}_{\mathcal{B}} \mathbb{1}_{\tilde{\mathbf{x}}, \tilde{\mathbf{x}}' \in \mathcal{T}(\mathcal{B})} \\ &= \frac{1}{Nb} \sum_{\tilde{\mathbf{x}} \in \mathcal{T}(\mathcal{S})} \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}} \rangle)^2 + \frac{b-1}{bN(N-1)} \sum_{\tilde{\mathbf{x}}, \tilde{\mathbf{x}}' \in \mathcal{T}(\mathcal{S}), \tilde{\mathbf{x}} \neq \tilde{\mathbf{x}}'} \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}} \rangle) \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}}' \rangle) \\ &\leq \frac{1}{Nb} \left( \sum_{\tilde{\mathbf{x}} \in \mathcal{T}(\mathcal{S})} \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}} \rangle) \right)^2. \end{aligned} \quad (15)$$

On the other hand,

$$\begin{aligned} \|\nabla\mathcal{L}(\mathbf{w})\| &= \frac{1}{N} \left\| \sum_{\tilde{\mathbf{x}} \in \mathcal{T}(\mathcal{S})} \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}} \rangle) \tilde{\mathbf{x}} \right\| \\ &\geq \frac{1}{N} \left\langle \sum_{\tilde{\mathbf{x}} \in \mathcal{T}(\mathcal{S})} \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}} \rangle) \tilde{\mathbf{x}}, -\frac{\hat{\mathbf{w}}}{\|\hat{\mathbf{w}}\|} \right\rangle \stackrel{(*)}{\geq} \frac{\gamma}{N} \sum_{\tilde{\mathbf{x}} \in \mathcal{T}(\mathcal{S})} \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}} \rangle) \end{aligned}$$

where Eq. (\*) is due to  $\forall \mathbf{z} \in \langle \tilde{\mathbf{x}}, -\hat{\mathbf{w}} \rangle \geq 1$  and  $\ell' < 0$ .

Therefore,

$$\|\nabla\mathcal{L}(\mathbf{w})\|^2 \geq \frac{\gamma^2}{N^2} \left( \sum_{\tilde{\mathbf{x}} \in \mathcal{T}(\mathcal{S})} \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}} \rangle) \right)^2. \quad (16)$$

The proof is completed by putting Eqs. (15) and (16) together.  $\square$

In the following lemma, we show the updates of GDM, Adam, and SGDM are all non-zero.

**Lemma 13.** *Regardless of GDM, Adam, or SGDM, the updates of all steps are non-zero, i.e.,*

$$\|\mathbf{w}(t+1) - \mathbf{w}(t)\| > 0, \forall t > 1.$$

*Proof.* We start with the alternative forms of the update rule of GDM, Adam, and SGDM using the gradients along the trajectory respectively. For GDM, by Eq. (5), the update rule can be written as

$$\mathbf{w}(t+1) - \mathbf{w}(t) = -\eta \left( \sum_{s=1}^t \beta^{t-s} \nabla \mathcal{L}(\mathbf{w}(s)) \right). \quad (17)$$

Similarly, the update rule of SGDM can be written as

$$\mathbf{w}(t+1) - \mathbf{w}(t) = -\eta \left( \sum_{s=1}^t \beta^{t-s} \nabla \mathcal{L}_{\mathbf{B}(s)}(\mathbf{w}(s)) \right), \quad (18)$$

while the update rule of Adam can be given as

$$\mathbf{w}(t+1) - \mathbf{w}(t) = -\eta \frac{\sum_{s=1}^t \frac{1-\beta_1}{1-\beta_1^s} \beta_1^{t-s} \nabla \mathcal{L}(\mathbf{w}(s))}{\sqrt{\varepsilon \mathbf{1}_d + \sum_{s=1}^t \frac{1-\beta_2}{1-\beta_2^s} \beta_2^{t-s} (\nabla \mathcal{L}(\mathbf{w}(s)))^2}}. \quad (19)$$

On the other hand, by the definition of empirical risk  $\mathcal{L}$ , the gradient of  $\mathcal{L}$  at point  $\mathbf{w}$  can be given as

$$\nabla \mathcal{L}(\mathbf{w}) = \frac{\sum_{i=1}^N \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}}_i \rangle) \tilde{\mathbf{x}}_i}{N}. \quad (20)$$

By Eq. (20) and Eq. (17), we further have for GDM,

$$\mathbf{w}(t+1) - \mathbf{w}(t) = -\eta \left( \sum_{s=1}^t \beta^{t-s} \frac{\sum_{i=1}^N \ell'(\langle \mathbf{w}(s), \tilde{\mathbf{x}}_i \rangle) \tilde{\mathbf{x}}_i}{N} \right). \quad (21)$$

By Assumption 1, there exists a non-zero parameter  $\hat{\mathbf{w}}$ , such that,  $\langle \hat{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle > 0, \forall i$ . Therefore, by executing inner product between Eq. (21) and  $\hat{\mathbf{w}}$ , we have

$$\|\mathbf{w}(t+1) - \mathbf{w}(t)\| \|\hat{\mathbf{w}}\| \geq \langle \mathbf{w}(t+1) - \mathbf{w}(t), \hat{\mathbf{w}} \rangle = -\eta \left( \sum_{s=1}^t \beta^{t-s} \frac{\sum_{i=1}^N \ell'(\langle \mathbf{w}(s), \tilde{\mathbf{x}}_i \rangle) \langle \tilde{\mathbf{x}}_i, \hat{\mathbf{w}} \rangle}{N} \right) \stackrel{(*)}{>} 0,$$

where Eq. (\*) is due to  $\ell' < 0$ . This complete the proof for GDM.

Similarly, for SGDM, we have

$$\|\mathbf{w}(t+1) - \mathbf{w}(t)\| \|\hat{\mathbf{w}}\| \geq -\eta \left( \sum_{s=1}^t \beta^{t-s} \frac{\sum_{(\mathbf{x}, \mathbf{y}) \in \mathbf{B}} \ell'(\langle \mathbf{w}(s), \mathbf{y} \mathbf{x} \rangle) \langle \mathbf{y} \mathbf{x}, \hat{\mathbf{w}} \rangle}{b} \right) > 0,$$

which completes the proof of SGDM.

For Adam, we have

$$\begin{aligned} & \|\mathbf{w}(t+1) - \mathbf{w}(t)\| \left\| \hat{\mathbf{w}} \odot \sqrt{\varepsilon \mathbf{1}_d + \sum_{s=1}^t \frac{1-\beta_2}{1-\beta_2^s} \beta_2^{t-s} (\nabla \mathcal{L}(\mathbf{w}(s)))^2} \right\| \\ & \geq - \left\langle \hat{\mathbf{w}} \odot \sqrt{\varepsilon \mathbf{1}_d + \sum_{s=1}^t \frac{1-\beta_2}{1-\beta_2^s} \beta_2^{t-s} (\nabla \mathcal{L}(\mathbf{w}(s)))^2}, \eta \frac{\sum_{s=1}^t \frac{1-\beta_1}{1-\beta_1^s} \beta_1^{t-s} \nabla \mathcal{L}(\mathbf{w}(s))}{\sqrt{\varepsilon \mathbf{1}_d + \sum_{s=1}^t \frac{1-\beta_2}{1-\beta_2^s} \beta_2^{t-s} (\nabla \mathcal{L}(\mathbf{w}(s)))^2}} \right\rangle \\ & = \left\langle \hat{\mathbf{w}}, \eta \sum_{s=1}^t \frac{1-\beta_1}{1-\beta_1^s} \beta_1^{t-s} \nabla \mathcal{L}(\mathbf{w}(s)) \right\rangle \\ & = -\eta \left( \sum_{s=1}^t \frac{1-\beta_1}{1-\beta_1^s} \beta_1^{t-s} \frac{\sum_{i=1}^N \ell'(\langle \mathbf{w}(s), \tilde{\mathbf{x}}_i \rangle) \langle \tilde{\mathbf{x}}_i, \hat{\mathbf{w}} \rangle}{N} \right) > 0, \end{aligned}$$

which completes the proof of Adam.

The proof is completed.  $\square$



## B IMPLICIT BIAS OF GD/SGD WITH MOMENTUM

This section collects the proof of the implicit bias of gradient descent with momentum and stochastic gradient descent with momentum.

### B.1 IMPLICIT BIAS OF GD WITH MOMENTUM

This section collects the proof of Theorem 2.

#### B.1.1 PROOF OF THE SUM OF SQUARED GRADIENTS CONVERGES

To begin with, we will prove the sum of squared norm of gradients along the trajectory is finite for gradient descent with momentum. To see this, we first define the continuous-time update rule as

$$\mathbf{w}(t + \alpha) - \mathbf{w}(t) = \alpha(\mathbf{w}(t + 1) - \mathbf{w}(t)), \forall t \in \mathbb{Z}^+, \forall \alpha \in [0, 1].$$

We then prove a generalized case of Lemma 1 for any  $\mathbf{w}(t + \alpha)$ .

**Lemma 14** (Lemma 1, extended). *Let all conditions in Theorem 2 hold. We then have*

$$\begin{aligned} \mathcal{L}(\mathbf{w}(t)) + \frac{\beta}{2\eta} \|\mathbf{w}(t) - \mathbf{w}(t-1)\|^2 \geq & \mathcal{L}(\mathbf{w}(t + \alpha)) + \frac{\beta}{2\eta} \alpha^2 \|\mathbf{w}(t + 1) - \mathbf{w}(t)\|^2 \\ & + \frac{(1 - \beta)(1 - C_1)\alpha^2}{\eta} \|\mathbf{w}(t + 1) - \mathbf{w}(t)\|^2, \end{aligned}$$

where  $C_1$  is a positive real such that  $\eta = 2\frac{1-\beta}{H_{s_0}}C_1$ .

*Proof of Lemma 14.* For brevity, we denote  $s_0 \triangleq \ell^{-1}(N\mathcal{L}(\mathbf{w}_1))$ . We prove this lemma by reduction to absurdity.

Concretely, let  $t^*$  be the smallest positive integer time such that there exists an  $\alpha \in [0, 1]$ , such that Eq. (7) doesn't hold. Let  $\alpha^* = \inf\{\alpha \in [0, 1] : \text{Eq. (7) doesn't hold for } (t^*, \alpha)\}$ . By continuity, Eq. (7) holds for  $(t^*, \alpha^*)$ .

We further divide the proof into two cases depending on the value of  $\alpha^*$ .

**Case 1:**  $\alpha^* = 0$ : For any  $t^* > t \geq 1$ , we have Eq. (7) holds for  $(t, 1)$ . Specifically, we have

$$\mathcal{L}(\mathbf{w}(t)) + \frac{\beta}{2\eta} \|\mathbf{w}(t) - \mathbf{w}(t-1)\|^2 \geq \mathcal{L}(\mathbf{w}(t+1)) + \frac{\beta}{2\eta} \|\mathbf{w}(t+1) - \mathbf{w}(t)\|^2,$$

which further leads to

$$\mathcal{L}(\mathbf{w}(1)) = \mathcal{L}(\mathbf{w}(1)) + \frac{\beta}{2\eta} \|\mathbf{w}(1) - \mathbf{w}(0)\|^2 \geq \mathcal{L}(\mathbf{w}(t^*)) + \frac{\beta}{2\eta} \|\mathbf{w}(t^*) - \mathbf{w}(t^* - 1)\|^2.$$

Since  $\frac{\beta}{2\eta} \|\mathbf{w}(t^*) - \mathbf{w}(t^* - 1)\|^2$  is non-negative, we have

$$\mathcal{L}(\mathbf{w}(1)) \geq \mathcal{L}(\mathbf{w}(t^*)).$$

By Lemma 10, we have  $\mathcal{L}$  is  $H_{s_0}$  smooth at  $\mathbf{w}(t^*)$ . Therefore, by Taylor's expansion for  $\mathcal{L}$  at point  $\mathbf{w}(t^*)$ , we have for small enough  $\alpha > 0$

$$\begin{aligned} & \mathcal{L}(\mathbf{w}(t^* + \alpha)) \\ \leq & \mathcal{L}(\mathbf{w}(t^*)) + \langle \nabla \mathcal{L}(\mathbf{w}(t^*)), \mathbf{w}(t^* + \alpha) - \mathbf{w}(t^*) \rangle + \frac{H_{s_0}}{2} \|\mathbf{w}(t^* + \alpha) - \mathbf{w}(t^*)\|^2 \\ = & \mathcal{L}(\mathbf{w}(t^*)) + \alpha \langle \nabla \mathcal{L}(\mathbf{w}(t^*)), \mathbf{w}(t^* + 1) - \mathbf{w}(t^*) \rangle + \frac{H_{s_0}\alpha^2}{2} \|\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)\|^2 \\ \stackrel{(*)}{=} & \mathcal{L}(\mathbf{w}(t^*)) + \alpha \left\langle \frac{1}{\eta} (\beta(\mathbf{w}(t^*) - \mathbf{w}(t^* - 1)) - (\mathbf{w}(t^* + 1) - \mathbf{w}(t^*))), \mathbf{w}(t^* + 1) - \mathbf{w}(t^*) \right\rangle \\ & + \frac{H_{s_0}\alpha^2}{2} \|\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)\|^2 \end{aligned}$$

$$\begin{aligned}
&= \mathcal{L}(\mathbf{w}(t^*)) + \frac{\alpha\beta}{\eta} \langle (\mathbf{w}(t^*) - \mathbf{w}(t^* - 1)), \mathbf{w}(t^* + 1) - \mathbf{w}(t^*) \rangle + \left( \frac{H_{s_0}\alpha^2}{2} - \frac{\alpha}{\eta} \right) \|\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)\|^2 \\
&\stackrel{(**)}{\leq} \mathcal{L}(\mathbf{w}(t^*)) + \frac{\alpha\beta}{2\eta} \|\mathbf{w}(t^*) - \mathbf{w}(t^* - 1)\|^2 + \frac{\alpha\beta}{2\eta} \|\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)\|^2 \\
&\quad + \left( \frac{H_{s_0}\alpha^2}{2} - \frac{\alpha}{\eta} \right) \|\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)\|^2 \\
&= \mathcal{L}(\mathbf{w}(t^*)) + \frac{\alpha\beta}{2\eta} \|\mathbf{w}(t^*) - \mathbf{w}(t^* - 1)\|^2 + \left( \frac{\alpha\beta}{2\eta} - \frac{\alpha}{\eta} + \frac{H_{s_0}\alpha^2}{2} \right) \|\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)\|^2 \\
&= \mathcal{L}(\mathbf{w}(t^*)) + \frac{\beta}{2\eta} \|\mathbf{w}(t^*) - \mathbf{w}(t^* - 1)\|^2 - \frac{(1-\alpha)\beta}{2\eta} \|\mathbf{w}(t^*) - \mathbf{w}(t^* - 1)\|^2 \\
&\quad + \left( \frac{\alpha\beta}{2\eta} - \frac{\alpha}{\eta} + \frac{H_{s_0}\alpha^2}{2} \right) \|\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)\|^2 \\
&\stackrel{(\diamond)}{\leq} \mathcal{L}(\mathbf{w}(t^*)) + \frac{\beta}{2\eta} \|\mathbf{w}(t^*) - \mathbf{w}(t^* - 1)\|^2 - \frac{\beta}{2\eta} \alpha^2 \|\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)\|^2 \\
&\quad - \frac{(1-\beta)(1-C_1)\alpha^2}{\eta} \|\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)\|^2, \tag{22}
\end{aligned}$$

where Eq. (\*) is due to a simple rearrangement of the update rule of gradient descent with momentum (Eq. (5)), i.e.,

$$\nabla \mathcal{L}(\mathbf{w}(t)) = \frac{1}{\eta} (\beta(\mathbf{w}(t) - \mathbf{w}(t-1)) - (\mathbf{w}(t+1) - \mathbf{w}(t))), \forall t \geq 1, \tag{23}$$

Inequality (\*\*) is due to Cauchy Schwarz's inequality and arithmetic-geometric average inequality, and Inequality ( $\diamond$ ) is due to

$$\begin{aligned}
&- \frac{(1-\alpha)\beta}{2\eta} \|\mathbf{w}(t^*) - \mathbf{w}(t^* - 1)\|^2 + \left( \frac{\alpha\beta}{2\eta} - \frac{\alpha}{\eta} + \frac{H_{s_0}\alpha^2}{2} \right) \|\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)\|^2 \\
&= - \frac{(1-\alpha)\beta}{2\eta} \|\mathbf{w}(t^*) - \mathbf{w}(t^* - 1)\|^2 + \mathcal{O}(\alpha) \\
&\leq \mathcal{O}(\alpha^2) \\
&= - \frac{\beta}{2\eta} \alpha^2 \|\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)\|^2 - \frac{(1-\beta)(1-C_1)\alpha^2}{\eta} \|\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)\|^2.
\end{aligned}$$

Here the inequality is due to that  $-\frac{(1-\alpha)\beta}{2\eta} \|\mathbf{w}(t^*) - \mathbf{w}(t^* - 1)\|^2$  tend to  $-\frac{\beta}{2\eta} \|\mathbf{w}(t^*) - \mathbf{w}(t^* - 1)\|^2$  as  $\alpha$  tend to zero, which is a negative constant by Lemma 13.

Eq. (22) indicates Eq. (7) holds at  $(t^*, \alpha)$  for  $\alpha > 0$  is small enough, which contradicts to  $\alpha^* = 0$ .

**Case 2:**  $\alpha^* \neq 0$ : Same as **Case 1**, we have for any  $1 \leq t < t^*$ ,

$$\mathcal{L}(\mathbf{w}(t)) + \frac{\beta}{2\eta} \|\mathbf{w}(t) - \mathbf{w}(t-1)\|^2 \geq \mathcal{L}(\mathbf{w}(t+1)) + \frac{\beta}{2\eta} \|\mathbf{w}(t+1) - \mathbf{w}(t)\|^2,$$

which further leads to

$$\mathcal{L}(\mathbf{w}(1)) \geq \mathcal{L}(\mathbf{w}(t^*)) + \frac{\beta}{2\eta} \|\mathbf{w}(t^*) - \mathbf{w}(t^* - 1)\|^2. \tag{24}$$

On the other hand, by the definition of  $\alpha^*$ , we have for any  $0 \leq \alpha < \alpha^*$ , we have Eq. (7) holds for  $(t^*, \alpha)$ , which by continuity further leads to Eq. (7) holds for  $(t^*, \alpha^*)$ . Therefore,  $\alpha^* < 1$ , otherwise, Eq. (7) holds for  $(t^*, \alpha)$ ,  $\forall \alpha \in [0, 1]$  which contradicts the definition of  $t^*$ .

Combining Eq. (7) with  $(t^*, \alpha)$  and Eq. (24), we further have

$$\mathcal{L}(\mathbf{w}(1)) \geq \mathcal{L}(\mathbf{w}(t^* + \alpha)) + \frac{\beta}{2\eta} \alpha^2 \|\mathbf{w}(t+1) - \mathbf{w}(t)\|^2 + \frac{1-C_1}{2C_1} H_{s_0} \alpha^2 \|\mathbf{w}(t+1) - \mathbf{w}(t)\|^2,$$

Consequently, for any  $\alpha \in [0, \alpha^*]$

$$\mathcal{L}(\mathbf{w}(1)) \geq \mathcal{L}(\mathbf{w}(t^* + \alpha)),$$

and by Lemma 10, we then have  $\mathcal{L}$  is  $H_{s_0}$  smooth at  $\mathbf{w}(t^* + \alpha)$ , which further by Taylor's expansion leads to

$$\begin{aligned} & \mathcal{L}(\mathbf{w}(t^* + \alpha^*)) \\ & \leq \mathcal{L}(\mathbf{w}(t^*)) + \langle \nabla \mathcal{L}(\mathbf{w}(t^*)), \mathbf{w}(t^* + \alpha^*) - \mathbf{w}(t^*) \rangle + \frac{H_{s_0}}{2} \|\mathbf{w}(t^* + \alpha^*) - \mathbf{w}(t^*)\|^2 \\ & \stackrel{(\circ)}{\leq} \mathcal{L}(\mathbf{w}(t^*)) + \frac{\alpha^* \beta}{2\eta} \|\mathbf{w}(t^*) - \mathbf{w}(t^* - 1)\|^2 + \frac{\alpha^* \beta}{2\eta} \|\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)\|^2 \\ & \quad + \left( \frac{H_{s_0}(\alpha^*)^2}{2} - \frac{\alpha^*}{\eta} \right) \|\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)\|^2 \\ & = \mathcal{L}(\mathbf{w}(t^*)) + \frac{\alpha^* \beta}{2\eta} \|\mathbf{w}(t^*) - \mathbf{w}(t^* - 1)\|^2 + \left( \frac{H_{s_0}(\alpha^*)^2}{2} - \frac{\alpha^*(2 - \beta)}{2\eta} \right) \|\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)\|^2 \\ & \stackrel{(\bullet)}{=} \mathcal{L}(\mathbf{w}(t^*)) + \frac{\alpha^* \beta}{2\eta} \|\mathbf{w}(t^*) - \mathbf{w}(t^* - 1)\|^2 + \left( \frac{(1 - \beta)C_1(\alpha^*)^2}{\eta} - \frac{\alpha^*(2 - \beta)}{2\eta} \right) \|\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)\|^2 \\ & \stackrel{(*)}{\leq} \mathcal{L}(\mathbf{w}(t^*)) + \frac{\beta}{2\eta} \|\mathbf{w}(t^*) - \mathbf{w}(t^* - 1)\|^2 - \frac{(\alpha^*)^2 \beta}{2\eta} \|\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)\|^2 \\ & \quad - \frac{(1 - \beta)(1 - C_1)(\alpha^*)^2}{\eta} \|\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)\|^2 \end{aligned}$$

where Eq. (◦) follows the same routine as **Case 1**, Eq. (•) is due to the definition of  $\eta$  and  $C_1$ , and Eq. (\*) is due to  $\alpha^* < 1$ , and  $\|\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)\|^2 > 0$  (given by Lemma 13).

By the continuity of  $\mathcal{L}$ , for any small enough  $\delta > 0$ , Eq. (7) holds for  $(t^*, \alpha^* + \delta)$ , which contradicts to the definition of  $\alpha^*$ .

The proof is completed.  $\square$

By Lemma 1, one can easily obtain the sum of the squared norms of the updates across the trajectory converges.

**Corollary 3.** *Let all conditions in Theorem 2 hold. We have*

$$\sum_{t=1}^{\infty} \|\mathbf{w}(t+1) - \mathbf{w}(t)\|^2 < \infty. \quad (25)$$

Consequently, we have

$$\|\mathbf{w}(t)\| = \mathcal{O}(\sqrt{t}).$$

*Proof.* By Lemma 1, we have

$$\begin{aligned} \mathcal{L}(\mathbf{w}(t)) + \frac{\beta}{2\eta} \|\mathbf{w}(t) - \mathbf{w}(t-1)\|^2 & - \left( \mathcal{L}(\mathbf{w}(t+1)) + \frac{\beta}{2\eta} \|\mathbf{w}(t+1) - \mathbf{w}(t)\|^2 \right) \\ & \geq \frac{(1 - C_1)(1 - \beta)}{\eta} \|\mathbf{w}(t+1) - \mathbf{w}(t)\|^2, \end{aligned}$$

which by summing over  $t$  further leads to

$$\mathcal{L}(\mathbf{w}(1)) \geq \mathcal{L}(\mathbf{w}(1)) - \left( \mathcal{L}(\mathbf{w}(t+1)) + \frac{\beta}{2\eta} \|\mathbf{w}(t+1) - \mathbf{w}(t)\|^2 \right) \geq \frac{(1 - C_1)(1 - \beta)}{\eta} \sum_{s=1}^t \|\mathbf{w}(s+1) - \mathbf{w}(s)\|^2.$$

Taking  $t \rightarrow \infty$  leads to

$$\sum_{s=1}^{\infty} \|\mathbf{w}(s+1) - \mathbf{w}(s)\|^2 < \infty.$$

By triangle inequality, we further have

$$\begin{aligned} \|\mathbf{w}(t)\| &\leq \sum_{s=1}^t \|\mathbf{w}(s+1) - \mathbf{w}(s)\| + \|\mathbf{w}(1)\| \\ &\stackrel{(*)}{\leq} \sqrt{t \left( \sum_{s=1}^t \|\mathbf{w}(s+1) - \mathbf{w}(s)\|^2 \right)} + \|\mathbf{w}(1)\| = \mathcal{O}(\sqrt{t}), \end{aligned}$$

where Eq. (\*) is due to Cauchy-Schwartz's inequality.

The proof is completed.  $\square$

By the negative derivative of the loss and the separable data, we can finally prove Corollary 1.

*Proof of Corollary 1.* By Eq. (21), we have

$$\begin{aligned} \|\mathbf{w}(t+1) - \mathbf{w}(t)\|^2 &= \eta^2 \left\| \sum_{s=1}^t \beta^{t-s} \frac{\sum_{i=1}^N \ell'(\langle \mathbf{w}(s), \tilde{\mathbf{x}}_i \rangle) \tilde{\mathbf{x}}_i}{N} \right\|^2 \\ &= \eta^2 \left\| \sum_{s=1}^t \beta^{t-s} \frac{\sum_{i=1}^N \ell'(\langle \mathbf{w}(s), \tilde{\mathbf{x}}_i \rangle) \tilde{\mathbf{x}}_i}{N} \right\|^2 \frac{\|\hat{\mathbf{w}}\|^2}{\|\hat{\mathbf{w}}\|^2} \\ &\stackrel{(*)}{\geq} \eta^2 \gamma^2 \left\langle \hat{\mathbf{w}}, \sum_{s=1}^t \beta^{t-s} \frac{\sum_{i=1}^N \ell'(\langle \mathbf{w}(s), \tilde{\mathbf{x}}_i \rangle) \tilde{\mathbf{x}}_i}{N} \right\rangle^2 \\ &\stackrel{(**)}{\geq} \eta^2 \gamma^2 \left( \sum_{s=1}^t \beta^{t-s} \frac{\sum_{i=1}^N \ell'(\langle \mathbf{w}(s), \tilde{\mathbf{x}}_i \rangle)}{N} \right)^2 \\ &\geq \eta^2 \gamma^2 \left( \frac{\sum_{i=1}^N \ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle)}{N} \right)^2 \\ &\stackrel{(\bullet)}{\geq} \eta^2 \gamma^2 \left( \frac{\sum_{i=1}^N \ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle) \|\tilde{\mathbf{x}}_i\|}{N} \right)^2 \\ &\geq \eta^2 \gamma^2 \left\| \frac{\sum_{i=1}^N \ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle) \tilde{\mathbf{x}}_i}{N} \right\|^2 \\ &= \eta^2 \gamma^2 \|\nabla \mathcal{L}(\mathbf{w}(t))\|^2, \end{aligned} \tag{26}$$

where Inequality (\*) is due to Cauchy-Schwartz's inequality, Inequality (\*\*) is due to  $\ell'(s) < 0$ ,  $\forall s \in \mathbb{R}$  and  $\langle \hat{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle \geq \gamma$ ,  $\forall i \in [N]$ , and Inequality (•) is due to  $\|\tilde{\mathbf{x}}_i\| \leq 1$ . By combining Eq. (25) and Eq. (26), we complete the proof.  $\square$

By the exponential-tailed assumption of the loss (Assumption 2), we further have the following corollary.

**Corollary 4.** *Let all conditions in Theorem 2 hold. Then,  $\lim_{t \rightarrow \infty} \|\nabla \mathcal{L}(\mathbf{w}(t))\| = 0$ , and*

$$\lim_{t \rightarrow \infty} \langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle = \infty, \forall i.$$

*Consequently, there exists an large enough time  $t_0$ , such that,  $\forall t > t_0$ ,  $\forall i$ , we have  $\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle > 0$ , and*

$$\begin{aligned} -\ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle) &\leq (1 + e^{-\mu + \langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle}) e^{-\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle}, \\ -\ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle) &\geq (1 - e^{-\mu - \langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle}) e^{-\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle}. \end{aligned}$$

## B.1.2 BOUNDING THE ORTHOGONAL PART

To prove Theorem 2, we only need to show  $\mathbf{w}(t) - \log(t)\hat{\mathbf{w}}$  ( $t \geq 1$ ) has bounded norm for any iteration  $t > 0$ . Letting  $\mathbf{C}_2 = \frac{\eta}{(1-\beta)N}$  in Corollary 2, we obtain a constant vector  $\tilde{\mathbf{w}}$  satisfying Eq. (12). Define

$$\mathbf{r}(t) \triangleq \mathbf{w}(t) - \log(t)\hat{\mathbf{w}} - \tilde{\mathbf{w}}. \quad (27)$$

As  $\tilde{\mathbf{w}}$  is a constant vector, that  $\mathbf{w}(t) - \log(t)\hat{\mathbf{w}}$  ( $t \geq 1$ ) has bounded norm is equivalent to  $\mathbf{r}(t)$  has bounded norm. Lemma 2 then propose an equivalent proposition of  $\|\mathbf{r}(t)\|$  is bounded, and further prove this proposition is fulfilled. As the proof is rather complex, we separate it into two sub-lemmas. We first prove  $\|\mathbf{r}(t)\|$  is bounded if and only if function  $g(t)$  is upper bounded.

**Lemma 15** (First argument in Lemma 2). *Let all conditions in Theorem 2 hold. Then,  $\|\mathbf{r}(t)\|$  is bounded if and only if function  $g(t)$  is upper bounded.*

*Proof.* We start the proof by showing that  $A_1(t) \triangleq \sum_{\tau=2}^t \langle \mathbf{r}(\tau) - \mathbf{r}(\tau-1), \mathbf{w}(\tau) - \mathbf{w}(\tau-1) \rangle$  has bounded absolute value.

By the definition of  $\mathbf{r}(t)$ , we have

$$\mathbf{r}(t) - \mathbf{r}(t-1) = \mathbf{w}(t) - \mathbf{w}(t-1) - \log\left(\frac{t}{t-1}\right)\hat{\mathbf{w}},$$

which further indicates

$$A_1(t) = \sum_{\tau=2}^t \left\langle \mathbf{w}(\tau) - \mathbf{w}(\tau-1) - \log\left(\frac{\tau}{\tau-1}\right)\hat{\mathbf{w}}, \mathbf{w}(\tau) - \mathbf{w}(\tau-1) \right\rangle.$$

Therefore, the absolute value of  $A_1(t)$  can be bounded as

$$\begin{aligned} |A_1(t)| &= \left| \sum_{\tau=2}^t \left\langle \mathbf{w}(\tau) - \mathbf{w}(\tau-1) - \log\left(\frac{\tau}{\tau-1}\right)\hat{\mathbf{w}}, \mathbf{w}(\tau) - \mathbf{w}(\tau-1) \right\rangle \right| \\ &\leq \sum_{\tau=2}^t \left| \left\langle \mathbf{w}(\tau) - \mathbf{w}(\tau-1) - \log\left(\frac{\tau}{\tau-1}\right)\hat{\mathbf{w}}, \mathbf{w}(\tau) - \mathbf{w}(\tau-1) \right\rangle \right| \\ &\leq \sum_{\tau=2}^t \|\mathbf{w}(\tau) - \mathbf{w}(\tau-1)\|^2 + \sum_{\tau=2}^t \left| \left\langle \log\left(\frac{\tau}{\tau-1}\right)\hat{\mathbf{w}}, \mathbf{w}(\tau) - \mathbf{w}(\tau-1) \right\rangle \right| \\ &\leq \sum_{\tau=2}^t \|\mathbf{w}(\tau) - \mathbf{w}(\tau-1)\|^2 + \sum_{\tau=2}^t \left\| \log\left(\frac{\tau}{\tau-1}\right)\hat{\mathbf{w}} \right\| \|\mathbf{w}(\tau) - \mathbf{w}(\tau-1)\| \\ &\stackrel{(*)}{\leq} \frac{3}{2} \sum_{\tau=2}^t \|\mathbf{w}(\tau) - \mathbf{w}(\tau-1)\|^2 + \frac{1}{2} \sum_{\tau=2}^t \left\| \log\left(\frac{\tau}{\tau-1}\right)\hat{\mathbf{w}} \right\|^2 \\ &\stackrel{(\circ)}{<} \infty, \end{aligned}$$

where Inequality  $(\star)$  is due to the Inequality of arithmetic and geometric means, and Inequality  $(\circ)$  is due to Corollary 3 and  $\log\frac{\tau}{\tau-1} = \mathcal{O}\left(\frac{1}{\tau}\right)$ .

Therefore,  $g(t)$  is upper bounded is then equivalent to  $\frac{1}{2}\|\mathbf{r}(t)\|^2 + \frac{\beta}{1-\beta}\langle \mathbf{r}(t), \mathbf{w}(t) - \mathbf{w}(t-1) \rangle$  is upper bounded. Now if  $\frac{1}{2}\|\mathbf{r}(t)\|^2 + \frac{\beta}{1-\beta}\langle \mathbf{r}(t), \mathbf{w}(t) - \mathbf{w}(t-1) \rangle$  is upper bounded, we will prove  $\|\mathbf{r}(t)\|$  is bounded by reduction to absurdity.

**Suppose that  $\|\mathbf{r}(t)\|$  has unbounded norm.** By Corollary 3, we have  $\lim_{t \rightarrow \infty} \|\mathbf{w}(t) - \mathbf{w}(t-1)\| = 0$ , and there exists a large enough time  $T$ , such that  $\|\mathbf{w}(t) - \mathbf{w}(t-1)\| < 1$  for any  $t \geq T$ . On the other hand, since  $\mathbf{r}(t)$  is unbounded from above, there exists an increasing time sequence  $k_i > T$ ,  $i \in \mathbb{Z}^+$ , such that

$$\lim_{i \rightarrow \infty} \|\mathbf{r}(k_i)\| = \infty.$$

Therefore, we have

$$\begin{aligned}
& \liminf_{i \rightarrow \infty} \frac{1}{2} \|\mathbf{r}(k_i)\|^2 + \frac{\beta}{1-\beta} \langle \mathbf{r}(k_i), \mathbf{w}(k_i) - \mathbf{w}(k_i - 1) \rangle \\
& \geq \liminf_{i \rightarrow \infty} \frac{1}{2} \|\mathbf{r}(k_i)\|^2 - \frac{\beta}{1-\beta} \|\mathbf{r}(k_i)\| \|\mathbf{w}(k_i) - \mathbf{w}(k_i - 1)\| \\
& \geq \liminf_{i \rightarrow \infty} \frac{1}{2} \|\mathbf{r}(k_i)\|^2 - \frac{\beta}{1-\beta} \|\mathbf{r}(k_i)\| = \infty,
\end{aligned}$$

which leads to contradictory, and completes the proof of necessity.

On the other hand, if  $\|\mathbf{r}(t)\|$  is upper bounded, since  $\|\mathbf{w}(t) - \mathbf{w}(t - 1)\|$  is also upper bounded, we have  $\frac{1}{2} \|\mathbf{r}(t)\|^2 + \frac{\beta}{1-\beta} \langle \mathbf{r}(t), \mathbf{w}(t) - \mathbf{w}(t - 1) \rangle$  is upper bounded, which completes the proof of sufficiency.

The proof is completed.  $\square$

Therefore, the last piece of this puzzle is to prove  $g(t)$  is upper bounded  $\forall t > 0$ .

**Lemma 16** (Second argument in Lemma 2). *Let all conditions in Theorem 2 hold. Then, we have  $g(t)$  is upper bounded.*

*Proof.* We start the proof by calculating  $g(t + 1) - g(t)$ . For any  $t \geq 2$ , we have

$$\begin{aligned}
g(t + 1) - g(t) &= \frac{1}{2} \|\mathbf{r}(t + 1) - \mathbf{r}(t)\|^2 + \langle \mathbf{r}(t), \mathbf{r}(t + 1) - \mathbf{r}(t) \rangle + \frac{\beta}{1-\beta} \langle \mathbf{r}(t + 1), \mathbf{w}(t + 1) - \mathbf{w}(t) \rangle \\
&\quad - \frac{\beta}{1-\beta} \langle \mathbf{r}(t), \mathbf{w}(t) - \mathbf{w}(t - 1) \rangle - \frac{\beta}{1-\beta} \langle \mathbf{r}(t + 1) - \mathbf{r}(t), \mathbf{w}(t + 1) - \mathbf{w}(t) \rangle \\
&= \frac{1}{2} \|\mathbf{r}(t + 1) - \mathbf{r}(t)\|^2 + \langle \mathbf{r}(t), \mathbf{r}(t + 1) - \mathbf{r}(t) \rangle + \frac{\beta}{1-\beta} \langle \mathbf{r}(t), \mathbf{w}(t + 1) + \mathbf{w}(t - 1) - 2\mathbf{w}(t) \rangle.
\end{aligned}$$

On the other hand, by simply rearranging the update rule Eq. (5), we have

$$\frac{\beta}{1-\beta} (\mathbf{w}(t + 1) + \mathbf{w}(t - 1) - 2\mathbf{w}(t)) = -\frac{\eta}{1-\beta} \nabla \mathcal{L}(\mathbf{w}(t)) - (\mathbf{w}(t + 1) - \mathbf{w}(t)), \quad (28)$$

which further indicates

$$\begin{aligned}
& g(t + 1) - g(t) \\
&= \frac{1}{2} \|\mathbf{r}(t + 1) - \mathbf{r}(t)\|^2 + \langle \mathbf{r}(t), \mathbf{r}(t + 1) - \mathbf{r}(t) \rangle + \left\langle \mathbf{r}(t), -\frac{\eta}{1-\beta} \nabla \mathcal{L}(\mathbf{w}(t)) - (\mathbf{w}(t + 1) - \mathbf{w}(t)) \right\rangle \\
&= \frac{1}{2} \|\mathbf{r}(t + 1) - \mathbf{r}(t)\|^2 + \left\langle \mathbf{r}(t), -\log \left( \frac{t+1}{t} \right) \hat{\mathbf{w}} - \frac{\eta}{1-\beta} \nabla \mathcal{L}(\mathbf{w}(t)) \right\rangle.
\end{aligned}$$

Denote  $A_2(t) = \|\mathbf{r}(t + 1) - \mathbf{r}(t)\|^2$ , and  $A_3(t) = \left\langle \mathbf{r}(t), -\log \left( \frac{t+1}{t} \right) \hat{\mathbf{w}} - \frac{\eta}{1-\beta} \nabla \mathcal{L}(\mathbf{w}(t)) \right\rangle$ . We then prove respectively  $\sum_{t=1}^{\infty} A_2(t)$  and  $\sum_{t=1}^{\infty} A_3(t)$  are upper bounded.

First of all, by definition of  $\mathbf{r}(t)$  Eq.(27), we have

$$\begin{aligned}
\sum_{t=1}^{\infty} A_2(t) &= \sum_{t=1}^{\infty} \left( \|\mathbf{w}(t + 1) - \mathbf{w}(t)\|^2 + \log \left( \frac{t+1}{t} \right)^2 \|\hat{\mathbf{w}}\|^2 - 2 \log \left( \frac{t+1}{t} \right) \langle \mathbf{w}(t + 1) - \mathbf{w}(t), \hat{\mathbf{w}} \rangle \right) \\
&\leq 2 \sum_{t=1}^{\infty} \left( \|\mathbf{w}(t + 1) - \mathbf{w}(t)\|^2 + \log \left( \frac{t+1}{t} \right)^2 \|\hat{\mathbf{w}}\|^2 \right) \stackrel{(\bullet)}{<} \infty, \quad (29)
\end{aligned}$$

where Eq. (•) is due to Lemma 3 and  $\log \left( \frac{t+1}{t} \right) = \mathcal{O}(\frac{1}{t})$ .

Then we only need to prove  $\sum_{t=1}^{\infty} A_3(t) < \infty$ .

To begin with, by adding one additional term  $\frac{1}{t}\hat{\mathbf{w}}$  into  $A_3$ , we have

$$A_3(t) = \left\langle \mathbf{r}(t), \frac{1}{t}\hat{\mathbf{w}} - \log\left(\frac{t+1}{t}\right)\hat{\mathbf{w}} \right\rangle + \left\langle \mathbf{r}(t), -\frac{1}{t}\hat{\mathbf{w}} - \frac{\eta}{1-\beta}\nabla\mathcal{L}(\mathbf{w}(t)) \right\rangle.$$

On the one hand, by Corollary 3,  $\|\mathbf{w}(t)\| = \mathcal{O}(\sqrt{t})$ , which further leads to

$$\|\mathbf{r}(t)\| = \|\mathbf{w}(t)\| + \log(t)\|\hat{\mathbf{w}}\| + \|\hat{\mathbf{w}}\| = \mathcal{O}(\sqrt{t})$$

By  $\frac{1}{t} - \log\frac{t+1}{t} = \mathcal{O}\left(\frac{1}{t^2}\right)$ , we have

$$\left\langle \mathbf{r}(t), \frac{1}{t}\hat{\mathbf{w}} - \log\left(\frac{t+1}{t}\right)\hat{\mathbf{w}} \right\rangle = \mathcal{O}\left(\frac{1}{t^{\frac{3}{2}}}\right). \quad (30)$$

On the other hand, by direct calculation of the gradient, we have

$$\begin{aligned} & \left\langle \mathbf{r}(t), -\frac{1}{t}\hat{\mathbf{w}} - \frac{\eta}{1-\beta}\nabla\mathcal{L}(\mathbf{w}(t)) \right\rangle \\ &= \left\langle \mathbf{r}(t), -\frac{1}{t}\hat{\mathbf{w}} - \frac{\eta}{(1-\beta)N}\sum_{i=1}^N \ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle)\tilde{\mathbf{x}}_i \right\rangle \\ &\stackrel{(*)}{=} \frac{1}{N} \left\langle \mathbf{r}(t), -\frac{1}{t}\frac{\eta}{1-\beta}\sum_{\tilde{\mathbf{x}}_i \in \mathcal{T}(\mathcal{S}_s)} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle}\tilde{\mathbf{x}}_i - \frac{\eta}{1-\beta}\sum_{i=1}^N \ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle)\tilde{\mathbf{x}}_i \right\rangle \\ &= \frac{1}{N} \left\langle \mathbf{r}(t), -\frac{\eta}{1-\beta}\sum_{\tilde{\mathbf{x}}_i \in \mathcal{T}(\mathcal{S}_s)} \left( \frac{1}{t}e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} + \ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle) \right)\tilde{\mathbf{x}}_i \right\rangle - \frac{1}{N} \left\langle \mathbf{r}(t), \frac{\eta}{1-\beta}\sum_{\tilde{\mathbf{x}}_i \notin \mathcal{T}(\mathcal{S}_s)} \ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle)\tilde{\mathbf{x}}_i \right\rangle, \end{aligned}$$

where Eq. (\*) is due to the definition of  $\tilde{\mathbf{w}}$  (Eq. (12) with  $C_2 = \frac{\eta}{1-\beta}$ ).

Denote

$$A_4(t) = - \left\langle \mathbf{r}(t), \frac{\eta}{1-\beta}\sum_{\tilde{\mathbf{x}}_i \notin \mathcal{T}(\mathcal{S}_s)} \ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle)\tilde{\mathbf{x}}_i \right\rangle,$$

and

$$A_5(t) = \left\langle \mathbf{r}(t), -\frac{\eta}{1-\beta}\sum_{\tilde{\mathbf{x}}_i \in \mathcal{T}(\mathcal{S}_s)} \left( \frac{1}{t}e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} + \ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle) \right)\tilde{\mathbf{x}}_i \right\rangle.$$

We then analysis these two terms respectively. As for  $A_4(t)$ , due to  $\ell' < 0$ , we have

$$A_4(t) \leq -\frac{\eta}{1-\beta} \left\langle \mathbf{r}(t), \sum_{\tilde{\mathbf{x}}_i \notin \mathcal{T}(\mathcal{S}_s), \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle > 0} \ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle)\tilde{\mathbf{x}}_i \right\rangle.$$

By Corollary 4, we further have  $\forall t > t_0$

$$-\ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle) \leq (1 + e^{-\mu + \langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle})e^{-\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle} \leq 2e^{-\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle},$$

which further indicates

$$\begin{aligned} A_4(t) &\leq -\frac{\eta}{1-\beta} \sum_{\tilde{\mathbf{x}}_i \notin \mathcal{T}(\mathcal{S}_s), \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle > 0} \ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\ &\leq \frac{\eta}{1-\beta} \sum_{\tilde{\mathbf{x}}_i \notin \mathcal{T}(\mathcal{S}_s), \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle > 0} 2e^{-\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle} \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\ &= \frac{\eta}{1-\beta} \sum_{\tilde{\mathbf{x}}_i \notin \mathcal{T}(\mathcal{S}_s), \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle > 0} 2e^{-\langle \mathbf{r}(t) + \log t \hat{\mathbf{w}} + \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\ &\leq \frac{\eta}{1-\beta} \left( \max_i e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \right) \sum_{\tilde{\mathbf{x}}_i \notin \mathcal{T}(\mathcal{S}_s), \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle > 0} 2e^{-\langle \mathbf{r}(t) + \log t \hat{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\ &\stackrel{(\circ)}{\leq} \frac{\eta}{1-\beta} \frac{(\max_i e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle})}{t^\theta} \sum_{\tilde{\mathbf{x}}_i \notin \mathcal{T}(\mathcal{S}_s), \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle > 0} 2e^{-\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle} \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\ &\stackrel{(\circ)}{\leq} \frac{\eta}{1-\beta} \frac{(\max_i e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle})}{t^\theta} 2N, \end{aligned}$$

where  $\theta$  in Eq. (◊) is defined as

$$\theta = \min_{\tilde{\mathbf{x}}_i \notin \mathcal{T}(\mathbf{S}_s)} \langle \tilde{\mathbf{x}}_i, \hat{\mathbf{w}} \rangle > 1. \quad (31)$$

As  $\sum_{t=1}^{\infty} \frac{1}{t^\theta} < \infty$ , we have

$$\sum_{t=1}^{\infty} A_4(t) < \infty^4. \quad (32)$$

For each term  $\left\langle \mathbf{r}(t), -\frac{\eta}{1-\beta} \left( \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} + \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}}_i \rangle) \right) \tilde{\mathbf{x}}_i \right\rangle$  ( $\tilde{\mathbf{x}}_i \notin \mathcal{T}(\mathbf{S}_s)$ ) in  $A_5(t)$ , we divide the analysis into two parts depending on the sign of  $\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle$ .

**Case 1:**  $\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \geq 0$ . By Corollary 4, we have

$$\begin{aligned} & \left\langle \mathbf{r}(t), -\frac{\eta}{1-\beta} \left( \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} + \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}}_i \rangle) \right) \tilde{\mathbf{x}}_i \right\rangle \\ &= -\frac{\eta}{1-\beta} \left( \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} + \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}}_i \rangle) \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\ &\leq \frac{\eta}{1-\beta} \left( -\frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} + (1 + e^{-\mu + \langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle}) e^{-\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\ &\stackrel{(\diamond)}{=} \frac{\eta}{1-\beta} \left( -\frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} + (1 + e^{-\mu + \langle \mathbf{r}(t) + \log t \hat{\mathbf{w}} + \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle}) e^{-\langle \mathbf{r}(t) + \log t \hat{\mathbf{w}} + \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle, \end{aligned}$$

where Eq. (◊) is due to the definition of  $\mathbf{r}(t)$  (Eq. (27)).

Since  $\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \geq 0$ , we further have

$$\begin{aligned} & \left\langle \mathbf{r}(t), -\frac{\eta}{1-\beta} \left( \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} + \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}}_i \rangle) \right) \tilde{\mathbf{x}}_i \right\rangle \\ &\leq \frac{\eta}{1-\beta} \left( -\frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} + (1 + e^{-\mu + \langle \log t \hat{\mathbf{w}} + \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle}) e^{-\langle \mathbf{r}(t) + \log t \hat{\mathbf{w}} + \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\ &\stackrel{(\square)}{=} \frac{\eta}{1-\beta} \left( -\frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} + \frac{1}{t} (1 + t^{-\mu + e^{-\mu + \langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle}}) e^{-\langle \mathbf{r}(t) + \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\ &= \frac{\eta}{1-\beta} \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \left( -1 + (1 + t^{-\mu + e^{-\mu + \langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle}}) e^{-\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle, \end{aligned}$$

where Eq. (□) is due to  $\langle \hat{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle = 1, \forall \tilde{\mathbf{x}}_i \in \mathcal{T}(\mathbf{S}_s)$ .

Specifically,

$$\begin{aligned} & -1 + (1 + t^{-\mu + e^{-\mu + \langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle}}) e^{-\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle} \\ &= -1 + e^{-\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle} + t^{-\mu + e^{-\mu + \langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle}} e^{-\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle} \\ &\leq t^{-\mu + e^{-\mu + \langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle}} e^{-\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle}. \end{aligned}$$

Therefore,

$$\begin{aligned} & \frac{\eta}{1-\beta} \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \left( -1 + (1 + t^{-\mu + e^{-\mu + \langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle}}) e^{-\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\ &\leq \frac{\eta}{1-\beta} \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \left( t^{-\mu + e^{-\mu + \langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle}} e^{-\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\ &\leq \frac{\eta}{(1-\beta)e} \frac{1}{t^{1+\mu_+}} e^{-(1+\mu_+) \langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} = \mathcal{O} \left( \frac{1}{t^{1+\mu_+}} \right). \end{aligned}$$

<sup>4</sup>In this paper, for a real series  $\{r_i\}_{i=1}^{\infty}$ , we use  $\sum_{i=1}^{\infty} r_i < \infty$  representing  $\sum_{i=1}^T r_i$  is uniformly upper bounded for any  $T$ .



**Case 2:**  $\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle < 0$ . Similar to **Case 1**, in this case we have

$$\begin{aligned}
& \left\langle \mathbf{r}(t), -\frac{\eta}{1-\beta} \left( \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} + \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}}_i \rangle) \right) \tilde{\mathbf{x}}_i \right\rangle \\
& \leq \frac{\eta}{1-\beta} \left( -\frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} + \left( 1 - e^{-\mu - \langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle} \right) e^{-\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\
& = \frac{\eta}{1-\beta} \left( -\frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} + \left( 1 - e^{-\mu - \langle \mathbf{r}(t) + \log t \tilde{\mathbf{w}} + \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \right) e^{-\langle \mathbf{r}(t) + \log t \tilde{\mathbf{w}} + \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\
& = \frac{\eta}{1-\beta} \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \left( -1 + \left( 1 - e^{-\mu - \langle \mathbf{r}(t) + \log t \tilde{\mathbf{w}} + \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \right) e^{-\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle.
\end{aligned}$$

Specifically, if  $\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \geq -t^{-0.5\mu_-}$ ,

$$\begin{aligned}
& \left| \frac{\eta}{1-\beta} \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \left( -1 + \left( 1 - e^{-\mu - \langle \mathbf{r}(t) + \log t \tilde{\mathbf{w}} + \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \right) e^{-\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \right| \\
& = \left| \frac{\eta}{1-\beta} \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \left( -1 + \left( 1 - t^{-\mu_-} e^{-\mu - \langle \mathbf{r}(t) + \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \right) e^{-\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \right| \\
& \leq \frac{\eta}{1-\beta} \frac{1}{t^{1+0.5\mu_-}} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \left| -1 + \left( 1 - t^{-\mu_-} e^{-\mu - \langle \mathbf{r}(t) + \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \right) e^{-\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle} \right| \\
& \stackrel{(\dagger)}{=} \mathcal{O} \left( \frac{1}{t^{1+0.5\mu_-}} \right),
\end{aligned}$$

where Eq. (†) is due to if  $\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \geq -t^{-0.5\mu_-}$ ,

$$\lim_{t \rightarrow \infty} \left| -1 + \left( 1 - t^{-\mu_-} e^{-\mu - \langle \mathbf{r}(t) + \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \right) e^{-\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle} \right| = 0.$$

If  $-2 \leq \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle < -t^{-0.5\mu_-}$ , we have

$$\begin{aligned}
& \frac{\eta}{1-\beta} \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \left( -1 + \left( 1 - e^{-\mu - \langle \mathbf{r}(t) + \log t \tilde{\mathbf{w}} + \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \right) e^{-\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\
& = \frac{\eta}{1-\beta} \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \left( -1 + \left( 1 - \frac{1}{t^{\mu_-}} e^{-\mu - \langle \mathbf{r}(t) + \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \right) e^{-\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\
& \leq \frac{\eta}{1-\beta} \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \left( -1 + \left( 1 - \frac{e^{2\mu_-}}{t^{\mu_-}} e^{-\mu - \langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \right) e^{-\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle.
\end{aligned}$$

Therefore, when  $t$  is large enough,  $1 - \frac{e^{2\mu_-}}{t^{\mu_-}} e^{-\mu - \langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} > 0$ , which by  $e^{-\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle} \geq 1 - \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle$  leads to

$$\begin{aligned}
& \frac{\eta}{1-\beta} \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \left( -1 + \left( 1 - \frac{e^{2\mu_-}}{t^{\mu_-}} e^{-\mu - \langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \right) e^{-\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\
& \leq \frac{\eta}{1-\beta} \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \left( -1 + \left( 1 - \frac{e^{2\mu_-}}{t^{\mu_-}} e^{-\mu - \langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \right) (1 - \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle) \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\
& \leq \frac{\eta}{1-\beta} \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \left( -1 + \left( 1 - \frac{e^{2\mu_-}}{t^{\mu_-}} e^{-\mu - \langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \right) \left( 1 + \frac{1}{t^{0.5\mu_-}} \right) \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\
& = \frac{\eta}{1-\beta} \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \left( \frac{1}{t^{0.5\mu_-}} + \mathcal{O} \left( \frac{1}{t^{0.5\mu_-}} \right) \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle < 0.
\end{aligned}$$

If  $-2 > \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle$ ,

$$\begin{aligned}
& \frac{\eta}{1-\beta} \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \left( -1 + \left( 1 - e^{-\mu - \langle \mathbf{r}(t) + \log t \tilde{\mathbf{w}} + \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \right) e^{-\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\
& = \frac{\eta}{1-\beta} \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \left( -1 + \left( 1 - e^{-\mu - \langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle} \right) e^{-\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle.
\end{aligned}$$

For large enough  $t$ ,  $1 - e^{-\mu - \langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle} > \frac{1}{2}$ , and

$$\begin{aligned} & \frac{\eta}{1-\beta} \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \left( -1 + \left( 1 - e^{-\mu - \langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle} \right) e^{-\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\ & \leq \frac{\eta}{1-\beta} \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \left( -1 + \left( 1 - e^{-\mu - \langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle} \right) e^2 \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\ & \leq \frac{\eta}{1-\beta} \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \left( -1 + \frac{e^2}{2} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle < 0. \end{aligned}$$

Therefore, in **Case 2.**, for large enough  $t$ , we have

$$\left\langle \mathbf{r}(t), -\frac{\eta}{1-\beta} \left( \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} + \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}}_i \rangle) \right) \tilde{\mathbf{x}}_i \right\rangle \leq \mathcal{O} \left( \frac{1}{t^{1+0.5\mu_-}} \right).$$

Combining **Case 1.** and **Case 2.**, we conclude that

$$A_5(t) \leq \mathcal{O} \left( \frac{1}{t^{1+0.5\mu_+}} \right),$$

which further yields

$$\sum_{t=1}^{\infty} A_5(t) < \infty. \quad (33)$$

Combining Eq. (32) and Eq. (33), we conclude that  $\sum_{t=1}^{\infty} A_3(t) < \infty$ , which together with Eq. (29) yields  $\sum_{t=2}^{\infty} g(t+1) - g(t) < \infty$ , and completes the proof.  $\square$

We are now ready to prove Theorem 2.

*Proof of Theorem 2.* By Lemma 16, we have  $g(t)$  is upper bounded. Therefore, by Lemma 2, we have  $\|\mathbf{r}(t)\|$  is bounded, which further indicates  $\|\mathbf{w}(t) - \log(t)\hat{\mathbf{w}}\|$  is bounded.

Therefore, the direction of  $\mathbf{w}(t)$  can be calculated as

$$\begin{aligned} \frac{\mathbf{w}(t)}{\|\mathbf{w}(t)\|} &= \frac{\log(t)\hat{\mathbf{w}}}{\|\mathbf{w}(t)\|} + \frac{\mathbf{w}(t) - \log(t)\hat{\mathbf{w}}}{\|\mathbf{w}(t)\|} = \frac{\log(t)\hat{\mathbf{w}}}{\|\log(t)\hat{\mathbf{w}} + \mathbf{w}(t) - \log(t)\hat{\mathbf{w}}\|} + \frac{\mathbf{w}(t) - \log(t)\hat{\mathbf{w}}}{\|\mathbf{w}(t)\|} \\ &= \frac{\hat{\mathbf{w}}}{\left\| \hat{\mathbf{w}} + \frac{\mathbf{w}(t) - \log(t)\hat{\mathbf{w}}}{\log t} \right\|} + \frac{\mathbf{w}(t) - \log(t)\hat{\mathbf{w}}}{\|\mathbf{w}(t)\|} \rightarrow \frac{\hat{\mathbf{w}}}{\|\hat{\mathbf{w}}\|} \quad (as \ t \rightarrow \infty). \end{aligned}$$

The proof is completed.  $\square$

## B.2 IMPLICIT BIAS OF SGD WITH MOMENTUM

This section collects the proof of Theorem 3. Following the same framework as Appendix B.1, we will first prove that the sum of the squared gradient norms along the trajectory is finite. One may expect  $\mathcal{L}(\mathbf{w}(t)) + \frac{\beta}{2\eta} \|\mathbf{w}(t) - \mathbf{w}(t-1)\|^2$  is a Lyapunov function of SGDM. However, due to the randomness of the update rule of SGDM,  $\mathcal{L}(\mathbf{w}(t)) + \frac{\beta}{2\eta} \|\mathbf{w}(t) - \mathbf{w}(t-1)\|^2$  may no longer decrease (we will show this in the end of Appendix B.2, please see Appendix B.2.3 for explanation).

### B.2.1 PROOF OF THE SUM OF GRADIENTS ALONG THE TRAJECTORY IS FINITE

We first provide a proof of Lemma 3.

*Proof of Lemma 3.* We denote the parameter of the  $t$ -th step in Eq. (9) as  $\tilde{\mathbf{w}}$ , while that in Eq. (1) as  $\mathbf{w}(t)$ . With  $\mathbf{w}(1) = \tilde{\mathbf{w}}(1)$ , we will prove  $\tilde{\mathbf{w}}(t) = \mathbf{w}(t)$ ,  $\forall t > 1$  iteratively. Specifically, suppose for any  $1 \leq k \leq t$   $\mathbf{w}(k) = \tilde{\mathbf{w}}(k)$ . Then, we have

$$\begin{aligned} \mathbf{u}(t+1) &= -\tilde{\eta} \nabla \mathcal{L}_{\mathbf{B}(t)}(\tilde{\mathbf{w}}(t)) + \mathbf{u}(t) = -\tilde{\eta} \nabla \mathcal{L}_{\mathbf{B}(t)}(\mathbf{w}(t)) + \mathbf{u}(t) \\ &= -\tilde{\eta} \nabla \mathcal{L}_{\mathbf{B}(t)}(\mathbf{w}(t)) - \tilde{\eta} \nabla \mathcal{L}_{\mathbf{B}(t-1)}(\mathbf{w}(t-1)) + \mathbf{u}(t-1) \\ &= \dots = -\sum_{k=1}^t \tilde{\eta} \nabla \mathcal{L}_{\mathbf{B}(k)}(\mathbf{w}(k)) + \mathbf{u}(1) = -\sum_{k=1}^t \tilde{\eta} \nabla \mathcal{L}_{\mathbf{B}(k)}(\mathbf{w}(k)) + \mathbf{w}(1). \end{aligned}$$

Therefore,

$$\tilde{\mathbf{w}}(t+1) = \beta \tilde{\mathbf{w}}(t) + (1-\beta) \mathbf{u}(t+1) = \beta \tilde{\mathbf{w}}(t) + (1-\beta) \left( -\sum_{k=1}^t \tilde{\eta} \nabla \mathcal{L}_{\mathbf{B}(k)}(\tilde{\mathbf{w}}(k)) + \tilde{\mathbf{w}}(1) \right),$$

which by iteration further indicates

$$\begin{aligned} \tilde{\mathbf{w}}(t+1) &= \tilde{\mathbf{w}}(1) - \tilde{\eta} \left( \sum_{k=1}^t (1 - \beta^{t+1-k}) \nabla \mathcal{L}_{\mathbf{B}(k)}(\tilde{\mathbf{w}}(k)) \right) = \mathbf{w}(1) - \tilde{\eta} \left( \sum_{k=1}^t (1 - \beta^{t+1-k}) \nabla \mathcal{L}_{\mathbf{B}(k)}(\mathbf{w}(k)) \right), \\ \mathbf{w}(t) &= \mathbf{w}(1) - \tilde{\eta} \left( \sum_{k=1}^{t-1} (1 - \beta^{t-k}) \nabla \mathcal{L}_{\mathbf{B}(k)}(\mathbf{w}(k)) \right), \end{aligned}$$

and

$$\mathbf{w}(t-1) = \mathbf{w}(1) - \tilde{\eta} \left( \sum_{k=1}^{t-2} (1 - \beta^{t-1-k}) \nabla \mathcal{L}_{\mathbf{B}(k)}(\mathbf{w}(k)) \right).$$

On the other hand, by Eq. (1), we have

$$\begin{aligned} &\mathbf{w}(t+1) \\ &= \mathbf{w}(t) - \eta \nabla \mathcal{L}_{\mathbf{B}(t)}(\mathbf{w}(t)) + \beta(\mathbf{w}(t) - \mathbf{w}(t-1)) \\ &= \mathbf{w}(1) - \tilde{\eta} \left( \sum_{k=1}^{t-1} (1 - \beta^{t-k}) \nabla \mathcal{L}_{\mathbf{B}(k)}(\mathbf{w}(k)) \right) - \eta \nabla \mathcal{L}_{\mathbf{B}(t)}(\mathbf{w}(t)) + \beta \tilde{\eta} \left( \sum_{k=1}^{t-1} (\beta^{t-1-k} - \beta^{t-k}) \nabla \mathcal{L}_{\mathbf{B}(k)}(\mathbf{w}(k)) \right) \\ &= \mathbf{w}(1) - \tilde{\eta} \left( \sum_{k=1}^t (1 - \beta^{t+1-k}) \nabla \mathcal{L}_{\mathbf{B}(k)}(\mathbf{w}(k)) \right) = \tilde{\mathbf{w}}(t+1). \end{aligned}$$

The proof is completed.  $\square$

Using the alternative form given by Lemma 3, we then prove Lemma 4, which indicates  $\mathcal{L}(\mathbf{u}(t))$  is a proper choice of Lyapunov function.

*Proof of Lemma 4.* We start the proof by applying the Taylor's expansion of  $\mathcal{L}$  at the point  $\mathbf{u}(t)$  to the point  $\mathbf{u}(t+1)$ . Concretely, by Assumption 3. (S), we have

$$\mathcal{L}(\mathbf{u}(t+1)) \leq \mathcal{L}(\mathbf{u}(t)) + \langle \mathbf{u}(t+1) - \mathbf{u}(t), \nabla \mathcal{L}(\mathbf{u}(t)) \rangle + \frac{H}{2} \|\mathbf{u}(t+1) - \mathbf{u}(t)\|^2,$$

which by Eq. (9) leads to

$$\mathcal{L}(\mathbf{u}(t+1)) \leq \mathcal{L}(\mathbf{u}(t)) - \tilde{\eta} \langle \nabla \mathcal{L}_{\mathbf{B}(t)}(\mathbf{w}(t)), \nabla \mathcal{L}(\mathbf{u}(t)) \rangle + \frac{H \tilde{\eta}^2}{2} \|\nabla \mathcal{L}_{\mathbf{B}(t)}(\mathbf{w}(t))\|^2. \quad (34)$$

Taking the expectation of Eq. (34) with respect to  $\mathbf{w}(t+1)$  conditioning on  $\{\mathbf{w}(s)\}_{s=1}^t$ , we have

$$\begin{aligned}
& \mathbb{E}_{\mathbf{w}(t+1)}[\mathcal{L}(\mathbf{u}(t+1))|\{\mathbf{w}(s)\}_{s=1}^t] \\
& \stackrel{(\star)}{=} \mathbb{E}_{\mathbf{B}(t)}[\mathcal{L}(\mathbf{u}(t+1))|\{\mathbf{w}(s)\}_{s=1}^t] \\
& \leq \mathbb{E}_{\mathbf{B}(t)} \left[ \mathcal{L}(\mathbf{u}(t)) - \tilde{\eta} \langle \nabla \mathcal{L}_{\mathbf{B}(t)}(\mathbf{w}(t)), \nabla \mathcal{L}(\mathbf{u}(t)) \rangle + \frac{H\tilde{\eta}^2}{2} \|\nabla \mathcal{L}_{\mathbf{B}(t)}(\mathbf{w}(t))\|^2 \middle| \{\mathbf{w}(s)\}_{s=1}^t \right] \\
& \stackrel{(\circ)}{=} \mathcal{L}(\mathbf{u}(t)) - \tilde{\eta} \langle \nabla \mathcal{L}(\mathbf{w}(t)), \nabla \mathcal{L}(\mathbf{u}(t)) \rangle + \frac{H\tilde{\eta}^2}{2} \mathbb{E}_{\mathbf{B}(t)} [\|\nabla \mathcal{L}_{\mathbf{B}(t)}(\mathbf{w}(t))\|^2] \\
& \stackrel{(\bullet)}{\leq} \mathcal{L}(\mathbf{u}(t)) - \tilde{\eta} \langle \nabla \mathcal{L}(\mathbf{w}(t)), \nabla \mathcal{L}(\mathbf{u}(t)) \rangle + \frac{H\tilde{\eta}^2 N}{2b\gamma^2} \|\nabla \mathcal{L}(\mathbf{w}(t))\|^2, \tag{35}
\end{aligned}$$

where Eq.  $(\star)$  is due to that  $\mathbf{w}(t+1)$  is uniquely determined by  $\mathbf{B}(t)$  given  $\{\mathbf{w}(s)\}_{s=1}^t$ , Eq.  $(\circ)$  is due to  $\mathbf{u}(t)$  is uniquely determined by  $\{\mathbf{w}(s)\}_{s=1}^t$ , and **Inequality**.  $(\bullet)$  is due to Lemma 12.

Therefore, we have

$$\begin{aligned}
& \mathbb{E}_{\mathbf{w}(t+1)}[\mathcal{L}(\mathbf{u}(t+1))|\{\mathbf{w}(s)\}_{s=1}^t] \\
& \leq \mathcal{L}(\mathbf{u}(t)) - \tilde{\eta} \langle \nabla \mathcal{L}(\mathbf{w}(t)), \nabla \mathcal{L}(\mathbf{u}(t)) \rangle + \frac{H\tilde{\eta}^2 N}{2b\gamma^2} \|\nabla \mathcal{L}(\mathbf{w}(t))\|^2 \\
& = \mathcal{L}(\mathbf{u}(t)) - \tilde{\eta} \langle \nabla \mathcal{L}(\mathbf{w}(t)), \nabla \mathcal{L}(\mathbf{w}(t)) \rangle + \tilde{\eta} \langle \nabla \mathcal{L}(\mathbf{w}(t)), \nabla \mathcal{L}(\mathbf{w}(t)) - \nabla \mathcal{L}(\mathbf{u}(t)) \rangle + \frac{H\tilde{\eta}^2 N}{2b\gamma^2} \|\nabla \mathcal{L}(\mathbf{w}(t))\|^2 \\
& = \mathcal{L}(\mathbf{u}(t)) - \tilde{\eta} \left( 1 - \frac{H\tilde{\eta}N}{2b\gamma^2} \right) \|\nabla \mathcal{L}(\mathbf{w}(t))\|^2 + \langle \tilde{\eta} \nabla \mathcal{L}(\mathbf{w}(t)), \nabla \mathcal{L}(\mathbf{w}(t)) - \nabla \mathcal{L}(\mathbf{u}(t)) \rangle \\
& \leq \mathcal{L}(\mathbf{u}(t)) - \tilde{\eta} \left( 1 - \frac{H\tilde{\eta}N}{2b\gamma^2} \right) \|\nabla \mathcal{L}(\mathbf{w}(t))\|^2 + \frac{1}{2} \|\tilde{\eta} \nabla \mathcal{L}(\mathbf{w}(t))\|^2 + \frac{1}{2} \|\nabla \mathcal{L}(\mathbf{w}(t)) - \nabla \mathcal{L}(\mathbf{u}(t))\|^2 \\
& = \mathcal{L}(\mathbf{u}(t)) - \tilde{\eta} \left( 1 - \left( \frac{1}{2} + \frac{HN}{2b\gamma^2} \right) \tilde{\eta} \right) \|\nabla \mathcal{L}(\mathbf{w}(t))\|^2 + \frac{1}{2} \|\nabla \mathcal{L}(\mathbf{w}(t)) - \nabla \mathcal{L}(\mathbf{u}(t))\|^2.
\end{aligned}$$

By Assumption 3. (S),  $\ell$  is  $H$ -smooth, which further leads to

$$\begin{aligned}
& \|\nabla \mathcal{L}(\mathbf{w}(t)) - \nabla \mathcal{L}(\mathbf{u}(t))\|^2 = \left\| \sum_{\tilde{\mathbf{x}} \in \tilde{\mathcal{S}}} (\ell'(\langle \tilde{\mathbf{x}}, \mathbf{w}(t) \rangle) - \ell'(\langle \tilde{\mathbf{x}}, \mathbf{u}(t) \rangle)) \tilde{\mathbf{x}} \right\|^2 \\
& \leq \left( \sum_{\tilde{\mathbf{x}} \in \tilde{\mathcal{S}}} |\ell'(\langle \tilde{\mathbf{x}}, \mathbf{w}(t) \rangle) - \ell'(\langle \tilde{\mathbf{x}}, \mathbf{u}(t) \rangle)| \right)^2 \leq \left( H \sum_{\tilde{\mathbf{x}} \in \tilde{\mathcal{S}}} |\langle \tilde{\mathbf{x}}, \mathbf{w}(t) \rangle - \langle \tilde{\mathbf{x}}, \mathbf{u}(t) \rangle| \right)^2 \\
& \leq H^2 N^2 \|\mathbf{w}(t) - \mathbf{u}(t)\|^2 \stackrel{(\square)}{\leq} \frac{H^2 N^2 \beta^2}{(1-\beta)^2} \|\mathbf{w}(t) - \mathbf{w}(t-1)\|^2 \\
& = \frac{H^2 N^2 \beta^2}{(1-\beta)^2} \left\| \sum_{s=1}^{t-1} \eta \beta^{t-1-s} \nabla \mathcal{L}_{\mathbf{B}(s)}(\mathbf{w}(s)) \right\|^2 = H^2 N^2 \tilde{\eta}^2 \beta^2 \left\| \sum_{s=1}^{t-1} \beta^{t-1-s} \nabla \mathcal{L}_{\mathbf{B}(s)}(\mathbf{w}(s)) \right\|^2 \\
& \stackrel{(\diamond)}{\leq} \frac{H^2 N^2 \tilde{\eta}^2 \beta^2}{\gamma^2} \left( \sum_{s=1}^{t-1} \beta^{t-1-s} \|\nabla \mathcal{L}_{\mathbf{B}(s)}(\mathbf{w}(s))\| \right)^2 \\
& \stackrel{(\clubsuit)}{\leq} \frac{H^2 N^2 \tilde{\eta}^2 \beta^2}{\gamma^2} \left( \sum_{s=1}^{t-1} \beta^{t-1-s} \|\nabla \mathcal{L}_{\mathbf{B}(s)}(\mathbf{w}(s))\|^2 \right) \left( \sum_{s=1}^{t-1} \beta^{t-1-s} \right) \\
& \leq \frac{H^2 N^2 \tilde{\eta}^2 \beta^2}{\gamma^2 (1-\beta)} \left( \sum_{s=1}^{t-1} \beta^{t-1-s} \|\nabla \mathcal{L}_{\mathbf{B}(s)}(\mathbf{w}(s))\|^2 \right), \tag{36}
\end{aligned}$$

where Inequality  $(\square)$  is due to  $\beta(\mathbf{w}(t) - \mathbf{w}(t-1)) = (1-\beta)(\mathbf{u}(t) - \mathbf{w}(t))$  by Eq. (9), Inequality  $(\diamond)$  is due to triangular inequality, and Inequality  $(\clubsuit)$  is due to Cauchy-Schwartz Inequality.

Combining Eqs. (35) and (36), we have

$$\begin{aligned} & \mathbb{E}_{\mathbf{w}(t+1)}[\mathcal{L}(\mathbf{u}(t+1)) | \{\mathbf{w}(s)\}_{s=1}^t] \\ & \leq \mathcal{L}(\mathbf{u}(t)) - \tilde{\eta} \left( 1 - \left( \frac{1}{2} + \frac{HN}{2b\gamma^2} \right) \tilde{\eta} \right) \|\nabla \mathcal{L}(\mathbf{w}(t))\|^2 + \frac{1}{2} \frac{H^2 N^2 \tilde{\eta}^2 \beta^2}{\gamma^2(1-\beta)} \left( \sum_{s=1}^{t-1} \beta^{t-1-s} \|\nabla \mathcal{L}_{\mathbf{B}(s)}(\mathbf{w}(s))\|^2 \right), \end{aligned}$$

which by taking expectation to  $\{\mathbf{w}(s)\}_{s=1}^t$  leads to

$$\begin{aligned} & \mathbb{E}_{\mathcal{F}_{t+1}}[\mathcal{L}(\mathbf{u}(t+1))] = \mathbb{E}_{\{\mathbf{w}(s)\}_{s=1}^{t+1}}[\mathcal{L}(\mathbf{u}(t+1))] \\ & \leq \mathbb{E}_{\mathcal{F}_t} \mathcal{L}(\mathbf{u}(t)) - \mathbb{E}_{\mathcal{F}_t} \tilde{\eta} \left( 1 - \left( \frac{1}{2} + \frac{HN}{2b\gamma^2} \right) \tilde{\eta} \right) \|\nabla \mathcal{L}(\mathbf{w}(t))\|^2 + \mathbb{E}_{\mathcal{F}_t} \frac{1}{2} \frac{H^2 N^2 \tilde{\eta}^2 \beta^2}{\gamma^2(1-\beta)} \left( \sum_{s=1}^{t-1} \beta^{t-1-s} \|\nabla \mathcal{L}_{\mathbf{B}(s)}(\mathbf{w}(s))\|^2 \right) \\ & = \mathbb{E}_{\mathcal{F}_t} \mathcal{L}(\mathbf{u}(t)) - \tilde{\eta} \left( 1 - \left( \frac{1}{2} + \frac{HN}{2b\gamma^2} \right) \tilde{\eta} \right) \mathbb{E}_{\mathcal{F}_t} \|\nabla \mathcal{L}(\mathbf{w}(t))\|^2 + \frac{1}{2} \frac{H^2 N^2 \tilde{\eta}^2 \beta^2}{\gamma^2(1-\beta)} \left( \sum_{s=1}^{t-1} \beta^{t-1-s} \mathbb{E}_{\mathcal{F}_{s+1}} \|\nabla \mathcal{L}_{\mathbf{B}(s)}(\mathbf{w}(s))\|^2 \right). \end{aligned}$$

By Lemma 12 and  $\eta < \min\{\frac{1-\beta}{1+\frac{HN}{b\gamma^2}}, \frac{(1-\beta)^3 \gamma^4 b}{2H^2 N^3 \beta^2}\}$ , we further have

$$\begin{aligned} & \mathbb{E}_{\mathcal{F}_{t+1}}[\mathcal{L}(\mathbf{u}(t+1))] = \mathbb{E}_{\{\mathbf{w}(s)\}_{s=1}^{t+1}}[\mathcal{L}(\mathbf{u}(t+1))] \\ & \leq \mathbb{E}_{\mathcal{F}_t} \mathcal{L}(\mathbf{u}(t)) - \tilde{\eta} \left( 1 - \left( \frac{1}{2} + \frac{HN}{2b\gamma^2} \right) \tilde{\eta} \right) \mathbb{E}_{\mathcal{F}_t} \|\nabla \mathcal{L}(\mathbf{w}(t))\|^2 \\ & \quad + \frac{1}{2} \frac{H^2 N^3 \tilde{\eta}^2 \beta^2}{\gamma^4 b(1-\beta)} \left( \sum_{s=1}^{t-1} \beta^{t-1-s} \mathbb{E}_{\mathcal{F}_s} \|\nabla \mathcal{L}(\mathbf{w}(s))\|^2 \right) \\ & \leq \mathbb{E}_{\mathcal{F}_t} \mathcal{L}(\mathbf{u}(t)) - \frac{\tilde{\eta}}{2} \mathbb{E}_{\mathcal{F}_t} \|\nabla \mathcal{L}(\mathbf{w}(t))\|^2 + \frac{1-\beta}{4} \tilde{\eta} \left( \sum_{s=1}^{t-1} \beta^{t-1-s} \mathbb{E}_{\mathcal{F}_s} \|\nabla \mathcal{L}(\mathbf{w}(s))\|^2 \right). \quad (37) \end{aligned}$$

Summing Eq. (37) with  $t$  from 1 to  $T$  leads to

$$\begin{aligned} & \mathbb{E}[\mathcal{L}(\mathbf{u}(T+1))] = \mathbb{E}_{\mathcal{F}_{T+1}}[\mathcal{L}(\mathbf{u}(T+1))] \\ & \leq \mathcal{L}(\mathbf{u}(1)) - \sum_{t=1}^T \frac{\tilde{\eta}}{2} \mathbb{E}_{\mathcal{F}_t} \|\nabla \mathcal{L}(\mathbf{w}(t))\|^2 + \sum_{t=1}^T \frac{1-\beta}{4} \tilde{\eta} \left( \sum_{s=1}^{t-1} \beta^{t-1-s} \mathbb{E}_{\mathcal{F}_s} \|\nabla \mathcal{L}(\mathbf{w}(s))\|^2 \right) \\ & = \mathcal{L}(\mathbf{u}(1)) - \sum_{t=1}^T \frac{\tilde{\eta}}{2} \mathbb{E}_{\mathcal{F}_t} \|\nabla \mathcal{L}(\mathbf{w}(t))\|^2 + \sum_{s=1}^{T-1} \frac{1-\beta}{4} \tilde{\eta} \left( \sum_{t=s+1}^T \beta^{t-1-s} \mathbb{E}_{\mathcal{F}_s} \|\nabla \mathcal{L}(\mathbf{w}(s))\|^2 \right) \\ & \leq \mathcal{L}(\mathbf{u}(1)) - \sum_{t=1}^T \frac{\tilde{\eta}}{2} \mathbb{E}_{\mathcal{F}_t} \|\nabla \mathcal{L}(\mathbf{w}(t))\|^2 + \sum_{s=1}^{T-1} \frac{1}{4} \tilde{\eta} \left( \mathbb{E}_{\mathcal{F}_s} \|\nabla \mathcal{L}(\mathbf{w}(s))\|^2 \right) \\ & \leq \mathcal{L}(\mathbf{u}(1)) - \sum_{t=1}^T \frac{\tilde{\eta}}{4} \mathbb{E}_{\mathcal{F}_t} \|\nabla \mathcal{L}(\mathbf{w}(t))\|^2 \\ & = \mathcal{L}(\mathbf{u}(1)) - \sum_{t=1}^T \frac{\tilde{\eta}}{4} \mathbb{E} \|\nabla \mathcal{L}(\mathbf{w}(t))\|^2. \end{aligned}$$

The proof is completed.  $\square$

As  $\mathcal{L}(\mathbf{u}(1))$  is upper bounded, we have the following corollary given by Lemma 4.

**Corollary 5.** *Let all conditions in Theorem 3 hold. Then, we have*

$$\sum_{t=1}^{\infty} \mathbb{E} \|\nabla \mathcal{L}(\mathbf{w}(t))\|^2 < \infty. \quad (38)$$

Consequently,

$$\sum_{t=1}^{\infty} \|\nabla \mathcal{L}(\mathbf{w}(t))\|^2 < \infty$$

and

$$\langle \mathbf{w}(t), \tilde{\mathbf{x}} \rangle \rightarrow \infty, \forall \tilde{\mathbf{x}} \in \tilde{\mathcal{S}}$$

hold almost surely.

*Proof.* By Lemma 4, we have for any  $T > 1$ ,

$$\sum_{t=1}^T \frac{\tilde{\eta}}{4} \mathbb{E} \|\nabla \mathcal{L}(\mathbf{w}(t))\|^2 \leq \mathcal{L}(\mathbf{u}(1)) - \mathbb{E}[\mathcal{L}(\mathbf{u}(T+1))] \leq \mathcal{L}(\mathbf{u}(1)) < \infty,$$

which completes the proof of Eq. (38). The rest of claims follows immediately by Fubini's Theorem and Assumption 2.

The proof is completed.  $\square$

## B.2.2 BOUNDING THE ORTHOGONAL PART

Similar to the case of GDM, we define  $\tilde{\mathbf{w}}$  as the solution of Eq. (12) with  $C_2 = \frac{\eta}{(1-\beta)N}$ . We also let  $\mathbf{n}(t)$  be given by Lemma 9, and define  $\mathbf{r}(t)$  in this case as

$$\mathbf{r}(t) \triangleq \mathbf{w}(t) - \log(t)\hat{\mathbf{w}} - \tilde{\mathbf{w}} - \mathbf{n}(t). \quad (39)$$

As  $\tilde{\mathbf{w}}$  is a constant vector, and  $\|\mathbf{n}(t)\| \rightarrow 0$  as  $t \rightarrow \infty$ , we have  $\mathbf{w}(t) - \log(t)\hat{\mathbf{w}}$  has bounded norm if and only if  $\|\mathbf{r}(t)\|$  is upper bounded. Similar to the GDM case, we have the following equivalent condition of that  $\|\mathbf{r}(t)\|$  is bounded.

**Lemma 17.** *Let all conditions in Theorem 3 hold. Then,  $\|\mathbf{r}(t)\|$  is bounded almost surely if and only if function  $g(t)$  is upper bounded almost surely, where  $g : \mathbb{Z}^+ \rightarrow \mathbb{R}$  is defined as*

$$g(t) \triangleq \frac{1}{2} \|\mathbf{r}(t)\|^2 + \frac{\beta}{1-\beta} \langle \mathbf{r}(t), \mathbf{w}(t) - \mathbf{w}(t-1) \rangle - \frac{\beta}{1-\beta} \sum_{\tau=2}^t \langle \mathbf{r}(\tau) - \mathbf{r}(\tau-1), \mathbf{w}(\tau) - \mathbf{w}(\tau-1) \rangle. \quad (40)$$

*Proof.* To begin with, we prove that almost surely  $|\sum_{\tau=2}^t \langle \mathbf{r}(\tau) - \mathbf{r}(\tau-1), \mathbf{w}(\tau) - \mathbf{w}(\tau-1) \rangle|$  is upper bounded for any  $t$ . By Corollary 5, we have almost surely

$$\sum_{t=1}^{\infty} \|\nabla \mathcal{L}(\mathbf{w}(t))\|^2 < \infty.$$

On the other hand, for any  $\mathbf{w}$ , we have

$$\begin{aligned} \|\nabla \mathcal{L}_{\mathcal{B}(t)}(\mathbf{w})\| &= \frac{1}{b} \left\| \sum_{\tilde{\mathbf{x}} \in \mathcal{T}(\mathcal{B}(t))} \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}} \rangle) \tilde{\mathbf{x}} \right\| \\ &\leq -\frac{1}{b} \sum_{\tilde{\mathbf{x}} \in \mathcal{T}(\mathcal{B}(t))} \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}} \rangle) < -\frac{1}{b} \sum_{\tilde{\mathbf{x}} \in \mathcal{T}(\mathcal{S})} \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}} \rangle) \\ &\leq -\frac{1}{b} \sum_{\tilde{\mathbf{x}} \in \mathcal{T}(\mathcal{S})} \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}} \rangle) \langle \hat{\mathbf{w}}, \tilde{\mathbf{x}} \rangle \leq \frac{N}{b} \left\| \frac{1}{N} \sum_{\tilde{\mathbf{x}} \in \mathcal{T}(\mathcal{S})} \ell'(\langle \mathbf{w}, \tilde{\mathbf{x}} \rangle) \tilde{\mathbf{x}} \right\| \|\hat{\mathbf{w}}\| \\ &= \frac{N}{b\gamma} \|\nabla \mathcal{L}(\mathbf{w})\|. \end{aligned}$$

Therefore, we have almost surely,

$$\sum_{t=1}^{\infty} \|\nabla \mathcal{L}_{\mathcal{B}(t)}(\mathbf{w}(t))\|^2 < \infty,$$

which further leads to almost surely

$$\begin{aligned}
\sum_{t=1}^{\infty} \|\mathbf{w}(t+1) - \mathbf{w}(t)\|^2 &\leq \eta^2 \sum_{t=1}^{\infty} \left\| \sum_{s=1}^t \beta^{t-s} \nabla \mathcal{L}_{\mathbf{B}(s)}(\mathbf{w}(s)) \right\|^2 \\
&\leq \eta^2 \sum_{t=1}^{\infty} \left( \sum_{s=1}^t \beta^{t-s} \|\nabla \mathcal{L}_{\mathbf{B}(s)}(\mathbf{w}(s))\| \right)^2 \\
&\leq \eta^2 \sum_{t=1}^{\infty} \left( \sum_{s=1}^t \beta^{t-s} \|\nabla \mathcal{L}_{\mathbf{B}(s)}(\mathbf{w}(s))\|^2 \right) \left( \sum_{s=1}^t \beta^{t-s} \right) \\
&\leq \frac{\eta^2}{1-\beta} \sum_{s=1}^{\infty} \|\nabla \mathcal{L}_{\mathbf{B}(s)}(\mathbf{w}(s))\|^2 \sum_{t=s}^{\infty} \beta^{t-s} \\
&= \frac{\eta^2}{1-\beta} \sum_{s=1}^{\infty} \|\nabla \mathcal{L}_{\mathbf{B}(s)}(\mathbf{w}(s))\|^2 < \infty.
\end{aligned}$$

By the definition of  $\mathbf{r}(t)$  (Eq. (39)), we further have

$$\begin{aligned}
&\left| \sum_{\tau=2}^t \langle \mathbf{r}(\tau) - \mathbf{r}(\tau-1), \mathbf{w}(\tau) - \mathbf{w}(\tau-1) \rangle \right| \\
&\leq \sum_{\tau=2}^t |\langle \mathbf{r}(\tau) - \mathbf{r}(\tau-1), \mathbf{w}(\tau) - \mathbf{w}(\tau-1) \rangle| \\
&= \sum_{\tau=2}^t \left| \left\langle \mathbf{w}(\tau) - \mathbf{w}(\tau-1) - \log \left( \frac{\tau+1}{\tau} \right) - (\mathbf{n}(\tau) - \mathbf{n}(\tau-1)), \mathbf{w}(\tau) - \mathbf{w}(\tau-1) \right\rangle \right| \\
&\leq \sum_{\tau=2}^t \|\mathbf{w}(\tau) - \mathbf{w}(\tau-1)\|^2 + \sum_{\tau=2}^t \left| \left\langle -\log \left( \frac{\tau+1}{\tau} \right) - (\mathbf{n}(\tau) - \mathbf{n}(\tau-1)), \mathbf{w}(\tau) - \mathbf{w}(\tau-1) \right\rangle \right| \\
&\leq \frac{3}{2} \sum_{\tau=2}^t \|\mathbf{w}(\tau) - \mathbf{w}(\tau-1)\|^2 + \frac{1}{2} \sum_{\tau=2}^t \left\| -\log \left( \frac{\tau+1}{\tau} \right) - (\mathbf{n}(\tau) - \mathbf{n}(\tau-1)) \right\|^2 \\
&\stackrel{(*)}{\leq} \frac{3}{2} \sum_{\tau=2}^t \|\mathbf{w}(\tau) - \mathbf{w}(\tau-1)\|^2 + \frac{1}{2} \sum_{\tau=2}^t \mathcal{O} \left( \frac{1}{\tau} \right)^2 < \infty,
\end{aligned}$$

where Inequality (\*) is due to  $\|\mathbf{n}(\tau) - \mathbf{n}(\tau-1)\| = \mathcal{O}(\frac{1}{\tau})$  and  $\log \frac{\tau+1}{\tau} = \mathcal{O}(\frac{1}{\tau})$ .

Therefore,  $g(t)$  is upper bounded almost surely is equivalent to  $\frac{1}{2} \|\mathbf{r}(t)\|^2 + \frac{\beta}{1-\beta} \langle \mathbf{r}(t), \mathbf{w}(t) - \mathbf{w}(t-1) \rangle$  is upper bounded, which can be shown to be equivalent with  $\|\mathbf{r}(t)\|$  is bounded following the same routine as Lemma 2.

The proof is completed.  $\square$

As the case of GDM, we only need to prove  $g(t)$  is upper bounded to complete the proof of Theorem 3.

**Lemma 18.** *Let all conditions in Theorem 3 hold. We have  $g(t)$  is upper bounded.*

*Proof.* Following the same routine as Lemma 2, we have

$$\begin{aligned}
&g(t+1) - g(t) \\
&= \frac{1}{2} \|\mathbf{r}(t+1) - \mathbf{r}(t)\|^2 + \langle \mathbf{r}(t), \mathbf{r}(t+1) - \mathbf{r}(t) \rangle + \left\langle \mathbf{r}(t), -\frac{\eta}{1-\beta} \nabla \mathcal{L}_{\mathbf{B}(t)}(\mathbf{w}(t)) - (\mathbf{w}(t+1) - \mathbf{w}(t)) \right\rangle,
\end{aligned}$$

where  $\sum_{t=1}^{\infty} \|\mathbf{r}(t+1) - \mathbf{r}(t)\|^2$  is upper bounded.

On the other hand, by the definition of  $\mathbf{r}(t)$  (Eq. (39)), we have

$$\begin{aligned} & \mathbf{r}(t+1) - \mathbf{r}(t) \\ &= \mathbf{w}(t+1) - \mathbf{w}(t) - \log\left(\frac{t+1}{t}\right) \hat{\mathbf{w}} - \mathbf{n}(t+1) + \mathbf{n}(t), \end{aligned}$$

while by Lemma 9,

$$\frac{N}{b} \frac{1}{t} \sum_{i: \tilde{\mathbf{x}}_i \in \mathcal{T}(\mathbf{B}(t) \cap \mathcal{S}_s)} \mathbf{v}_i \tilde{\mathbf{x}}_i = \log\left(\frac{t+1}{t}\right) \hat{\mathbf{w}} + \mathbf{n}(t+1) - \mathbf{n}(t).$$

Combining the above two equations, we further have

$$\mathbf{r}(t+1) - \mathbf{r}(t) = \mathbf{w}(t+1) - \mathbf{w}(t) - \frac{N}{bt} \sum_{i: \tilde{\mathbf{x}}_i \in \mathcal{T}(\mathbf{B}(t) \cap \mathcal{S}_s)} \mathbf{v}_i \tilde{\mathbf{x}}_i,$$

which further indicates

$$g(t+1) - g(t) = \frac{1}{2} \|\mathbf{r}(t+1) - \mathbf{r}(t)\|^2 + \left\langle \mathbf{r}(t), -\frac{\eta}{1-\beta} \nabla \mathcal{L}_{\mathbf{B}(t)}(\mathbf{w}(t)) - \frac{N}{bt} \sum_{i: \tilde{\mathbf{x}}_i \in \mathcal{T}(\mathbf{B}(t) \cap \mathcal{S}_s)} \mathbf{v}_i \tilde{\mathbf{x}}_i \right\rangle.$$

Therefore, we only need to prove  $\sum_{t=1}^{\infty} \langle \mathbf{r}(t), -\frac{\eta}{1-\beta} \nabla \mathcal{L}(\mathbf{w}(t)) - \frac{N}{bt} \sum_{i: \tilde{\mathbf{x}}_i \in \mathcal{T}(\mathbf{B}(t) \cap \mathcal{S}_s)} \mathbf{v}_i \tilde{\mathbf{x}}_i \rangle < \infty$ . By directly applying the form of  $\nabla \mathcal{L}(\mathbf{w}(t))$ , we have

$$\begin{aligned} & \left\langle \mathbf{r}(t), -\frac{\eta}{1-\beta} \nabla \mathcal{L}_{\mathbf{B}(t)}(\mathbf{w}(t)) - \frac{N}{bt} \sum_{i: \tilde{\mathbf{x}}_i \in \mathcal{T}(\mathbf{B}(s) \cap \mathcal{S}_s)} \mathbf{v}_i \tilde{\mathbf{x}}_i \right\rangle \\ &= \left\langle \mathbf{r}(t), -\frac{\eta}{(1-\beta)b} \sum_{i: \tilde{\mathbf{x}}_i \in \mathcal{T}(\mathbf{B}(s))} \ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle) \tilde{\mathbf{x}}_i - \frac{N}{bt} \sum_{i: \tilde{\mathbf{x}}_i \in \mathcal{T}(\mathbf{B}(s) \cap \mathcal{S}_s)} \mathbf{v}_i \tilde{\mathbf{x}}_i \right\rangle \\ &= \frac{\eta}{(1-\beta)b} \left\langle \mathbf{r}(t), -\sum_{i: \tilde{\mathbf{x}}_i \in \mathcal{T}(\mathbf{B}(s))} \ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle) \tilde{\mathbf{x}}_i - \frac{N(1-\beta)}{\eta} \frac{1}{t} \sum_{i: \tilde{\mathbf{x}}_i \in \mathcal{T}(\mathbf{B}(s) \cap \mathcal{S}_s)} \mathbf{v}_i \tilde{\mathbf{x}}_i \right\rangle \\ &= \frac{\eta}{(1-\beta)b} \left\langle \mathbf{r}(t), -\sum_{i: \tilde{\mathbf{x}}_i \in \mathcal{T}(\mathbf{B}(s))} \ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle) \tilde{\mathbf{x}}_i - \frac{1}{t} \sum_{i: \tilde{\mathbf{x}}_i \in \mathcal{T}(\mathbf{B}(s) \cap \mathcal{S}_s)} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \tilde{\mathbf{x}}_i \right\rangle \\ &= \frac{\eta}{(1-\beta)b} \sum_{i: \tilde{\mathbf{x}}_i \in \mathcal{T}(\mathbf{B}(s) \cap \mathcal{S}_s)} \left( -\ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle) - \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\ &+ \frac{\eta}{(1-\beta)b} \sum_{i: \tilde{\mathbf{x}}_i \in \mathcal{T}(\mathbf{B}(s) \cap \mathcal{S}_s^c)} \langle \mathbf{r}(t), -\ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle) \tilde{\mathbf{x}}_i \rangle. \end{aligned}$$

Let  $A_6(t) = \sum_{i: \tilde{\mathbf{x}}_i \in \mathcal{T}(\mathbf{B}(s) \cap \mathcal{S}_s)} \left( -\ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle) - \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle$ , and  $A_7(t) = \sum_{i: \tilde{\mathbf{x}}_i \in \mathcal{T}(\mathbf{B}(s) \cap \mathcal{S}_s^c)} \langle \mathbf{r}(t), -\ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle) \tilde{\mathbf{x}}_i \rangle$ . We will investigate these two terms respectively.

As  $\langle \mathbf{w}(t), \tilde{\mathbf{x}} \rangle \rightarrow \infty, \forall \tilde{\mathbf{x}} \in \mathcal{T}(\mathcal{S})$ , a.s., we have a.s., there exists a large enough time  $t_0$ , s.t.,  $\forall t \geq t_0, \forall \tilde{\mathbf{x}} \in \mathcal{T}(\mathcal{S})$ ,

$$\begin{aligned} -\ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}} \rangle) &\leq (1 + e^{-\mu_+ \langle \mathbf{w}(t), \tilde{\mathbf{x}} \rangle}) e^{-\langle \mathbf{w}(t), \tilde{\mathbf{x}} \rangle}, \\ -\ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle) &\geq (1 - e^{-\mu_- \langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle}) e^{-\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle}, \\ &\langle \tilde{\mathbf{x}}, \mathbf{w}(t) \rangle > 0. \end{aligned}$$



Therefore,

$$\begin{aligned}
A_7(t) &\leq \sum_{i:\tilde{\mathbf{x}}_i \in \mathcal{T}(\mathbf{B}(s) \cap \mathcal{S}_s^c)} -\ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \mathbb{1}_{\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \geq 0} \\
&\leq \sum_{i:\tilde{\mathbf{x}}_i \in \mathcal{T}(\mathbf{B}(s) \cap \mathcal{S}_s^c)} (1 + e^{-\mu_+ \langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle}) e^{-\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle} \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \mathbb{1}_{\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \geq 0} \\
&\leq 2 \sum_{i:\tilde{\mathbf{x}}_i \in \mathcal{T}(\mathbf{B}(s) \cap \mathcal{S}_s^c)} e^{-\langle \mathbf{r}(t) + \log(t)\hat{\mathbf{w}} + \tilde{\mathbf{w}} + \mathbf{n}(t), \tilde{\mathbf{x}}_i \rangle} \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \mathbb{1}_{\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \geq 0} \\
&\stackrel{(\star)}{\leq} 2 \sum_{i:\tilde{\mathbf{x}}_i \in \mathcal{T}(\mathbf{B}(s) \cap \mathcal{S}_s^c)} \frac{1}{t^\theta} e^{-\langle \tilde{\mathbf{w}} + \mathbf{n}(t), \tilde{\mathbf{x}}_i \rangle} e^{-\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle} \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \mathbb{1}_{\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \geq 0} \\
&\stackrel{(\dagger)}{\leq} \frac{2}{e} \frac{1}{t^\theta} \sum_{i:\tilde{\mathbf{x}}_i \in \mathcal{T}(\mathbf{B}(s) \cap \mathcal{S}_s^c)} e^{-\langle \tilde{\mathbf{w}} + \mathbf{n}(t), \tilde{\mathbf{x}}_i \rangle} \mathbb{1}_{\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \geq 0} \\
&\stackrel{(\circ)}{=} \mathcal{O}\left(\frac{1}{t^\theta}\right),
\end{aligned}$$

where Inequality.  $(\star)$  is due the definition of  $\theta$  (Eq. (31)), Inequality.  $(\dagger)$  is due to  $e^{-\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle} \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \leq e^{-1}$ , and Eq.  $(\circ)$  is due to  $\lim_{t \rightarrow \infty} e^{-\langle \tilde{\mathbf{w}} + \mathbf{n}(t), \tilde{\mathbf{x}}_i \rangle} = e^{-\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle}$ . Thus,

$$\sum_{t=1}^{\infty} A_7(t) < \infty.$$

On the other hand,  $A_6(t)$  can be rewritten as

$$\begin{aligned}
A_6(t) &= \sum_{i:\tilde{\mathbf{x}}_i \in \mathcal{T}(\mathbf{B}(s) \cap \mathcal{S}_s)} \left( -\ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle) - \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \mathbb{1}_{\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \geq 0} \\
&\quad + \sum_{i:\tilde{\mathbf{x}}_i \in \mathcal{T}(\mathbf{B}(s) \cap \mathcal{S}_s)} \left( -\ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle) - \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \mathbb{1}_{\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle < 0}.
\end{aligned}$$

If  $\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \geq 0$ , we have for  $\varepsilon < 0.5$ ,

$$\begin{aligned}
&\left( -\ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle) - \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\
&\leq \left( \left( 1 + e^{-\mu_+ \langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle} \right) e^{-\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle} - \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\
&= \left( \left( 1 + e^{-\mu_+ \langle \mathbf{r}(t) + \log(t)\hat{\mathbf{w}} + \tilde{\mathbf{w}} + \mathbf{n}(t), \tilde{\mathbf{x}}_i \rangle} \right) e^{-\langle \mathbf{r}(t) + \log(t)\hat{\mathbf{w}} + \tilde{\mathbf{w}} + \mathbf{n}(t), \tilde{\mathbf{x}}_i \rangle} - \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\
&= \left( \left( 1 + e^{-\mu_+ \langle \mathbf{r}(t) + \log(t)\hat{\mathbf{w}} + \tilde{\mathbf{w}} + \mathbf{n}(t), \tilde{\mathbf{x}}_i \rangle} \right) e^{-\langle \mathbf{r}(t) + \mathbf{n}(t), \tilde{\mathbf{x}}_i \rangle} - 1 \right) \frac{1}{t} \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle e^{-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle} \\
&\leq \left( \left( 1 + \frac{1}{t^{\mu_+}} e^{-\mu_+ \langle \tilde{\mathbf{w}} + \mathbf{n}(t), \tilde{\mathbf{x}}_i \rangle} \right) e^{-\langle \mathbf{r}(t) + \mathbf{n}(t), \tilde{\mathbf{x}}_i \rangle} - 1 \right) \frac{1}{t} \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle e^{-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle} \\
&\stackrel{(\bullet)}{=} \left( \left( 1 + \mathcal{O}\left(\frac{1}{t^{\mu_+}}\right) \right) \left( 1 + \mathcal{O}\left(\frac{1}{t^{0.5-\varepsilon}}\right) \right) e^{-\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle} - 1 \right) \frac{1}{t} \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle e^{-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle} \\
&= \left( e^{-\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle} - 1 \right) \frac{1}{t} \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle e^{-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle} + \frac{1}{t} \mathcal{O}\left(\frac{1}{t^{\min\{\mu_+, 0.5-\varepsilon\}}}\right) e^{-\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle} \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle e^{-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle} \\
&\leq \frac{1}{t} \mathcal{O}\left(\frac{1}{t^{\min\{\mu_+, 0.5-\varepsilon\}}}\right) e^{-\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle} \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle e^{-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle} \\
&\stackrel{(\circ)}{=} \mathcal{O}\left(\frac{1}{t^{\min\{1+\mu_+, 1.5-\varepsilon\}}}\right),
\end{aligned}$$

where Eq.  $(\bullet)$  is due to  $\mathbf{n}(t) = \mathcal{O}\left(\frac{1}{t^{0.5-\varepsilon}}\right)$ , and Eq.  $(\circ)$  is due to  $e^{-\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle} \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \leq \frac{1}{e}$ .

On the other hand, if  $\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle < 0$ , we have

$$\begin{aligned} & \left( -\ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle) - \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\ & \leq \left( \left( 1 - e^{-\mu_- \langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle} \right) e^{-\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle} - \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\ & = \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle} \left( -1 + \left( 1 - e^{-\mu_- \langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle} \right) e^{-\langle \mathbf{r}(t) + \mathbf{n}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \end{aligned}$$

Specifically, if  $\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \geq -t^{-0.5 \min\{\mu_-, 0.5\}}$ ,

$$\begin{aligned} & \left| \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle} \left( -1 + \left( 1 - e^{-\mu_- \langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle} \right) e^{-\langle \mathbf{r}(t) + \mathbf{n}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \right| \\ & \leq \frac{1}{t^{1+0.5 \min\{\mu_-, 0.5\}}} e^{-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle} \left| -1 + \left( 1 - e^{-\mu_- \langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle} \right) e^{-\langle \mathbf{r}(t) + \mathbf{n}(t), \tilde{\mathbf{x}}_i \rangle} \right| \\ & \stackrel{(\square)}{=} \mathcal{O} \left( \frac{1}{t^{1+0.5 \min\{\mu_-, 0.5\}}} \right), \end{aligned}$$

where Eq. (□) is due to that as  $\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle \rightarrow \infty$  and  $t^{-0.5 \min\{\mu_-, 0.5\}} \rightarrow 0$  as  $t \rightarrow \infty$ , there exists a large enough time  $T$ , s.t.,  $\forall t > T$ , under the circumstance  $0 > \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \geq -t^{-0.5 \min\{\mu_-, 0.5\}}$ ,  $e^{-\langle \mathbf{r}(t) + \mathbf{n}(t), \tilde{\mathbf{x}}_i \rangle} < 1$  and  $e^{-\mu_- \langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle} < 1$ .

If  $-2 \leq \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle < -t^{-0.5 \min\{\mu_-, 0.5\}}$ , then, for large enough  $t$ ,  $|\langle \tilde{\mathbf{x}}_i, \mathbf{n}(t) \rangle| < 2$ ,  $1 - \frac{e^{\mu_- (-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle + 4)}}{t^{\mu_-}} > 0$ , and

$$\begin{aligned} & \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle} \left( -1 + \left( 1 - e^{-\mu_- \langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle} \right) e^{-\langle \mathbf{r}(t) + \mathbf{n}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\ & = \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle} \left( -1 + \left( 1 - e^{-\mu_- \langle \mathbf{r}(t) + \log(t) \tilde{\mathbf{w}} + \tilde{\mathbf{w}} + \mathbf{n}(t), \tilde{\mathbf{x}}_i \rangle} \right) e^{-\langle \mathbf{r}(t) + \mathbf{n}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\ & = \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle} \left( -1 + \left( 1 - \frac{e^{-\mu_- \langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle}}{t^{\mu_-}} e^{-\mu_- \langle \mathbf{r}(t) + \mathbf{n}(t), \tilde{\mathbf{x}}_i \rangle} \right) e^{-\langle \mathbf{r}(t) + \mathbf{n}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\ & \leq \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle} \left( -1 + \left( 1 - \frac{e^{\mu_- (-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle + 4)}}{t^{\mu_-}} \right) e^{-\mu_- \langle \mathbf{r}(t) + \mathbf{n}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\ & \leq \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle} \left( -1 + \left( 1 - \frac{e^{\mu_- (-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle + 4)}}{t^{\mu_-}} \right) (1 - \mu_- \langle \mathbf{r}(t) + \mathbf{n}(t), \tilde{\mathbf{x}}_i \rangle) \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\ & \leq \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle} \left( -1 + \left( 1 - \frac{e^{\mu_- (-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle + 4)}}{t^{\mu_-}} \right) \left( 1 + \mu_- t^{-0.5 \min\{\mu_-, 0.5\}} - \mu_- \langle \mathbf{n}(t), \tilde{\mathbf{x}}_i \rangle \right) \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\ & = \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle} \left( -1 + \left( 1 - \frac{e^{\mu_- (-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle + 4)}}{t^{\mu_-}} \right) \left( 1 + \mu_- t^{-0.5 \min\{\mu_-, 0.5\}} + \mathcal{O} \left( t^{-0.5 \min\{\mu_-, 0.5\}} \right) \right) \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\ & = \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle} \left( -1 + 1 + \mu_- t^{-0.5 \min\{\mu_-, 0.5\}} + \mathcal{O} \left( t^{-0.5 \min\{\mu_-, 0.5\}} \right) \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle < 0. \end{aligned}$$

If  $-2 > \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle$ , then for large enough time  $t$ ,  $e^{-\langle \mathbf{r}(t) + \mathbf{n}(t), \tilde{\mathbf{x}}_i \rangle} \geq e^{\frac{3}{2}}$ ,  $1 - e^{-\mu_- \langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle} \geq e^{-\frac{1}{2}}$ , and

$$\begin{aligned} & \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle} \left( -1 + \left( 1 - e^{-\mu_- \langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle} \right) e^{-\langle \mathbf{r}(t) + \mathbf{n}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \\ & \leq \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle} (-1 + e) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle < 0. \end{aligned}$$

Conclusively, if  $\langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle < 0$ , for large enough  $t$ , we have

$$\left( -\ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}}_i \rangle) - \frac{1}{t} e^{-\langle \tilde{\mathbf{w}}(t), \tilde{\mathbf{x}}_i \rangle} \right) \langle \mathbf{r}(t), \tilde{\mathbf{x}}_i \rangle \leq \mathcal{O} \left( \frac{1}{t^{1+0.5 \min\{\mu_-, 0.5\}}} \right),$$

which further indicates, for large enough  $t$ , we have

$$A_6(t) \leq \max \left\{ \mathcal{O} \left( \frac{1}{t^{1+0.5 \min\{\mu_-, 0.5\}}} \right), \mathcal{O} \left( \frac{1}{t^{\min\{1+\mu_+, 1.5-\varepsilon\}}} \right) \right\},$$

which indicates

$$\sum_{t=1}^{\infty} A_6(t) < \infty.$$

Therefore,

$$\begin{aligned} & \sum_{t=1}^{\infty} (g(t+1) - g(t)) \\ &= \sum_{t=1}^{\infty} \left( \frac{1}{2} \|\mathbf{r}(t+1) - \mathbf{r}(t)\|^2 + \left\langle \mathbf{r}(t), -\frac{\eta}{1-\beta} \nabla \mathcal{L}_{\mathbf{B}(t)}(\mathbf{w}(t)) - \frac{N}{bt} \sum_{i: \tilde{\mathbf{x}}_i \in \mathcal{T}(\mathbf{B}(t) \cap \mathcal{S}_s)} \mathbf{v}_i \tilde{\mathbf{x}}_i \right\rangle \right) \\ &= \sum_{t=1}^{\infty} \left( \frac{1}{2} \|\mathbf{r}(t+1) - \mathbf{r}(t)\|^2 + \frac{\eta}{1-\beta} A_6(t) + \frac{\eta}{1-\beta} A_7(t) \right) \\ &< \infty. \end{aligned}$$

The proof is completed.  $\square$

### B.2.3 EXPLANATION FOR PROPER LYAPUNOV FUNCTION

Based on the success of applying Lyapunov function  $\mathcal{L}(\mathbf{w}(t)) + \frac{\beta}{2\eta} \|\mathbf{w}(t) - \mathbf{w}(t-1)\|^2$  to analyze gradient descent with momentum, it is natural to try to extend this routine to analyze stochastic gradient descent with momentum. However, in this section, we will show such Lyapunov function is not proper to analyze SGDM as this will put constraints on the range of the momentum rate  $\beta$ . Specifically, at any step  $t$ , since the loss  $\mathcal{L}$  is  $H$  smooth at  $\mathbf{w}(t)$ , we can expand the loss  $\mathcal{L}$  in the same way as the GDM case:

$$\mathcal{L}(\mathbf{w}(t+1)) \leq \mathcal{L}(\mathbf{w}(t)) + \langle \mathbf{w}(t+1) - \mathbf{w}(t), \nabla \mathcal{L}(\mathbf{w}(t)) \rangle + \frac{H}{2} \|\mathbf{w}(t+1) - \mathbf{w}(t)\|^2.$$

By taking expectation with respect to  $\mathbf{w}(t+1)$  conditioning on  $\{\mathbf{w}(s)\}_{s=1}^t$  for both sides, we further obtain

$$\begin{aligned} & \mathbb{E}_{\mathbf{w}(t+1)} [\mathcal{L}(\mathbf{w}(t+1)) | \{\mathbf{w}(s)\}_{s=1}^t] \\ & \leq \mathcal{L}(\mathbf{w}(t)) + \langle \mathbb{E}_{\mathbf{w}(t+1)} [\mathbf{w}(t+1) - \mathbf{w}(t) | \{\mathbf{w}(s)\}_{s=1}^t], \nabla \mathcal{L}(\mathbf{w}(t)) \rangle + \frac{H}{2} \mathbb{E}_{\mathbf{w}(t+1)} [\|\mathbf{w}(t+1) - \mathbf{w}(t)\|^2 | \{\mathbf{w}(s)\}_{s=1}^t] \\ & \stackrel{(*)}{=} \mathcal{L}(\mathbf{w}(t)) + \frac{1}{\eta} \langle \mathbb{E}_{\mathbf{w}(t+1)} [\mathbf{w}(t+1) - \mathbf{w}(t) | \{\mathbf{w}(s)\}_{s=1}^t], \beta(\mathbf{w}(t) - \mathbf{w}(t-1)) - \mathbb{E}_{\mathbf{w}(t+1)} [\mathbf{w}(t+1) - \mathbf{w}(t) | \{\mathbf{w}(s)\}_{s=1}^t] \rangle \\ & \quad + \frac{H}{2} \mathbb{E}_{\mathbf{w}(t+1)} [\|\mathbf{w}(t+1) - \mathbf{w}(t)\|^2 | \{\mathbf{w}(s)\}_{s=1}^t] \\ & = \mathcal{L}(\mathbf{w}(t)) + \frac{\beta}{\eta} \langle (\mathbf{w}(t) - \mathbf{w}(t-1)), \mathbb{E}_{\mathbf{w}(t+1)} [\mathbf{w}(t+1) - \mathbf{w}(t) | \{\mathbf{w}(s)\}_{s=1}^t] \rangle \\ & \quad + \frac{H}{2} \mathbb{E}_{\mathbf{w}(t+1)} [\|\mathbf{w}(t+1) - \mathbf{w}(t)\|^2 | \{\mathbf{w}(s)\}_{s=1}^t] - \frac{1}{\eta} \|\mathbb{E}_{\mathbf{w}(t+1)} [\mathbf{w}(t+1) - \mathbf{w}(t) | \{\mathbf{w}(s)\}_{s=1}^t]\|^2 \\ & \leq \mathcal{L}(\mathbf{w}(t)) + \frac{\beta}{2\eta} \|\mathbf{w}(t) - \mathbf{w}(t-1)\|^2 + \frac{\beta}{2\eta} \|\mathbb{E}_{\mathbf{w}(t+1)} [\mathbf{w}(t+1) - \mathbf{w}(t) | \{\mathbf{w}(s)\}_{s=1}^t]\|^2 \\ & \quad + \frac{H}{2} \mathbb{E}_{\mathbf{w}(t+1)} [\|\mathbf{w}(t+1) - \mathbf{w}(t)\|^2 | \{\mathbf{w}(s)\}_{s=1}^t] - \frac{1}{\eta} \|\mathbb{E}_{\mathbf{w}(t+1)} [\mathbf{w}(t+1) - \mathbf{w}(t) | \{\mathbf{w}(s)\}_{s=1}^t]\|^2, \end{aligned}$$

where Eq. (\*) is because  $\mathbb{E}_{\mathbf{w}(t+1)} [\mathbf{w}(t+1) - \mathbf{w}(t) | \{\mathbf{w}(s)\}_{s=1}^t] = -\eta \nabla \mathcal{L}(\mathbf{w}(t)) + \beta(\mathbf{w}(t) - \mathbf{w}(t-1))$  due to the definition of SGDM (Eq. (1)). Rearranging the above inequal-

ity and taking expectations of both sides with respect to  $\{\mathbf{w}(s)\}_{s=1}^t$  leads to

$$\begin{aligned} & \mathbb{E}_{\mathbf{w}(t+1)} [\mathcal{L}(\mathbf{w}(t+1))] + \frac{2-\beta}{2\eta} \mathbb{E}_{\{\mathbf{w}(s)\}_{s=1}^t} \left\| \mathbb{E}_{\mathbf{w}(t+1)} [\mathbf{w}(t+1) - \mathbf{w}(t) | \{\mathbf{w}(s)\}_{s=1}^t] \right\|^2 \\ & - \frac{H}{2} \mathbb{E}_{\{\mathbf{w}(s)\}_{s=1}^{t+1}} [\|\mathbf{w}(t+1) - \mathbf{w}(t)\|^2] \\ & \leq \mathbb{E}_{\mathbf{w}(t)} \mathcal{L}(\mathbf{w}(t)) + \frac{\beta}{2\eta} \mathbb{E}_{\{\mathbf{w}(s)\}_{s=1}^t} \|\mathbf{w}(t) - \mathbf{w}(t-1)\|^2. \end{aligned} \quad (41)$$

On the other hand, we wish to obtain some positive constant  $\alpha$  from Eq. (41), such that (at least),

$$\begin{aligned} & \mathbb{E}_{\mathbf{w}(t+1)} [\mathcal{L}(\mathbf{w}(t+1))] + \alpha \mathbb{E}_{\{\mathbf{w}(s)\}_{s=1}^t} \|\mathbf{w}(t+1) - \mathbf{w}(t)\|^2 \\ & \leq \mathbb{E}_{\mathbf{w}(t)} \mathcal{L}(\mathbf{w}(t)) + \alpha \mathbb{E}_{\{\mathbf{w}(s)\}_{s=1}^t} \|\mathbf{w}(t) - \mathbf{w}(t-1)\|^2, \end{aligned} \quad (42)$$

which requires to lower bound  $\mathbb{E}_{\{\mathbf{w}(s)\}_{s=1}^t} \left\| \mathbb{E}_{\mathbf{w}(t+1)} [\mathbf{w}(t+1) - \mathbf{w}(t) | \{\mathbf{w}(s)\}_{s=1}^t] \right\|^2$  by  $\mathbb{E}_{\{\mathbf{w}(s)\}_{s=1}^{t+1}} \|\mathbf{w}(t+1) - \mathbf{w}(t)\|^2$ . However, in general cases,  $\mathbb{E}_{\{\mathbf{w}(s)\}_{s=1}^t} \left\| \mathbb{E}_{\mathbf{w}(t+1)} [\mathbf{w}(t+1) - \mathbf{w}(t) | \{\mathbf{w}(s)\}_{s=1}^t] \right\|^2$  is only upper bounded by  $\mathbb{E}_{\{\mathbf{w}(s)\}_{s=1}^{t+1}} \|\mathbf{w}(t+1) - \mathbf{w}(t)\|^2$  (Holder's Inequality), although in our case,  $\left\| \mathbb{E}_{\mathbf{w}(t+1)} [\mathbf{w}(t+1) - \mathbf{w}(t) | \{\mathbf{w}(s)\}_{s=1}^t] \right\|^2$  can be bounded as

$$\begin{aligned} & \left\| \mathbb{E}_{\mathbf{w}(t+1)} [\mathbf{w}(t+1) - \mathbf{w}(t) | \{\mathbf{w}(s)\}_{s=1}^t] \right\|^2 \\ & = \left\| -\eta \nabla \mathcal{L}(\mathbf{w}(t)) + \beta (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\|^2 \\ & = \left\| -\eta \nabla \mathcal{L}(\mathbf{w}(t)) \right\|^2 + \left\| \beta (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\|^2 + 2\beta\eta \langle \mathbf{w}(t) - \mathbf{w}(t-1), -\nabla \mathcal{L}(\mathbf{w}(t)) \rangle, \end{aligned}$$

while  $\mathbb{E}_{\mathbf{w}(t+1)} \left[ \|\mathbf{w}(t+1) - \mathbf{w}(t)\|^2 | \{\mathbf{w}(s)\}_{s=1}^t \right]$  can be calculated as

$$\begin{aligned} & \mathbb{E}_{\mathbf{w}(t+1)} \left[ \|\mathbf{w}(t+1) - \mathbf{w}(t)\|^2 | \{\mathbf{w}(s)\}_{s=1}^t \right] \\ & = \mathbb{E}_{\mathbf{B}(t)} \left\| -\eta \nabla \mathcal{L}_{\mathbf{B}(t)}(\mathbf{w}(t)) + \beta (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\|^2 \\ & = \mathbb{E}_{\mathbf{B}(t)} \left\| -\eta \nabla \mathcal{L}_{\mathbf{B}(t)}(\mathbf{w}(t)) \right\|^2 + \left\| \beta (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\|^2 + 2\eta\beta \mathbb{E}_{\mathbf{B}(t)} \langle -\nabla \mathcal{L}_{\mathbf{B}(t)}(\mathbf{w}(t)), \mathbf{w}(t) - \mathbf{w}(t-1) \rangle \\ & = \mathbb{E}_{\mathbf{B}(t)} \left\| -\eta \nabla \mathcal{L}_{\mathbf{B}(t)}(\mathbf{w}(t)) \right\|^2 + \left\| \beta (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\|^2 + 2\eta\beta \langle -\nabla \mathcal{L}(\mathbf{w}(t)), \mathbf{w}(t) - \mathbf{w}(t-1) \rangle \\ & \leq \frac{N}{b\gamma^2} \left\| -\eta \nabla \mathcal{L}(\mathbf{w}(t)) \right\|^2 + \left\| \beta (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\|^2 + 2\eta\beta \langle -\nabla \mathcal{L}(\mathbf{w}(t)), \mathbf{w}(t) - \mathbf{w}(t-1) \rangle \\ & \leq \frac{N}{b\gamma^2} \left( \left\| -\eta \nabla \mathcal{L}(\mathbf{w}(t)) \right\|^2 + \left\| \beta (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\|^2 + 2\eta\beta \langle -\nabla \mathcal{L}(\mathbf{w}(t)), \mathbf{w}(t) - \mathbf{w}(t-1) \rangle \right) \\ & = \frac{N}{b\gamma^2} \left\| \mathbb{E}_{\mathbf{w}(t+1)} [\mathbf{w}(t+1) - \mathbf{w}(t) | \{\mathbf{w}(s)\}_{s=1}^t] \right\|^2. \end{aligned} \quad (43)$$

By Eqs. (41) and (43), we have that to ensure Eq. (42), it is required that

$$\frac{2-\beta}{2\eta} \frac{b\gamma^2}{N} - \frac{H}{2} \leq \frac{\beta}{2\eta},$$

which puts constraint on  $\beta$  that

$$\beta \leq \frac{2\frac{b}{N}\gamma^2 - H\eta}{1 + \frac{b}{N}\gamma^2}.$$

Specifically,  $\beta < \frac{2\frac{b}{N}\gamma^2}{1 + \frac{b}{N}\gamma^2}$ , which approaches 0 when  $\gamma$  approaches 0, and constrains  $\beta$  in a small range.

## C IMPLICIT BIAS OF ADAM

This section collects the proof of the convergent direction of Adam, i.e., Theorem 4. The methodology of this section bears great similarity with GDM, although the preconditioner of Adam requires specific treatment for analysis. The proof is still divided into two stages: (1). we first prove the sum of squared gradients along the trajectory is finite. Additionally, we prove the convergent rate of loss is  $\mathcal{O}(\frac{1}{t})$ ; (2). we prove  $\mathbf{w}(t) - \log(t)\hat{\mathbf{w}}$  has bounded norm. Before we present these two stages of proof, we will first give the required range of  $\eta$  for which Theorem 3 holds.

### C.1 CHOICE OF LEARNING RATE

Let  $H_{s_0}$  be the smooth parameter over  $[s_0, \infty)$  given by Assumption 3. (D). Let  $\beta_2 = (c\beta_1)^4$  ( $c > 1$ ). The "sufficiently small learning rate" in Theorem 3 means

$$\eta \leq \frac{\sqrt{\varepsilon} \inf_{t \geq 2} \left( \frac{1-\beta_1^t}{1-\beta_1} - \frac{1-\beta_1^{t-1}}{c(1-\beta_1)} \frac{1-(c\beta_1)^t}{1-(c\beta_1)^{t-1}} \right)}{H_{\ell^{-1}((1-c\beta_1)^{-1}N\mathcal{L}(\mathbf{w}(1)))}}.$$

To ensure  $\eta$  is well-defined, we need to prove

$$\inf_{t \geq 2} \left( \frac{1-\beta_1^t}{1-\beta_1} - \frac{1-\beta_1^{t-1}}{c(1-\beta_1)} \frac{1-(c\beta_1)^t}{1-(c\beta_1)^{t-1}} \right) > 0,$$

and we introduce the following technical lemma:

**Lemma 19.** Define  $f_t(x) = \frac{1-x^t}{x(1-x^{t-1})}$ ,  $\forall t \in \mathbb{Z}, t \geq 2$ . We have  $f_t(x)$  is decreasing with respect to  $x$ . Furthermore, for any  $x \in [0, 1)$ , we have

$$f(x) \geq \sqrt[4]{f(x^4)}. \quad (44)$$

*Proof.* First of all, by definition,

$$f(x) = \frac{1-x^t}{x-x^t} = 1 + \frac{1-x}{x-x^t} = 1 + \frac{1-x}{x(1-x^{t-1})} = 1 + \frac{1}{x(1+x+\dots+x^{t-2})}$$

is monotonously decreasing as  $0 \leq x < 1$ . Secondly, Eq. (44) is equivalent to

$$\begin{aligned} \frac{(1-x^t)^4}{\beta_1^4(1-x^{t-1})^4} &\geq \frac{(1-x^{4t})}{x^4(1-x^{4(t-1)})} \\ \iff \frac{(1-x^t)^3}{(1-x^{t-1})^3} &\geq \frac{(1+x^t)(1+x^{2t})}{(1+x^{t-1})(1+x^{2(t-1)})}. \end{aligned}$$

The left side of the above inequality is no smaller than 1, while the right side is no larger than 1, which completes the proof.  $\square$

We are now ready to prove  $\eta$  is well-defined. First of all, for every  $t$ , we have

$$\begin{aligned} &\frac{1-\beta_1^t}{1-\beta_1} - \frac{1-\beta_1^{t-1}}{c(1-\beta_1)} \frac{1-(c\beta_1)^t}{1-(c\beta_1)^{t-1}} \\ &= \frac{\beta_1(1-\beta_1^{t-1})}{1-\beta_1} \left( \frac{1-\beta_1^t}{\beta_1(1-\beta_1^{t-1})} - \frac{1-(c\beta_1)^t}{(c\beta_1)(1-(c\beta_1)^{t-1})} \right) \\ &\stackrel{(*)}{=} \frac{\beta_1(1-\beta_1^{t-1})}{1-\beta_1} (f_t(\beta_1) - f_t(c\beta_1)) > 0, \end{aligned} \quad (45)$$

where Eq. (\*) is by Lemma 19 and  $c\beta_1 = \sqrt[4]{\beta_2} < 1$ .

On the other hand, we have

$$\lim_{t \rightarrow \infty} \left( \frac{1-\beta_1^t}{1-\beta_1} - \frac{1-\beta_1^{t-1}}{c(1-\beta_1)} \frac{1-(c\beta_1)^t}{1-(c\beta_1)^{t-1}} \right) = \left(1 - \frac{1}{c}\right) \frac{1}{1-\beta_1}. \quad (46)$$

By Eq. (45) and Eq. (46), we obtain  $\frac{1-\beta_1^t}{1-\beta_1} - \frac{1-\beta_1^{t-1}}{c(1-\beta_1)} \frac{1-(c\beta_1)^t}{1-(c\beta_1)^{t-1}}$  is lower bounded by some positive constant across  $t$ , and  $\eta$  is well defined.

## C.2 SUM OF GRADIENTS ALONG THE TRAJECTORY IS BOUNDED

We start with the following lemma, which indicates  $\mathcal{L}(\mathbf{w}(t)) + \|\sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t) - \mathbf{w}(t-1))\|^2$  is a proper Lyapunov function for Adam.

*Proof of Lemma 5.* We start with the case  $t = 1$ . To begin with, we have  $\mathcal{L}$  is  $H_{\ell^{-1}(N\mathcal{L}(\mathbf{w}(1)))}$  smooth around  $\mathbf{w}(1)$ . By definition  $H_x$  is non-increasing with respect to  $x$ , and since  $\ell^{-1}$  is also non-increasing, we have

$$H_{\ell^{-1}(N\mathcal{L}(\mathbf{w}(1)))} \leq H_{\ell^{-1}(\frac{1}{1-c\beta_1}N\mathcal{L}(\mathbf{w}(1)))},$$

which further indicates when  $\alpha$  is small enough,

$$\begin{aligned} \mathcal{L}(\mathbf{w}(1+\alpha)) &\stackrel{(*)}{\leq} \mathcal{L}(\mathbf{w}(1)) + \alpha \langle \nabla \mathcal{L}(\mathbf{w}(1)), \mathbf{w}(2) - \mathbf{w}(1) \rangle + \frac{L}{2} \alpha^2 \|\mathbf{w}(2) - \mathbf{w}(1)\|^2 \\ &= \mathcal{L}(\mathbf{w}(1)) - \alpha \left\langle \nabla \mathcal{L}(\mathbf{w}(1)), \eta \frac{1}{\sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(1)}} \odot \nabla \mathcal{L}(\mathbf{w}(1)) \right\rangle + \mathbf{o}(\alpha^2) \\ &\leq \mathcal{L}(\mathbf{w}(1)) - \frac{1}{2\eta} \alpha^2 \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(1)} \odot (\mathbf{w}(2) - \mathbf{w}(1)) \right\|^2, \end{aligned}$$

where in Eq. (★) we denote  $L \triangleq H_{\ell^{-1}(\frac{1}{1-c\beta_1}N\mathcal{L}(\mathbf{w}(1)))}$ , and the last inequality is due to  $\frac{1}{2\eta} \alpha^2 \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(1)} \odot (\mathbf{w}(2) - \mathbf{w}(1)) \right\|^2 = \mathbf{o}(\alpha^2)$ , and  $\left\langle \nabla \mathcal{L}(\mathbf{w}(1)), \eta \frac{1}{\sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(1)}} \odot \nabla \mathcal{L}(\mathbf{w}(1)) \right\rangle$  is positive.

Now if there exists an  $\alpha \in (0, 1)$ , such that Eq. (10) fails, we denote  $\alpha^* = \inf\{\alpha : \text{Eq. (10) fails for } 1 + \alpha\}$ . We have  $\alpha^* > 0$ , and the equality in Eq. (10) holds for  $1 + \alpha^*$ . Therefore, we have for any  $\alpha \in (0, \alpha^*)$ ,

$$\mathcal{L}(\mathbf{w}(1+\alpha)) \leq \mathcal{L}(\mathbf{w}(1+\alpha)) + \frac{1}{2\eta} \alpha^2 \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(1)} \odot (\mathbf{w}(2) - \mathbf{w}(1)) \right\|^2 \leq \mathcal{L}(\mathbf{w}(1)),$$

which by Lemma 10 leads to  $\mathcal{L}$  is  $H_{\ell^{-1}(N\mathcal{L}(\mathbf{w}(1)))}$  smooth (thus  $L$  smooth) over the set  $\{\mathbf{w}(1+\alpha) : \alpha \in [0, \alpha^*]\}$ , and

$$\begin{aligned} &\mathcal{L}(\mathbf{w}(1+\alpha^*)) \\ &\leq \mathcal{L}(\mathbf{w}(1)) + \alpha^* \langle \nabla \mathcal{L}(\mathbf{w}(1)), \mathbf{w}(2) - \mathbf{w}(1) \rangle + \frac{L}{2} (\alpha^*)^2 \|\mathbf{w}(2) - \mathbf{w}(1)\|^2 \\ &= \mathcal{L}(\mathbf{w}(1)) - \alpha^* \left\langle \frac{1}{\eta} \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(1)} \odot (\mathbf{w}(2) - \mathbf{w}(1)), \mathbf{w}(2) - \mathbf{w}(1) \right\rangle + \frac{L}{2} (\alpha^*)^2 \|\mathbf{w}(2) - \mathbf{w}(1)\|^2 \\ &= \mathcal{L}(\mathbf{w}(1)) - \alpha^* \frac{1}{\eta} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(1)} \odot (\mathbf{w}(2) - \mathbf{w}(1)) \right\|^2 + \frac{L}{2} (\alpha^*)^2 \left\| \frac{1}{\sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(1)}} \odot \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(1)} \odot (\mathbf{w}(2) - \mathbf{w}(1)) \right\|^2 \\ &\leq \mathcal{L}(\mathbf{w}(1)) - \alpha^* \frac{1}{\eta} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(1)} \odot (\mathbf{w}(2) - \mathbf{w}(1)) \right\|^2 + \frac{L}{2\sqrt{\varepsilon}} (\alpha^*)^2 \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(1)} \odot (\mathbf{w}(2) - \mathbf{w}(1)) \right\|^2 \\ &< \mathcal{L}(\mathbf{w}(1)) - (\alpha^*)^2 \frac{1}{\eta} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(1)} \odot (\mathbf{w}(2) - \mathbf{w}(1)) \right\|^2 + \frac{L}{2\sqrt{\varepsilon}} (\alpha^*)^2 \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(1)} \odot (\mathbf{w}(2) - \mathbf{w}(1)) \right\|^2 \\ &\leq \mathcal{L}(\mathbf{w}(1)) - (\alpha^*)^2 \frac{1}{2\eta} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(1)} \odot (\mathbf{w}(2) - \mathbf{w}(1)) \right\|^2, \tag{47} \end{aligned}$$

where the second-to-last inequality is due to  $\|\mathbf{w}(2) - \mathbf{w}(1)\| > 0$  (by Lemma 13) and  $\alpha^* > (\alpha^*)^2$ , while the last inequality is due to

$$\begin{aligned} \eta &\leq \frac{\sqrt{\varepsilon} \inf_{t \geq 2} \left( \frac{1-\beta_1^t}{1-\beta_1} - \frac{1-\beta_1^{t-1}}{c(1-\beta_1)} \frac{1-(c\beta_1)^t}{1-(c\beta_1)^{t-1}} \right)}{L} \leq \frac{\sqrt{\varepsilon} \left( \frac{1-\beta_1^2}{1-\beta_1} - \frac{1-\beta_1}{c(1-\beta_1)} \frac{1-(c\beta_1)^2}{1-(c\beta_1)} \right)}{L} \\ &= \frac{\sqrt{\varepsilon} \left( 1 + \beta_1 - \frac{1+(c\beta_1)}{c} \right)}{L} = \frac{\sqrt{\varepsilon} (1 - \frac{1}{c})}{L} < \frac{\sqrt{\varepsilon}}{L}. \end{aligned}$$

Eq. (47) contradicts the fact that the equality in Eq. (10) holds for  $1 + \alpha^*$ , which completes the proof of  $t = 1$ .

If  $t \geq 2$ , following the similar routine as  $t = 1$ , we also prove Eq. (10) by reduction to absurdity. If there exist  $t$  and  $\alpha$  such that Eq. (10) fails. Denote  $t^*$  as the smallest time such that there exists an  $\alpha \in [0, 1)$  such that Eq. (10) fails for  $t^*$  and  $\alpha$ . By Lemma 13,  $\left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^* - 1)} \odot (\mathbf{w}(t^*) - \mathbf{w}(t^* - 1)) \right\|^2$  is positive, and strict inequality in Eq. (10) holds for  $t$  and  $\alpha = 0$ , which by continuity leads to

$$1 > \alpha^* \triangleq \inf\{\alpha \in [0, 1] : \text{Eq. (10) fails for } 1 + \alpha\} > 0.$$

Then, for any  $\alpha \in [0, \alpha^*]$ , we have

$$\begin{aligned} & \mathcal{L}(\mathbf{w}(t^* + \alpha)) \\ & \leq \mathcal{L}(\mathbf{w}(t^* + \alpha)) + \frac{1}{2} \alpha^2 \frac{1 - \beta_1^{t^*}}{\eta(1 - \beta_1)} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^*)} \odot (\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)) \right\|^2 \\ & \leq \mathcal{L}(\mathbf{w}(t^*)) + \frac{1 - \beta_1^{t^* - 1}}{2c\eta(1 - \beta_1)} \frac{1 - (c\beta_1)^{t^*}}{1 - (c\beta_1)^{t^* - 1}} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^* - 1)} \odot (\mathbf{w}(t^*) - \mathbf{w}(t^* - 1)) \right\|^2. \end{aligned}$$

On the other hand, for any time  $2 \leq s \leq t^* - 1$ , we have

$$\begin{aligned} & \mathcal{L}(\mathbf{w}(s + 1)) + \frac{1}{2} \frac{1 - \beta_1^s}{\eta(1 - \beta_1)} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(s)} \odot (\mathbf{w}(s + 1) - \mathbf{w}(s)) \right\|^2 \\ & \leq \mathcal{L}(\mathbf{w}(s)) + \frac{\beta_1(1 - \beta_1^{s-1})}{2c\eta(1 - \beta_1)} \frac{1 - (c\beta_1)^s}{1 - (c\beta_1)^{s-1}} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(s-1)} \odot (\mathbf{w}(s) - \mathbf{w}(s-1)) \right\|^2. \quad (48) \end{aligned}$$

By Eq. (46), we have

$$\frac{1 - \beta_1^s}{\eta(1 - \beta_1)} > \frac{1 - \beta_1^s}{c\eta(1 - \beta_1)} = \frac{1 - \beta_1^s}{c\eta(1 - \beta_1)} \frac{1 - (c\beta_1)^{s+1}}{1 - (c\beta_1)^s} \frac{1 - (c\beta_1)^s}{1 - (c\beta_1)^{s+1}},$$

which by  $\frac{(1 - \beta_1^{s-1})}{(1 - \beta_1^s)}$  further leads to

$$\begin{aligned} & \mathcal{L}(\mathbf{w}(s)) + \frac{1 - \beta_1^{s-1}}{2c\eta(1 - \beta_1)} \frac{1 - (c\beta_1)^s}{1 - (c\beta_1)^{s-1}} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(s-1)} \odot (\mathbf{w}(s) - \mathbf{w}(s-1)) \right\|^2 \\ & \geq \mathcal{L}(\mathbf{w}(s+1)) + \frac{1}{2} \frac{1 - \beta_1^s}{\eta(1 - \beta_1)} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(s)} \odot (\mathbf{w}(s+1) - \mathbf{w}(s)) \right\|^2 \\ & > \mathcal{L}(\mathbf{w}(s+1)) + \frac{1 - \beta_1^s}{2c\eta(1 - \beta_1)} \frac{1 - (c\beta_1)^{s+1}}{1 - (c\beta_1)^s} \frac{1 - (c\beta_1)^s}{1 - (c\beta_1)^{s+1}} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(s)} \odot (\mathbf{w}(s+1) - \mathbf{w}(s)) \right\|^2 \\ & > \frac{1 - (c\beta_1)^s}{1 - (c\beta_1)^{s+1}} \left( \mathcal{L}(\mathbf{w}(s+1)) + \frac{1 - \beta_1^s}{2c\eta(1 - \beta_1)} \frac{1 - (c\beta_1)^{s+1}}{1 - (c\beta_1)^s} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(s)} \odot (\mathbf{w}(s+1) - \mathbf{w}(s)) \right\|^2 \right). \quad (49) \end{aligned}$$

On the other hand, for  $s = 1$ , we have

$$\begin{aligned} \mathcal{L}(\mathbf{w}(1)) & \geq \mathcal{L}(\mathbf{w}(2)) + \frac{1}{2} \frac{1 - \beta_1}{\eta(1 - \beta_1)} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(1)} \odot (\mathbf{w}(2) - \mathbf{w}(1)) \right\|^2 \\ & \geq \frac{1 - (c\beta_1)}{1 - (c\beta_1)^2} \left( \mathcal{L}(\mathbf{w}(2)) + \frac{1 - \beta_1}{2c\eta(1 - \beta_1)} \frac{1 - (c\beta_1)^2}{1 - (c\beta_1)} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(1)} \odot (\mathbf{w}(2) - \mathbf{w}(1)) \right\|^2 \right). \quad (50) \end{aligned}$$

Combining Eqs. (48), (49), and (50), we have

$$\begin{aligned}
& \mathcal{L}(\mathbf{w}(t^* + \alpha)) \\
& \leq \mathcal{L}(\mathbf{w}(t^*)) + \frac{1 - \beta_1^{t^*-1}}{2c\eta(1 - \beta_1)} \frac{1 - (c\beta_1)^{t^*}}{1 - (c\beta_1)^{t^*-1}} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^* - 1)} \odot (\mathbf{w}(t^*) - \mathbf{w}(t^* - 1)) \right\|^2 \\
& < \frac{1 - (c\beta_1)^{t^*}}{1 - (c\beta_2)^{t^*-1}} \left( \mathcal{L}(\mathbf{w}(t^* - 1)) + \frac{1 - \beta_1^{t^*-2}}{2c\eta(1 - \beta_1)} \frac{1 - (c\beta_1)^{t^*-1}}{1 - (c\beta_1)^{t^*-2}} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^* - 2)} \odot (\mathbf{w}(t^* - 1) - \mathbf{w}(t^* - 2)) \right\|^2 \right) \\
& < \dots \\
& < \frac{1 - (c\beta_1)^{t^*}}{1 - (c\beta_1)^2} \left( \mathcal{L}(\mathbf{w}(2)) + \frac{1 - \beta_1}{2c\eta(1 - \beta_1)} \frac{1 - (c\beta_1)^2}{1 - (c\beta_1)} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(1)} \odot (\mathbf{w}(2) - \mathbf{w}(1)) \right\|^2 \right) \\
& \leq \frac{1 - (c\beta_1)^{t^*}}{1 - c\beta_1} \mathcal{L}(\mathbf{w}(1)) < \frac{1}{1 - c\beta_1} \mathcal{L}(\mathbf{w}(1)).
\end{aligned}$$

Therefore, by Lemma 10,  $\mathcal{L}$  is  $H_{\ell^{-1}(\frac{1}{1-c\beta_1}N\mathcal{L}(\mathbf{w}(1)))}$  smooth (thus  $L$  smooth) over the set  $\{\mathbf{w}(t^* + \alpha) : \alpha \in [0, \alpha^*]\}$ , which further leads to

$$\begin{aligned}
& \mathcal{L}(\mathbf{w}(t^* + \alpha^*)) \\
& \leq \mathcal{L}(\mathbf{w}(t^*)) + \alpha^* \langle \nabla \mathcal{L}(\mathbf{w}(t^*)), \mathbf{w}(t^* + 1) - \mathbf{w}(t^*) \rangle + \frac{L}{2} (\alpha^*)^2 \|\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)\|^2 \\
& \stackrel{(\bullet)}{=} - \frac{\alpha^*}{\eta(1 - \beta_1)} \left\langle \mathbf{w}(t^* + 1) - \mathbf{w}(t^*), (1 - \beta_1^{t^*}) \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t^*)} \odot (\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)) \right. \\
& \quad \left. - \beta_1 (1 - \beta_1^{t^*-1}) \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t^*)} \odot (\mathbf{w}(t^*) - \mathbf{w}(t^* - 1)) \right\rangle \\
& + \mathcal{L}(\mathbf{w}(t^*)) + \frac{L}{2} (\alpha^*)^2 \|\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)\|^2 \\
& = \mathcal{L}(\mathbf{w}(t^*)) + \frac{L}{2} (\alpha^*)^2 \|\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)\|^2 - \frac{\alpha^* (1 - \beta_1^{t^*})}{\eta(1 - \beta_1)} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^*)} \odot (\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)) \right\|^2 \\
& + \beta_1 \frac{\alpha^* (1 - \beta_1^{t^*-1})}{\eta(1 - \beta_1)} \left\langle \mathbf{w}(t^* + 1) - \mathbf{w}(t^*), \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t^*)} \odot (\mathbf{w}(t^*) - \mathbf{w}(t^* - 1)) \right\rangle \\
& = \mathcal{L}(\mathbf{w}(t^*)) + \frac{L}{2} (\alpha^*)^2 \|\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)\|^2 - \frac{\alpha^* (1 - \beta_1^{t^*})}{\eta(1 - \beta_1)} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^*)} \odot (\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)) \right\|^2 \\
& + \beta_1 \frac{\alpha^* (1 - \beta_1^{t^*-1})}{\eta(1 - \beta_1)} \left\langle \frac{\sqrt[8]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^* - 1)}}{\sqrt[8]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^*)}} \odot \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^*)} \odot (\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)), \right. \\
& \quad \left. \frac{\sqrt[8]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^* - 1)}}{\sqrt[8]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^*)}} \odot \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^* - 1)} \odot (\mathbf{w}(t^*) - \mathbf{w}(t^* - 1)) \right\rangle \\
& \leq \mathcal{L}(\mathbf{w}(t^*)) + \frac{L}{2} (\alpha^*)^2 \|\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)\|^2 - \frac{\alpha^* (1 - \beta_1^{t^*})}{\eta(1 - \beta_1)} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^*)} \odot (\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)) \right\|^2 \\
& + \beta_1 \frac{(\alpha^*)^2 (1 - \beta_1^{t^*-1})}{2\eta(1 - \beta_1)} \left\| \frac{\sqrt[8]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^* - 1)}}{\sqrt[8]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^*)}} \odot \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^*)} \odot (\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)) \right\|^2 \\
& + \beta_1 \frac{(1 - \beta_1^{t^*-1})}{2\eta(1 - \beta_1)} \left\| \frac{\sqrt[8]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^* - 1)}}{\sqrt[8]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^*)}} \odot \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^* - 1)} \odot (\mathbf{w}(t^*) - \mathbf{w}(t^* - 1)) \right\|^2
\end{aligned}$$



$$\begin{aligned}
&\stackrel{(\diamond)}{\leq} \mathcal{L}(\mathbf{w}(t^*)) + \frac{L}{2}(\alpha^*)^2 \|\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)\|^2 - \frac{\alpha^*(1 - \beta_1^{t^*})}{\eta(1 - \beta_1)} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^*)} \odot (\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)) \right\|^2 \\
&\quad + \beta_1 \frac{(\alpha^*)^2(1 - \beta_1^{t^*-1})}{2\eta(1 - \beta_1)} \frac{1 - (c\beta_1)^{t^*}}{c\beta_1(1 - (c\beta_1)^{t^*-1})} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^*)} \odot (\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)) \right\|^2 \\
&\quad + \beta_1 \frac{(1 - \beta_1^{t^*-1})}{2\eta(1 - \beta_1)} \frac{1 - (c\beta_1)^{t^*}}{c\beta_1(1 - (c\beta_1)^{t^*-1})} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^* - 1)} \odot (\mathbf{w}(t^*) - \mathbf{w}(t^* - 1)) \right\|^2 \\
&\leq \mathcal{L}(\mathbf{w}(t^*)) + \frac{L}{2\sqrt{\varepsilon}}(\alpha^*)^2 \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^*)} \odot (\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)) \right\|^2 \\
&\quad - \frac{\alpha^*(1 - \beta_1^{t^*})}{\eta(1 - \beta_1)} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^*)} \odot (\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)) \right\|^2 \\
&\quad + \beta_1 \frac{(\alpha^*)^2(1 - \beta_1^{t^*})}{2\eta(1 - \beta_1)} \frac{1 - (c\beta_1)^{t^*}}{c\beta_1(1 - (c\beta_1)^{t^*-1})} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^*)} \odot (\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)) \right\|^2 \\
&\quad + \beta_1 \frac{(1 - \beta_1^{t^*})}{2\eta(1 - \beta_1)} \frac{1 - (c\beta_1)^{t^*}}{c\beta_1(1 - (c\beta_1)^{t^*-1})} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^* - 1)} \odot (\mathbf{w}(t^*) - \mathbf{w}(t^* - 1)) \right\|^2 \\
&\stackrel{(\square)}{<} \mathcal{L}(\mathbf{w}(t^*)) - \frac{(\alpha^*)^2(1 - \beta_1^{t^*})}{2\eta(1 - \beta_1)} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^*)} \odot (\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)) \right\|^2 \\
&\quad + \frac{(1 - \beta_1^{t^*-1})}{2\eta(1 - \beta_1)} \frac{1 - (c\beta_1)^{t^*}}{c(1 - (c\beta_1)^{t^*-1})} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^* - 1)} \odot (\mathbf{w}(t^*) - \mathbf{w}(t^* - 1)) \right\|^2,
\end{aligned}$$

where Eq. (•) is due to an alternative form of the Adam's update rule:

$$\begin{aligned}
&(1 - \beta_1^{t^*})\sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t^*)} \odot (\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)) - \beta_1(1 - \beta_1^{t^*-1})\sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t^* - 1)} \odot (\mathbf{w}(t^*) - \mathbf{w}(t^* - 1)) \\
&= -\eta(1 - \beta_1)\nabla \mathcal{L}(\mathbf{w}(t^*)), \tag{51}
\end{aligned}$$

Inequality (◇) is due to

$$\begin{aligned}
&\frac{\sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^* - 1)}}{\sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^*)}} = \sqrt[4]{\frac{\varepsilon \mathbf{1}_d + \hat{\nu}(t^* - 1)}{\varepsilon \mathbf{1}_d + \hat{\nu}(t^*)}} = \sqrt[4]{\frac{\varepsilon \mathbf{1}_d + \hat{\nu}(t^* - 1)}{\varepsilon \mathbf{1}_d + \frac{\beta_2 \nu(t^* - 1) + (1 - \beta_2)\nabla \mathcal{L}(\mathbf{w}(t^*))^2}{1 - \beta_2^{t^*}}}} \\
&\leq \sqrt[4]{\frac{\varepsilon \mathbf{1}_d + \hat{\nu}(t^* - 1)}{\varepsilon \mathbf{1}_d + \frac{\beta_2 \nu(t^* - 1)}{1 - \beta_2^{t^*}}}} = \sqrt[4]{\frac{\varepsilon \mathbf{1}_d + \hat{\nu}(t^* - 1)}{\varepsilon \mathbf{1}_d + \frac{\beta_2(1 - \beta_2^{t^*-1})\hat{\nu}(t^* - 1)}{1 - \beta_2^{t^*}}} \leq \sqrt[4]{\frac{\varepsilon \mathbf{1}_d + \hat{\nu}(t^* - 1)}{\frac{\beta_2(1 - \beta_2^{t^*-1})\hat{\nu}(t^* - 1)}{1 - \beta_2^{t^*}} \varepsilon \mathbf{1}_d + \frac{\beta_2(1 - \beta_2^{t^*-1})\hat{\nu}(t^* - 1)}{1 - \beta_2^{t^*}}} \\
&= \sqrt[4]{\frac{1 - \beta_2^{t^*}}{\beta_2(1 - \beta_2^{t^*-1})}} \mathbf{1}_d \text{ (all the computations are component-wisely)},
\end{aligned}$$

and  $f(c\beta_1) \geq \sqrt[4]{f((c\beta_1)^4)}$ , and Inequality (□) is due to

$$\frac{L}{2\sqrt{\varepsilon}} \leq \frac{\inf_{t \geq 2} \left( \frac{1 - \beta_1^t}{1 - \beta_1} - \frac{1 - \beta_1^{t-1}}{c(1 - \beta_1)} \frac{1 - (c\beta_1)^t}{1 - (c\beta_1)^{t-1}} \right)}{2\eta} \leq \frac{\left( \frac{1 - \beta_1^{t^*}}{1 - \beta_1} - \frac{1 - \beta_1^{t^*-1}}{c(1 - \beta_1)} \frac{1 - (c\beta_1)^{t^*}}{1 - (c\beta_1)^{t^*-1}} \right)}{2\eta},$$

$\alpha^* > (\alpha^*)^2$ , and  $\left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t^*)} \odot (\mathbf{w}(t^* + 1) - \mathbf{w}(t^*)) \right\|^2 > 0$ .

This contradicts to that the equality in Eq. (10) holds for  $t^* + \alpha^*$ .

The proof is completed.  $\square$

As  $\lim_{t \rightarrow \infty} \beta_1^t = 0$  and  $\lim_{t \rightarrow \infty} (c\beta_1)^t = 0$ , we have the following corollary based on Lemma 1.

**Corollary 6.** *Let all assumptions in Theorem 4 hold. Then, for large enough  $t$ , we have*

$$\begin{aligned}
&\mathcal{L}(\mathbf{w}(t + 1)) + \frac{1}{2\sqrt[4]{c\eta}(1 - \beta_1)} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t + 1) - \mathbf{w}(t)) \right\|^2 \\
&\leq \mathcal{L}(\mathbf{w}(t)) + \frac{1}{2\sqrt[4]{c\eta}(1 - \beta_1)} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t - 1)} \odot (\mathbf{w}(t) - \mathbf{w}(t - 1)) \right\|^2. \tag{52}
\end{aligned}$$

Consequently, we have

$$\sum_{t=1}^{\infty} \|\nabla \mathcal{L}(\mathbf{w}(t))\|^2 < \infty. \quad (53)$$

The proof of Corollary 6 requires the following classical lemma on the equivalence between the convergence of two non-negative sequence. The proof is omitted here and can be found in Wang et al. (2021).

**Lemma 20** (c.f. Lemma 27, Wang et al. (2021)). *Let  $\{a_i\}_{i=1}^{\infty}$  be a series of non-negative reals, and  $\varepsilon$  be a positive real. Then,  $\sum_{i=1}^{\infty} a_i < \infty$  is equivalent to  $\sum_{i=1}^{\infty} \frac{a_i}{\sqrt{\varepsilon + \sum_{s=1}^i a_s}} < \infty$ .*

*Proof of Corollary 6.* We have

$$\begin{aligned} \lim_{t \rightarrow \infty} \frac{1 - \beta_1^{t-1}}{2c\eta(1 - \beta_1)} \frac{1 - (c\beta_1)^t}{1 - (c\beta_1)^{t-1}} &= \frac{1}{2c\eta(1 - \beta_1)} < \frac{1}{2\sqrt[2]{c\eta(1 - \beta_1)}}, \\ \lim_{t \rightarrow \infty} \frac{1 - \beta_1^t}{2\eta(1 - \beta_1)} &= \frac{1}{2\eta(1 - \beta_1)} > \frac{1}{2\sqrt[4]{c\eta(1 - \beta_1)}}, \end{aligned}$$

which completes the proof of Eq. (52). Rearranging Eq. (52) leads to

$$\begin{aligned} \frac{\sqrt[4]{c} - 1}{2\sqrt[2]{c\eta(1 - \beta_1)}} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) \right\|^2 &\leq \frac{1}{2\sqrt[4]{c\eta(1 - \beta_1)}} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t-1)} \odot (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\|^2 \\ + \mathcal{L}(\mathbf{w}(t)) - \left( \mathcal{L}(\mathbf{w}(t+1)) + \frac{1}{2\sqrt[4]{c\eta(1 - \beta_1)}} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) \right\|^2 \right), \end{aligned}$$

which by iteration further leads to that for a large enough time  $T_1$

$$\begin{aligned} &\sum_{t=T_1}^{T_2} \frac{\sqrt[4]{c} - 1}{2\sqrt[2]{c\eta(1 - \beta_1)}} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) \right\|^2 \\ &\leq \mathcal{L}(\mathbf{w}(T_1)) + \frac{1}{2\sqrt[4]{c\eta(1 - \beta_1)}} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(T_1 - 1)} \odot (\mathbf{w}(T_1) - \mathbf{w}(T_1 - 1)) \right\|^2 \\ &\quad - \mathcal{L}(\mathbf{w}(T_2 + 1)) + \frac{1}{2\sqrt[4]{c\eta(1 - \beta_1)}} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(T_2)} \odot (\mathbf{w}(T_2 + 1) - \mathbf{w}(T_2 + 1)) \right\|^2 \\ &< \mathcal{L}(\mathbf{w}(T_1)) + \frac{1}{2\sqrt[4]{c\eta(1 - \beta_1)}} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(T_1 - 1)} \odot (\mathbf{w}(T_1) - \mathbf{w}(T_1 - 1)) \right\|^2. \end{aligned}$$

Consequently, we obtain

$$\sum_{t=1}^{\infty} \frac{\sqrt[4]{c} - 1}{2\sqrt[2]{c\eta(1 - \beta_1)}} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) \right\|^2 < \infty. \quad (54)$$

On the other hand, for any  $t$ , we have

$$\begin{aligned} &\left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) \right\| \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot \hat{\mathbf{w}} \right\| \\ &\geq \left\langle \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)), \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot \hat{\mathbf{w}} \right\rangle = \left\langle \sqrt[2]{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)), \hat{\mathbf{w}} \right\rangle \\ &= \langle -\eta \hat{\mathbf{m}}(t), \hat{\mathbf{w}} \rangle = -\frac{\eta(1 - \beta_1)}{1 - \beta_1^t} \left\langle \sum_{s=1}^t \beta_1^{t-s} \nabla \mathcal{L}(\mathbf{w}(s)), \hat{\mathbf{w}} \right\rangle \\ &= -\frac{\eta(1 - \beta_1)}{1 - \beta_1^t} \frac{1}{N} \left\langle \sum_{s=1}^t \beta_1^{t-s} \sum_{\tilde{\mathbf{x}}_i \in \mathcal{T}(\mathbf{S})} \ell'(\langle \tilde{\mathbf{x}}_i, \mathbf{w}(s) \rangle) \tilde{\mathbf{x}}_i, \hat{\mathbf{w}} \right\rangle \geq -\frac{\eta(1 - \beta_1)}{1 - \beta_1^t} \frac{1}{N} \sum_{s=1}^t \beta_1^{t-s} \sum_{\tilde{\mathbf{x}}_i \in \mathcal{T}(\mathbf{S})} \ell'(\langle \tilde{\mathbf{x}}_i, \mathbf{w}(s) \rangle) \\ &\geq -\frac{\eta(1 - \beta_1)}{1 - \beta_1^t} \frac{1}{N} \sum_{\tilde{\mathbf{x}}_i \in \mathcal{T}(\mathbf{S})} \ell'(\langle \tilde{\mathbf{x}}_i, \mathbf{w}(t) \rangle) \geq \frac{\eta(1 - \beta_1)}{1 - \beta_1^t} \|\nabla \mathcal{L}(\mathbf{w}(t))\|, \end{aligned}$$

which by Eq. (54) indicates

$$\sum_{t=1}^{\infty} \left( \frac{\eta(1-\beta_1)}{1-\beta_1^t} \right)^2 \frac{\|\nabla \mathcal{L}(\mathbf{w}(t))\|^2}{\left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot \hat{\mathbf{w}} \right\|^2} < \infty.$$

As  $\lim_{t \rightarrow \infty} \left( \frac{\eta(1-\beta_1)}{1-\beta_1^t} \right)^2 = \eta^2(1-\beta_1)^2$ , we then obtain

$$\begin{aligned} & \sum_{t=1}^{\infty} \frac{\|\nabla \mathcal{L}(\mathbf{w}(t))\|^2}{\sqrt[2]{\varepsilon + \sum_{s=1}^t \|\nabla \mathcal{L}(\mathbf{w}(s))\|^2}} \leq \sum_{t=1}^{\infty} \frac{\|\nabla \mathcal{L}(\mathbf{w}(t))\|^2}{\sqrt[2]{\varepsilon + \sum_{s=1}^t (1-\beta)\beta^{t-s} \|\nabla \mathcal{L}(\mathbf{w}(s))\|^2}} \\ & \leq \sqrt{\frac{1}{1-\beta}} \sum_{t=1}^{\infty} \frac{\|\nabla \mathcal{L}(\mathbf{w}(t))\|^2}{\sqrt[2]{\varepsilon + \frac{\sum_{s=1}^t (1-\beta)\beta^{t-s} \|\nabla \mathcal{L}(\mathbf{w}(s))\|^2}{1-\beta^t}}} \leq d \sqrt{\frac{1}{1-\beta}} \sum_{t=1}^{\infty} \frac{\|\nabla \mathcal{L}(\mathbf{w}(t))\|^2}{\left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \frac{\sum_{s=1}^t (1-\beta)\beta^{t-s} \|\nabla \mathcal{L}(\mathbf{w}(s))\|^2}{1-\beta^t}} \right\|^2} \\ & = d \sqrt{\frac{1}{1-\beta}} \sum_{t=1}^{\infty} \frac{\|\nabla \mathcal{L}(\mathbf{w}(t))\|^2}{\left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \right\|^2} \leq d \|\hat{\mathbf{w}}\|_{\infty}^2 \sqrt{\frac{1}{1-\beta}} \sum_{t=1}^{\infty} \frac{\|\nabla \mathcal{L}(\mathbf{w}(t))\|^2}{\left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot \hat{\mathbf{w}} \right\|^2} < \infty, \end{aligned}$$

which by Lemma 20 completes the proof.  $\square$

Based on Corollary 6, we can further prove Lemma 6, characterizing the convergent rate of loss  $\mathcal{L}$  directly.

*Proof of Lemma 6.* To begin with, Eq. (51) indicates

$$\begin{aligned} & \|\eta(1-\beta_1)\nabla \mathcal{L}(\mathbf{w}(t))\|^2 \\ & = \left\| (1-\beta_1^t)\sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) - \beta_1(1-\beta_1^{t-1})\sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t-1)} \odot (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\|^2 \\ & \leq \left\| (1-\beta_1^t)\sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) \right\| + \left\| \beta_1(1-\beta_1^{t-1})\sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t-1)} \odot (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\|^2 \\ & \leq \left( \left\| \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) \right\| + \left\| \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t-1)} \odot (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\| \right)^2 \\ & \leq 2 \left( \left\| \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) \right\|^2 + \left\| \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t-1)} \odot (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\|^2 \right) \end{aligned} \quad (55)$$

On the other hand, by Corollary 6,

$$\sum_{s=1}^{\infty} \|\nabla \mathcal{L}(\mathbf{w}(s))\|^2 < \infty,$$

which following the same routine as Corollary 5 leads to

$$\langle \mathbf{w}(t), \tilde{\mathbf{x}} \rangle \rightarrow \infty, \forall \tilde{\mathbf{x}} \in \tilde{\mathcal{S}}.$$

Therefore, by Lemma 11, there exists a large enough time  $T_1$ , such that  $\forall t \geq T_1$ ,

$$\frac{1}{K} \ell(\langle \mathbf{w}(t), \tilde{\mathbf{x}} \rangle) \leq -\ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}} \rangle) \leq K \ell(\langle \mathbf{w}(t), \tilde{\mathbf{x}} \rangle), \forall \tilde{\mathbf{x}} \in \mathcal{T}(\mathcal{S}),$$

which by the separable assumption further leads to

$$\begin{aligned} \frac{\gamma}{K} \mathcal{L}(\mathbf{w}(t)) & \leq -\frac{\gamma}{N} \sum_{\tilde{\mathbf{x}} \in \mathcal{T}(\mathcal{S})} \ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}} \rangle) \leq \frac{1}{N} \left\langle - \sum_{\tilde{\mathbf{x}} \in \mathcal{T}(\mathcal{S})} \ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}} \rangle) \tilde{\mathbf{x}}, \gamma \hat{\mathbf{w}} \right\rangle \\ & \leq \frac{1}{N} \left\| \sum_{\tilde{\mathbf{x}} \in \mathcal{T}(\mathcal{S})} \ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}} \rangle) \tilde{\mathbf{x}} \right\| \|\gamma \hat{\mathbf{w}}\| = \|\nabla \mathcal{L}(\mathbf{w}(t))\| \\ & \leq -\frac{1}{N} \sum_{\tilde{\mathbf{x}} \in \mathcal{T}(\mathcal{S})} \ell'(\langle \mathbf{w}(t), \tilde{\mathbf{x}} \rangle) \leq K \mathcal{L}(\mathbf{w}(t)). \end{aligned} \quad (56)$$

Combining Eq. (26) and the above inequality, we have

$$\begin{aligned} \left(\frac{\eta(1-\beta_1)\gamma}{K}\right)^2 \mathcal{L}(\mathbf{w}(t))^2 &\leq 2 \left( \left\| \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) \right\|^2 \right. \\ &\quad \left. + \left\| \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t-1)} \odot (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\|^2 \right). \end{aligned} \quad (57)$$

On the other hand, by Eq. (54), we have

$$\sum_{t=1}^{\infty} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) \right\|^2 < \infty.$$

Therefore, there exists large enough time  $T_2$ , such that  $\forall t > T_2$ ,

$$\left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) \right\|^2 < 1,$$

and thus,

$$\left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) \right\|^4 < \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) \right\|^2. \quad (58)$$

Combining Eq. (57) and Eq. (58), there exists a positive real constant  $C$ , such that

$$\begin{aligned} \mathcal{L}(\mathbf{w}(t))^2 &\leq C \left( \left\| \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) \right\|^2 \right. \\ &\quad \left. + \left\| \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t-1)} \odot (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\|^2 \right), \\ \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) \right\|^4 &\leq C \left( \left\| \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) \right\|^2 \right. \\ &\quad \left. + \left\| \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t-1)} \odot (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\|^2 \right). \end{aligned}$$

Rearranging Eq. (52) leads to

$$\begin{aligned} &\frac{\sqrt[4]{c} - 1}{4\sqrt[4]{c}\eta(1-\beta_1)} \left( \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t-1)} \odot (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\|^2 \right. \\ &\quad \left. + \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t-1)} \odot (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\|^2 \right) \\ &\leq \mathcal{L}(\mathbf{w}(t)) + \frac{\sqrt[4]{c} + 1}{4\sqrt[4]{c}\eta(1-\beta_1)} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t-1)} \odot (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\|^2 \\ &\quad - \left( \mathcal{L}(\mathbf{w}(t+1)) + \frac{\sqrt[4]{c} + 1}{4\sqrt[4]{c}\eta(1-\beta_1)} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) \right\|^2 \right), \end{aligned}$$

which further indicates

$$\begin{aligned}
& \left( \mathcal{L}(\mathbf{w}(t)) + \frac{\sqrt[4]{c} + 1}{4\sqrt[2]{c}\eta(1 - \beta_1)} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t-1)} \odot (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\|^2 \right)^2 \\
& \leq 2 \left( \mathcal{L}(\mathbf{w}(t))^2 + \frac{\sqrt[4]{c} + 1}{4\sqrt[2]{c}\eta(1 - \beta_1)} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t-1)} \odot (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\|^4 \right) \\
& \leq 2C \left( 1 + \frac{\sqrt[4]{c} + 1}{4\sqrt[2]{c}\eta(1 - \beta_1)} \right) \left( \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) \right\|^2 \right. \\
& \quad \left. + \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t-1)} \odot (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\|^2 \right) \\
& \leq 2C \left( 1 + \frac{\sqrt[4]{c} + 1}{4\sqrt[2]{c}\eta(1 - \beta_1)} \right) \frac{4\sqrt[2]{c}\eta(1 - \beta_1)}{\sqrt[4]{c} - 1} \left( \mathcal{L}(\mathbf{w}(t)) + \frac{\sqrt[4]{c} + 1}{4\sqrt[2]{c}\eta(1 - \beta_1)} \right. \\
& \quad \cdot \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t-1)} \odot (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\|^2 - (\mathcal{L}(\mathbf{w}(t+1))) \\
& \quad \left. + \frac{\sqrt[4]{c} + 1}{4\sqrt[2]{c}\eta(1 - \beta_1)} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) \right\|^2 \right).
\end{aligned}$$

Denote  $\xi(t)$  as

$$\xi(t) \triangleq \mathcal{L}(\mathbf{w}(t)) + \frac{\sqrt[4]{c} + 1}{4\sqrt[2]{c}\eta(1 - \beta_1)} \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t-1)} \odot (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\|^2.$$

We then have

$$\xi(t)^2 \leq 2C \left( 1 + \frac{\sqrt[4]{c} + 1}{4\sqrt[2]{c}\eta(1 - \beta_1)} \right) \frac{4\sqrt[2]{c}\eta(1 - \beta_1)}{\sqrt[4]{c} - 1} (\xi(t) - \xi(t+1)),$$

which leads to

$$\begin{aligned}
& \xi(t) = \mathcal{O}\left(\frac{1}{t}\right), \text{ i.e., } \mathcal{L}(\mathbf{w}(t)) = \mathcal{O}\left(\frac{1}{t}\right), \\
& \text{and } \left\| \sqrt[4]{\varepsilon \mathbf{1}_d + \hat{\nu}(t-1)} \odot (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\|^2 = \mathcal{O}\left(\frac{1}{t}\right).
\end{aligned}$$

Due to Eq. (56), we further have  $\|\nabla \mathcal{L}(\mathbf{w}(t))\| = \mathcal{O}(t^{-1})$ , which indicates

$$\begin{aligned}
\|\mathbf{w}(t)\| & \leq \|\mathbf{w}(1)\| + \sum_{s=1}^t \|\mathbf{w}(s+1) - \mathbf{w}(s)\| = \|\mathbf{w}(1)\| + \eta \sum_{s=1}^t \left\| \frac{\hat{\mathbf{m}}(s)}{\sqrt{\hat{\nu}(s) + \varepsilon \mathbf{1}_d}} \right\| \\
& \leq \|\mathbf{w}(1)\| + \frac{\eta}{\sqrt{\varepsilon}} \sum_{s=1}^t \|\hat{\mathbf{m}}(s)\| = \|\mathbf{w}(1)\| + \frac{\eta}{\sqrt{\varepsilon}} \sum_{s=1}^t \frac{1}{(1 - \beta^s)} \left\| \sum_{i=1}^s \beta^{s-i} \nabla \mathcal{L}(\mathbf{w}(i)) \right\| \\
& \leq \|\mathbf{w}(1)\| + \frac{\eta}{\sqrt{\varepsilon}(1 - \beta)} \sum_{s=1}^t \sum_{i=1}^s \beta^{s-i} \|\nabla \mathcal{L}(\mathbf{w}(i))\| \\
& \leq \|\mathbf{w}(1)\| + \frac{\eta}{\sqrt{\varepsilon}(1 - \beta)^2} \sum_{s=1}^t \|\nabla \mathcal{L}(\mathbf{w}(s))\| = \mathcal{O}(\log(t)).
\end{aligned}$$

Therefore, for any  $\tilde{\mathbf{x}} \in \mathcal{T}(\mathbf{S})$ , we have  $\langle \mathbf{w}(t), \tilde{\mathbf{x}} \rangle = \mathcal{O}(\log(t))$ , which by  $\ell$  is exponential-tailed leads to  $\ell(\langle \mathbf{w}(t), \tilde{\mathbf{x}} \rangle) = \Omega(t^{-1})$ , and thus  $\mathcal{L}(\mathbf{w}(t)) = \Theta(t^{-1})$ . Also, since  $\mathcal{L}(\mathbf{w}(t)) = \mathcal{O}(t^{-1})$ , we have  $\langle \mathbf{w}(t), \tilde{\mathbf{x}} \rangle = \Omega(\log(t))$ , which further leads to  $\|\mathbf{w}(t)\| = \Omega(\log(t))$ , and thus  $\|\mathbf{w}(t)\| = \Theta(\log(t))$ .

Finally, we have

$$\begin{aligned}
\gamma \sum_{s=1}^t \beta^{t-s} \|\nabla \mathcal{L}(\mathbf{w}(s))\| &= \frac{\gamma}{N} \sum_{s=1}^t \beta^{t-s} \left\| \sum_{\tilde{\mathbf{x}} \in \mathcal{T}(s)} \ell'(\langle \mathbf{w}(s), \tilde{\mathbf{x}} \rangle) \tilde{\mathbf{x}} \right\| \\
&\leq -\frac{\gamma}{N} \sum_{s=1}^t \beta^{t-s} \sum_{\tilde{\mathbf{x}} \in \mathcal{T}(s)} \ell'(\langle \mathbf{w}(s), \tilde{\mathbf{x}} \rangle) \leq -\frac{1}{N} \sum_{s=1}^t \beta^{t-s} \left\langle \sum_{\tilde{\mathbf{x}} \in \mathcal{T}(s)} \ell'(\langle \mathbf{w}(s), \tilde{\mathbf{x}} \rangle) \tilde{\mathbf{x}}, \gamma \hat{\mathbf{w}} \right\rangle \\
&\leq \left\| \sum_{s=1}^t \beta^{t-s} \nabla \mathcal{L}(\mathbf{w}(s)) \right\| \|\gamma \hat{\mathbf{w}}\| = \|\mathbf{m}(t)\| \leq \sum_{s=1}^t \beta^{t-s} \|\nabla \mathcal{L}(\mathbf{w}(s))\|,
\end{aligned}$$

which leads to  $\|\mathbf{m}(t)\| = \Theta(t^{-1})$ . Similarly, we have  $\nu(t) = \mathcal{O}(t^{-2})$ , component-wisely. As  $\lim_{t \rightarrow \infty} \beta_1^t = 0$  and  $\lim_{t \rightarrow \infty} \beta_2^t = 0$ , we have

$$\|\mathbf{w}(t) - \mathbf{w}(t-1)\| = \left\| \frac{\hat{\mathbf{m}}(t)}{\sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t)}} \right\| = \Theta(t^{-1}).$$

The proof is completed.  $\square$

### C.3 BOUNDING THE ORTHOGONAL PART

By Lemma 8, there exists a solution  $\tilde{\mathbf{w}}$  as the solution of Eq. (12) with  $\mathbf{C}_2 = \frac{\eta}{(1-\beta)\sqrt{\varepsilon}}$ . Define  $\mathbf{r}(t)$  as

$$\mathbf{r}(t) \triangleq \mathbf{w}(t) - \log(t) \hat{\mathbf{w}} - \tilde{\mathbf{w}}, \quad (59)$$

and we only need to prove  $\|\mathbf{r}(t)\|$  is bounded over time. We then prove Lemma 7, providing an equivalent condition of  $\|\mathbf{r}(t)\|$  being bounded. As the GDM and SGDM case, we separate the proof into two sub-lemmas.

**Lemma 21.** *Let all conditions in Theorem 4 hold. Then,  $\|\mathbf{r}(t)\|$  is bounded if and only if  $g(t)$  is upper bounded.*

*Proof.* Following the same routine as Lemma 2 and Lemma 17, we only need to prove

$$\lim_{t \rightarrow \infty} \left\| (1 - \beta_1^{t-1}) \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t-1)} \odot (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\| = 0, \quad (60)$$

and

$$\sum_{\tau=2}^{\infty} \left| \langle \mathbf{r}(\tau) - \mathbf{r}(\tau-1), (1 - \beta_1^{\tau-1}) \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(\tau-1)} \odot (\mathbf{w}(\tau) - \mathbf{w}(\tau-1)) \rangle \right| < \infty. \quad (61)$$

As for Eq. (60), by Lemma 6, we have

$$\begin{aligned}
&\left\| (1 - \beta_1^{t-1}) \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t-1)} \odot (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\| \\
&= \mathcal{O}(t^{-1}) = \mathbf{o}(1).
\end{aligned}$$

As for Eq. (60), we have

$$\begin{aligned}
&\left| \langle \mathbf{r}(\tau) - \mathbf{r}(\tau-1), (1 - \beta_1^{\tau-1}) \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(\tau-1)} \odot (\mathbf{w}(\tau) - \mathbf{w}(\tau-1)) \rangle \right| \\
&= \left| \left\langle \mathbf{w}(\tau) - \mathbf{w}(\tau-1) - \log \frac{\tau}{\tau-1} \hat{\mathbf{w}}, (1 - \beta_1^{\tau-1}) \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(\tau-1)} \odot (\mathbf{w}(\tau) - \mathbf{w}(\tau-1)) \right\rangle \right| \\
&\leq \left| \left\langle \log \frac{\tau}{\tau-1} \hat{\mathbf{w}}, (1 - \beta_1^{\tau-1}) \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(\tau-1)} \odot (\mathbf{w}(\tau) - \mathbf{w}(\tau-1)) \right\rangle \right| \\
&\quad + (1 - \beta_1^{\tau-1}) \left\| \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(\tau-1)} \odot (\mathbf{w}(\tau) - \mathbf{w}(\tau-1)) \right\|^2 \\
&\stackrel{(\star)}{=} \mathcal{O}(\tau^{-2}),
\end{aligned}$$

where Eq.  $(\star)$  is due to Lemma 6 and  $\log(\frac{\tau}{\tau-1}) = \Theta(\tau^{-1})$ .

The proof is completed.  $\square$

We conclude the proof of Theorem 4 by showing  $g(t)$  is upper bounded.

**Lemma 22.** *Let all conditions in Theorem 4 hold. Then,  $g(t)$  is upper bounded.*

*Proof.*  $g(t)$  is upper bounded is equivalent to  $\sum_{t=1}^{\infty} g(t+1) - g(t) < \infty$ . We then prove this lemma by calculating  $g(t+1) - g(t)$  directly.

$$\begin{aligned}
& g(t+1) - g(t) \\
&= \frac{\sqrt{\varepsilon}}{2} \|\mathbf{r}(t+1)\|^2 + \frac{\beta_1}{1-\beta_1} \left\langle \mathbf{r}(t+1), (1-\beta_1^t) \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) \right\rangle \\
&\quad - \left( \frac{\sqrt{\varepsilon}}{2} \|\mathbf{r}(t)\|^2 + \frac{\beta_1}{1-\beta_1} \left\langle \mathbf{r}(t), (1-\beta_1^{t-1}) \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t-1)} \odot (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\rangle \right) \\
&\quad - \frac{\beta_1}{1-\beta_1} \langle \mathbf{r}(t+1) - \mathbf{r}(t), (1-\beta_1^t) \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) \rangle \\
&= \frac{\sqrt{\varepsilon}}{2} \|\mathbf{r}(t+1) - \mathbf{r}(t)\|^2 + \frac{\beta_1}{1-\beta_1} \left\langle \mathbf{r}(t), (1-\beta_1^t) \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) \right. \\
&\quad \left. - (1-\beta_1^{t-1}) \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t-1)} \odot (\mathbf{w}(t) - \mathbf{w}(t-1)) \right\rangle + \sqrt{\varepsilon} \langle \mathbf{r}(t+1) - \mathbf{r}(t), \mathbf{r}(t) \rangle \\
&\stackrel{(*)}{=} \left\langle \mathbf{r}(t), -(1-\beta_1^t) \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) - \frac{\eta}{1-\beta} \nabla \mathcal{L}(\mathbf{w}(t)) \right\rangle \\
&\quad + \frac{\sqrt{\varepsilon}}{2} \|\mathbf{r}(t+1) - \mathbf{r}(t)\|^2 + \sqrt{\varepsilon} \langle \mathbf{r}(t+1) - \mathbf{r}(t), \mathbf{r}(t) \rangle,
\end{aligned}$$

where Eq. (\*) is due to a simple rearranging of the update rule of Adam, i.e.,

$$\begin{aligned}
& \frac{\beta_1}{1-\beta_1} \left( (1-\beta_1^t) \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) - (1-\beta_1^{t-1}) \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t-1)} \odot (\mathbf{w}(t) - \mathbf{w}(t-1)) \right) \\
&= -\frac{\eta}{1-\beta_1} \nabla \mathcal{L}(\mathbf{w}(t)) - (1-\beta_1^t) \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)).
\end{aligned}$$

On the one hand, as  $\|\mathbf{r}(t+1) - \mathbf{r}(t)\| = \|\mathbf{w}(t+1) - \mathbf{w}(t) - \log \frac{t+1}{t} \hat{\mathbf{w}}\| = \mathcal{O}(t^{-1})$ ,

$$\sum_{t=1}^{\infty} \frac{\sqrt{\varepsilon}}{2} \|\mathbf{r}(t+1) - \mathbf{r}(t)\|^2 < \infty.$$

On the other hand,

$$\begin{aligned}
& \left\langle \mathbf{r}(t), -(1-\beta_1^t) \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) - \frac{\eta}{1-\beta} \nabla \mathcal{L}(\mathbf{w}(t)) \right\rangle \\
&\quad + \sqrt{\varepsilon} \langle \mathbf{r}(t+1) - \mathbf{r}(t), \mathbf{r}(t) \rangle \\
&= \left\langle \mathbf{r}(t), -(1-\beta_1^t) \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) - \frac{\eta}{1-\beta} \nabla \mathcal{L}(\mathbf{w}(t)) \right\rangle \\
&\quad + \sqrt{\varepsilon} \left\langle \mathbf{w}(t+1) - \mathbf{w}(t) - \log \left( \frac{t+1}{t} \right) \hat{\mathbf{w}}, \mathbf{r}(t) \right\rangle \\
&= \left\langle \mathbf{r}(t), -(1-\beta_1^t) \sqrt{\varepsilon \mathbf{1}_d + \hat{\nu}(t)} \odot (\mathbf{w}(t+1) - \mathbf{w}(t)) + \sqrt{\varepsilon} (\mathbf{w}(t+1) - \mathbf{w}(t)) \right\rangle \\
&\quad + \left\langle \mathbf{r}(t), -\sqrt{\varepsilon} \log \left( \frac{t+1}{t} \right) \hat{\mathbf{w}} - \frac{\eta}{1-\beta} \nabla \mathcal{L}(\mathbf{w}(t)) \right\rangle \\
&\stackrel{(\bullet)}{=} \mathcal{O}(\beta_1^t + t^{-2}) + \left\langle \mathbf{r}(t), -\sqrt{\varepsilon} \log \left( \frac{t+1}{t} \right) \hat{\mathbf{w}} - \frac{\eta}{1-\beta} \nabla \mathcal{L}(\mathbf{w}(t)) \right\rangle,
\end{aligned}$$

where Eq. (•) is due to  $\hat{\nu}(t) = \mathcal{O}(t^{-2})$ .

Furthermore, following exactly the same routine as Lemma 16, we have

$$\sum_{t=1}^{\infty} \left\langle \mathbf{r}(t), -\sqrt{\varepsilon} \log \left( \frac{t+1}{t} \right) \hat{\mathbf{w}} - \frac{\eta}{1-\beta} \nabla \mathcal{L}(\mathbf{w}(t)) \right\rangle < \infty.$$

The proof is completed.  $\square$

## D APPLICATIONS&EXTENSIONS

### D.1 APPLICATION TO THE MINI-BATCH SGDM

This section provides formal description of the implicit bias of mini-batch SGDM and its corresponding proof. To begin with, we would like to provide a formal definition of mini-batch SGDM. Mini-batch SGDM differs from SGDM by applying sampling without replacement to obtain  $\mathbf{B}(t)$  in Eq. (1). Specifically, let  $K = \frac{N}{b}$ . For any  $T \geq 0$ , we call time series  $\{KT + 1, \dots, KT + K\}$  the  $(T + 1)$ -th epoch, and during the  $T + 1$ -th epoch, the dataset  $\mathcal{S}$  is randomly uniformly divided into  $K$  parts  $\{\mathbf{B}(KT + 1), \dots, \mathbf{B}(KT + K)\}$ , with  $\bigcup_{t=KT+1}^{KT+K} \mathbf{B}(t) = \mathcal{S}$ . The implicit bias of mini-batch SGDM is then stated as the following theorem:

**Theorem 6.** *Let Assumptions 1, 2, and 3. (S) hold. Let learning rate  $\eta$  be small enough, and  $\beta \in [0, 1)$ . Then, for almost every dataset  $\mathcal{S}$ , mini-batch SGDM satisfies  $\mathbf{w}(t) - \log(t)\hat{\mathbf{w}}$  is bounded as  $t \rightarrow \infty$ , and  $\lim_{t \rightarrow \infty} \frac{\mathbf{w}(t)}{\|\mathbf{w}(t)\|} = \frac{\hat{\mathbf{w}}}{\|\hat{\mathbf{w}}\|}$ .*

The without-replacement sampling method leads to the direction of every trajectory of mini-SGDM converge to the max-margin solution, compared to the same conclusion holds for SGDM a.s.. We prove the theorem following the same framework of GDM, by proceeding with two stages.

**Stage I.** The following lemma proves  $\mathcal{L}(\mathbf{u}(t))$  is an Lyapunov function for mini-batch SGDM and without the a.s. condition.

**Lemma 23.** *Let all conditions in Theorem 6 hold. Then, we have*

$$\mathcal{L}(\mathbf{u}(t + 1)) \leq \mathcal{L}(\mathbf{u}(1)) - \Omega(\eta) \sum_{s=1}^t \|\nabla \mathcal{L}(\mathbf{w}(s))\|^2.$$

*Proof.* By the Taylor Expansion of  $\mathcal{L}(\mathbf{u}(t + 1))$  at  $\mathbf{u}(t)$ , we have

$$\begin{aligned} & \mathcal{L}(\mathbf{u}(KT + T + 1)) \\ & \leq \mathcal{L}(\mathbf{u}(KT + 1)) - \tilde{\eta} \left\langle \nabla \mathcal{L}(\mathbf{u}(KT + 1)), \sum_{t=1}^K \nabla \mathcal{L}_{\mathbf{B}(t+KT)}(\mathbf{w}(t + KT)) \right\rangle \\ & \quad + \frac{H\tilde{\eta}^2}{2} \left\| \sum_{t=1}^K \nabla \mathcal{L}_{\mathbf{B}(t+KT)}(\mathbf{w}(t + KT)) \right\|^2. \end{aligned} \tag{62}$$

On the other hand, for any  $t \in \{2, \dots, K\}$ , we have

$$\begin{aligned} \mathbf{w}(KT + t) - \mathbf{w}(KT + 1) &= \eta \sum_{s=1}^t \left( \sum_{\ell=1}^{KT+s} \beta^{KT+s-\ell} \nabla \mathcal{L}_{\mathbf{B}(\ell)}(\mathbf{w}(\ell)) \right) \\ &= \eta \sum_{s=1}^t \left( \sum_{\ell=KT+1}^{KT+s} \beta^{KT+s-\ell} \nabla \mathcal{L}_{\mathbf{B}(\ell)}(\mathbf{w}(\ell)) \right) + \eta \sum_{s=1}^t \left( \sum_{\ell=1}^{KT} \beta^{KT+s-\ell} \nabla \mathcal{L}_{\mathbf{B}(\ell)}(\mathbf{w}(\ell)) \right) \\ &= \eta \sum_{\ell=1}^t \frac{1 - \beta^{t-\ell+1}}{1 - \beta} \nabla \mathcal{L}_{\mathbf{B}(KT+\ell)}(\mathbf{w}(KT + \ell)) + \eta \frac{\beta(1 - \beta^t)}{1 - \beta} \sum_{\ell=1}^{KT} \beta^{KT-\ell} \nabla \mathcal{L}_{\mathbf{B}(\ell)}(\mathbf{w}(\ell)) \\ &= \eta \sum_{\ell=1}^t \frac{1 - \beta^{t-\ell+1}}{1 - \beta} \nabla \mathcal{L}_{\mathbf{B}(KT+\ell)}(\mathbf{w}(KT + \ell)) - \eta \sum_{\ell=1}^t \frac{1 - \beta^{t-\ell+1}}{1 - \beta} \nabla \mathcal{L}_{\mathbf{B}(KT+\ell)}(\mathbf{w}(KT + 1)) \\ & \quad + \eta \frac{\beta(1 - \beta^t)}{1 - \beta} \sum_{\ell=1}^{KT} \beta^{KT-\ell} \nabla \mathcal{L}_{\mathbf{B}(\ell)}(\mathbf{w}(\ell)) + \eta \sum_{\ell=1}^t \frac{1 - \beta^{t-\ell+1}}{1 - \beta} \nabla \mathcal{L}_{\mathbf{B}(KT+\ell)}(\mathbf{w}(KT + 1)), \end{aligned}$$



which by  $\eta$  is small enough further indicates

$$\begin{aligned}
& \|\mathbf{w}(KT+t) - \mathbf{w}(KT+1)\| \\
& \leq \eta \left\| \sum_{\ell=1}^t \frac{1-\beta^{t-\ell+1}}{1-\beta} \nabla \mathcal{L}_{\mathbf{B}(KT+\ell)}(\mathbf{w}(KT+\ell)) - \sum_{\ell=1}^t \frac{1-\beta^{t-\ell+1}}{1-\beta} \nabla \mathcal{L}_{\mathbf{B}(KT+\ell)}(\mathbf{w}(KT+1)) \right\| \\
& + \eta \left\| \frac{\beta(1-\beta^t)}{1-\beta} \sum_{\ell=1}^{KT} \beta^{KT-\ell} \nabla \mathcal{L}_{\mathbf{B}(\ell)}(\mathbf{w}(\ell)) \right\| + \eta \left\| \sum_{\ell=1}^t \frac{1-\beta^{t-\ell+1}}{1-\beta} \nabla \mathcal{L}_{\mathbf{B}(KT+\ell)}(\mathbf{w}(KT+1)) \right\| \\
& = \mathcal{O}(\eta) \sum_{\ell=2}^t \|\mathbf{w}(KT+\ell) - \mathbf{w}(KT+1)\| + \mathcal{O}(\eta) \left( \sum_{\ell=1}^{KT} \beta^{KT-\ell} \|\nabla \mathcal{L}_{\mathbf{B}(\ell)}(\mathbf{w}(\ell))\| \right) \\
& + \mathcal{O}(\eta) \|\nabla \mathcal{L}(\mathbf{w}(KT+1))\|.
\end{aligned}$$

Applying the same analysis to  $\|\mathbf{w}(KT+t-1) - \mathbf{w}(KT+1)\|$  recursively, we finally obtain

$$\begin{aligned}
& \|\mathbf{w}(KT+t) - \mathbf{w}(KT+1)\| \\
& \leq \mathcal{O}(\eta) \left( \sum_{\ell=1}^{KT} \beta^{KT-\ell} \|\nabla \mathcal{L}_{\mathbf{B}(\ell)}(\mathbf{w}(\ell))\| \right) + \mathcal{O}(\eta) \|\nabla \mathcal{L}(\mathbf{w}(KT+1))\|. \quad (63)
\end{aligned}$$

Applying Eq. (63) to the  $\|\nabla \mathcal{L}_{\mathbf{B}(\ell)}(\mathbf{w}(\ell))\|$  in Eq. (63) ( $\forall \ell \in [1, KT]$ ) iterative and choosing  $\eta$  to be small enough, we further have

$$\begin{aligned}
& \|\mathbf{w}(KT+t) - \mathbf{w}(KT+1)\| \\
& \leq \mathcal{O}(\eta) \left( \sum_{\ell=0}^{T-1} \sqrt{\beta^{K(T-\ell)}} \|\nabla \mathcal{L}_{\mathbf{B}(K\ell+1)}(\mathbf{w}(K\ell+1))\| \right) + \mathcal{O}(\eta) \|\nabla \mathcal{L}(\mathbf{w}(KT+1))\| \\
& = \mathcal{O}(\eta) \left( \sum_{\ell=0}^T \sqrt{\beta^{K(T-\ell)}} \|\nabla \mathcal{L}_{\mathbf{B}(K\ell+1)}(\mathbf{w}(K\ell+1))\| \right).
\end{aligned}$$

Therefore,

$$\begin{aligned}
& \sum_{t=1}^K \nabla \mathcal{L}_{\mathbf{B}(t+KT)}(\mathbf{w}(t+KT)) \\
& = \sum_{t=1}^K \nabla \mathcal{L}_{\mathbf{B}(t+KT)}(\mathbf{w}(t)) + \mathcal{O} \left( \eta \left( \sum_{\ell=0}^T \sqrt{\beta^{K(T-\ell)}} \|\nabla \mathcal{L}_{\mathbf{B}(K\ell+1)}(\mathbf{w}(K\ell+1))\| \right) \right) \\
& = K \nabla \mathcal{L}(\mathbf{w}(t)) + \mathcal{O} \left( \eta \left( \sum_{\ell=0}^T \sqrt{\beta^{K(T-\ell)}} \|\nabla \mathcal{L}_{\mathbf{B}(K\ell+1)}(\mathbf{w}(K\ell+1))\| \right) \right). \quad (64)
\end{aligned}$$

Similarly, one can obtain

$$\begin{aligned}
& \nabla \mathcal{L}(\mathbf{u}(KT+1)) \\
& = \nabla \mathcal{L}(\mathbf{w}(KT+1)) + \mathcal{O}(\|\mathbf{w}(KT+1) - \mathbf{w}(KT)\|) \\
& = \nabla \mathcal{L}(\mathbf{w}(KT+1)) + \mathcal{O} \left( \eta \left( \sum_{\ell=0}^T \sqrt{\beta^{K(T-\ell)}} \|\nabla \mathcal{L}_{\mathbf{B}(K\ell+1)}(\mathbf{w}(K\ell+1))\| \right) \right). \quad (65)
\end{aligned}$$

Applying Eq. (64) and Eq. (65) back to the Taylor Expansion (Eq. (62)), we have

$$\begin{aligned} & \mathcal{L}(\mathbf{u}(KT + T + 1)) \\ & \leq \mathcal{L}(\mathbf{u}(KT + 1)) - \Omega(\eta) \langle \nabla \mathcal{L}(\mathbf{w}(KT + 1)), \nabla \mathcal{L}(\mathbf{w}(KT + 1)) \rangle \\ & \quad + \mathcal{O} \left( \eta^2 \left( \sum_{\ell=0}^T \sqrt{\beta^{K(T-\ell)}} \|\nabla \mathcal{L}_{\mathbf{B}(K\ell+1)}(\mathbf{w}(K\ell + 1))\| \right)^2 \right) \\ & \leq \mathcal{L}(\mathbf{u}(KT + 1)) - \Omega(\eta) \langle \nabla \mathcal{L}(\mathbf{w}(KT + 1)), \nabla \mathcal{L}(\mathbf{w}(KT + 1)) \rangle \\ & \quad + \mathcal{O} \left( \eta^2 \left( \sum_{\ell=0}^T \sqrt{\beta^{K(T-\ell)}} \|\nabla \mathcal{L}_{\mathbf{B}(K\ell+1)}(\mathbf{w}(K\ell + 1))\|^2 \right) \right). \end{aligned}$$

Summing the above inequality over  $T$  and setting  $\eta$  small enough leads to the conclusion.

The proof is completed.  $\square$

## D.2 EXTENSION TO THE MULTI-CLASS CLASSIFICATION PROBLEM

Here we use several notations and lemmas from (Soudry et al., 2018). We define  $\mathbf{w} = \text{vec}(\mathbf{W})$ ,  $\hat{\mathbf{w}} = \text{vec}(\hat{\mathbf{W}})$ ,  $\mathbf{e}_i \in \mathbb{R}^C$  ( $i \in \{1, \dots, C\}$ ) satisfying  $(\mathbf{e}_i)_j = \delta_{ij}$ , and  $\mathbf{A}_i = \mathbf{e}_i \otimes \mathbb{I}_{d_X}$ , where  $\mathbb{I}_{d_X}$  is the identity matrix with dimension  $d_X$ . We still consider the normalized data, i.e.,  $\|\mathbf{x}\| \leq 1$ ,  $\forall (\mathbf{x}, \mathbf{y}) \in \mathcal{S}$ . Then, the individual loss of sample  $(\mathbf{x}, \mathbf{y})$  can be then represented as

$$\ell(\mathbf{y}, \mathbf{W}\mathbf{x}) = \log \frac{e^{\langle \mathbf{w}, \mathbf{A}_y \mathbf{x} \rangle}}{\sum_{i=1}^C e^{\langle \mathbf{w}, \mathbf{A}_i \mathbf{x} \rangle}}.$$

Furthermore, the gradient of training error at  $\mathbf{W}$  has the form

$$\nabla \mathcal{L}(\mathbf{w}) = \frac{1}{N} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{S}} \sum_{i=1}^C \frac{1}{\sum_{j=1}^C e^{\langle \mathbf{w}, (\mathbf{A}_j - \mathbf{A}_i) \mathbf{x} \rangle}} (\mathbf{A}_i - \mathbf{A}_y) \mathbf{x}.$$

and the Hessian matrix of  $\mathcal{L}$  can be represented as

$$\mathcal{H}\mathcal{L}(\mathbf{w}) = \frac{1}{N} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{S}} \sum_{i=1}^C \frac{\sum_{j=1}^C e^{\langle \mathbf{w}, (\mathbf{A}_j - \mathbf{A}_i) \mathbf{x} \rangle}}{\left( \sum_{j=1}^C e^{\langle \mathbf{w}, (\mathbf{A}_j - \mathbf{A}_i) \mathbf{x} \rangle} \right)^2} (\mathbf{A}_i - \mathbf{A}_y) \mathbf{x} ((\mathbf{A}_j - \mathbf{A}_i) \mathbf{x})^\top,$$

one can then easily verify all absolute value of the eigenvalues of  $\mathcal{H}\mathcal{L}(\mathbf{w})$  is no larger than 2, which indicates  $\mathcal{L}$  is 2-globally smooth.

On the other hand, the separable assumption leads to  $\langle \hat{\mathbf{w}}, (\mathbf{A}_y - \mathbf{A}_i) \mathbf{x} \rangle > 0$ ,  $\forall \mathbf{y} \neq i$ , which further indicates

$$\langle \nabla \mathcal{L}(\mathbf{w}), \hat{\mathbf{w}} \rangle > 0.$$

Let  $\gamma = \frac{1}{\|\hat{\mathbf{w}}\|}$ , following the similar routine as the binary case, we have for a random subset of  $\mathcal{S}$  sampled uniformly without replacement with size  $b$ , we have

$$\|\nabla \mathcal{L}(\mathbf{w})\|^2 \leq \mathbb{E}_{\mathbf{B}(t)} \|\nabla \mathcal{L}_{\mathbf{B}(t)}(\mathbf{w})\|^2 \leq \frac{2N}{\gamma b^2} \|\nabla \mathcal{L}(\mathbf{w})\|^2. \quad (66)$$

Similarly, we have for any positive real series  $\{a_t\}_{t=t_1}^{t_2}$ ,

$$\gamma \sum_{t=t_1}^{t_2} a(t) \|\nabla \mathcal{L}(\mathbf{w}(t))\| \leq \left\| \sum_{t=t_1}^{t_2} a(t) \nabla \mathcal{L}(\mathbf{w}(t)) \right\| \leq \sum_{t=t_1}^{t_2} a(t) \|\nabla \mathcal{L}(\mathbf{w}(t))\|. \quad (67)$$

The proofs of Stage I can then be obtained with Lyapunov functions unchanged and by replacing the corresponding lemmas using Eq. (66) and Eq. (67).

As for the proofs of Stage II, the Lyapunov functions are still the same, while we only need to prove the sum of  $\langle \mathbf{r}(t), -\frac{\eta}{1-\beta} \nabla \mathcal{L}(\mathbf{w}(t)) - \log \frac{t+1}{t} \hat{\mathbf{w}} \rangle$  (for GDM,  $\langle \mathbf{r}(t), -\frac{\eta}{1-\beta} \nabla \mathcal{L}_{\mathbf{B}(t)}(\mathbf{w}(t)) -$

$\frac{N}{bt} \sum_{i: \tilde{\mathbf{x}}_i \in \mathcal{T}(\mathcal{B}(t) \cap \mathcal{S}_s)} \langle \mathbf{v}_i, \tilde{\mathbf{x}}_i \rangle$  for SGDM,  $\langle \mathbf{r}(t), -\sqrt{\varepsilon} \log\left(\frac{t+1}{t}\right) \hat{\mathbf{w}} - \frac{\eta}{1-\beta} \nabla \mathcal{L}(\mathbf{w}(t)) \rangle$  for Adam). For the multi-class case using GDM, We present the following lemma from (Soudry et al., 2018), while the other two cases can be proved similarly:

**Lemma 24** (Part of the proof of Lemma 20, (Soudry et al., 2018)). *If  $\langle \mathbf{w}(t), (\mathbf{A}_y - \mathbf{A}_i)\mathbf{x} \rangle \rightarrow \infty$  as  $t \rightarrow \infty$ ,  $\forall (\mathbf{x}, \mathbf{y}) \in \mathcal{S}$  and  $\forall i \neq y$ , we have the sum of  $\langle \mathbf{r}(t), -\frac{\eta}{1-\beta} \nabla \mathcal{L}(\mathbf{w}(t)) - \log\left(\frac{t+1}{t}\right) \hat{\mathbf{w}} \rangle$  is upper bounded.*

The proof of Theorem 5 is then completed.

## E EXPERIMENTS

This section collects several experiments supporting our theoretical results.

### E.1 EXPERIMENTS OF LINEAR MODEL

#### E.1.1 COMPARING THE TRAINING BEHAVIOR OF GD, GDM, AND ADAM (W/S)

The experiments in this section is designed to verify Theorem 2, i.e., with proper learning rates, gradient descent with momentum converges to the max margin solution, which is the same as gradient descent. We use the synthetic dataset as (Figure 1, (Soudry et al., 2018)) and run GD, GDM and Adam (w/s) over it with different learning rates  $\eta = 0.1, 0.01$  and different random seeds (for random initialization and random samples despite the support sets  $\{((1.5, 0.5), 1), ((0.5, 1.5), 1), ((-1.5, -0.5), -1), ((-1.5, -0.5), -1)\}$ ). Both the angle between the output parameter and max margin solution and the training accuracy are plotted in Figure 1. Specifically, we plot the results with learning rate (1).  $\eta_{GD} = \eta_{GDM} = \eta_{Adam} = 0.001$ , (2).  $\eta_{GD} = \eta_{GDM} = \eta_{Adam} = 0.1$ . We plot the training accuracy and the angle between the output parameter and the max margin solution for each setting respectively. The results are shown in Figure 1. The observations can be summarized as follows:

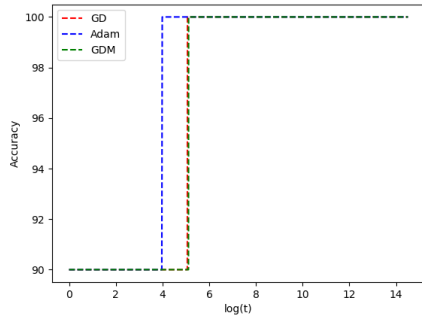
- When learning rate is small enough ( $\eta = 0.1, 0.001$ ), both GD, GDM, and Adam converge to the max margin solution, which supports our theoretical results;
- (Similarity between GD and GDM) The training behaviors of GD and GDM are highly similar.
- (The acceleration effect of Adam) Adam achieves smaller angle with the max margin solution under the same number of iterations.

#### E.1.2 COMPARING THE TRAINING BEHAVIOR OF SGD AND SGDM

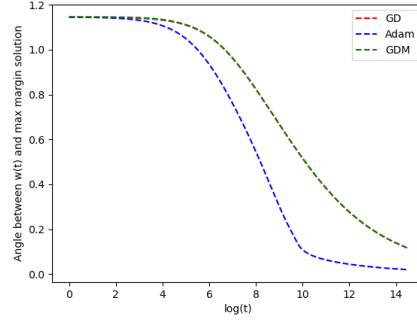
We also run the stochastic optimizers (SGD and SGDM) on the same synthetic dataset as Figure 1. The learning rate is same as Figure 1, namely 0.001 and 0.1. The results of the experiment are plotted in Figure 2. It can be observed that when learning rate is small enough, SGD and SGDM have the similar training behavior, both converging to the max margin solution.

#### E.1.3 ADAM ON ILL-POSED DATASET IN (SOUDRY ET AL., 2018)

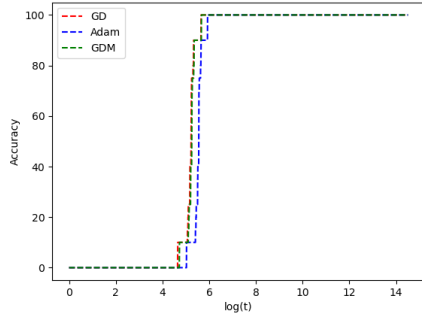
In Figure 3 of (Soudry et al., 2018), an ill-posed synthetic dataset is proposed to support the argument "Adam does not converge to max margin solution", which contradicts to the theoretical results of this paper. We re-conduct the experiment of Figure 3 in (Soudry et al., 2018) with the same ill-posed synthetic dataset with different learning rates and different random seeds as Figure 3. Figure 3. (f) is similar to Figure 3 in (Soudry et al., 2018), where with learning rate  $\eta = 0.1$  and random seed 1, the angle of GD to the max margin solution is smaller than Adam all the time. However, it can be observed from the amplified figure that the angle of GD keeps still and above 0 all the time, meaning that GD doesn't converge to the max margin solution under this setting. However, the angle of Adam to the max margin solution still keeps decreasing and it's unreasonable to claim "Adam doesn't converge to the max margin solution" in this case (the same issue exists in Figure 3 in (Soudry et al., 2018)). Also, as we mentioned at the beginning of this section, this dataset is ill-posed, which is due to the imbalance between the two components of the data (for all data  $((x_1, x_2), y)$  in the dataset,



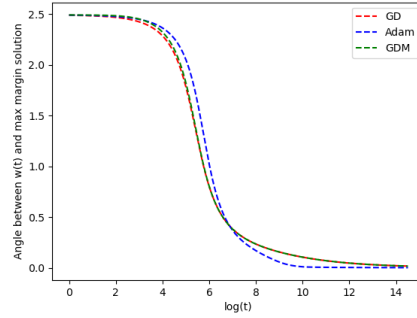
(a) Comparison of Accuracy:  $\eta = 0.001$ , random seed = 1



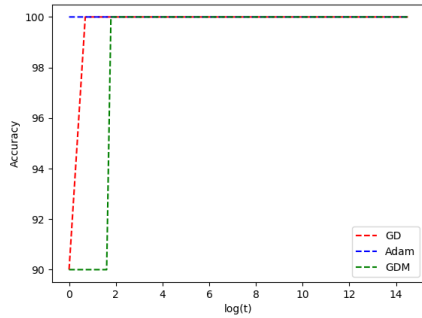
(b) Comparison of Angle:  $\eta = 0.001$ , random seed = 1



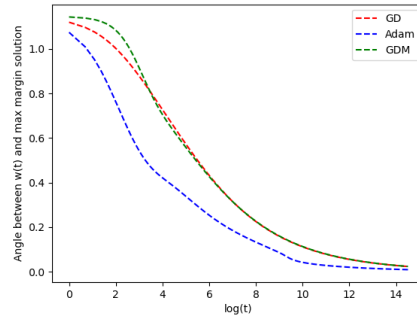
(c) Comparison of Accuracy:  $\eta = 0.001$ , random seed = 2



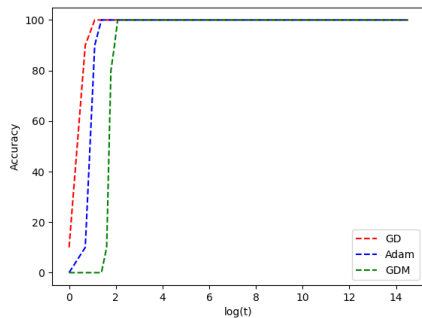
(d) Comparison of Angle:  $\eta = 0.001$ , random seed = 2



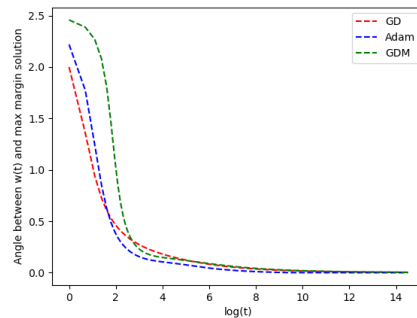
(e) Comparison of Accuracy:  $\eta = 0.1$ , random seed = 1



(f) Comparison of Angle:  $\eta = 0.1$ , random seed = 1

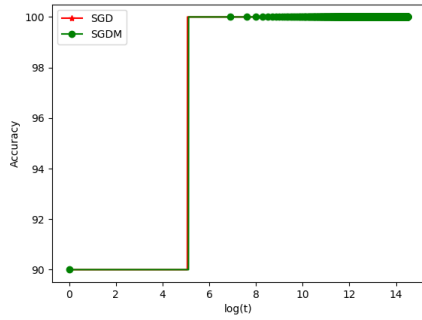


(g) Comparison of Accuracy:  $\eta = 0.1$ , random seed = 2

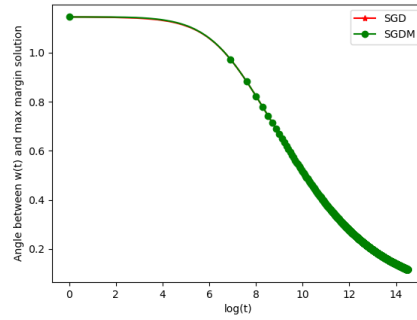


(h) Comparison of Angle:  $\eta = 0.1$ , random seed = 2

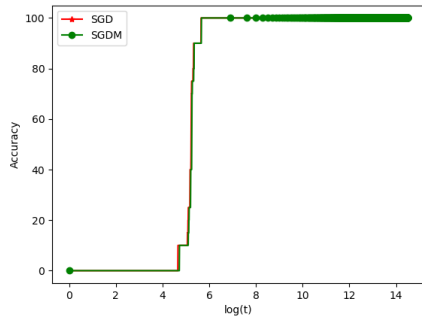
Figure 1: Comparison of GD, GDM, and Adam on the synthetic dataset in (Soudry et al., 2018). In (a-g), the GD curve coincides with the GDM curve.



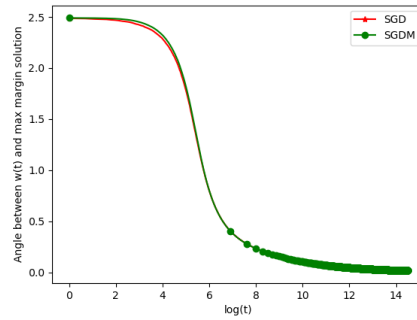
(a) Comparison of Accuracy:  $\eta = 0.001$ , random seed = 1



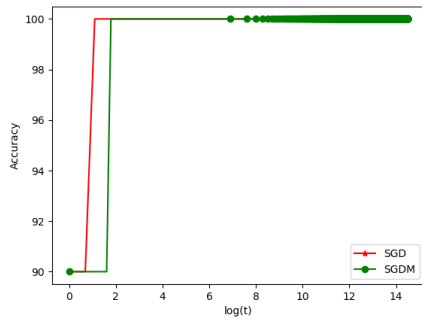
(b) Comparison of Angle:  $\eta = 0.001$ , random seed = 1



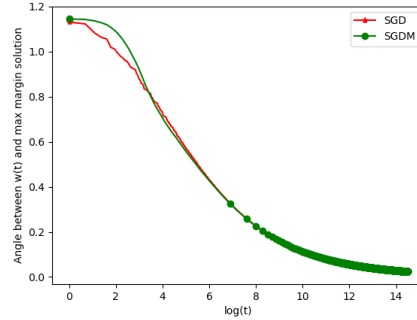
(c) Comparison of Accuracy:  $\eta = 0.001$ , random seed = 2



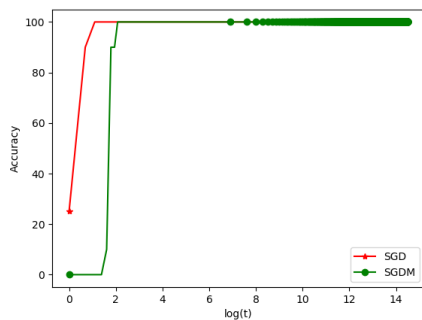
(d) Comparison of Angle:  $\eta = 0.001$ , random seed = 2



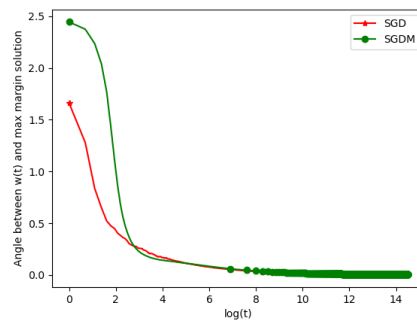
(e) Comparison of Accuracy:  $\eta = 0.1$ , random seed = 1



(f) Comparison of Angle:  $\eta = 0.1$ , random seed = 1



(g) Comparison of Accuracy:  $\eta = 0.1$ , random seed = 2



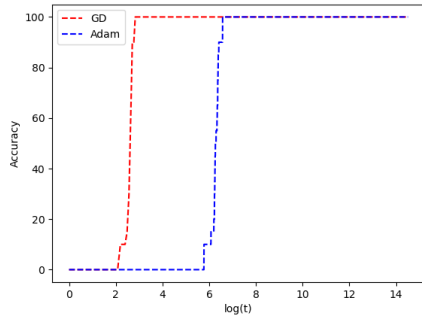
(h) Comparison of Angle:  $\eta = 0.1$ , random seed = 2

Figure 2: Comparison of SGD and SGDM on the synthetic dataset in (Soudry et al., 2018).

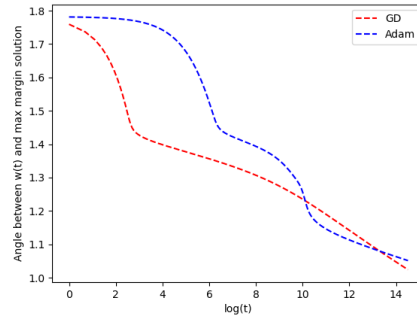
$|x_1|$  is always smaller than 2, while  $|x_2|$  is larger than 10 (and even larger than 30 despite two data in the dataset)), which requires smaller learning rate. To tackle this problem, we need to tune down the learning rate. By Figure 3. (b),(d) and Figure 4. (b), after scaling down the learning rate, both GD's angle and Adam's angle keep decreasing.

## E.2 EVIDENCE IN THE DEEP NEURAL NETWORKS

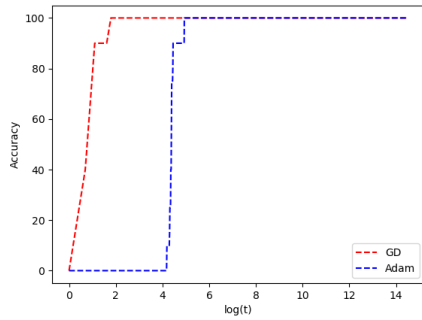
We conduct an experiment on the MNIST dataset using the four layer convolutional networks used in (Lyu & Li, 2019; Wang et al., 2021) (first proposed by (Madry et al., 2018)) to verify whether SGD and SGDM still behave similarly in (homogeneous) deep neural networks. The learning rates of the optimizers are all set to be the default in Pytorch. The results can be seen in Figure 5. It can be observed that (1). SGDM achieves similar test accuracy compared to SGD while (2). SGDM converges faster than SGD.



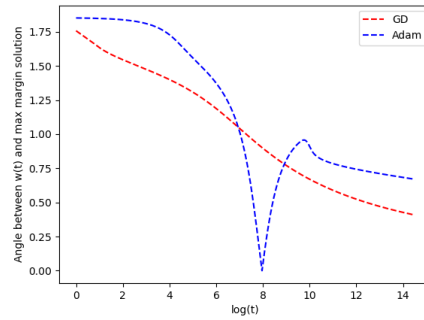
(a) Comparison of Accuracy:  $\eta = 0.1$ , random seed = 1



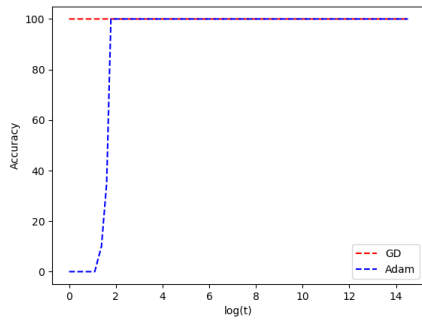
(b) Comparison of Angle:  $\eta = 0.001$ , random seed = 1



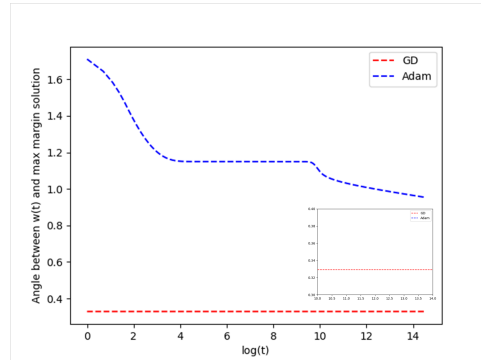
(c) Comparison of Accuracy:  $\eta = 0.001$ , random seed = 2



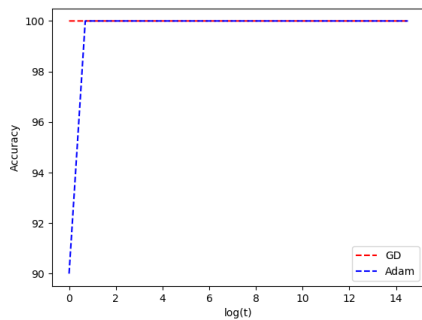
(d) Comparison of Angle:  $\eta = 0.001$ , random seed = 2



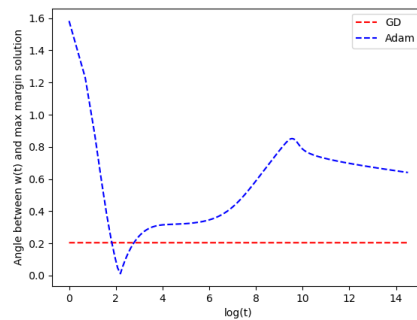
(e) Comparison of Accuracy:  $\eta = 0.1$ , random seed = 1



(f) Comparison of Angle:  $\eta = 0.1$ , random seed = 1

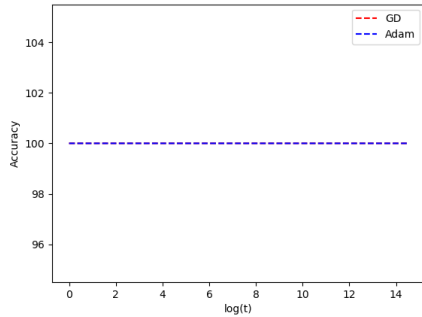


(g) Comparison of Accuracy:  $\eta = 0.1$ , random seed = 2

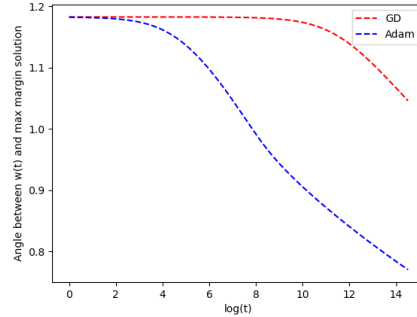


(h) Comparison of Angle:  $\eta = 0.1$ , random seed = 2

Figure 3: Comparison of GD and Adam on the ill-posed synthetic dataset in (Soudry et al., 2018).

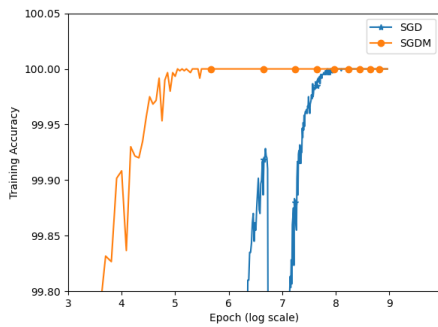


(a) Comparison of Accuracy:  $\eta = 0.1$ , random seed = 3

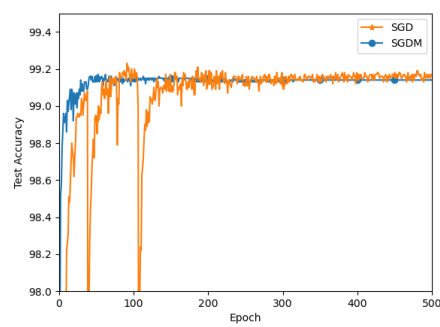


(b) Comparison of Angle:  $\eta = 0.001$ , random seed = 3

Figure 4: Comparison of GD and Adam on the ill-posed synthetic dataset in (Soudry et al., 2018) (continue).



(a) Comparison of Training Accuracy



(b) Comparison of Test Accuracy

Figure 5: Comparison of SGD and SGDM on MNIST with four layer CNN.