

---

# A Modular Analysis of Provable Acceleration via Polyak’s Momentum: Training a Wide ReLU Network and a Deep Linear Network

---

Jun-Kun Wang<sup>1</sup> Chi-Heng Lin<sup>2</sup> Jacob Abernethy<sup>1</sup>

## Abstract

Incorporating a so-called “momentum” dynamic in gradient descent methods is widely used in neural net training as it has been broadly observed that, at least empirically, it often leads to significantly faster convergence. At the same time, there are very few theoretical guarantees in the literature to explain this apparent acceleration effect. Even for the classical strongly convex quadratic problems, several existing results only show Polyak’s momentum has an accelerated linear rate asymptotically. In this paper, we first revisit the quadratic problems and show a non-asymptotic accelerated linear rate of Polyak’s momentum. Then, we provably show that Polyak’s momentum achieves acceleration for training a one-layer wide ReLU network and a deep linear network, which are perhaps the two most popular canonical models for studying optimization and deep learning in the literature. Prior work (Du et al., 2019b; Wu et al., 2019c) showed that using vanilla gradient descent, and with an use of over-parameterization, the error decays as  $(1 - \Theta(\frac{1}{\kappa'}))^t$  after  $t$  iterations, where  $\kappa'$  is the condition number of a Gram Matrix. Our result shows that with the appropriate choice of parameters Polyak’s momentum has a rate of  $(1 - \Theta(\frac{1}{\sqrt{\kappa'}}))^t$ . For the deep linear network, prior work (Hu et al., 2020b) showed that vanilla gradient descent has a rate of  $(1 - \Theta(\frac{1}{\kappa}))^t$ , where  $\kappa$  is the condition number of a data matrix. Our result shows an acceleration rate  $(1 - \Theta(\frac{1}{\sqrt{\kappa}}))^t$  is achievable by Polyak’s momentum. This work establishes that momentum does indeed speed up neural net training.

## 1. Introduction

Momentum methods are very popular for training neural networks in various applications (e.g. He et al. (2016); Vaswani et al. (2017); Krizhevsky et al. (2012)). It has been widely observed that the use of momentum helps faster training in deep learning (e.g. Loshchilov & Hutter (2019); Wilson et al. (2017); Cutkosky & Orabona (2019)). Among all the momentum methods, the most popular one seems to be Polyak’s momentum (a.k.a. Heavy Ball momentum) (Polyak, 1964), which is the default choice of momentum in PyTorch and Tensorflow. The success of Polyak’s momentum in deep learning is widely appreciated and almost all of the recently developed adaptive gradient methods like Adam (Kingma & Ba, 2015), AMSGrad (Reddi et al., 2018), and AdaBound (Luo et al., 2019) adopt the use of Polyak’s momentum, instead of Nesterov’s momentum.

However, despite its popularity, little is known in theory about why Polyak’s momentum helps to accelerate training neural networks. Even for convex optimization, problems like strongly convex quadratic problems seem to be one of the few cases that discrete-time Polyak’s momentum method provably achieves faster convergence than standard gradient descent (e.g. Lessard et al. (2016); Goh (2017); Ghadimi et al. (2015); Gitman et al. (2019); Loizou & Richtárik (2017; 2018); Can et al. (2019); Scieur & Pedregosa (2020); Flammarion & Bach (2015); Wilson et al. (2021); Franca et al. (2020); Diakonikolas & Jordan (2019); Shi et al. (2018); Hu (2020)). On the other hand, the theoretical guarantees of Adam, AMSGrad, or AdaBound are only worse if the momentum parameter  $\beta$  is non-zero and the guarantees deteriorate as the momentum parameter increases, which do not show any advantage of the use of momentum (Alacaoglu et al., 2020). Moreover, the convergence rates that have been established for Polyak’s momentum in several related works (Gadat et al., 2016; Sun et al., 2019; Yang et al., 2018; Liu et al., 2020c; Mai & Johansson, 2020) do not improve upon those for vanilla gradient descent or vanilla SGD in the worst case. Lessard et al. (2016); Ghadimi et al. (2015) even show negative cases in *convex* optimization that the use of Polyak’s momentum results in divergence. Furthermore, Kidambi et al. (2018) construct a problem instance for which the momentum method under its

---

<sup>1</sup>School of Computer Science, Georgia Institute of Technology <sup>2</sup>School of Electrical and Computer Engineering, Georgia Institute of Technology. Correspondence to: Jun-Kun Wang <jimwang@gatech.edu>, Chi-Heng Lin <cl3385@gatech.edu>, Jacob Abernethy <prof@gatech.edu>.

---

**Algorithm 1** Gradient descent with Polyak’s momentum (Polyak, 1964) (Equivalent Version 1)

---

- 1: Required: Step size  $\eta$  and momentum parameter  $\beta$ .
  - 2: Init:  $w_0 \in \mathbb{R}^d$  and  $M_{-1} = 0 \in \mathbb{R}^d$ .
  - 3: **for**  $t = 0$  to  $T$  **do**
  - 4:   Given current iterate  $w_t$ , obtain gradient  $\nabla\ell(w_t)$ .
  - 5:   Update momentum  $M_t = \beta M_{t-1} + \nabla\ell(w_t)$ .
  - 6:   Update iterate  $w_{t+1} = w_t - \eta M_t$ .
  - 7: **end for**
- 

optimal tuning is outperformed by other algorithms. Wang et al. (2020) show that Polyak’s momentum helps escape saddle points faster compared with the case without momentum, which is the only provable advantage of Polyak’s momentum in non-convex optimization that we are aware of. A solid understanding of the empirical success of Polyak’s momentum in deep learning has eluded researchers for some time.

We begin this paper by first revisiting the use of Polyak’s momentum for the class of strongly convex quadratic problems,

$$\min_{w \in \mathbb{R}^d} \frac{1}{2} w^\top \Gamma w + b^\top w, \quad (1)$$

where  $\Gamma \in \mathbb{R}^{d \times d}$  is a PSD matrix such that  $\lambda_{\max}(\Gamma) = \alpha$ ,  $\lambda_{\min}(\Gamma) = \mu > 0$ . This is one of the few<sup>1</sup> known examples that Polyak’s momentum has a provable *globally accelerated* linear rate in the *discrete-time* setting. Yet even for this class of problems existing results only establish an accelerated linear rate in an asymptotic sense and several of them do not have an explicit rate in the non-asymptotic regime (e.g. Polyak (1964); Lessard et al. (2016); Mitliagkas (2019); Recht (2018)). Is it possible to prove a non-asymptotic accelerated linear rate in this case? We will return to this question soon.

For general  $\mu$ -strongly convex,  $\alpha$ -smooth, and twice differentiable functions (not necessarily quadratic), denoted as  $F_{\mu, \alpha}^2$ , Theorem 9 in Polyak (1964) shows an asymptotic accelerated linear rate when the iterate is *sufficiently* close to the minimizer so that the landscape can be well approximated by that of a quadratic function. However, the definition of the neighborhood was not very precise in the paper. In this work, we show a locally accelerated linear rate under a quantifiable definition of the neighborhood.

Furthermore, we provably show that Polyak’s momentum helps to achieve a faster convergence for training two neural networks, compared to vanilla GD. The first is training a one-layer ReLU network. Over the past few years there have appeared an enormous number of works considering training a one-layer ReLU network, provably showing con-

---

<sup>1</sup>In Section 2 and Appendix A, we will provide more discussions about this point.

---

**Algorithm 2** Gradient descent with Polyak’s momentum (Polyak, 1964) (Equivalent Version 2)

---

- 1: Required: step size  $\eta$  and momentum parameter  $\beta$ .
  - 2: Init:  $w_0 = w_{-1} \in \mathbb{R}^d$
  - 3: **for**  $t = 0$  to  $T$  **do**
  - 4:   Given current iterate  $w_t$ , obtain gradient  $\nabla\ell(w_t)$ .
  - 5:   Update iterate  $w_{t+1} = w_t - \eta \nabla\ell(w_t) + \beta(w_t - w_{t-1})$ .
  - 6: **end for**
- 

vergence results for vanilla (stochastic) gradient descent (e.g. Li & Liang (2018); Ji & Telgarsky (2020); Li & Yuan (2017); Du et al. (2019b;a); Allen-Zhu et al. (2019); Song & Yang (2019); Zou et al. (2019); Arora et al. (2019c); Jacot et al. (2018); Lee et al. (2019); Chizat et al. (2019); Oymak & Soltanolkotabi (2019); Brutzkus & Globerson (2017); Chen et al. (2020a); Tian (2017); Soltanolkotabi (2017); Bai & Lee (2020); Ghorbani et al. (2019); Li et al. (2020); Hanin & Nica (2020); Daniely (2017); Zou & Gu (2019); Dukler et al. (2020); Daniely (2020); Wei et al. (2019); Yehudai & Shamir (2020); Fang et al. (2019); Su & Yang (2019); Chen et al. (2020b)), as well as for other algorithms (e.g. Zhang et al. (2019); Wu et al. (2019b); Cai et al. (2019); Zhong et al. (2017); Ge et al. (2019); van den Brand et al. (2020); Lee et al. (2020); Pilanci & Ergen (2020)). However, we are not aware of any theoretical works that study the momentum method in neural net training except the work Krichene et al. (2020). These authors show that SGD with Polyak’s momentum (a.k.a. stochastic Heavy Ball) with infinitesimal step size, i.e.  $\eta \rightarrow 0$ , for training a one-hidden-layer network with an infinite number of neurons, i.e.  $m \rightarrow \infty$ , converges to a stationary solution. However, the theoretical result does not show a faster convergence by momentum. In this paper we consider the discrete-time setting and nets with finitely many neurons. We provide a non-asymptotic convergence rate of Polyak’s momentum, establishing a concrete improvement relative to the best-known rates for vanilla gradient descent.

Our setting of training a ReLU network follows the same framework as previous results, including Du et al. (2019b); Arora et al. (2019c); Song & Yang (2019). Specifically, we study training a one-hidden-layer ReLU neural net of the form,

$$\mathcal{N}_W^{\text{ReLU}}(x) := \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \sigma(\langle w^{(r)}, x \rangle), \quad (2)$$

where  $\sigma(z) := z \cdot \mathbb{1}\{z \geq 0\}$  is the ReLU activation,  $w^{(1)}, \dots, w^{(m)} \in \mathbb{R}^d$  are the weights of  $m$  neurons on the first layer,  $a_1, \dots, a_m \in \mathbb{R}$  are weights on the second layer, and  $\mathcal{N}_W^{\text{ReLU}}(x) \in \mathbb{R}$  is the output predicted on input  $x$ . Assume  $n$  number of samples  $\{x_i \in \mathbb{R}^d\}_{i=1}^n$  is given. Following Du et al. (2019b); Arora et al. (2019c); Song &

Yang (2019), we define a Gram matrix  $H \in \mathbb{R}^{n \times n}$  for the weights  $W$  and its expectation  $\bar{H} \in \mathbb{R}^{n \times n}$  over the random draws of  $w^{(r)} \sim N(0, I_d) \in \mathbb{R}^d$  whose  $(i, j)$  entries are defined as follows,

$$H(W)_{i,j} = \sum_{r=1}^m \frac{x_i^\top x_j}{m} \mathbb{1}\{\langle w^{(r)}, x_i \rangle \geq 0 \ \& \ \langle w^{(r)}, x_j \rangle \geq 0\}$$

$$\bar{H}_{i,j} := \mathbb{E}_{w^{(r)}} [x_i^\top x_j \mathbb{1}\{\langle w^{(r)}, x_i \rangle \geq 0 \ \& \ \langle w^{(r)}, x_j \rangle \geq 0\}]. \quad (3)$$

The matrix  $\bar{H}$  is also called a neural tangent kernel (NTK) matrix in the literature (e.g. Jacot et al. (2018); Yang (2019); Bietti & Mairal (2019)). Assume that the smallest eigenvalue  $\lambda_{\min}(\bar{H})$  is strictly positive and certain conditions about the step size and the number of neurons are satisfied. Previous works (Du et al., 2019b; Song & Yang, 2019) show a linear rate of vanilla gradient descent, while we show an accelerated linear rate<sup>2</sup> of gradient descent with Polyak’s momentum. As far as we are aware, our result is the first acceleration result of training an over-parametrized ReLU network.

The second result is training a deep linear network. The deep linear network is a canonical model for studying optimization and deep learning, and in particular for understanding gradient descent (e.g. Shamir (2019); Saxe et al. (2014); Hu et al. (2020b)), studying the optimization landscape (e.g. Kawaguchi (2016); Laurent & von Brecht (2018)), and establishing the effect of implicit regularization (e.g. Moroshko et al. (2020); Ji & Telgarsky (2019); Li et al. (2018); Razin & Cohen (2020); Arora et al. (2019b); Gidel et al. (2019); Gunasekar et al. (2017); Lyu & Li (2020)). In this paper, following (Du & Hu, 2019; Hu et al., 2020b), we study training a  $L$ -layer linear network of the form,

$$\mathcal{N}_W^{L\text{-linear}}(x) := \frac{1}{\sqrt{m^{L-1}d_y}} W^{(L)} W^{(L-1)} \dots W^{(1)} x, \quad (4)$$

where  $W^{(l)} \in \mathbb{R}^{d_l \times d_{l-1}}$  is the weight matrix of the layer  $l \in [L]$ , and  $d_0 = d$ ,  $d_L = d_y$  and  $d_l = m$  for  $l \neq 1, L$ . Therefore, except the first layer  $W^{(1)} \in \mathbb{R}^{m \times d}$  and the last layer  $W^{(L)} \in \mathbb{R}^{d_y \times m}$ , all the intermediate layers are  $m \times m$  square matrices. The scaling  $\frac{1}{\sqrt{m^{L-1}d_y}}$  is necessary to ensure that the network’s output at the initialization  $\mathcal{N}_{W_0}^{L\text{-linear}}(x)$  has the same size as that of the input  $x$ , in the sense that  $\mathbb{E}[\|\mathcal{N}_{W_0}^{L\text{-linear}}(x)\|^2] = \|x\|^2$ , where the expectation is taken over some appropriate random initialization of the network (see e.g. Du & Hu (2019); Hu et al. (2020b)). Hu et al. (2020b) show vanilla gradient descent with orthogonal initialization converges linearly and the required width of the network  $m$  is independent of the depth  $L$ , while we

<sup>2</sup>We borrow the term “accelerated linear rate” from the convex optimization literature (Nesterov, 2013), because the result here has a resemblance to those results in convex optimization, even though the neural network training is a non-convex problem.

show an accelerated linear rate of Polyak’s momentum and the width  $m$  is also independent of  $L$ . To our knowledge, this is the first acceleration result of training a deep linear network.

A careful reader may be tempted by the following line of reasoning: a deep linear network (without activation) is effectively a simple linear model, and we already know that a linear model with the squared loss gives a quadratic objective for which Polyak’s momentum exhibits an accelerated convergence rate. But this intuition, while natural, is not quite right: it is indeed nontrivial even to show that vanilla gradient descent provides a linear rate on deep linear networks (Hu et al., 2020b; Du & Hu, 2019; Shamir, 2019; Arora et al., 2019a; Hardt & Ma, 2016; Wu et al., 2019a; Zou et al., 2020), as the optimization landscape is non-convex. Existing works show that under certain assumptions, all the local minimum are global (Kawaguchi, 2016; Laurent & von Brecht, 2018; Yun et al., 2018; Lu & Kawaguchi, 2017; Zhou & Liang, 2018; Hardt & Ma, 2016). These results are not sufficient to explain the linear convergence of momentum, let alone the acceleration; see Section H in the appendix for an empirical result.

Similarly, it is known that under the NTK regime the output of the ReLU network trained by gradient descent can be approximated by a linear model (e.g. Hu et al. (2020a)). However, this result alone neither implies a global convergence of any algorithm nor characterizes the optimization landscape. While (Liu et al., 2020a) attempt to derive an algorithm-independent equivalence of a class of linear models and a family of wide networks, their result requires the activation function to be differentiable which does not hold for the most prevalent networks like ReLU. Also, their work heavily depends on the regularity of Hessian, making it hard to generalize beyond differentiable networks. Hence, while there has been some progress understanding training of wide networks through linear models, there remains a significant gap in applying this to the momentum dynamics of a non-differentiable networks. Liu et al. (2020b) establish an interesting connection between solving an over-parametrized non-linear system of equations and solving the classical linear system. They show that for smooth and twice differentiable activation, the optimization landscape of an over-parametrized network satisfies a (non-convex) notion called the Polyak-Lokasiewicz (PL) condition (Polyak, 1963), i.e.  $\frac{1}{2} \|\nabla \ell(w)\|^2 \geq \mu (\ell(w) - \ell(w_*))$ , where  $w_*$  is a global minimizer and  $\mu > 0$ . It is not clear whether their result can be extended to ReLU activation, however, and the existing result of Danilova et al. (2018) for the discrete-time Polyak’s momentum under the PL condition does not give an accelerated rate nor is it better than that of vanilla GD. Aujol et al. (2020) show a *variant* of Polyak’s momentum method having an accelerated rate in a *continuous-time* limit for a problem that satisfies PL and has a unique global

minimizer. It is unclear if their result is applicable to our problem. Therefore, showing the advantage of training the ReLU network and the deep linear network by using existing results of Polyak’s momentum can be difficult.

To summarize, our contributions in the present work include

- In convex optimization, we show an accelerated linear rate in the non-asymptotic sense for solving the class of the strongly convex quadratic problems via Polyak’s momentum (Theorem 7). We also provide an analysis of the accelerated local convergence for the class of functions in  $F_{\mu,\alpha}^2$  (Theorem 8). We establish a technical result (Theorem 5) that helps to obtain these non-asymptotic rates.
- In non-convex optimization, we show accelerated linear rates of the discrete-time Polyak’s momentum for training an over-parametrized ReLU network and a deep linear network (Theorems 9 and 10).

Furthermore, we will develop a modular analysis to show all the results in this paper. We identify conditions and propose a meta theorem of acceleration when the momentum method exhibits a certain dynamic, which can be of independent interest. We show that when applying Polyak’s momentum for these problems, the induced dynamics exhibit a form where we can directly apply our meta theorem.

## 2. Preliminaries

Throughout this paper,  $\|\cdot\|_F$  represents the Frobenius norm and  $\|\cdot\|_2$  represents the spectral norm of a matrix, while  $\|\cdot\|$  represents  $l_2$  norm of a vector. We also denote  $\otimes$  the Kronecker product,  $\sigma_{\max}(\cdot) = \|\cdot\|_2$  and  $\sigma_{\min}(\cdot)$  the largest and the smallest singular value of a matrix respectively.

For the case of training neural networks, we will consider minimizing the squared loss

$$\ell(W) := \frac{1}{2} \sum_{i=1}^n (y_i - \mathcal{N}_W(x_i))^2, \quad (5)$$

where  $x_i \in \mathbb{R}^d$  is the feature vector,  $y_i \in \mathbb{R}^{d_y}$  is the label of sample  $i$ , and there are  $n$  number of samples. For training the ReLU network, we have  $\mathcal{N}_W(\cdot) := \mathcal{N}_W^{\text{ReLU}}(\cdot)$ ,  $d_y = 1$ , and  $W := \{w^{(r)}\}_{r=1}^m$ , while for the deep linear network, we have  $\mathcal{N}_W(\cdot) := \mathcal{N}_W^{\text{linear}}(\cdot)$ , and  $W$  represents the set of all the weight matrices, i.e.  $W := \{W^{(l)}\}_{l=1}^L$ . The notation  $A^k$  represents the  $k$ th matrix power of  $A$ .

### 2.1. Prior result of Polyak’s momentum

Algorithm 1 and Algorithm 2 show two equivalent presentations of gradient descent with Polyak’s momentum. Given the same initialization, one can show that Algorithm 1 and Algorithm 2 generate exactly the same iterates during optimization.

Let us briefly describe a prior acceleration result of Polyak’s momentum. The recursive dynamics of Polyak’s momentum for solving the strongly convex quadratic problems (1) can be written as

$$\begin{bmatrix} w_{t+1} - w_* \\ w_t - w_* \end{bmatrix} = \underbrace{\begin{bmatrix} I_d - \eta\Gamma + \beta I_d & -\beta I_d \\ I_d & 0_d \end{bmatrix}}_{:=A} \begin{bmatrix} w_t - w_* \\ w_{t-1} - w_* \end{bmatrix}, \quad (6)$$

where  $w_*$  is the unique minimizer. By a recursive expansion, one can get

$$\left\| \begin{bmatrix} w_t - w_* \\ w_{t-1} - w_* \end{bmatrix} \right\| \leq \|A^t\|_2 \left\| \begin{bmatrix} w_0 - w_* \\ w_{-1} - w_* \end{bmatrix} \right\|. \quad (7)$$

Hence, it suffices to control the spectral norm of the matrix power  $\|A^t\|_2$  for obtaining a convergence rate. In the literature, this is achieved by using Gelfand’s formula.

**Theorem 1.** (Gelfand (1941); see also Foucart (2018)) (Gelfand’s formula) Let  $A$  be a  $d \times d$  matrix. Define the spectral radius  $\rho(A) := \max_{i \in [d]} |\lambda_i(A)|$ , where  $\lambda_i(\cdot)$  is the  $i$ th eigenvalue. Then, there exists a non-negative sequence  $\{\epsilon_t\}$  such that  $\|A^t\|_2 = (\rho(A) + \epsilon_t)^t$  and  $\lim_{t \rightarrow \infty} \epsilon_t = 0$ .

We remark that there is a lack of the convergence rate of  $\epsilon_t$  in Gelfand’s formula in general.

Denote  $\kappa := \alpha/\mu$  the condition number. One can control the spectral radius  $\rho(A)$  as  $\rho(A) \leq 1 - \frac{2}{\sqrt{\kappa+1}}$  by choosing  $\eta$  and  $\beta$  appropriately, which leads to the following result.

**Theorem 2.** (Polyak (1964); see also Lessard et al. (2016); Recht (2018); Mitliagkas (2019)) Gradient descent with Polyak’s momentum with the step size  $\eta = \frac{4}{(\sqrt{\mu} + \sqrt{\alpha})^2}$  and the momentum parameter  $\beta = \left(1 - \frac{2}{\sqrt{\kappa+1}}\right)^2$  has

$$\left\| \begin{bmatrix} w_{t+1} - w_* \\ w_t - w_* \end{bmatrix} \right\| \leq \left(1 - \frac{2}{\sqrt{\kappa+1}} + \epsilon_t\right)^{t+1} \left\| \begin{bmatrix} w_0 - w_* \\ w_{-1} - w_* \end{bmatrix} \right\|,$$

where  $\epsilon_t$  is a non-negative sequence that goes to zero.

That is, when  $t \rightarrow \infty$ , Polyak’s momentum has the  $\left(1 - \frac{2}{\sqrt{\kappa+1}}\right)$  rate, which has a better dependency on the condition number  $\kappa$  than the  $1 - \Theta\left(\frac{1}{\kappa}\right)$  rate of vanilla gradient descent. A concern is that the bound is not quantifiable for a finite  $t$ . On the other hand, we are aware of a different analysis that leverages Chebyshev polynomials instead of Gelfand’s formula (e.g. Liu & Belkin (2018)), which manages to obtain a  $t(1 - \Theta\left(\frac{1}{\sqrt{\kappa}}\right))^t$  convergence rate. So the accelerated linear rate is still obtained in an asymptotic sense. Theorem 9 in Can et al. (2019) shows a rate  $\max\{\bar{C}_1, t\bar{C}_2\}(1 - \Theta\left(\frac{1}{\sqrt{\kappa}}\right))^t$  for some constants  $\bar{C}_1$  and  $\bar{C}_2$  under the same choice of the momentum parameter and the step size as Theorem 2. However, for a large  $t$ , the dominant term could be  $t(1 - \Theta\left(\frac{1}{\sqrt{\kappa}}\right))^t$ . In this paper, we aim at

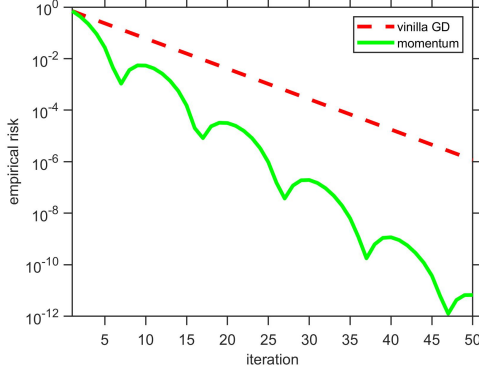


Figure 1. Empirical risk  $\ell(W_t)$  vs. iteration  $t$ . Polyak’s momentum accelerates the optimization process of training an over-parametrized one-layer ReLU network. Experimental details are available in Appendix H.

obtaining a bound that (I) holds for a wide range of values of the parameters, (II) has a dependency on the squared root of the condition number  $\sqrt{\kappa}$ , (III) is quantifiable in each iteration and is better than the rate  $t(1 - \Theta(\frac{1}{\sqrt{\kappa}}))^t$ .

## 2.2. (One-layer ReLU network) Settings and Assumptions

The ReLU activation is not differentiable at zero. So for solving (5), we will replace the notion of gradient in Algorithm 1 and 2 with subgradient  $\frac{\partial \ell(W_t)}{\partial w_t^{(r)}} := \frac{1}{\sqrt{m}} \sum_{i=1}^n (\mathcal{N}_{W_t}(x_i) - y_i) a_r \cdot \mathbb{I}[\langle w_t^{(r)}, x_i \rangle \geq 0] x_i$  and update the neuron  $r$  as  $w_{t+1}^{(r)} = w_t^{(r)} - \eta \frac{\partial \ell(W_t)}{\partial w_t^{(r)}} + \beta(w_t^{(r)} - w_{t-1}^{(r)})$ . As described in the introduction, we assume that the smallest eigenvalue of the Gram matrix  $\bar{H} \in \mathbb{R}^{n \times n}$  is strictly positive, i.e.  $\lambda_{\min}(\bar{H}) > 0$ . We will also denote the largest eigenvalue of the Gram matrix  $\bar{H}$  as  $\lambda_{\max}(\bar{H})$  and denote the condition number of the Gram matrix as  $\kappa := \frac{\lambda_{\max}(\bar{H})}{\lambda_{\min}(\bar{H})}$ . Du et al. (2019b) show that the strict positiveness assumption is indeed mild. Specifically, they show that if no two inputs are parallel, then the least eigenvalue is strictly positive. Panigrahi et al. (2020) were able to provide a quantitative lower bound under certain conditions. Following the same framework of Du et al. (2019b), we consider that each weight vector  $w^{(r)} \in \mathbb{R}^d$  is initialized according to the normal distribution, i.e.  $w^{(r)} \sim N(0, I_d)$ , and each  $a_r \in R$  is sampled from the Rademacher distribution, i.e.  $a_r = 1$  with probability 0.5; and  $a_r = -1$  with probability 0.5. We also assume  $\|x_i\| \leq 1$  for all samples  $i$ . As the previous works (e.g. Li & Liang (2018); Ji & Telgarsky (2020); Du et al. (2019b)), we consider only training the first layer  $\{w^{(r)}\}$  and the second layer  $\{a_r\}$  is fixed throughout the iterations. We will denote  $u_t \in \mathbb{R}^n$  whose  $i_{th}$  entry is the network’s prediction for sample  $i$ , i.e.  $u_t[i] := \mathcal{N}_{W_t}^{\text{ReLU}}(x_i)$  in iteration

$t$  and denote  $y \in \mathbb{R}^n$  the vector whose  $i_{th}$  element is the label of sample  $i$ . The following theorem is a prior result due to Du et al. (2019b).

**Theorem 3.** (Theorem 4.1 in Du et al. (2019b)) Assume that  $\lambda := \lambda_{\min}(\bar{H})/2 > 0$  and that  $w_0^{(r)} \sim N(0, I_d)$  and  $a_r$  uniformly sampled from  $\{-1, 1\}$ . Set the number of nodes  $m = \Omega(\lambda^{-4} n^6 \delta^{-3})$  and the constant step size  $\eta = O(\frac{\lambda}{n^2})$ . Then, with probability at least  $1 - \delta$  over the random initialization, vanilla gradient descent, i.e. Algorithm 1 & 2 with  $\beta = 0$ , has  $\|u_t - y\|^2 \leq (1 - \eta\lambda)^t \cdot \|u_0 - y\|^2$ .

Later Song & Yang (2019) improve the network size  $m$  to  $m = \Omega(\lambda^{-4} n^4 \log^3(n/\delta))$ . Wu et al. (2019c) provide an improved analysis over Du et al. (2019b), which shows that the step size  $\eta$  of vanilla gradient descent can be set as  $\eta = \frac{1}{c_1 \lambda_{\max}(\bar{H})}$  for some quantity  $c_1 > 0$ . The result in turn leads to a convergence rate  $(1 - \frac{1}{c_2 \kappa})$  for some quantity  $c_2 > 0$ . However, the quantities  $c_1$  and  $c_2$  are not universal constants and actually depend on the problem parameters  $\lambda_{\min}(\bar{H})$ ,  $n$ , and  $\delta$ . A question that we will answer in this paper is “Can Polyak’s momentum achieve an accelerated linear rate  $(1 - \Theta(\frac{1}{\sqrt{\kappa}}))$ , where the factor  $\Theta(\frac{1}{\sqrt{\kappa}})$  does not depend on any other problem parameter?”.

## 2.3. (Deep Linear network) Settings and Assumptions

For the case of deep linear networks, we will denote  $X := [x_1, \dots, x_n] \in \mathbb{R}^{d \times n}$  the data matrix and  $Y := [y_1, \dots, y_n] \in \mathbb{R}^{d_y \times n}$  the corresponding label matrix. We will also denote  $\bar{r} := \text{rank}(X)$  and the condition number  $\kappa := \frac{\lambda_{\max}(X^T X)}{\lambda_{\bar{r}}(X^T X)}$ . Following Hu et al. (2020b), we will assume that the linear network is initialized by the orthogonal initialization, which is conducted by sampling uniformly from (scaled) orthogonal matrices such that  $(W_0^{(1)})^T W_0^{(1)} = mI_d$ ,  $W_0^{(L)} (W_0^{(L)})^T = mI_{d_y}$ , and  $(W_0^{(l)})^T W_0^{(l)} = W_0^{(l)} (W_0^{(l)})^T = mI_m$  for layer  $2 \leq l \leq L - 1$ . We will denote  $W^{(j:i)} := W_j W_{j-1} \dots W_i = \prod_{l=i}^j W_l$ , where  $1 \leq i \leq j \leq L$  and  $W^{(i-1:i)} = I$ . We also denote the network’s output  $U := \frac{1}{\sqrt{m^{L-1} d_y}} W^{(L:1)} X \in \mathbb{R}^{d_y \times n}$ .

In our analysis, following Du & Hu (2019); Hu et al. (2020b), we will further assume that (A1) there exists a  $W^*$  such that  $Y = W^* X$ ,  $X \in \mathbb{R}^{d \times \bar{r}}$ , and  $\bar{r} = \text{rank}(X)$ , which is actually without loss of generality (see e.g. the discussion in Appendix B of Du & Hu (2019)).

**Theorem 4.** (Theorem 4.1 in Hu et al. (2020b)) Assume (A1) and the use of the orthogonal initialization. Suppose the width of the deep linear network satisfies  $m \geq C \frac{\|X\|_F^2}{\sigma_{\max}^2(X)} \kappa^2 (d_y (1 + \|W^*\|_2^2) + \log(\bar{r}/\delta))$  and  $m \geq \max\{d_x, d_y\}$  for some  $\delta \in (0, 1)$  and a sufficiently large constant  $C > 0$ . Set the constant step size  $\eta = \frac{d_y}{2L\sigma_{\max}^2(X)}$ .

Then, with probability at least  $1 - \delta$  over the random initialization, vanilla gradient descent, i.e. Algorithm 1 & 2 with  $\beta = 0$ , has  $\|U_t - Y\|_F^2 \leq (1 - \Theta(\frac{1}{\kappa}))^t \cdot \|U_0 - Y\|_F^2$ .

### 3. Modular Analysis

In this section, we will provide a meta theorem for the following dynamics of the residual vector  $\xi_t \in \mathbb{R}^{n_0}$ ,

$$\begin{bmatrix} \xi_{t+1} \\ \xi_t \end{bmatrix} = \begin{bmatrix} I_{n_0} - \eta H + \beta I_{n_0} & -\beta I_{n_0} \\ I_{n_0} & 0_{n_0} \end{bmatrix} \begin{bmatrix} \xi_t \\ \xi_{t-1} \end{bmatrix} + \begin{bmatrix} \varphi_t \\ 0_{n_0} \end{bmatrix}, \quad (8)$$

where  $\eta$  is the step size,  $\beta$  is the momentum parameter,  $H \in \mathbb{R}^{n_0 \times n_0}$  is a PSD matrix,  $\varphi_t \in \mathbb{R}^{n_0}$  is some vector, and  $I_{n_0}$  is the  $n_0 \times n_0$ -dimensional identity matrix. Note that  $\xi_t$  and  $\varphi_t$  depend on the underlying model learned at iteration  $t$ , i.e. depend on  $W_t$ .

We first show that the residual dynamics of Polyak's momentum for solving all the four problems in this paper are in the form of (8). The proof of the following lemmas (Lemma 2, 3, and 4) are available in Appendix B.

#### 3.1. Realization: Strongly convex quadratic problems

One can easily see that the dynamics of Polyak's momentum (6) for solving the strongly convex quadratic problem (1) is in the form of (8). We thus have the following lemma.

**Lemma 1.** *Applying Algorithm 1 or Algorithm 2 to solving the class of strongly convex quadratic problems (1) induces a residual dynamics in the form of (8), where  $\xi_t = w_t - w_*$  (and hence  $n_0 = d$ ),  $H = \Gamma$ ,  $\varphi_t = 0_d$ .*

#### 3.2. Realization: Solving $F_{\mu,\alpha}^2$

A similar result holds for optimizing functions in  $F_{\mu,\alpha}^2$ .

**Lemma 2.** *Applying Algorithm 1 or Algorithm 2 to minimizing a function  $f(w) \in F_{\mu,\alpha}^2$  induces a residual dynamics in the form of (8), where  $\xi_t = w_t - w_*$ ,  $H = \int_0^1 \nabla^2 f((1-\tau)w_0 + \tau w_*) d\tau$ ,  $\varphi_t = \eta(\int_0^1 \nabla^2 f((1-\tau)w_0 + \tau w_*) d\tau - \int_0^1 \nabla^2 f((1-\tau)w_t + \tau w_*) d\tau)(w_t - w_*)$ , where  $w_* := \arg \min_w f(w)$ .*

#### 3.3. Realization: One-layer ReLU network

**More notations:** For the analysis, let us define the event  $A_{ir} := \{\exists w \in \mathbb{R}^d : \|w - w_0^{(r)}\| \leq R^{\text{ReLU}}, \mathbb{1}\{x_i^\top w_0^{(r)}\} \neq \mathbb{1}\{x_i^\top w \geq 0\}\}$ , where  $R^{\text{ReLU}} > 0$  is a number to be determined later. The event  $A_{ir}$  means that there exists a  $w \in \mathbb{R}^d$  which is within the  $R^{\text{ReLU}}$ -ball centered at the initial point  $w_0^{(r)}$  such that its activation pattern of sample  $i$  is different from that of  $w_0^{(r)}$ . We also denote a random set  $S_i := \{r \in [m] : \mathbb{1}\{A_{ir}\} = 0\}$  and its complementary set  $S_i^\perp := [m] \setminus S_i$ .

Lemma 3 below shows that training the ReLU network  $\mathcal{N}_W^{\text{ReLU}}(\cdot)$  via momentum induces the residual dynamics in the form of (8).

**Lemma 3.** *(Residual dynamics of training the ReLU network  $\mathcal{N}_W^{\text{ReLU}}(\cdot)$ ) Denote*

$$(H_t)_{i,j} := H(W_t)_{i,j} = \frac{1}{m} \sum_{r=1}^m x_i^\top x_j \times \mathbb{1}\{\langle w_t^{(r)}, x_i \rangle \geq 0 \& \langle w_t^{(r)}, x_j \rangle \geq 0\}.$$

*Applying Algorithm 1 or Algorithm 2 to (5) for training the ReLU network  $\mathcal{N}_W^{\text{ReLU}}(x)$  induces a residual dynamics in the form of (8) such that  $\xi_t[i] = \mathcal{N}_{W_t}^{\text{ReLU}}(x_i) - y_i$  (and hence  $n_0 = d$ ),  $H = H_0$ , and  $\varphi_t = \phi_t + \iota_t$ , where each element  $i$  of  $\xi_t \in \mathbb{R}^n$  is the residual error of the sample  $i$ , and the  $i_{th}$ -element of  $\phi_t \in \mathbb{R}^n$  satisfies*

$$|\phi_t[i]| \leq \frac{2\eta\sqrt{n}|S_i^\perp|}{m} (\|u_t - y\| + \beta \sum_{s=0}^{t-1} \beta^{t-1-s} \|u_s - y\|),$$

and  $\iota_t = \eta(H_0 - H_t)\xi_t \in \mathbb{R}^n$ .

#### 3.4. Realization: Deep Linear network

Lemma 4 below shows that the residual dynamics due to Polyak's momentum for training the deep linear network is indeed in the form of (8). In the lemma, "vec" stands for the vectorization of the underlying matrix in column-first order.

**Lemma 4.** *(Residual dynamics of training  $\mathcal{N}_W^{L\text{-linear}}(\cdot)$ ) Denote  $M_{t,l}$  the momentum term of layer  $l$  at iteration  $t$ , which is recursively defined as  $M_{t,l} = \beta M_{t,l-1} + \frac{\partial \ell(W_t^{(L:1)})}{\partial W_t^{(l)}}$ . Denote*

$$H_t := \frac{1}{m^{L-1}d_y} \sum_{l=1}^L [(W_t^{(l-1:1)} X)^\top (W_t^{(l-1:1)} X) \otimes W_t^{(L:l+1)} (W_t^{(L:l+1)})^\top] \in \mathbb{R}^{d_y n \times d_y n}.$$

*Applying Algorithm 1 or Algorithm 2 to (5) for training the deep linear network  $\mathcal{N}_W^{L\text{-linear}}(x)$  induces a residual dynamics in the form of (8) such that  $\xi_t = \text{vec}(U_t - Y) \in \mathbb{R}^{d_y n}$  (and hence  $n_0 = d_y n$ ),  $H = H_0$ , and  $\varphi_t = \phi_t + \psi_t + \iota_t \in \mathbb{R}^{d_y n}$ , where the vector  $\phi_t = \frac{1}{\sqrt{m^{L-1}d_y}} \text{vec}(\Phi_t X)$  with*

$$\Phi_t = \Pi_l \left( W_t^{(l)} - \eta M_{t,l} \right) - W_t^{(L:1)} + \eta \sum_{l=1}^L W_t^{(L:l+1)} M_{t,l} W_t^{(l-1:1)},$$

and the vector  $\psi_t$  is

$$\psi_t = \frac{1}{\sqrt{m^{L-1}d_y}} \text{vec}((L-1)\beta W_t^{(L:1)} X + \beta W_{t-1}^{(L:1)} X - \beta \sum_{l=1}^L W_t^{(L:l+1)} W_{t-1}^{(l)} W_t^{(l-1:1)} X),$$

and  $\iota_t = \eta(H_0 - H_t)\xi_t$ .

### 3.5. A key theorem of bounding a matrix-vector product

Our meta theorem of acceleration will be based on Theorem 5 in the following, which upper-bounds the size of the matrix-vector product of a matrix power  $A^k$  and a vector  $v_0$ . Compared to Gelfand's formula (Theorem 1), Theorem 5 below provides a better control of the size of the matrix-vector product, since it avoids the dependency on the unknown sequence  $\{\epsilon_t\}$ . The result can be of independent interest and might be useful for analyzing Polyak's momentum for other problems in future research.

**Theorem 5.** Let  $A := \begin{bmatrix} (1+\beta)I_n - \eta H & -\beta I_n \\ I_n & 0 \end{bmatrix} \in \mathbb{R}^{2n \times 2n}$ . Suppose that  $H \in \mathbb{R}^{n \times n}$  is a positive semidefinite matrix. Fix a vector  $v_0 \in \mathbb{R}^n$ . If  $\beta$  is chosen to satisfy  $1 \geq \beta > \max\left\{\left(1 - \sqrt{\eta\lambda_{\min}(H)}\right)^2, \left(1 - \sqrt{\eta\lambda_{\max}(H)}\right)^2\right\}$ , then

$$\|A^k v_0\| \leq (\sqrt{\beta})^k C_0 \|v_0\|, \quad (9)$$

where the constant

$$C_0 := \frac{\sqrt{2}(\beta+1)}{\sqrt{\min\{h(\beta, \eta\lambda_{\min}(H)), h(\beta, \eta\lambda_{\max}(H))\}}} \geq 1, \quad (10)$$

and the function  $h(\beta, z)$  is defined as  $h(\beta, z) := -\left(\beta - (1 - \sqrt{z})^2\right)\left(\beta - (1 + \sqrt{z})^2\right)$ .

Note that the constant  $C_0$  in Theorem 5 depends on  $\beta$  and  $\eta H$ . It should be written as  $C_0(\beta, \eta H)$  to be precise. However, for the brevity, we will simply denote it as  $C_0$  when the underlying choice of  $\beta$  and  $\eta H$  is clear from the context. The proof of Theorem 5 is available in Appendix C. Theorem 5 allows us to derive a concrete upper bound of the residual errors in each iteration of momentum, and consequently allows us to show an accelerated linear rate in the non-asymptotic sense. The favorable property of the bound will also help to analyze Polyak's momentum for training the neural networks. As shown later in this paper, we will need to guarantee the progress of Polyak's momentum in each iteration, which is not possible if we only have a quantifiable bound in the limit. Based on Theorem 5, we have the following corollary. The proof is in Appendix C.1.

**Corollary 1.** Assume that  $\lambda_{\min}(H) > 0$ . Denote  $\kappa := \lambda_{\max}(H)/\lambda_{\min}(H)$ . Set  $\eta = 1/\lambda_{\max}(H)$  and set  $\beta = \left(1 - \frac{1}{2}\sqrt{\eta\lambda_{\min}(H)}\right)^2 = \left(1 - \frac{1}{2\sqrt{\kappa}}\right)^2$ . Then,  $C_0 \leq 4\sqrt{\kappa}$ .

### 3.6. Meta theorem

Let  $\lambda > 0$  be the smallest eigenvalue of the matrix  $H$  that appears on the residual dynamics (8). Our goal is to show

that the residual errors satisfy

$$\left\| \begin{bmatrix} \xi_s \\ \xi_{s-1} \end{bmatrix} \right\| \leq (\sqrt{\beta} + \mathbb{1}_\varphi C_2)^s (C_0 + \mathbb{1}_\varphi C_1) \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|, \quad (11)$$

where  $C_0$  is the constant defined on (10), and  $C_1, C_2 \geq 0$  are some constants,  $\mathbb{1}_\varphi$  is an indicator if any  $\varphi_t$  on the residual dynamics (8) is a non-zero vector. For the case of training the neural networks, we have  $\mathbb{1}_\varphi = 1$ .

**Theorem 6.** (Meta theorem for the residual dynamics (8)) Assume that the step size  $\eta$  and the momentum parameter  $\beta$  satisfying  $1 \geq \beta > \max\left\{\left(1 - \sqrt{\eta\lambda_{\min}(H)}\right)^2, \left(1 - \sqrt{\eta\lambda_{\max}(H)}\right)^2\right\}$ , are set appropriately so that (11) holds at iteration  $s = 0, 1, \dots, t-1$  implies that

$$\left\| \sum_{s=0}^{t-1} A^{t-s-1} \begin{bmatrix} \varphi_s \\ 0 \end{bmatrix} \right\| \leq (\sqrt{\beta} + \mathbb{1}_\varphi C_2)^t C_3 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|. \quad (12)$$

Then, we have

$$\left\| \begin{bmatrix} \xi_t \\ \xi_{t-1} \end{bmatrix} \right\| \leq (\sqrt{\beta} + \mathbb{1}_\varphi C_2)^t (C_0 + \mathbb{1}_\varphi C_1) \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|, \quad (13)$$

holds for all  $t$ , where  $C_0$  is defined on (10) and  $C_1, C_2, C_3 \geq 0$  are some constants satisfying:

$$\begin{aligned} (\sqrt{\beta})^t C_0 + (\sqrt{\beta} + \mathbb{1}_\varphi C_2)^t \mathbb{1}_\varphi C_3 \leq \\ (\sqrt{\beta} + \mathbb{1}_\varphi C_2)^t (C_0 + \mathbb{1}_\varphi C_1). \end{aligned} \quad (14)$$

*Proof.* The proof is by induction. At  $s = 0$ , (11) holds since  $C_0 \geq 1$  by Theorem 5. Now assume that the inequality holds at  $s = 0, 1, \dots, t-1$ . Consider iteration  $t$ . Recursively expanding the dynamics (8), we have

$$\begin{bmatrix} \xi_t \\ \xi_{t-1} \end{bmatrix} = A^t \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} + \sum_{s=0}^{t-1} A^{t-s-1} \begin{bmatrix} \varphi_s \\ 0 \end{bmatrix}. \quad (15)$$

By Theorem 5, the first term on the r.h.s. of (15) can be bounded by

$$\|A^t \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix}\| \leq (\sqrt{\beta})^t C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\| \quad (16)$$

By assumption, given (11) holds at  $s = 0, 1, \dots, t-1$ , we have (12). Combining (12), (14), (15), and (16), we have (13) and hence the proof is completed.  $\square$

**Remark:** As shown in the proof, we need the residual errors be tightly bounded as (11) in each iteration. Theorem 5 is critical for establishing the desired result. On the other hand, it would become tricky if instead we use Gelfand's formula or other techniques in the related works that lead to a convergence rate in the form of  $O(t\theta^t)$ .

## 4. Main results

The important lemmas and theorems in the previous section help to show our main results in the following subsections. The high-level idea to obtain the results is by using the meta theorem (i.e. Theorem 6). Specifically, we will need to show that if the underlying residual dynamics satisfy (11) for all the previous iterations, then the terms  $\{\varphi_s\}$  in the dynamics satisfy (12). This condition trivially holds for the case of the quadratic problems, since there is no such term. On the other hand, for solving the other problems, we need to carefully show that the condition holds. For example, according to Lemma 3, showing acceleration for the ReLU network will require bounding terms like  $\|(H_0 - H_s)\xi_s\|$  (and other terms as well), where  $H_0 - H_s$  corresponds to the difference of the kernel matrix at two different time steps. By controlling the width of the network, we can guarantee that the change is not too much. A similar result can be obtained for the problem of the deep linear network. The high-level idea is simple but the analysis of the problems of the neural networks can be tedious.

### 4.1. Non-asymptotic accelerated linear rate for solving strongly convex quadratic problems

**Theorem 7.** *Assume the momentum parameter  $\beta$  satisfies  $1 \geq \beta > \max\{(1 - \sqrt{\eta\mu})^2, (1 - \sqrt{\eta\alpha})^2\}$ . Gradient descent with Polyak’s momentum for solving (1) has*

$$\left\| \begin{bmatrix} w_t - w_* \\ w_{t-1} - w_* \end{bmatrix} \right\| \leq \left(\sqrt{\beta}\right)^t C_0 \left\| \begin{bmatrix} w_0 - w_* \\ w_{-1} - w_* \end{bmatrix} \right\|, \quad (17)$$

where the constant  $C_0$  is defined as

$$C_0 := \frac{\sqrt{2(\beta+1)}}{\sqrt{\min\{h(\beta, \eta\lambda_{\min}(\Gamma)), h(\beta, \eta\lambda_{\max}(\Gamma))\}}} \geq 1, \quad (18)$$

and  $h(\beta, z) = -\left(\beta - (1 - \sqrt{z})^2\right)\left(\beta - (1 + \sqrt{z})^2\right)$ . Consequently, if the step size  $\eta = \frac{1}{\alpha}$  and the momentum parameter  $\beta = \left(1 - \frac{1}{2\sqrt{\kappa}}\right)^2$ , then it has

$$\left\| \begin{bmatrix} w_t - w_* \\ w_{t-1} - w_* \end{bmatrix} \right\| \leq \left(1 - \frac{1}{2\sqrt{\kappa}}\right)^t 4\sqrt{\kappa} \left\| \begin{bmatrix} w_0 - w_* \\ w_{-1} - w_* \end{bmatrix} \right\|. \quad (19)$$

Furthermore, if  $\eta = \frac{4}{(\sqrt{\mu} + \sqrt{\alpha})^2}$  and  $\beta$  approaches  $\beta \rightarrow \left(1 - \frac{2}{\sqrt{\kappa+1}}\right)^2$  from above, then it has a convergence rate approximately  $\left(1 - \frac{2}{\sqrt{\kappa+1}}\right)$  as  $t \rightarrow \infty$ .

The convergence rates shown in the above theorem do not depend on the unknown sequence  $\{\epsilon_t\}$ . Moreover, the rates depend on the squared root of the condition number  $\sqrt{\kappa}$ . We have hence established a non-asymptotic accelerated linear rate of Polyak’s momentum, which helps to show the advantage of Polyak’s momentum over vanilla gradient

descent in the finite  $t$  regime. Our result also recovers the rate  $\left(1 - \frac{2}{\sqrt{\kappa+1}}\right)$  asymptotically under the same choices of the parameters as the previous works. The detailed proof can be found in Appendix D, which is actually a trivial application of Lemma 1, Theorem 6, and Corollary 1 with  $C_1 = C_2 = C_3 = 0$ .

### 4.2. Non-asymptotic accelerated linear rate of the local convergence for solving $f(\cdot) \in F_{\mu, \alpha}^2$

Here we provide a local acceleration result of the discrete-time Polyak’s momentum for general smooth strongly convex and twice differentiable function  $F_{\mu, \alpha}^2$ . Compared to Theorem 9 of (Polyak, 1964), Theorem 8 clearly indicates the required distance that ensures an acceleration when the iterate is in the neighborhood of the global minimizer. Furthermore, the rate is in the non-asymptotic sense instead of the asymptotic one. We defer the proof of Theorem 8 to Appendix E.

**Theorem 8.** *Assume that the function  $f(\cdot) \in F_{\mu, \alpha}^2$  and its Hessian is  $\alpha$ -Lipschitz. Denote the condition number  $\kappa := \frac{\alpha}{\mu}$ . Suppose that the initial point satisfies*

*$\left\| \begin{bmatrix} w_0 - w_* \\ w_{-1} - w_* \end{bmatrix} \right\| \leq \frac{1}{683\kappa^{3/2}}$ . Then, Gradient descent with Polyak’s momentum with the step size  $\eta = \frac{1}{\alpha}$  and the momentum parameter  $\beta = \left(1 - \frac{1}{2\sqrt{\kappa}}\right)^2$  for solving  $\min_w f(w)$  has*

$$\left\| \begin{bmatrix} w_{t+1} - w_* \\ w_t - w_* \end{bmatrix} \right\| \leq \left(1 - \frac{1}{4\sqrt{\kappa}}\right)^{t+1} 8\sqrt{\kappa} \left\| \begin{bmatrix} w_0 - w_* \\ w_{-1} - w_* \end{bmatrix} \right\|, \quad (20)$$

where  $w_* = \arg \min_w f(w)$ .

### 4.3. Acceleration for training $\mathcal{N}_W^{\text{ReLU}}(x)$

Before introducing our result of training the ReLU network, we need the following lemma.

**Lemma 5.** *[Lemma 3.1 in Du et al. (2019b) and Song & Yang (2019)] Set  $m = \Omega(\lambda^{-2}n^2 \log(n/\delta))$ . Suppose that the neurons  $w_0^{(1)}, \dots, w_0^{(m)}$  are i.i.d. generated by  $N(0, I_d)$  initially. Then, with probability at least  $1 - \delta$ , it holds that*

$$\|H_0 - \bar{H}\|_F \leq \frac{\lambda_{\min}(\bar{H})}{4}, \quad \lambda_{\min}(H_0) \geq \frac{3}{4}\lambda_{\min}(\bar{H}),$$

$$\text{and} \quad \lambda_{\max}(H_0) \leq \lambda_{\max}(\bar{H}) + \frac{\lambda_{\min}(\bar{H})}{4}.$$

Lemma 5 shows that by the random initialization, with probability  $1 - \delta$ , the least eigenvalue of the Gram matrix  $H := H_0$  defined in Lemma 3 is lower-bounded and the largest eigenvalue is close to  $\lambda_{\max}(\bar{H})$ . Furthermore, Lemma 5 implies that the condition number of the Gram matrix  $H_0$  at the initialization  $\hat{\kappa} := \frac{\lambda_{\max}(H_0)}{\lambda_{\min}(H_0)}$  satisfies



$\hat{\kappa} \leq \frac{4}{3}\kappa + \frac{1}{3}$ , where  $\kappa := \frac{\lambda_{\max}(\bar{H})}{\lambda_{\min}(\bar{H})}$ .

**Theorem 9.** (One-layer ReLU network  $\mathcal{N}_W^{\text{ReLU}}(x)$ ) Assume that  $\lambda := \frac{3\lambda_{\min}(\bar{H})}{4} > 0$  and that  $w_0^{(r)} \sim N(0, I_d)$  and  $a_r$  uniformly sampled from  $\{-1, 1\}$ . Denote  $\lambda_{\max} := \lambda_{\max}(\bar{H}) + \frac{\lambda_{\min}(\bar{H})}{4}$  and denote  $\hat{\kappa} := \lambda_{\max}/\lambda = (4\kappa + 1)/3$ . Set a constant step size  $\eta = \frac{1}{\lambda_{\max}}$ , fix momentum parameter  $\beta = \left(1 - \frac{1}{2\hat{\kappa}}\right)^2$ , and finally set the number of network nodes  $m = \Omega(\lambda^{-4}n^4\kappa^2 \log^3(n/\delta))$ . Then, with probability at least  $1 - \delta$  over the random initialization, gradient descent with Polyak’s momentum satisfies for any  $t$ ,

$$\left\| \begin{bmatrix} \xi_t \\ \xi_{t-1} \end{bmatrix} \right\| \leq \left(1 - \frac{1}{4\sqrt{\hat{\kappa}}}\right)^t \cdot 8\sqrt{\hat{\kappa}} \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|. \quad (21)$$

We remark that  $\hat{\kappa}$ , which is the condition number of the Gram matrix  $H_0$ , is within a constant factor of the condition number of  $\bar{H}$ . Therefore, Theorem 9 essentially shows an accelerated linear rate  $\left(1 - \Theta\left(\frac{1}{\sqrt{\kappa}}\right)\right)$ . The rate has an improved dependency on the condition number, i.e.  $\sqrt{\kappa}$  instead of  $\kappa$ , which shows the advantage of Polyak’s momentum over vanilla GD when the condition number is large. We believe this is an interesting result, as the acceleration is akin to that in convex optimization, e.g. Nesterov (2013); Shi et al. (2018).

Our result also implies that over-parametrization helps acceleration in optimization. To our knowledge, in the literature, there is little theory of understanding why over-parametrization can help training a neural network faster. The only exception that we are aware of is Arora et al. (2018), which shows that the dynamic of vanilla gradient descent for an over-parametrized objective function exhibits some momentum terms, although their message is very different from ours. The proof of Theorem 9 is in Appendix F.

#### 4.4. Acceleration for training $\mathcal{N}_W^{L\text{-linear}}(x)$

**Theorem 10.** (Deep linear network  $\mathcal{N}_W^{L\text{-linear}}(x)$ ) Denote  $\lambda := \frac{L\sigma_{\min}^2(X)}{d_y}$  and  $\kappa := \frac{\sigma_{\max}^2(X)}{\sigma_{\min}^2(X)}$ . Set a constant step size  $\eta = \frac{d_y}{L\sigma_{\max}^2(X)}$ , fix momentum parameter  $\beta = \left(1 - \frac{1}{2\sqrt{\kappa}}\right)^2$ , and finally set a parameter  $m$  that controls the width  $m \geq C \frac{\kappa^5}{\sigma_{\max}^2(X)} (d_y(1 + \|W^*\|_2^2) + \log(\bar{r}/\delta))$  and  $m \geq \max\{d_x, d_y\}$  for some constant  $C > 0$ . Then, with probability at least  $1 - \delta$  over the random orthogonal initialization, gradient descent with Polyak’s momentum satisfies for any  $t$ ,

$$\left\| \begin{bmatrix} \xi_t \\ \xi_{t-1} \end{bmatrix} \right\| \leq \left(1 - \frac{1}{4\sqrt{\kappa}}\right)^t \cdot 8\sqrt{\kappa} \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|. \quad (22)$$

Compared with Theorem 4 of Hu et al. (2020b) for vanilla GD, our result clearly shows the acceleration via Polyak’s

momentum. Furthermore, the result suggests that the depth does not hurt optimization. Acceleration is achieved for any depth  $L$  and the required width  $m$  is independent of the depth  $L$  as Hu et al. (2020b); Zou et al. (2020) (of vanilla GD). The proof of Theorem 10 is in Appendix G.

## 5. Conclusion

We show some non-asymptotic acceleration results of the discrete-time Polyak’s momentum in this paper. The results not only improve the previous results in convex optimization but also establish the first time that Polyak’s momentum has provable acceleration for training certain neural networks. We analyze all the acceleration results from a modular framework. We hope the framework can serve as a building block towards understanding Polyak’s momentum in a more unified way.

## 6. Acknowledgment

The authors thank Daniel Pozo for catching a typo. The authors acknowledge support of NSF IIS Award 1910077. JW also thanks IDEaS-TRIAD Research Scholarship 03GR10000818.

## References

- Alacaoglu, A., Malitsky, Y., Mertikopoulos, P., and Cevher, V. A new regret analysis for adam-type algorithms. *ICML*, 2020.
- Allen-Zhu, Z., Li, Y., and Song, Z. A convergence theory for deep learning via overparameterization. *ICML*, 2019.
- Arora, S., Cohen, N., and Hazan, E. On the optimization of deep networks: Implicit acceleration by overparameterization. *ICML*, 2018.
- Arora, S., Cohen, N., Golowich, N., and Hu, W. A convergence analysis of gradient descent for deep linear neural networks. *ICLR*, 2019a.
- Arora, S., Cohen, N., Hu, W., and Luo, Y. Implicit regularization in deep matrix factorization. *NeurIPS*, 2019b.
- Arora, S., Du, S. S., Hu, W., Li, Z., and Wang, R. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. *NeurIPS*, 2019c.
- Aujol, J.-F., Dossal, C., and Rondepierre, A. Convergence rates of the heavy-ball method with lojasiewicz property. *hal-02928958*, 2020.
- Bai, Y. and Lee, J. D. Beyond linearization: On quadratic and higher-order approximation of wide neural networks. *ICLR*, 2020.

- Bietti, A. and Mairal, J. On the inductive bias of neural tangent kernels. *NeurIPS*, 2019.
- Brutzkus, A. and Globerson, A. Globally optimal gradient descent for a convnet with gaussian inputs. *ICML*, 2017.
- Cai, T., Gao, R., Hou, J., Chen, S., Wang, D., He, D., Zhang, Z., and Wang, L. A gram-gauss-newton method learning overparameterized deep neural networks for regression problems. *arXiv.org:1905.11675*, 2019.
- Can, B., Gürbüzbalaban, M., and Zhu, L. Accelerated linear convergence of stochastic momentum methods in wasserstein distances. *ICML*, 2019.
- Chen, S., He, H., and Su, W. J. Label-aware neural tangent kernel: Toward better generalization and local elasticity. *NeurIPS*, 2020a.
- Chen, Z., Cao, Y., Gu, Q., and Zhang, T. A generalized neural tangent kernel analysis for two-layer neural network. *NeurIPS*, 2020b.
- Chizat, L., Oyallon, E., and Bach, F. On lazy training in differentiable programming. *NeurIPS*, 2019.
- Cutkosky, A. and Orabona, F. Momentum-based variance reduction in non-convex sgd. *NeurIPS*, 2019.
- Daniely, A. Sgd learns the conjugate kernel class of the network. *NeurIPS*, 2017.
- Daniely, A. Memorizing gaussians with no over-parameterization via gradient decent on neural networks. *arXiv:1909.11837*, 2020.
- Danilova, M., Kulakova, A., and Polyak, B. Non-monotone behavior of the heavy ball method. *arXiv:1811.00658*, 2018.
- Diakonikolas, J. and Jordan, M. I. Generalized momentum-based methods: A hamiltonian perspective. *arXiv:1906.00436*, 2019.
- Du, S. S. and Hu, W. Width provably matters in optimization for deep linear neural networks. *ICML*, 2019.
- Du, S. S., Lee, J. D., Li, H., Wang, L., and Zhai, X. Gradient descent finds global minima of deep neural networks. *ICML*, 2019a.
- Du, S. S., Zhai, X., Póczos, B., and Singh, A. Gradient descent provably optimizes over-parameterized neural networks. *ICLR*, 2019b.
- Dukler, Y., Gu, Q., and Montufar, G. Optimization theory for relu neural networks trained with normalization layers. *ICML*, 2020.
- Fang, C., Dong, H., and Zhang, T. Over parameterized two-level neural networks can learn near optimal feature representations. *arXiv:1910.11508*, 2019.
- Flammarion, N. and Bach, F. From averaging to acceleration, there is only a step-size. *COLT*, 2015.
- Foucart, S. Matrix norms and spectral radii. *Online lecture note*, 2018.
- Franca, G., Sulam, J., Robinson, D. P., and Vidal, R. Conformal symplectic and relativistic optimization. *Journal of Statistical Mechanics: Theory and Experiment*, 2020.
- Gadat, S., Panloup, F., and Saadane, S. Stochastic heavy ball. *arXiv:1609.04228*, 2016.
- Ge, R., Kuditipudi, R., Li, Z., and Wang, X. Learning two-layer neural networks with symmetric inputs. *ICLR*, 2019.
- Gelfand, I. Normierte ringe. *Mat. Sbornik*, 1941.
- Ghadimi, E., Feyzmahdavian, H. R., and Johansson, M. Global convergence of the heavy-ball method for convex optimization. *ECC*, 2015.
- Ghorbani, B., Mei, S., Misiakiewicz, T., , and Montanari, A. Linearized two-layers neural networks in high dimension. *arXiv:1904.12191*, 2019.
- Gidel, G., Bach, F., and Lacoste-Julien, S. Implicit regularization of discrete gradient dynamics in linear neural networks. *NeurIPS*, 2019.
- Gitman, I., Lang, H., Zhang, P., and Xiao, L. Understanding the role of momentum in stochastic gradient methods. *NeurIPS*, 2019.
- Goh, G. Why momentum really works. *Distill*, 2017.
- Gunasekar, S., Woodworth, B., Bhojanapalli, S., Neyshabur, B., and Srebro, N. Implicit regularization in matrix factorization. *NeurIPS*, 2017.
- Hanin, B. and Nica, M. Finite depth and width corrections to the neural tangent kernel. *ICLR*, 2020.
- Hardt, M. and Ma, T. Identity matters in deep learning. *ICLR*, 2016.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- Hu, B. Unifying the analysis in control and optimization via semidefinite programs. *Lecture Note*, 2020.
- Hu, W., Xiao, L., Adlam, B., and Pennington, J. The surprising simplicity of the early-time learning dynamics of neural networks. *NeurIPS*, 2020a.

- Hu, W., Xiao, L., and Pennington, J. Provable benefit of orthogonal initialization in optimizing deep linear networks. *ICLR*, 2020b.
- Jacot, A., Gabriel, F., and Hongler, C. Neural tangent kernel: Convergence and generalization in neural networks. *NeurIPS*, 2018.
- Ji, Z. and Telgarsky, M. Gradient descent aligns the layers of deep linear networks. *ICLR*, 2019.
- Ji, Z. and Telgarsky, M. Polylogarithmic width suffices for gradient descent to achieve arbitrarily small test error with shallow relu networks. *ICLR*, 2020.
- Kawaguchi, K. Deep learning without poor local minima. *NeurIPS*, 2016.
- Kidambi, R., Netrapalli, P., Jain, P., and Kakade, S. M. On the insufficiency of existing momentum schemes for stochastic optimization. *ICLR*, 2018.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *ICLR*, 2015.
- Krichene, W., Caluyay, K. F., and Halder, A. Global convergence of second-order dynamics in two-layer neural networks. *arXiv:2006.07867*, 2020.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. Imagenet classification with deep convolutional neural networks. *NeurIPS*, 2012.
- Laurent, T. and von Brecht, J. Deep linear networks with arbitrary loss: All local minima are global. *ICML*, 2018.
- Lee, J., Xiao, L., Schoenholz, S. S., Bahri, Y., Sohl-Dickstein, J., and Pennington, J. Wide neural networks of any depth evolve as linear models under gradient descent. *NeurIPS*, 2019.
- Lee, J. D., Shen, R., Song, Z., Wang, M., and Yu, Z. Generalized leverage score sampling for neural networks. *arXiv:2009.09829*, 2020.
- Lessard, L., Recht, B., and Packard, A. Analysis and design of optimization algorithms via integral quadratic constraints. *SIAM Journal on Optimization*, 2016.
- Li, Y. and Liang, Y. Learning overparameterized neural networks via stochastic gradient descent on structured data. *NeurIPS*, 2018.
- Li, Y. and Yuan, Y. Convergence analysis of two-layer neural networks with relu activation. *NeurIPS*, 2017.
- Li, Y., Ma, T., and Zhang, H. Algorithmic regularization in over-parameterized matrix sensing and neural networks with quadratic activations. *COLT*, 2018.
- Li, Y., Ma, T., and Zhang, H. Learning over-parametrized two-layer relu neural networks beyond ntk. *COLT*, 2020.
- Liu, C. and Belkin, M. Parametrized accelerated methods free of condition number. *arXiv:1802.10235*, 2018.
- Liu, C., Zhu, L., and Belkin, M. On the linearity of large non-linear models: when and why the tangent kernel is constant. *arXiv:2010.01092*, 2020a.
- Liu, C., Zhu, L., and Belkin, M. Toward a theory of optimization for over-parameterized systems of non-linear equations: the lessons of deep learning. *arXiv:2003.00307*, 2020b.
- Liu, Y., Gao, Y., and Yin, W. An improved analysis of stochastic gradient descent with momentum. *NeurIPS*, 2020c.
- Loizou, N. and Richtárik, P. Momentum and stochastic momentum for stochastic gradient, newton, proximal point and subspace descent methods. *arXiv:1712.09677*, 2017.
- Loizou, N. and Richtárik, P. Accelerated gossip via stochastic heavy ball method. *Allerton*, 2018.
- Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. *ICLR*, 2019.
- Lu, H. and Kawaguchi, K. Depth creates no bad local minima. *arXiv:1702.08580*, 2017.
- Luo, L., Xiong, Y., Liu, Y., and Sun, X. Adaptive gradient methods with dynamic bound of learning rate. *ICLR*, 2019.
- Lyu, K. and Li, J. Gradient descent maximizes the margin of homogeneous neural networks. *ICLR*, 2020.
- Mai, V. V. and Johansson, M. Convergence of a stochastic gradient method with momentum for non-smooth non-convex optimization. *ICML*, 2020.
- Mitliagkas, I. Accelerated methods - polyak's momentum (heavy ball method). *Online Lecture Note*, 2019.
- Moroshko, E., Gunasekar, S., Woodworth, B., Lee, J. D., Srebro, N., and Soudry, D. Implicit bias in deep linear classification: Initialization scale vs training accuracy. *NeurIPS*, 2020.
- Nesterov, Y. Introductory lectures on convex optimization: a basic course. *Springer*, 2013.
- Oymak, S. and Soltanolkotabi, M. Towards moderate overparameterization: global convergence guarantees for training shallow neural networks. *arXiv:1902.04674*, 2019.

- Panigrahi, A., Shetty, A., and Goyal, N. Effect of activation functions on the training of overparametrized neural nets. *ICLR*, 2020.
- Pilanci, M. and Ergen, T. Neural networks are convex regularizers: Exact polynomial-time convex optimization formulations for two-layer networks. *ICML*, 2020.
- Polyak, B. Gradient methods for minimizing functionals. *Zhurnal Vychislitel'noi Matematiki i Matematicheskoi Fiziki*, 1963.
- Polyak, B. Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics*, 1964.
- Razin, N. and Cohen, N. Implicit regularization in deep learning may not be explainable by norms. *NeurIPS2020*, 2020.
- Recht, B. Lyapunov analysis and the heavy ball method. *Lecture note*, 2018.
- Reddi, S. J., Kale, S., and Kumar, S. On the convergence of adam and beyond. *ICLR*, 2018.
- Saxe, A. M., McClelland, J. L., and Ganguli, S. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *ICLR*, 2014.
- Scieur, D. and Pedregosa, F. Universal average-case optimality of polyak momentum. *ICML*, 2020.
- Shamir, O. Exponential convergence time of gradient descent for one-dimensional deep linear neural networks. *COLT*, 2019.
- Shi, B., Du, S. S., Jordan, M. I., and Su, W. J. Understanding the acceleration phenomenon via high-resolution differential equations. *arXiv:1810.08907*, 2018.
- Soltanolkotabi, M. Learning relus via gradient descent. *NeurIPS*, 2017.
- Song, Z. and Yang, X. Quadratic suffices for over-parametrization via matrix chernoff bound. *arXiv:1906.03593*, 2019.
- Su, L. and Yang, P. On learning over-parameterized neural networks: A functional approximation perspective. *NeurIPS*, 2019.
- Sun, T., Yin, P., Li, D., Huang, C., Guan, L., and Jiang, H. Non-ergodic convergence analysis of heavy-ball algorithms. *AAAI*, 2019.
- Tian, Y. An analytical formula of population gradient for two-layered relu network and its applications in convergence and critical point analysis. *ICML*, 2017.
- van den Brand, J., Peng, B., Song, Z., and Weinstein, O. Training (overparametrized) neural networks in near-linear time. *arXiv:2006.11648*, 2020.
- Vaswani, A., Shazeer, N., Parmar, N., and et al. Attention is all you need. *NeurIPS*, 2017.
- Wang, J.-K., Lin, C.-H., and Abernethy, J. Escaping saddle points faster with stochastic momentum. *ICLR*, 2020.
- Wei, C., Lee, J. D., Liu, Q., and Ma, T. Regularization matters: Generalization and optimization of neural nets v.s. their induced kernel. *NeurIPS*, 2019.
- Wilson, A. C., Roelofs, R., Stern, M., Srebro, N., , and Recht, B. The marginal value of adaptive gradient methods in machine learning. *NeurIPS*, 2017.
- Wilson, A. C., Jordan, M., and Recht, B. A lyapunov analysis of momentum methods in optimization. *JMLR*, 2021.
- Wu, L., Wang, Q., and Ma, C. Global convergence of gradient descent for deep linear residual networks. *NeurIPS*, 2019a.
- Wu, S., Dimakis, A. G., and Sanghavi, S. Learning distributions generated by one-layer relu networks. *NeurIPS*, 2019b.
- Wu, X., Du, S. S., and Ward, R. Global convergence of adaptive gradient methods for an over-parameterized neural network. *arXiv:1902.07111*, 2019c.
- Yang, G. Scaling limits of wide neural networks with weight sharing: Gaussian process behavior, gradient independence, and neural tangent kernel derivation. *arXiv:1902.04760*, 2019.
- Yang, T., Lin, Q., and Li, Z. Unified convergence analysis of stochastic momentum methods for convex and non-convex optimization. *IJCAI*, 2018.
- Yehudai, G. and Shamir, O. Learning a single neuron with gradient methods. *COLT*, 2020.
- Yun, C., Sra, S., and Jadbabaie, A. Global optimality conditions for deep neural networks. *ICLR*, 2018.
- Zhang, G., Martens, J., and Grosse, R. B. Fast convergence of natural gradient descent for over-parameterized neural networks. *NeurIPS*, 2019.
- Zhong, K., Song, Z., Jain, P., Bartlett, P. L., and Dhillon, I. S. Recovery guarantees for one-hidden-layer neural networks. *ICML*, 2017.
- Zhou, Y. and Liang, Y. Critical points of linear neural networks: Analytical forms and landscape. *ICLR*, 2018.

Zou, D. and Gu, Q. An improved analysis of training overparameterized deep neural networks. *NeurIPS*, 2019.

Zou, D., Cao, Y., Zhou, D., and Gu, Q. Stochastic gradient descent optimizes overparameterized deep relu networks. *Machine Learning, Springer*, 2019.

Zou, D., Long, P. M., and Gu, Q. On the global convergence of training deep linear resnets. *ICLR*, 2020.

## A. Linear-rate results of the discrete-time Polyak's momentum

In the discrete-time setting, for general smooth, strongly convex, and differentiable functions, a linear rate of the global convergence is shown by Ghadimi et al. (2015) and Shi et al. (2018). However, the rate is not an accelerated rate and is not better than that of the vanilla gradient descent. To our knowledge, the class of the strongly convex quadratic problems is the only known example that Polyak's momentum has a provable *accelerated linear rate* in terms of the *global convergence* in the *discrete-time* setting.

## B. Proof of Lemma 2, Lemma 3, and Lemma 4

**Lemma 2:** *Applying Algorithm 1 or Algorithm 2 to minimizing a function  $f(w) \in F_{\mu,\alpha}^2$  induces a residual dynamics in the form of (8), where*

$$\begin{aligned}\xi_t &= w_t - w_* \\ H &= \int_0^1 \nabla^2 f((1-\tau)w_0 + \tau w_*) d\tau \\ \varphi_t &= \eta \left( \int_0^1 \nabla^2 f((1-\tau)w_0 + \tau w_*) d\tau - \int_0^1 \nabla^2 f((1-\tau)w_t + \tau w_*) d\tau \right) (w_t - w_*),\end{aligned}$$

where  $w_* := \arg \min_w f(w)$ .

*Proof.* We have

$$\begin{aligned}\begin{bmatrix} w_{t+1} - w_* \\ w_t - w_* \end{bmatrix} &= \begin{bmatrix} I_d + \beta I_d & -\beta I_d \\ I_d & 0_d \end{bmatrix} \cdot \begin{bmatrix} w_t - w_* \\ w_{t-1} - w_* \end{bmatrix} + \begin{bmatrix} -\eta \nabla f(w_t) \\ 0 \end{bmatrix} \\ &= \begin{bmatrix} I_d - \eta \int_0^1 \nabla^2 f((1-\tau)w_t + \tau w_*) d\tau + \beta I_d & -\beta I_d \\ I_d & 0_d \end{bmatrix} \cdot \begin{bmatrix} w_t - w_* \\ w_{t-1} - w_* \end{bmatrix} \\ &= \begin{bmatrix} I_d - \eta \int_0^1 \nabla^2 f((1-\tau)w_0 + \tau w_*) d\tau + \beta I_d & -\beta I_d \\ I_d & 0_d \end{bmatrix} \cdot \begin{bmatrix} w_t - w_* \\ w_{t-1} - w_* \end{bmatrix} \\ &\quad + \eta \left( \int_0^1 \nabla^2 f((1-\tau)w_0 + \tau w_*) d\tau - \int_0^1 \nabla^2 f((1-\tau)w_t + \tau w_*) d\tau \right) (w_t - w_*),\end{aligned}\tag{23}$$

where the second equality is by the fundamental theorem of calculus.

$$\nabla f(w_t) - \nabla f(w_*) = \left( \int_0^1 \nabla^2 f((1-\tau)w_t + \tau w_*) d\tau \right) (w_t - w_*),\tag{24}$$

and that  $\nabla f(w_*) = 0$ . □

**Lemma 3:** *(Residual dynamics of training the ReLU network  $\mathcal{N}_W^{\text{ReLU}}(\cdot)$ ) Denote*

$$(H_t)_{i,j} := H(W_t)_{i,j} = \frac{1}{m} \sum_{r=1}^m x_i^\top x_j \mathbb{1}\{\langle w_t^{(r)}, x_i \rangle \geq 0 \ \& \ \langle w_t^{(r)}, x_j \rangle \geq 0\}.$$

*Applying Algorithm 1 or Algorithm 2 to (5) for training the ReLU network  $\mathcal{N}_W^{\text{ReLU}}(x)$  induces a residual dynamics in the form of (8) such that*

$$\begin{aligned}\xi_t[i] &= \mathcal{N}_{W_t}^{\text{ReLU}}(x_i) - y_i \quad \text{and hence } n_0 = d \\ H &= H_0 \\ \varphi_t &= \phi_t + \iota_t,\end{aligned}$$

where each element  $i$  of  $\xi_t \in \mathbb{R}^n$  is the residual error of the sample  $i$ , the  $i_{th}$ -element of  $\phi_t \in \mathbb{R}^n$  satisfies

$$|\phi_t[i]| \leq \frac{2\eta\sqrt{n}|S_i^{\perp}|}{m} (\|u_t - y\| + \beta \sum_{s=0}^{t-1} \beta^{t-1-s} \|u_s - y\|),$$

and  $\iota_t = \eta(H_0 - H_t)\xi_t \in \mathbb{R}^n$ .

*Proof.* For each sample  $i$ , we will divide the contribution to  $\mathcal{N}(x_i)$  into two groups.

$$\begin{aligned}\mathcal{N}(x_i) &= \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \sigma(\langle w^{(r)}, x_i \rangle) \\ &= \frac{1}{\sqrt{m}} \sum_{r \in S_i} a_r \sigma(\langle w^{(r)}, x_i \rangle) + \frac{1}{\sqrt{m}} \sum_{r \in S_i^\perp} a_r \sigma(\langle w^{(r)}, x_i \rangle).\end{aligned}\quad (25)$$

To continue, let us recall some notations; the subgradient with respect to  $w^{(r)} \in \mathbb{R}^d$  is

$$\frac{\partial L(W)}{\partial w^{(r)}} := \frac{1}{\sqrt{m}} \sum_{i=1}^n (\mathcal{N}(x_i) - y_i) a_r x_i \mathbb{1}\{\langle w^{(r)}, x \rangle \geq 0\}, \quad (26)$$

and the Gram matrix  $H_t$  whose  $(i, j)$  element is

$$H_t[i, j] := \frac{1}{m} x_i^\top x_j \sum_{r=1}^m \mathbb{1}\{\langle w_t^{(r)}, x_i \rangle \geq 0 \ \& \ \langle w_t^{(r)}, x_j \rangle \geq 0\}. \quad (27)$$

Let us also denote

$$H_t^\perp[i, j] := \frac{1}{m} x_i^\top x_j \sum_{r \in S_i^\perp} \mathbb{1}\{\langle w_t^{(r)}, x_i \rangle \geq 0 \ \& \ \langle w_t^{(r)}, x_j \rangle \geq 0\}. \quad (28)$$

We have that

$$\begin{aligned}\xi_{t+1}[i] &= \mathcal{N}_{t+1}(x_i) - y_i \\ &\stackrel{(25)}{=} \underbrace{\frac{1}{\sqrt{m}} \sum_{r \in S_i} a_r \sigma(\langle w_{t+1}^{(r)}, x_i \rangle)}_{\text{first term}} + \frac{1}{\sqrt{m}} \sum_{r \in S_i^\perp} a_r \sigma(\langle w_{t+1}^{(r)}, x_i \rangle) - y_i.\end{aligned}\quad (29)$$

For the first term above, we have that

$$\begin{aligned}&\underbrace{\frac{1}{\sqrt{m}} \sum_{r \in S_i} a_r \sigma(\langle w_{t+1}^{(r)}, x_i \rangle)}_{\text{first term}} = \frac{1}{\sqrt{m}} \sum_{r \in S_i} a_r \sigma(\langle w_t^{(r)} - \eta \frac{\partial L(W_t)}{\partial w_t^{(r)}} + \beta(w_t^{(r)} - w_{t-1}^{(r)}), x_i \rangle) \\ &= \frac{1}{\sqrt{m}} \sum_{r \in S_i} a_r \langle w_t^{(r)} - \eta \frac{\partial L(W_t)}{\partial w_t^{(r)}} + \beta(w_t^{(r)} - w_{t-1}^{(r)}), x_i \rangle \cdot \mathbb{1}\{\langle w_{t+1}^{(r)}, x_i \rangle \geq 0\} \\ &\stackrel{(a)}{=} \frac{1}{\sqrt{m}} \sum_{r \in S_i} a_r \langle w_t^{(r)}, x_i \rangle \cdot \mathbb{1}\{\langle w_t^{(r)}, x_i \rangle \geq 0\} + \frac{\beta}{\sqrt{m}} \sum_{r \in S_i} a_r \langle w_t^{(r)}, x_i \rangle \cdot \mathbb{1}\{\langle w_t^{(r)}, x_i \rangle \geq 0\} \\ &\quad - \frac{\beta}{\sqrt{m}} \sum_{r \in S_i} a_r \langle w_{t-1}^{(r)}, x_i \rangle \cdot \mathbb{1}\{\langle w_{t-1}^{(r)}, x_i \rangle \geq 0\} - \eta \frac{1}{\sqrt{m}} \sum_{r \in S_i} a_r \langle \frac{\partial L(W_t)}{\partial w_t^{(r)}}, x_i \rangle \mathbb{1}\{\langle w_t^{(r)}, x_i \rangle \geq 0\} \\ &= \mathcal{N}_t(x_i) + \beta(\mathcal{N}_t(x_i) - \mathcal{N}_{t-1}(x_i)) - \frac{1}{\sqrt{m}} \sum_{r \in S_i^\perp} a_r \langle w_t^{(r)}, x_i \rangle \mathbb{1}\{\langle w_t^{(r)}, x_i \rangle \geq 0\} \\ &\quad - \frac{\beta}{\sqrt{m}} \sum_{r \in S_i^\perp} a_r \langle w_t^{(r)}, x_i \rangle \mathbb{1}\{\langle w_t^{(r)}, x_i \rangle \geq 0\} + \frac{\beta}{\sqrt{m}} \sum_{r \in S_i^\perp} a_r \langle w_{t-1}^{(r)}, x_i \rangle \mathbb{1}\{\langle w_{t-1}^{(r)}, x_i \rangle \geq 0\} \\ &\quad - \underbrace{\eta \frac{1}{\sqrt{m}} \sum_{r \in S_i} a_r \langle \frac{\partial L(W_t)}{\partial w_t^{(r)}}, x_i \rangle \mathbb{1}\{\langle w_t^{(r)}, x_i \rangle \geq 0\}}_{\text{last term}},\end{aligned}\quad (30)$$

where (a) uses that for  $r \in S_i$ ,  $\mathbb{1}\{\langle w_{t+1}^{(r)}, x_i \rangle \geq 0\} = \mathbb{1}\{\langle w_t^{(r)}, x_i \rangle \geq 0\} = \mathbb{1}\{\langle w_{t-1}^{(r)}, x_i \rangle \geq 0\}$  as the neurons in  $S_i$  do not

change their activation patterns. We can further bound (30) as

$$\begin{aligned}
 & \stackrel{(b)}{=} \mathcal{N}_t(x_i) + \beta(\mathcal{N}_t(x_i) - \mathcal{N}_{t-1}(x_i)) - \eta \sum_{j=1}^n (\mathcal{N}_t(x_j) - y_j) H(W_t)_{i,j} \\
 & - \frac{\eta}{m} \sum_{j=1}^n x_i^\top x_j (\mathcal{N}_t(x_j) - y_j) \sum_{r \in S_i^\perp} \mathbb{1}\{\langle w_t^{(r)}, x_i \rangle \geq 0 \ \& \ \langle w_t^{(r)}, x_j \rangle \geq 0\} \\
 & - \frac{1}{\sqrt{m}} \sum_{r \in S_i^\perp} a_r \langle w_t^{(r)}, x_i \rangle \mathbb{1}\{\langle w_t^{(r)}, x_i \rangle \geq 0\} - \frac{\beta}{\sqrt{m}} \sum_{r \in S_i^\perp} a_r \langle w_t^{(r)}, x_i \rangle \mathbb{1}\{\langle w_t^{(r)}, x_i \rangle \geq 0\} \\
 & + \frac{\beta}{\sqrt{m}} \sum_{r \in S_i^\perp} a_r \langle w_{t-1}^{(r)}, x_i \rangle \mathbb{1}\{\langle w_{t-1}^{(r)}, x_i \rangle \geq 0\}, \tag{31}
 \end{aligned}$$

where (b) is due to that

$$\begin{aligned}
 & \underbrace{\frac{1}{\sqrt{m}} \sum_{r \in S_i} a_r \langle \frac{\partial L(W_t)}{\partial w_t^{(r)}}, x_i \rangle \mathbb{1}\{\langle w_t^{(r)}, x_i \rangle \geq 0\}}_{\text{last term}} \\
 & = \frac{1}{m} \sum_{j=1}^n x_i^\top x_j (\mathcal{N}_t(x_j) - y_j) \sum_{r \in S_i} \mathbb{1}\{\langle w_t^{(r)}, x_i \rangle \geq 0 \ \& \ \langle w_t^{(r)}, x_j \rangle \geq 0\} \\
 & = \sum_{j=1}^n (\mathcal{N}_t(x_j) - y_j) H(W_t)_{i,j} - \frac{1}{m} \sum_{j=1}^n x_i^\top x_j (\mathcal{N}_t(x_j) - y_j) \sum_{r \in S_i^\perp} \mathbb{1}\{\langle w_t^{(r)}, x_i \rangle \geq 0 \ \& \ \langle w_t^{(r)}, x_j \rangle \geq 0\}. \tag{32}
 \end{aligned}$$

Combining (29) and (31), we have that

$$\begin{aligned}
 \xi_{t+1}[i] & = \xi_t[i] + \beta(\xi_t[i] - \xi_{t-1}[i]) - \eta \sum_{j=1}^n H_t[i, j] \xi_t[j] \\
 & - \frac{\eta}{m} \sum_{j=1}^n x_i^\top x_j (\mathcal{N}_t(x_j) - y_j) \sum_{r \in S_i^\perp} \mathbb{1}\{\langle w_t^{(r)}, x_i \rangle \geq 0 \ \& \ \langle w_t^{(r)}, x_j \rangle \geq 0\} \\
 & + \frac{1}{\sqrt{m}} \sum_{r \in S_i^\perp} a_r \sigma(\langle w_{t+1}^{(r)}, x_i \rangle) - a_r \sigma(\langle w_t^{(r)}, x_i \rangle) - \beta a_r \sigma(\langle w_t^{(r)}, x_i \rangle) + \beta a_r \sigma(\langle w_{t-1}^{(r)}, x_i \rangle). \tag{33}
 \end{aligned}$$

So we can write the above into a matrix form.

$$\begin{aligned}
 \xi_{t+1} & = (I_n - \eta H_t) \xi_t + \beta(\xi_t - \xi_{t-1}) + \phi_t \\
 & = (I_n - \eta H_0) \xi_t + \beta(\xi_t - \xi_{t-1}) + \phi_t + \iota_t, \tag{34}
 \end{aligned}$$

where the  $i$  element of  $\phi_t \in \mathbb{R}^n$  is defined as

$$\begin{aligned}
 \phi_t[i] & = -\frac{\eta}{m} \sum_{j=1}^n x_i^\top x_j (\mathcal{N}_t(x_j) - y_j) \sum_{r \in S_i^\perp} \mathbb{1}\{\langle w_t^{(r)}, x_i \rangle \geq 0 \ \& \ \langle w_t^{(r)}, x_j \rangle \geq 0\} \\
 & + \frac{1}{\sqrt{m}} \sum_{r \in S_i^\perp} \{a_r \sigma(\langle w_{t+1}^{(r)}, x_i \rangle) - a_r \sigma(\langle w_t^{(r)}, x_i \rangle) - \beta a_r \sigma(\langle w_t^{(r)}, x_i \rangle) + \beta a_r \sigma(\langle w_{t-1}^{(r)}, x_i \rangle)\}. \tag{35}
 \end{aligned}$$



Now let us bound  $\phi_t[i]$  as follows.

$$\begin{aligned}
 \phi_t[i] &= -\frac{\eta}{m} \sum_{j=1}^n x_i^\top x_j (\mathcal{N}_t(x_j) - y_j) \sum_{r \in S_i^\perp} \mathbb{1}\{\langle w_t^{(r)}, x_i \rangle \geq 0 \ \& \ \langle w_t^{(r)}, x_j \rangle \geq 0\} \\
 &\quad + \frac{1}{\sqrt{m}} \sum_{r \in S_i^\perp} \{a_r \sigma(\langle w_{t+1}^{(r)}, x_i \rangle) - a_r \sigma(\langle w_t^{(r)}, x_i \rangle) - \beta a_r \sigma(\langle w_t^{(r)}, x_i \rangle) + \beta a_r \sigma(\langle w_{t-1}^{(r)}, x_i \rangle)\} \\
 &\stackrel{(a)}{\leq} \frac{\eta \sqrt{n} |S_i^\perp|}{m} \|u_t - y\| + \frac{1}{\sqrt{m}} \sum_{r \in S_i^\perp} (\|w_{t+1}^{(r)} - w_t^{(r)}\| + \beta \|w_t^{(r)} - w_{t-1}^{(r)}\|) \\
 &\stackrel{(b)}{=} \frac{\eta \sqrt{n} |S_i^\perp|}{m} \|u_t - y\| + \frac{\eta}{\sqrt{m}} \sum_{r \in S_i^\perp} \left( \left\| \sum_{s=0}^t \beta^{t-s} \frac{\partial L(W_s)}{\partial w_s^{(r)}} \right\| + \beta \left\| \sum_{s=0}^{t-1} \beta^{t-1-s} \frac{\partial L(W_s)}{\partial w_s^{(r)}} \right\| \right) \\
 &\stackrel{(c)}{\leq} \frac{\eta \sqrt{n} |S_i^\perp|}{m} \|u_t - y\| + \frac{\eta}{\sqrt{m}} \sum_{r \in S_i^\perp} \left( \sum_{s=0}^t \beta^{t-s} \left\| \frac{\partial L(W_s)}{\partial w_s^{(r)}} \right\| + \beta \sum_{s=0}^{t-1} \beta^{t-1-s} \left\| \frac{\partial L(W_s)}{\partial w_s^{(r)}} \right\| \right) \\
 &\stackrel{(d)}{\leq} \frac{\eta \sqrt{n} |S_i^\perp|}{m} \|u_t - y\| + \frac{\eta \sqrt{n} |S_i^\perp|}{m} \left( \sum_{s=0}^t \beta^{t-s} \|u_s - y\| + \beta \sum_{s=0}^{t-1} \beta^{t-1-s} \|u_s - y\| \right) \\
 &= \frac{2\eta \sqrt{n} |S_i^\perp|}{m} \left( \|u_t - y\| + \beta \sum_{s=0}^{t-1} \beta^{t-1-s} \|u_s - y\| \right), \tag{36}
 \end{aligned}$$

where (a) is because  $-\frac{\eta}{m} \sum_{j=1}^n x_i^\top x_j (\mathcal{N}_t(x_j) - y_j) \sum_{r \in S_i^\perp} \mathbb{1}\{\langle w_t^{(r)}, x_i \rangle \geq 0 \ \& \ \langle w_t^{(r)}, x_j \rangle \geq 0\} \leq \frac{\eta |S_i^\perp|}{m} \sum_{j=1}^n |\mathcal{N}_t(x_j) - y_j| \leq \frac{\eta \sqrt{n} |S_i^\perp|}{m} \|u_t - y\|$ , and that  $\sigma(\cdot)$  is 1-Lipschitz so that

$$\begin{aligned}
 &\frac{1}{\sqrt{m}} \sum_{r \in S_i^\perp} (a_r \sigma(\langle w_{t+1}^{(r)}, x_i \rangle) - a_r \sigma(\langle w_t^{(r)}, x_i \rangle)) \leq \frac{1}{\sqrt{m}} \sum_{r \in S_i^\perp} |\langle w_{t+1}^{(r)}, x_i \rangle - \langle w_t^{(r)}, x_i \rangle| \\
 &\leq \frac{1}{\sqrt{m}} \sum_{r \in S_i^\perp} \|w_{t+1}^{(r)} - w_t^{(r)}\| \|x_i\| \leq \frac{1}{\sqrt{m}} \sum_{r \in S_i^\perp} \|w_{t+1}^{(r)} - w_t^{(r)}\|,
 \end{aligned}$$

similarly,  $\frac{-\beta}{\sqrt{m}} \sum_{r \in S_i^\perp} (a_r \sigma(\langle w_t^{(r)}, x_i \rangle) - a_r \sigma(\langle w_{t-1}^{(r)}, x_i \rangle)) \leq \beta \frac{1}{\sqrt{m}} \sum_{r \in S_i^\perp} \|w_t^{(r)} - w_{t-1}^{(r)}\|$ , (b) is by the update rule (Algorithm 1), (c) is by Jensen's inequality, (d) is because  $|\frac{\partial L(W_s)}{\partial w_s^{(r)}}| = |\frac{1}{\sqrt{m}} \sum_{i=1}^n (u_s[i] - y_i) a_r x_i \mathbb{1}\{x^\top w_t^{(r)} \geq 0\}| \leq \frac{\sqrt{n}}{m} \|u_s - y\|$ .

□

**Lemma: 4** (Residual dynamics of training  $\mathcal{N}_W^{L\text{-linear}}(\cdot)$ ) Denote  $M_{t,l}$  the momentum term of layer  $l$  at iteration  $t$ , which is recursively defined as  $M_{t,l} = \beta M_{t,l-1} + \frac{\partial \ell(W_t^{(L:1)})}{\partial W_t^{(l)}}$ . Denote

$$H_t := \frac{1}{m^{L-1} d_y} \sum_{l=1}^L [(W_t^{(l-1:1)} X)^\top (W_t^{(l-1:1)} X) \otimes W_t^{(L:l+1)} (W_t^{(L:l+1)})^\top] \in \mathbb{R}^{d_y n \times d_y n}.$$

Applying Algorithm 1 or Algorithm 2 to (5) for training the deep linear network  $\mathcal{N}_W^{L\text{-linear}}(x)$  induces a residual dynamics in the form of (8) such that

$$\begin{aligned}
 \xi_t &= \text{vec}(U_t - Y) \in \mathbb{R}^{d_y n}, \text{ and hence } n_0 = d_y n \\
 H &= H_0 \\
 \varphi_t &= \phi_t + \psi_t + \iota_t \in \mathbb{R}^{d_y n},
 \end{aligned}$$

where

$$\begin{aligned}\phi_t &= \frac{1}{\sqrt{m^{L-1}d_y}} \text{vec}(\Phi_t X) \text{ with } \Phi_t = \Pi_l(W_t^{(l)} - \eta M_{t,l}) - W_t^{(L:1)} + \eta \sum_{l=1}^L W_t^{(L:l+1)} M_{t,l} W_t^{(l-1:1)} \\ \psi_t &= \frac{1}{\sqrt{m^{L-1}d_y}} \text{vec} \left( (L-1)\beta W_t^{(L:1)} X + \beta W_{t-1}^{(L:1)} X - \beta \sum_{l=1}^L W_t^{(L:l+1)} W_{t-1}^{(l)} W_t^{(l-1:1)} X \right) \\ \iota_t &= \eta(H_0 - H_t)\xi_t.\end{aligned}$$

*Proof.* According to the update rule of gradient descent with Polyak's momentum, we have

$$W_{t+1}^{(L:1)} = \Pi_l \left( W_t^{(l)} - \eta M_{t,l} \right) = W_t^{(L:1)} - \eta \sum_{l=1}^L W_t^{(L:l+1)} M_{t,l} W_t^{(l-1:1)} + \Phi_t, \quad (37)$$

where  $M_{t,l}$  stands for the momentum term of layer  $l$ , which is  $M_{t,l} = \beta M_{t,l-1} + \frac{\partial \ell(W_t^{(L:1)})}{\partial W_t^{(l)}} = \sum_{s=0}^t \beta^{t-s} \frac{\partial \ell(W_s^{(L:1)})}{\partial W_s^{(l)}}$ , and  $\Phi_t$  contains all the high-order terms (in terms of  $\eta$ ), e.g. those with  $\eta M_{t,i}$  and  $\eta M_{t,j}$ ,  $i \neq j \in [L]$ , or higher. Based on the equivalent update expression of gradient descent with Polyak's momentum  $-\eta M_{t,l} = -\eta \frac{\partial \ell(W_t^{(L:1)})}{\partial W_t^{(l)}} + \beta(W_t^{(l)} - W_{t-1}^{(l)})$ , we can rewrite (37) as

$$\begin{aligned}W_{t+1}^{(L:1)} &= W_t^{(L:1)} - \eta \sum_{l=1}^L W_t^{(L:l+1)} \frac{\partial \ell(W_t^{(L:1)})}{\partial W_t^{(l)}} W_t^{(l-1:1)} + \sum_{l=1}^L W_t^{(L:l+1)} \beta(W_t^{(l)} - W_{t-1}^{(l)}) W_t^{(l-1:1)} + \Phi_t \\ &= W_t^{(L:1)} - \eta \sum_{l=1}^L W_t^{(L:l+1)} \frac{\partial \ell(W_t^{(L:1)})}{\partial W_t^{(l)}} W_t^{(l-1:1)} + \beta(W_t^{(L:1)} - W_{t-1}^{(L:1)}) + \Phi_t \\ &\quad + (L-1)\beta W_t^{(L:1)} + \beta W_{t-1}^{(L:1)} - \beta \sum_{l=1}^L W_t^{(L:l+1)} W_{t-1}^{(l)} W_t^{(l-1:1)}.\end{aligned} \quad (38)$$

Multiplying the above equality with  $\frac{1}{\sqrt{m^{L-1}d_y}} X$ , we get

$$\begin{aligned}U_{t+1} &= U_t - \eta \frac{1}{m^{L-1}d_y} \sum_{l=1}^L W_t^{(L:l+1)} (W_t^{(L:l+1)})^\top (U_t - Y) (W_t^{(l-1:1)} X)^\top W_t^{(l-1:1)} X + \beta(U_t - U_{t-1}) \\ &\quad + \frac{1}{\sqrt{m^{L-1}d_y}} \left( (L-1)\beta W_t^{(L:1)} + \beta W_{t-1}^{(L:1)} - \beta \sum_{l=1}^L W_t^{(L:l+1)} W_{t-1}^{(l)} W_t^{(l-1:1)} \right) X + \frac{1}{\sqrt{m^{L-1}d_y}} \Phi_t X.\end{aligned} \quad (39)$$

Using  $\text{vec}(ACB) = (B^\top \otimes A)\text{vec}(C)$ , where  $\otimes$  stands for the Kronecker product, we can apply a vectorization of the above equation and obtain

$$\begin{aligned}\text{vec}(U_{t+1}) - \text{vec}(U_t) &= -\eta H_t \text{vec}(U_t - Y) + \beta(\text{vec}(U_t) - \text{vec}(U_{t-1})) \\ &\quad + \text{vec} \left( \frac{1}{\sqrt{m^{L-1}d_y}} \left( (L-1)\beta W_t^{(L:1)} + \beta W_{t-1}^{(L:1)} - \beta \sum_{l=1}^L W_t^{(L:l+1)} W_{t-1}^{(l)} W_t^{(l-1:1)} \right) X \right) \\ &\quad + \frac{1}{\sqrt{m^{L-1}d_y}} \text{vec}(\Phi_t X),\end{aligned} \quad (40)$$

where

$$H_t = \frac{1}{m^{L-1}d_y} \sum_{l=1}^L \left[ \left( (W_t^{(l-1:1)} X)^\top (W_t^{(l-1:1)} X) \right) \otimes W_t^{(L:l+1)} (W_t^{(L:l+1)})^\top \right], \quad (41)$$

which is a positive semi-definite matrix.

In the following, we will denote  $\xi_t := \text{vec}(U_t - Y)$  as the vector of the residual errors. Also, we denote  $\phi_t := \frac{1}{\sqrt{m^{L-1}d_y}} \text{vec}(\Phi_t X)$  with  $\Phi_t = \Pi_l(W_t^{(l)} - \eta M_{t,l}) - W_t^{(L:1)} + \eta \sum_{l=1}^L W_t^{(L:l+1)} M_{t,l} W_t^{(l-1:1)}$ , and  $\psi_t := \text{vec}\left(\frac{1}{\sqrt{m^{L-1}d_y}} \left((L-1)\beta W_t^{(L:1)} + \beta W_{t-1}^{(L:1)} - \beta \sum_{l=1}^L W_t^{(L:l+1)} W_{t-1}^{(l)} W_t^{(l-1:1)}\right) X\right)$ . Using the notations, we can rewrite (40) as

$$\begin{aligned} \begin{bmatrix} \xi_{t+1} \\ \xi_t \end{bmatrix} &= \begin{bmatrix} I_{d_y n} - \eta H_t + \beta I_{d_y n} & -\beta I_{d_y n} \\ & I_{d_y n} & 0_{d_y n} \end{bmatrix} \begin{bmatrix} \xi_t \\ \xi_{t-1} \end{bmatrix} + \begin{bmatrix} \phi_t + \psi_t \\ 0_{d_y n} \end{bmatrix} \\ &= \begin{bmatrix} I_{d_y n} - \eta H_0 + \beta I_{d_y n} & -\beta I_{d_y n} \\ & I_{d_y n} & 0_{d_y n} \end{bmatrix} \begin{bmatrix} \xi_t \\ \xi_{t-1} \end{bmatrix} + \begin{bmatrix} \varphi_t \\ 0_{d_y n} \end{bmatrix}, \end{aligned} \quad (42)$$

where  $\varphi_t = \phi_t + \psi_t + \iota_t \in \mathbb{R}^{d_y n}$  and  $I_{d_y n}$  is the  $d_y n \times d_y n$ -dimensional identity matrix. □

### C. Proof of Theorem 5

**Theorem 5** Let  $A := \begin{bmatrix} (1+\beta)I_n - \eta H & -\beta I_n \\ I_n & 0 \end{bmatrix} \in \mathbb{R}^{2n \times 2n}$ . Suppose that  $H \in \mathbb{R}^{n \times n}$  is a positive semidefinite matrix.

Fix a vector  $v_0 \in \mathbb{R}^n$ . If  $\beta$  is chosen to satisfy  $1 \geq \beta > \max\left\{\left(1 - \sqrt{\eta \lambda_{\min}(H)}\right)^2, \left(1 - \sqrt{\eta \lambda_{\max}(H)}\right)^2\right\}$ , then

$$\|A^k v_0\| \leq (\sqrt{\beta})^k C_0 \|v_0\|, \quad (43)$$

where the constant

$$C_0 := \frac{\sqrt{2}(\beta+1)}{\sqrt{\min\{h(\beta, \eta \lambda_{\min}(H)), h(\beta, \eta \lambda_{\max}(H))\}}} \geq 1, \quad (44)$$

and the function  $h(\beta, z)$  is defined as

$$h(\beta, z) := -\left(\beta - (1 - \sqrt{z})^2\right) \left(\beta - (1 + \sqrt{z})^2\right). \quad (45)$$

We would first prove some lemmas for the analysis.

**Lemma 6.** Under the assumption of Theorem 5,  $A$  is diagonalizable with respect to complex field  $\mathbb{C}$  in  $\mathbb{C}^n$ , i.e.,  $\exists P$  such that  $A = PDP^{-1}$  for some diagonal matrix  $D$ . Furthermore, the diagonal elements of  $D$  all have magnitudes bounded by  $\sqrt{\beta}$ .

*Proof.* In the following, we will use the notation/operation  $\text{Diag}(\dots)$  to represents a block-diagonal matrix that has the arguments on its main diagonal. Let  $U \text{Diag}([\lambda_1, \dots, \lambda_n]) U^*$  be the singular-value-decomposition of  $H$ , then

$$A = \begin{bmatrix} U & 0 \\ 0 & U \end{bmatrix} \begin{bmatrix} (1+\beta)I_n - \eta \text{Diag}([\lambda_1, \dots, \lambda_n]) & -\beta I_n \\ I_n & 0 \end{bmatrix} \begin{bmatrix} U^* & 0 \\ 0 & U^* \end{bmatrix}. \quad (46)$$

Let  $\tilde{U} = \begin{bmatrix} U & 0 \\ 0 & U \end{bmatrix}$ . Then, after applying some permutation matrix  $\tilde{P}$ ,  $A$  can be further simplified into

$$A = \tilde{U} \tilde{P} \tilde{\Sigma} \tilde{P}^T \tilde{U}^*, \quad (47)$$

where  $\tilde{\Sigma}$  is a block diagonal matrix consisting of  $n$  2-by-2 matrices  $\tilde{\Sigma}_i := \begin{bmatrix} 1 + \beta - \eta \lambda_i & -\beta \\ 1 & 0 \end{bmatrix}$ . The characteristic polynomial of  $\tilde{\Sigma}_i$  is  $x^2 - (1 + \beta - \eta \lambda_i)x + \beta$ . Hence it can be shown that when  $\beta > (1 - \sqrt{\eta \lambda_i})^2$  then the roots of polynomial are conjugate and have magnitude  $\sqrt{\beta}$ . These roots are exactly the eigenvalues of  $\tilde{\Sigma}_i \in \mathbb{R}^{2 \times 2}$ . On the other hand, the

corresponding eigenvectors  $q_i, \bar{q}_i$  are also conjugate to each other as  $\tilde{\Sigma}_i \in \mathbb{R}^{2 \times 2}$  is a real matrix. As a result,  $\Sigma \in \mathbb{R}^{2n \times 2n}$  admits a block eigen-decomposition as follows,

$$\begin{aligned} \Sigma &= \text{Diag}(\tilde{\Sigma}_1, \dots, \tilde{\Sigma}_n) \\ &= \text{Diag}(Q_1, \dots, Q_n) \text{Diag} \left( \begin{bmatrix} z_1 & 0 \\ 0 & \bar{z}_1 \end{bmatrix}, \dots, \begin{bmatrix} z_n & 0 \\ 0 & \bar{z}_n \end{bmatrix} \right) \text{Diag}(Q_1^{-1}, \dots, Q_n^{-1}), \end{aligned} \quad (48)$$

where  $Q_i = [q_i, \bar{q}_i]$  and  $z_i, \bar{z}_i$  are eigenvalues of  $\tilde{\Sigma}_i$  (they are conjugate by the condition on  $\beta$ ). Denote  $Q := \text{Diag}(Q_1, \dots, Q_n)$  and

$$D := \text{Diag} \left( \begin{bmatrix} z_1 & 0 \\ 0 & \bar{z}_1 \end{bmatrix}, \dots, \begin{bmatrix} z_n & 0 \\ 0 & \bar{z}_n \end{bmatrix} \right). \quad (49)$$

By combining (47) and (48), we have

$$A = P \text{Diag} \left( \begin{bmatrix} z_1 & 0 \\ 0 & \bar{z}_1 \end{bmatrix}, \dots, \begin{bmatrix} z_n & 0 \\ 0 & \bar{z}_n \end{bmatrix} \right) P^{-1} = P D P^{-1}, \quad (50)$$

where

$$P = \tilde{U} \tilde{P} Q, \quad (51)$$

by the fact that  $\tilde{P}^{-1} = \tilde{P}^T$  and  $\tilde{U}^{-1} = \tilde{U}^*$ .  $\square$

*Proof.* (of Theorem 5) Now we proceed the proof of Theorem 5. In the following, we denote  $v_k := A^k v_0$  (so  $v_k = A v_{k-1}$ ). Let  $P$  be the matrix in Lemma 6, and  $u_k := P^{-1} v_k$ , the dynamic can be rewritten as  $u_k = P^{-1} A v_{k-1} = P^{-1} A P u_{k-1} = D u_{k-1}$ . As  $D$  is diagonal, we immediately have

$$\begin{aligned} \|u_k\| &\leq \max_{i \in [n]} |D_{ii}|^k \|u_0\| \\ \Rightarrow \|P^{-1} v_k\| &\leq \max_{i \in [n]} |D_{ii}|^k \|P^{-1} v_0\| \\ \Rightarrow \sigma_{\min}(P^{-1}) \|v_k\| &\leq \sqrt{\beta}^k \sigma_{\max}(P^{-1}) \|v_0\| \quad (\text{Lemma 6.}) \\ \Rightarrow \sigma_{\max}^{-1}(P) \|v_k\| &\leq \sqrt{\beta}^k \sigma_{\min}^{-1}(P) \|v_0\| \\ \Rightarrow \|v_k\| &\leq \sqrt{\beta}^k \frac{\sigma_{\max}(P)}{\sigma_{\min}(P)} \|v_0\| \\ \Rightarrow \|v_k\| &\leq \sqrt{\beta}^k \sqrt{\frac{\lambda_{\max}(PP^*)}{\lambda_{\min}(PP^*)}} \|v_0\|. \end{aligned} \quad (52)$$

Hence, now it suffices to prove upper bound and lower bound of  $\lambda_{\max}$  and  $\lambda_{\min}$ , respectively. By using Lemma 7 in the following, we obtain the inequality of (43). We remark that as  $C_0$  is an upper-bound of the squared root of the condition number  $\sqrt{\frac{\lambda_{\max}(PP^*)}{\lambda_{\min}(PP^*)}}$ , it is lower bounded by 1.  $\square$

**Lemma 7.** *Let  $P$  be the matrix in Lemma 6, then we have  $\lambda_{\max}(PP^*) \leq 2(\beta + 1)$  and  $\lambda_{\min}(PP^*) \geq \min\{h(\beta, \eta \lambda_{\min}(H)), h(\beta, \eta \lambda_{\max}(H))\} / (1 + \beta)$ , where*

$$h(\beta, z) = - \left( \beta - (1 - \sqrt{z})^2 \right) \left( \beta - (1 + \sqrt{z})^2 \right). \quad (53)$$

*Proof.* As (51) in the proof of Lemma 2,  $P = \tilde{U} \tilde{P} \text{Diag}(Q_1, \dots, Q_n)$ . Since  $\tilde{U} \tilde{P}$  is unitary, it does not affect the spectrum of  $P$ , therefore, it suffices to analyze the eigenvalues of  $Q Q^*$ , where  $Q = \text{Diag}(Q_1, \dots, Q_n)$ . Observe that  $Q Q^*$  is a block diagonal matrix with blocks  $Q_i Q_i^*$ , the eigenvalues of it are exactly that of  $Q_i Q_i^*$ , i.e.,  $\lambda_{\max}(Q Q^*) = \max_{i \in [n]} \lambda_{\max}(Q_i Q_i^*)$

and likewise for the minimum. Recall  $Q_i = [q_i, \bar{q}_i]$  consisting of eigenvectors of  $\tilde{\Sigma}_i := \begin{bmatrix} 1 + \beta - \eta\lambda_i & -\beta \\ 1 & 0 \end{bmatrix}$  with corresponding eigenvalues  $z_i, \bar{z}_i$ . The eigenvalues satisfy

$$z_i + \bar{z}_i = 2\Re z_i = 1 + \beta - \eta\lambda_i, \quad (54)$$

$$z_i \bar{z}_i = |z_i|^2 = \beta. \quad (55)$$

On the other hand, the eigenvalue equation  $\tilde{\Sigma}_i q_i = z_i q_i$  together with (54) implies  $q_i = [z_i, 1]^T$ . Furthermore,  $Q_i Q_i^* = q_i q_i^* + \bar{q}_i \bar{q}_i^* = 2\Re q_i q_i^* = 2\Re q_i \Re q_i^T + 2\Im q_i \Im q_i^T$ . Thus,

$$\begin{aligned} Q_i Q_i^* &= 2\Re q_i \Re q_i^T + 2\Im q_i \Im q_i^T \\ &= 2 \left( \begin{bmatrix} \Re z_i \\ 1 \end{bmatrix} [\Re z_i \ 1] + \begin{bmatrix} \Im z_i \\ 0 \end{bmatrix} [\Im z_i \ 0] \right) \\ &= 2 \begin{bmatrix} |z_i|^2 & \Re z_i \\ \Re z_i & 1 \end{bmatrix}. \end{aligned} \quad (56)$$

Let the eigenvalues of  $Q_i Q_i^*$  be  $\theta_1, \theta_2$ , then by (54)-(56) we must have

$$\theta_1 + \theta_2 = 2(\beta + 1), \quad (57)$$

$$\begin{aligned} \theta_1 \theta_2 &= 4 \left( \beta - \left( \frac{1 + \beta - \eta\lambda_i}{2} \right)^2 \right) \\ &= - \left( \beta - \left( 1 - \sqrt{\eta\lambda_i} \right)^2 \right) \left( \beta - \left( 1 + \sqrt{\eta\lambda_i} \right)^2 \right) \geq 0. \end{aligned} \quad (58)$$

From (57), as both eigenvalues are nonnegative, we deduce that

$$2(1 + \beta) \geq \max\{\theta_1, \theta_2\} \geq \beta + 1. \quad (59)$$

On the other hand, from (57) we also have

$$\begin{aligned} \min\{\theta_1, \theta_2\} &= \theta_1 \theta_2 / \max\{\theta_1, \theta_2\} \\ &\geq - \left( \beta - \left( 1 - \sqrt{\eta\lambda_i} \right)^2 \right) \left( \beta - \left( 1 + \sqrt{\eta\lambda_i} \right)^2 \right) / (1 + \beta) \\ &:= h(\beta, \eta\lambda_i) / (1 + \beta). \end{aligned} \quad (60)$$

Finally, as the eigenvalues of  $Q Q^*$  are composed of exactly that of  $Q_i Q_i^*$ , applying the bound of (60) to each  $i$  we have

$$\begin{aligned} \lambda_{\min}(P P^*) &\geq \min_{i \in [n]} h(\beta, \eta\lambda_i) / (1 + \beta) \\ &\geq \min\{h(\beta, \eta\lambda_{\min}(H)), h(\beta, \eta\lambda_{\max}(H))\} / (1 + \beta), \end{aligned} \quad (61)$$

where the last inequality follows from the facts that  $\lambda_{\min}(H) \leq \lambda_i \leq \lambda_{\max}(H)$  and  $h$  is concave quadratic function of  $\lambda$  in which the minimum must occur at the boundary.  $\square$

### C.1. Proof of Corollary 1

**Corollary 1** Assume that  $\lambda_{\min}(H) > 0$ . Denote  $\kappa := \lambda_{\max}(H) / \lambda_{\min}(H)$ . Set  $\eta = 1 / \lambda_{\max}(H)$  and set  $\beta = \left( 1 - \frac{1}{2} \sqrt{\eta \lambda_{\min}(H)} \right)^2 = \left( 1 - \frac{1}{2\sqrt{\kappa}} \right)^2$ . Then,  $C_0 \leq \max\{4, 2\sqrt{\kappa}\} \leq 4\sqrt{\kappa}$ .

*Proof.* For notation brevity, in the following, we let  $\mu := \lambda_{\min}(H)$  and  $\alpha := \lambda_{\max}(H)$ . Recall that  $h(\beta, z) = - \left( \beta - \left( 1 - \sqrt{z} \right)^2 \right) \left( \beta - \left( 1 + \sqrt{z} \right)^2 \right)$ . We have

$$\begin{aligned} h(\beta, \eta\mu) &= - \left( \left( 1 - \frac{1}{2} \sqrt{\eta\mu} \right)^2 - \left( 1 - \sqrt{\eta\mu} \right)^2 \right) \left( \left( 1 - \frac{1}{2} \sqrt{\eta\mu} \right)^2 - \left( 1 + \sqrt{\eta\mu} \right)^2 \right) \\ &= 3 \left( \sqrt{\eta\mu} - \frac{3}{4} \eta\mu \right) \left( \sqrt{\eta\mu} + \frac{1}{4} \eta\mu \right) = 3 \left( \frac{1}{\sqrt{\kappa}} - \frac{3}{4\kappa} \right) \left( \frac{1}{\sqrt{\kappa}} + \frac{1}{4\kappa} \right) \end{aligned} \quad (62)$$

and

$$\begin{aligned}
 h(\beta, \eta\alpha) &= - \left( \left(1 - \frac{1}{2}\sqrt{\eta\mu}\right)^2 - (1 - \sqrt{\eta\alpha})^2 \right) \left( \left(1 - \frac{1}{2}\sqrt{\eta\mu}\right)^2 - (1 + \sqrt{\eta\alpha})^2 \right) \\
 &= \left( 2\sqrt{\eta\alpha} - \sqrt{\eta\mu} - \eta\alpha + \frac{1}{4}\eta\mu \right) \left( \sqrt{\eta\mu} + 2\sqrt{\eta\alpha} + \eta\alpha - \frac{1}{4}\eta\mu \right) \\
 &= \left( 1 - \frac{1}{\sqrt{\kappa}} + \frac{1}{4\kappa} \right) \left( 3 + \frac{1}{\sqrt{\kappa}} - \frac{1}{4\kappa} \right). \tag{63}
 \end{aligned}$$

We can simplify it to get that  $h(\beta, \eta\alpha) = 3 - \frac{2}{\sqrt{\kappa}} - \frac{1}{2\kappa} + \frac{1}{2\kappa^{3/2}} - \frac{1}{16\kappa^2} \geq 0.5$ .

Therefore, we have

$$\frac{\sqrt{2}(\beta+1)}{\sqrt{h(\beta, \eta\mu)}} = \frac{\sqrt{2}(\beta+1)}{\sqrt{3\eta\mu(1 - \frac{1}{2}\sqrt{\eta\mu} - \frac{3}{16}\eta\mu)}} = \frac{\sqrt{2}(\beta+1)}{\sqrt{3(1 - \frac{1}{2}\sqrt{\eta\mu} - \frac{3}{16}\eta\mu)}} \sqrt{\kappa} \leq \frac{1}{\sqrt{(1 - \frac{1}{2} - \frac{3}{16})}} \sqrt{\kappa} \leq 2\sqrt{\kappa}, \tag{64}$$

where we use  $\eta\mu = \frac{1}{\kappa}$ . On the other hand,  $\frac{\sqrt{2}(\beta+1)}{\sqrt{h(\beta, \eta\alpha)}} \leq 4$ . We conclude that

$$C_0 = \frac{\sqrt{2}(\beta+1)}{\sqrt{\min\{h(\beta, \eta\nu), h(\beta, \eta\alpha)\}}} \leq \max\{4, 2\sqrt{\kappa}\} \leq 4\sqrt{\kappa}. \tag{65}$$

□

## D. Proof of Theorem 7

**Theorem 7** Assume the momentum parameter  $\beta$  satisfies  $1 \geq \beta > \max\{(1 - \sqrt{\eta\mu})^2, (1 - \sqrt{\eta\alpha})^2\}$ . Gradient descent with Polyak's momentum has

$$\left\| \begin{bmatrix} w_t - w_* \\ w_{t-1} - w_* \end{bmatrix} \right\| \leq \left( \sqrt{\beta} \right)^t C_0 \left\| \begin{bmatrix} w_0 - w_* \\ w_{-1} - w_* \end{bmatrix} \right\|, \tag{66}$$

where the constant

$$C_0 := \frac{\sqrt{2}(\beta+1)}{\sqrt{\min\{h(\beta, \eta\lambda_{\min}(\Gamma)), h(\beta, \eta\lambda_{\max}(\Gamma))\}}}, \tag{67}$$

and  $h(\beta, z) = - \left( \beta - (1 - \sqrt{z})^2 \right) \left( \beta - (1 + \sqrt{z})^2 \right)$ . Consequently, if the step size  $\eta = \frac{1}{\alpha}$  and the momentum parameter  $\beta = (1 - \sqrt{\eta\mu})^2$ , then it has

$$\left\| \begin{bmatrix} w_t - w_* \\ w_{t-1} - w_* \end{bmatrix} \right\| \leq \left( 1 - \frac{1}{2\sqrt{\kappa}} \right)^t 4\sqrt{\kappa} \left\| \begin{bmatrix} w_0 - w_* \\ w_{-1} - w_* \end{bmatrix} \right\|. \tag{68}$$

Furthermore, if  $\eta = \frac{4}{(\sqrt{\mu} + \sqrt{\alpha})^2}$  and  $\beta$  approaches  $\beta \rightarrow \left(1 - \frac{2}{\sqrt{\kappa+1}}\right)^2$  from above, then it has a convergence rate approximately  $\left(1 - \frac{2}{\sqrt{\kappa+1}}\right)$  as  $t \rightarrow \infty$ .

*Proof.* The result (66) and (68) is due to a trivial combination of Lemma 1, Theorem 6, and Corollary 1.

On the other hand, set  $\eta = \frac{4}{(\sqrt{\mu} + \sqrt{\alpha})^2}$ , the lower bound on  $\beta$  becomes  $\max\{(1 - \sqrt{\eta\mu})^2, (1 - \sqrt{\eta\alpha})^2\} = \left(1 - \frac{2}{\sqrt{\kappa+1}}\right)^2$ .

Since the rate is  $r = \lim_{t \rightarrow \infty} \frac{1}{t} \log(\sqrt{\beta}^{t+1} C_0) = \sqrt{\beta}$ , setting  $\beta \downarrow \left(1 - \frac{2}{\sqrt{\kappa+1}}\right)^2$  from above leads to the rate of  $\left(1 - \frac{2}{\sqrt{\kappa+1}}\right)$ . Formally, it is straightforward to show that  $C_0 = \Theta\left(\frac{1}{\sqrt{\beta - (1 - \frac{2}{1 + \sqrt{\kappa}})^2}}\right)$ , hence, for any  $\beta$  converges to  $\left(1 - \frac{2}{\sqrt{\kappa+1}}\right)^2$  slower than inverse exponential of  $\kappa$ , i.e.,  $\beta = \left(1 - \frac{2}{\sqrt{\kappa+1}}\right)^2 + \left(\frac{1}{\kappa}\right)^{o(t)}$ , we have  $r = 1 - \frac{2}{\sqrt{\kappa+1}}$ .

□

## E. Proof of Theorem 8

*Proof.* (of Theorem 8) In the following, we denote  $\xi_t := w_t - w_*$  and denote  $\lambda := \mu > 0$ , which is a lower bound of  $\lambda_{\min}(H)$  of the matrix  $H := \int_0^1 \nabla^2 f((1-\tau)w_0 + \tau w_*) d\tau$  defined in Lemma 2, i.e.  $\lambda_{\min}(H) \geq \lambda$ . Also, denote  $\beta_* := 1 - \frac{1}{2}\sqrt{\eta\lambda}$  and  $\theta := \beta_* + \frac{1}{4}\sqrt{\eta\lambda} = 1 - \frac{1}{4}\sqrt{\eta\lambda}$ . Suppose  $\eta = \frac{1}{\alpha}$ , where  $\alpha$  is the smoothness constant. Denote  $C_0 := \frac{\sqrt{2}(\beta+1)}{\sqrt{\min\{h(\beta, \eta\lambda_{\min}(H)), h(\beta, \eta\lambda_{\max}(H))\}}} \leq 4\sqrt{\kappa}$  by Corollary 1. Let  $C_1 = C_3 = C_0$  and  $C_2 = \frac{1}{4}\sqrt{\eta\lambda}$  in Theorem 6. The goal is to show that  $\left\| \begin{bmatrix} \xi_t \\ \xi_{t-1} \end{bmatrix} \right\| \leq \theta^t 2C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|$  for all  $t$  by induction. To achieve this, we will also use induction to show that for all iterations  $s$ ,

$$\|w_s - w_*\| \leq R := \frac{3}{64\sqrt{\kappa}C_0}. \quad (69)$$

A sufficient condition for the base case  $s = 0$  of (69) to hold is

$$\left\| \begin{bmatrix} w_0 - w_* \\ w_{-1} - w_* \end{bmatrix} \right\| \leq \frac{R}{2C_0} = \frac{3}{128\sqrt{\kappa}C_0^2}, \quad (70)$$

as  $C_0 \geq 1$  by Theorem 5, which in turn can be guaranteed if  $\left\| \begin{bmatrix} w_0 - w_* \\ w_{-1} - w_* \end{bmatrix} \right\| \leq \frac{1}{683\kappa^{3/2}}$  by using the upper bound  $C_0 \leq 4\sqrt{\kappa}$  of Corollary 1.

From Lemma 2, we have

$$\begin{aligned} \|\phi_s\| &\leq \eta \left\| \int_0^1 \nabla^2 f((1-\tau)w_s + \tau w_*) d\tau - \int_0^1 \nabla^2 f((1-\tau)w_0 + \tau w_*) d\tau \right\| \|\xi_s\| \\ &\stackrel{(a)}{\leq} \eta\alpha \left( \int_0^1 (1-\tau) \|w_s - w_0\| d\tau \right) \|\xi_s\| \leq \eta\alpha \|w_s - w_0\| \|\xi_s\| \\ &\stackrel{(b)}{\leq} \eta\alpha (\|w_s - w_*\| + \|w_0 - w_*\|) \|\xi_s\|, \end{aligned} \quad (71)$$

where (a) is by  $\alpha$ -Lipschitzness of the Hessian and (b) is by the triangle inequality. By (69), (71), Lemma 2, Theorem 6, and Corollary 1, it suffices to show that given  $\left\| \begin{bmatrix} \xi_s \\ \xi_{s-1} \end{bmatrix} \right\| \leq \theta^s 2C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|$  and  $\|w_s - w_*\| \leq R := \frac{3}{64\sqrt{\kappa}C_0}$  hold at  $s = 0, 1, \dots, t-1$ , one has

$$\left\| \sum_{s=0}^{t-1} A^{t-s-1} \begin{bmatrix} \varphi_s \\ 0 \end{bmatrix} \right\| \leq \theta^t C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\| \quad (72)$$

$$\|w_t - w_*\| \leq R := \frac{3}{64\sqrt{\kappa}C_0}, \quad (73)$$

where  $A := \begin{bmatrix} (1+\beta)I_n - \eta \int_0^1 \nabla^2 f((1-\tau)w_0 + \tau w_*) d\tau & -\beta I_n \\ I_n & 0 \end{bmatrix}$ .

We have

$$\begin{aligned} \left\| \sum_{s=0}^{t-1} A^{t-s-1} \begin{bmatrix} \varphi_s \\ 0 \end{bmatrix} \right\| &\leq \sum_{s=0}^{t-1} \|A^{t-s-1} \begin{bmatrix} \varphi_s \\ 0 \end{bmatrix}\| \\ &\stackrel{(a)}{\leq} \sum_{s=0}^{t-1} \beta_*^{t-s-1} C_0 \|\varphi_s\| \\ &\stackrel{(b)}{\leq} 4\eta\alpha RC_0^2 \sum_{s=0}^{t-1} \beta_*^{t-s-1} \theta^s \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\| \\ &\stackrel{(c)}{\leq} RC_0^2 \frac{64}{3\sqrt{\eta\lambda}} \theta^t \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\| \\ &\stackrel{(d)}{\leq} C_0 \theta^t \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|, \end{aligned} \quad (74)$$

where (a) uses Theorem 5 with  $\beta = \beta_*^2$ , (b) is by (71), (69), and the induction that  $\|\xi_s\| \leq \theta^s 2C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|$ , (c) is because  $\sum_{s=0}^{t-1} \beta_*^{t-1-s} \theta^s = \theta^{t-1} \sum_{s=0}^{t-1} \left(\frac{\beta_*}{\theta}\right)^{t-1-s} \leq \theta^{t-1} \sum_{s=0}^{t-1} \theta^{t-1-s} \leq \theta^{t-1} \frac{4}{\sqrt{\eta\lambda}} \leq \theta^t \frac{16}{3\sqrt{\eta\lambda}}$ , and (d) is due to the definition of  $R := \frac{3}{64\sqrt{\kappa}C_0}$ . Therefore, by Theorem 6, we have  $\left\| \begin{bmatrix} \xi_t \\ \xi_{t-1} \end{bmatrix} \right\| \leq \theta^t 2C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|$ .

Now let us switch to show (73). We have

$$\|\xi_t\| := \|w_t - w_*\| \stackrel{\text{induction}}{\leq} \theta^t 2C_0 \left\| \begin{bmatrix} w_0 - w_* \\ w_{-1} - w_* \end{bmatrix} \right\| \leq R, \quad (75)$$

where the last inequality uses the constraint  $\left\| \begin{bmatrix} w_0 - w_* \\ w_{-1} - w_* \end{bmatrix} \right\| \leq \frac{R}{2C_0}$  by (70). □



## F. Proof of Theorem 9

We will need some supporting lemmas in the following for the proof. In the following analysis, we denote  $C_0 := \frac{\sqrt{2(\beta+1)}}{\sqrt{\min\{h(\beta, \eta\lambda_{\min}(H)), h(\beta, \eta\lambda_{\max}(H))\}}}$ , where  $h(\beta, \cdot)$  is defined in Theorem 5 and  $H = H_0$  whose  $(i, j)$  entry is  $(H_0)_{i,j} := H(W_0)_{i,j} = \frac{1}{m} \sum_{r=1}^m x_i^\top x_j \mathbb{1}\{\langle w_0^{(r)}, x_i \rangle \geq 0 \ \& \ \langle w_0^{(r)}, x_j \rangle \geq 0\}$ , as defined in Lemma 3. In the following, we also denote  $\beta = (1 - \frac{1}{2}\sqrt{\eta\lambda})^2 := \beta_*^2$ . We summarize the notations in Table 1.

Notation	definition (or value)	meaning
$\mathcal{N}_W^{\text{ReLU}}(x)$	$\mathcal{N}_W^{\text{ReLU}}(x) := \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \sigma(\langle w^{(r)}, x \rangle)$	the ReLU network's output given $x$
$\bar{H}$	$\bar{H}_{i,j} := \mathbb{E}_{w^{(r)}} [x_i^\top x_j \mathbb{1}\{\langle w^{(r)}, x_i \rangle \geq 0 \ \& \ \langle w^{(r)}, x_j \rangle \geq 0\}]$ .	the expectation of the Gram matrix
$H_0$	$H(W_0)_{i,j} = \frac{1}{m} \sum_{r=1}^m x_i^\top x_j \mathbb{1}\{\langle w_0^{(r)}, x_i \rangle \geq 0 \ \& \ \langle w_0^{(r)}, x_j \rangle \geq 0\}$	the Gram matrix at the initialization
$\lambda_{\min}(\bar{H})$	$\lambda_{\min}(\bar{H}) > 0$ (by assumption)	the least eigenvalue of $\bar{H}$ .
$\lambda_{\max}(\bar{H})$		the largest eigenvalue of $\bar{H}$
$\kappa$	$\kappa := \lambda_{\max}(\bar{H})/\lambda_{\min}(\bar{H})$	the condition number of $\bar{H}$
$\lambda$	$\lambda := \frac{3}{4} \lambda_{\min}(\bar{H})$	(a lower bound of) the least eigenvalue of $H_0$ .
$\lambda_{\max}$	$\lambda_{\max} := \lambda_{\max}(\bar{H}) + \frac{\lambda_{\min}(\bar{H})}{4}$	(an upper bound of) the largest eigenvalue of $H_0$ .
$\hat{\kappa}$	$\hat{\kappa} := \frac{\lambda_{\max}}{\lambda} = \frac{4}{3}\kappa + \frac{1}{3}$	the condition number of $H_0$ .
$\eta$	$\eta = 1/\lambda_{\max}$	step size
$\beta$	$\beta = (1 - \frac{1}{2}\sqrt{\eta\lambda})^2 = (1 - \frac{1}{2\sqrt{\hat{\kappa}}})^2 := \beta_*^2$	momentum parameter
$\beta_*$	$\beta_* = \sqrt{\beta} = 1 - \frac{1}{2}\sqrt{\eta\lambda}$	squared root of $\beta$
$\theta$	$\theta = \beta_* + \frac{1}{4}\sqrt{\eta\lambda} = 1 - \frac{1}{4}\sqrt{\eta\lambda} = 1 - \frac{1}{4\sqrt{\hat{\kappa}}}$	the convergence rate
$C_0$	$C_0 := \frac{\sqrt{2(\beta+1)}}{\sqrt{\min\{h(\beta, \eta\lambda_{\min}(H_0)), h(\beta, \eta\lambda_{\max}(H_0))\}}}$	the constant used in Theorem 5

Table 1. Summary of the notations for proving Theorem 9.

**Lemma 8.** Suppose that the neurons  $w_0^{(1)}, \dots, w_0^{(m)}$  are i.i.d. generated by  $N(0, I_d)$  initially. Then, for any set of weight vectors  $W_t := \{w_t^{(1)}, \dots, w_t^{(m)}\}$  that satisfy for any  $r \in [m]$ ,  $\|w_t^{(r)} - w_0^{(r)}\| \leq R^{\text{ReLU}} := \frac{\lambda}{1024nC_0}$ , it holds that

$$\|H_t - H_0\|_F \leq 2nR^{\text{ReLU}} = \frac{\lambda}{512C_0},$$

with probability at least  $1 - n^2 \cdot \exp(-mR^{\text{ReLU}}/10)$ .

*Proof.* This is an application of Lemma 3.2 in (Song & Yang, 2019).  $\square$

Lemma 8 shows that if the distance between the current iterate  $W_t$  and its initialization  $W_0$  is small, then the distance between the Gram matrix  $H(W_t)$  and  $H(W_0)$  should also be small. Lemma 8 allows us to obtain the following lemma, which bounds the size of  $\varphi_t$  (defined in Lemma 3) in the residual dynamics.

**Lemma 9.** Following the setting as Theorem 9, denote  $\theta := \beta_* + \frac{1}{4}\sqrt{\eta\lambda} = 1 - \frac{1}{4}\sqrt{\eta\lambda}$ . Suppose that  $\forall i \in [n], |S_i^\perp| \leq 4mR^{\text{ReLU}}$  for some constant  $R^{\text{ReLU}} := \frac{\lambda}{1024nC_0} > 0$ . If we have (I) for any  $s \leq t$ , the residual dynamics satisfies  $\| \begin{bmatrix} \xi_s \\ \xi_{s-1} \end{bmatrix} \| \leq \theta^s \cdot \nu C_0 \| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \|$ , for some constant  $\nu > 0$ , and (II) for any  $r \in [m]$  and any  $s \leq t$ ,  $\|w_s^{(r)} - w_0^{(r)}\| \leq R^{\text{ReLU}}$ , then  $\phi_t$  and  $\iota_t$  in Lemma 3 satisfies

$$\|\phi_t\| \leq \frac{\sqrt{\eta\lambda}}{16} \theta^t \nu \| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \|, \text{ and } \|\iota_t\| \leq \frac{\eta\lambda}{512} \theta^t \nu \| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \|.$$

Consequently,  $\varphi_t$  in Lemma 3 satisfies

$$\|\varphi_t\| \leq \left( \frac{\sqrt{\eta\lambda}}{16} + \frac{\eta\lambda}{512} \right) \theta^t \nu \| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \|.$$

*Proof.* Denote  $\beta_* := 1 - \frac{1}{2}\sqrt{\eta\lambda}$  and  $\theta := \beta_* + \frac{1}{4}\sqrt{\eta\lambda} = 1 - \frac{1}{4}\sqrt{\eta\lambda}$ . We have by Lemma 3

$$\begin{aligned}
 \|\phi_t\| &= \sqrt{\sum_{i=1}^n \phi_t[i]^2} \leq \sqrt{\sum_{i=1}^n \left( \frac{2\eta\sqrt{n}|S_i^\perp|}{m} (\|\xi_t\| + \beta \sum_{\tau=0}^{t-1} \beta^{t-1-\tau} \|\xi_\tau\|) \right)^2} \\
 &\stackrel{(a)}{\leq} 8\eta m R^{\text{ReLU}} (\|\xi_t\| + \beta \sum_{\tau=0}^{t-1} \beta^{t-1-\tau} \|\xi_\tau\|) \\
 &\stackrel{(b)}{\leq} 8\eta m R^{\text{ReLU}} \left( \theta^t \nu C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\| + \beta \sum_{\tau=0}^{t-1} \beta^{t-1-\tau} \theta^\tau \nu C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\| \right) \\
 &\stackrel{(c)}{=} 8\eta m R^{\text{ReLU}} \left( \theta^t \nu C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\| + \beta_*^2 \nu C_0 \sum_{\tau=0}^{t-1} \beta_*^{2(t-1-\tau)} \theta^\tau \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\| \right) \\
 &\stackrel{(d)}{\leq} 8\eta m R^{\text{ReLU}} \left( \theta^t \nu C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\| + \beta_*^2 \nu C_0 \theta^{t-1} \sum_{\tau=0}^{t-1} \theta^{t-1-\tau} \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\| \right) \\
 &\leq 8\eta m R^{\text{ReLU}} \theta^t (1 + \beta_* \sum_{\tau=0}^{t-1} \theta^\tau) \nu C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\| \\
 &\leq 8\eta m R^{\text{ReLU}} \theta^t (1 + \frac{\beta_*}{1-\theta}) \nu C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\| \\
 &\stackrel{(e)}{\leq} \frac{\sqrt{\eta\lambda}}{16} \theta^t \nu \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|, \tag{76}
 \end{aligned}$$

where (a) is by  $|S_i^\perp| \leq 4mR^{\text{ReLU}}$ , (b) is by induction that  $\|\xi_t\| \leq \theta^t \nu C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|$  as  $u_0 = u_{-1}$ , (c) uses that  $\beta = \beta_*^2$ , (d) uses  $\beta_* \leq \theta$ , (e) uses  $1 + \frac{\beta_*}{1-\theta} \leq \frac{2}{1-\theta} \leq \frac{8}{\sqrt{\eta\lambda}}$  and  $R^{\text{ReLU}} := \frac{\lambda}{1024nC_0}$ .

Now let us switch to bound  $\|\iota_t\|$ .

$$\|\iota_t\| \leq \eta \|H_0 - H_t\|_2 \|\xi_t\| \leq \frac{\eta\lambda}{512C_0} \theta^t \nu C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|, \tag{77}$$

where we uses Lemma 8 that  $\|H_0 - H_t\|_2 \leq \frac{\lambda}{512C_0}$  and the induction that  $\left\| \begin{bmatrix} \xi_t \\ \xi_{t-1} \end{bmatrix} \right\| \leq \theta^t \nu C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|$ . □

The assumption of Lemma 9,  $\forall i \in [n], |S_i^\perp| \leq 4mR^{\text{ReLU}}$  only depends on the initialization. Lemma 11 shows that it holds with probability at least  $1 - n \cdot \exp(-mR^{\text{ReLU}})$ .

**Lemma 10.** *Following the setting as Theorem 9, denote  $\theta := \beta_* + \frac{1}{4}\sqrt{\eta\lambda} = 1 - \frac{1}{4}\sqrt{\eta\lambda}$ . Suppose that the initial error satisfies  $\|\xi_0\|^2 = O(n \log(m/\delta) \log^2(n/\delta))$ . If for any  $s < t$ , the residual dynamics satisfies  $\left\| \begin{bmatrix} \xi_s \\ \xi_{s-1} \end{bmatrix} \right\| \leq \theta^s \cdot \nu C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|$ , for some constant  $\nu > 0$ , then*

$$\|w_t^{(r)} - w_0^{(r)}\| \leq R^{\text{ReLU}} := \frac{\lambda}{1024nC_0}.$$

*Proof.* We have

$$\begin{aligned}
 \|w_{t+1}^{(r)} - w_0^{(r)}\| &\stackrel{(a)}{\leq} \eta \sum_{s=0}^t \|M_s^{(r)}\| \stackrel{(b)}{=} \eta \sum_{s=0}^t \left\| \sum_{\tau=0}^s \beta^{s-\tau} \frac{\partial L(W_\tau)}{\partial w_\tau^{(r)}} \right\| \leq \eta \sum_{s=0}^t \sum_{\tau=0}^s \beta^{s-\tau} \left\| \frac{\partial L(W_\tau)}{\partial w_\tau^{(r)}} \right\| \\
 &\stackrel{(c)}{\leq} \eta \sum_{s=0}^t \sum_{\tau=0}^s \beta^{s-\tau} \frac{\sqrt{n}}{\sqrt{m}} \|y - u_\tau\| \\
 &\stackrel{(d)}{\leq} \eta \sum_{s=0}^t \sum_{\tau=0}^s \beta^{s-\tau} \frac{\sqrt{2n}}{\sqrt{m}} \theta^\tau \nu C_0 \|y - u_0\| \\
 &\stackrel{(e)}{\leq} \frac{\eta\sqrt{2n}}{\sqrt{m}} \sum_{s=0}^t \frac{\theta^s}{1-\theta} \nu C_0 \|y - u_0\| \leq \frac{\eta\sqrt{2n}}{\sqrt{m}} \left( \frac{\nu C_0}{(1-\theta)^2} \right) \|y - u_0\| \\
 &\stackrel{(f)}{=} \frac{\eta\sqrt{2n}}{\sqrt{m}} \left( \frac{16\nu C_0}{\eta\lambda} \right) \|y - u_0\| \\
 &\stackrel{(g)}{\leq} \frac{\eta\sqrt{2n}}{\sqrt{m}} \left( \frac{16\nu C_0}{\eta\lambda} \right) O(\sqrt{n \log(m/\delta) \log^2(n/\delta)}) \\
 &\stackrel{(h)}{\leq} \frac{\lambda}{1024n C_0}, \tag{78}
 \end{aligned}$$

where (a), (b) is by the update rule of momentum, which is  $w_{t+1}^{(r)} - w_t^{(r)} = -\eta M_t^{(r)}$ , where  $M_t^{(r)} := \sum_{s=0}^t \beta^{t-s} \frac{\partial L(W_s)}{\partial w_s^{(r)}}$ , (c) is because  $\left\| \frac{\partial L(W_s)}{\partial w_s^{(r)}} \right\| = \left\| \sum_{i=1}^n (y_i - u_s[i]) \frac{1}{\sqrt{m}} a_r x_i \cdot \mathbb{1}\{\langle w_s^{(r)}, x \rangle \geq 0\} \right\| \leq \frac{1}{\sqrt{m}} \sum_{i=1}^n |y_i - u_s[i]| \leq \frac{\sqrt{n}}{\sqrt{m}} \|y - u_s\|$ , (d) is by  $\left\| \begin{bmatrix} \xi_s \\ \xi_{s-1} \end{bmatrix} \right\| \leq \theta^s \nu C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|$  (e) is because that  $\beta = \beta_*^2 \leq \theta^2$ , (f) we use  $\theta := (1 - \frac{1}{4}\sqrt{\eta\lambda})$ , so that  $\frac{1}{(1-\theta)^2} = \frac{16}{\eta\lambda}$ , (g) is by that the initial error satisfies  $\|y - u_0\|^2 = O(n \log(m/\delta) \log^2(n/\delta))$ , and (h) is by the choice of the number of neurons  $m = \Omega(\lambda^{-4} n^4 C_0^4 \log^3(n/\delta)) = \Omega(\lambda^{-4} n^4 \kappa^2 \log^3(n/\delta))$ , as  $C_0 = \Theta(\sqrt{\kappa})$  by Corollary 1.

The proof is completed.  $\square$

Lemma 10 basically says that if the size of the residual errors is bounded and decays over iterations, then the distance between the current iterate  $W_t$  and its initialization  $W_0$  is well-controlled. The lemma will allow us to invoke Lemma 8 and Lemma 9 when proving Theorem 9. The proof of Lemma 10 is in Appendix F. The assumption of Lemma 10,  $\|\xi_0\|^2 = O(n \log(m/\delta) \log^2(n/\delta))$ , is satisfied by the random initialization with probability at least  $1 - \delta/3$  according to Lemma 12.

**Lemma 11.** (Claim 3.12 of (Song & Yang, 2019)) Fix a number  $R_1 \in (0, 1)$ . Recall that  $S_i^\perp$  is a random set defined in Subsection 3.3. With probability at least  $1 - n \cdot \exp(-mR_1)$ , we have that for all  $i \in [n]$ ,

$$|S_i^\perp| \leq 4mR_1.$$

A similar lemma also appears in (Du et al., 2019b). Lemma 11 says that the number of neurons whose activation patterns for a sample  $i$  could change during the execution is only a small fraction of  $m$  if  $R_1$  is a small number, i.e.  $|S_i^\perp| \leq 4mR_1 \ll m$ .

**Lemma 12.** (Claim 3.10 in (Song & Yang, 2019)) Assume that  $w_0^{(r)} \sim N(0, I_d)$  and  $a_r$  uniformly sampled from  $\{-1, 1\}$ . For  $0 < \delta < 1$ , we have that

$$\|y - u_0\|^2 = O(n \log(m/\delta) \log^2(n/\delta)),$$

with probability at least  $1 - \delta$ .

### F.1. Proof of Theorem 9

*Proof.* (of Theorem 9) Denote  $\lambda := \frac{3}{4} \lambda_{\min}(\bar{H}) > 0$ . Lemma 5 shows that  $\lambda$  is a lower bound of  $\lambda_{\min}(H)$  of the matrix  $H$  defined in Lemma 3. Also, denote  $\beta_* := 1 - \frac{1}{2}\sqrt{\eta\lambda}$  (note that  $\beta = \beta_*^2$ ) and  $\theta := \beta_* + \frac{1}{4}\sqrt{\eta\lambda} = 1 - \frac{1}{4}\sqrt{\eta\lambda}$ . In the

following, we let  $\nu = 2$  in Lemma 9, 10, and let  $C_1 = C_3 = C_0$  and  $C_2 = \frac{1}{4}\sqrt{\eta\lambda}$  in Theorem 6. The goal is to show that

$$\left\| \begin{bmatrix} \xi_t \\ \xi_{t-1} \end{bmatrix} \right\| \leq \theta^t 2C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\| \text{ for all } t \text{ by induction. To achieve this, we will also use induction to show that for all iterations } s,$$

$$\forall r \in [m], \|w_s^{(r)} - w_0^{(r)}\| \leq R^{\text{ReLU}} := \frac{\lambda}{1024nC_0}, \quad (79)$$

which is clear true in the base case  $s = 0$ .

By Lemma 3, 5, 8, 9, Theorem 6, and Corollary 1, it suffices to show that given  $\left\| \begin{bmatrix} \xi_s \\ \xi_{s-1} \end{bmatrix} \right\| \leq \theta^s 2C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|$  and (79) hold at  $s = 0, 1, \dots, t-1$ , one has

$$\left\| \sum_{s=0}^{t-1} A^{t-s-1} \begin{bmatrix} \varphi_s \\ 0 \end{bmatrix} \right\| \leq \theta^t C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|, \quad (80)$$

$$\forall r \in [m], \|w_t^{(r)} - w_0^{(r)}\| \leq R^{\text{ReLU}} := \frac{\lambda}{1024nC_0}, \quad (81)$$

where the matrix  $A$  and the vector  $\varphi_t$  are defined in Lemma 3. The inequality (80) is the required condition for using the result of Theorem 6, while the inequality (81) helps us to show (80) through invoking Lemma 9 to bound the terms  $\{\varphi_s\}$  as shown in the following.

We have

$$\begin{aligned} \left\| \sum_{s=0}^{t-1} A^{t-s-1} \begin{bmatrix} \varphi_s \\ 0 \end{bmatrix} \right\| &\stackrel{(a)}{\leq} \sum_{s=0}^{t-1} \beta_*^{t-s-1} C_0 \|\varphi_s\| \\ &\stackrel{(b)}{\leq} \left( \frac{\sqrt{\eta\lambda}}{16} + \frac{\eta\lambda}{512} \right) 2C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\| \left( \sum_{s=0}^{t-1} \beta_*^{t-1-s} \theta^s \right) \\ &\stackrel{(c)}{\leq} \left( \frac{1}{2} + \frac{1}{64} \sqrt{\eta\lambda} \right) \theta^{t-1} C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\| \stackrel{(d)}{\leq} \theta^t C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|, \end{aligned} \quad (82)$$

where (a) uses Theorem 5, (b) is due to Lemma 9, Lemma 11, (c) is because  $\sum_{s=0}^{t-1} \beta_*^{t-1-s} \theta^s = \theta^{t-1} \sum_{s=0}^{t-1} \left( \frac{\beta_*}{\theta} \right)^{t-1-s} \leq \theta^{t-1} \sum_{s=0}^{t-1} \theta^{t-1-s} \leq \theta^{t-1} \frac{4}{\sqrt{\eta\lambda}}$ , (d) uses that  $\theta \geq \frac{3}{4}$  and  $\eta\lambda \leq 1$ . Hence, we have shown (80). Therefore, by Theorem 6, we have  $\left\| \begin{bmatrix} \xi_t \\ \xi_{t-1} \end{bmatrix} \right\| \leq \theta^t 2C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|$ .

By Lemma 10 and Lemma 12, we have (81). Furthermore, with the choice of  $m$ , we have  $3n^2 \exp(-mR^{\text{ReLU}}/10) \leq \delta$ . Thus, we have completed the proof.  $\square$

## G. Proof of Theorem 10

We will need some supporting lemmas in the following for the proof. In the following analysis, we denote  $C_0 := \frac{\sqrt{2(\beta+1)}}{\sqrt{\min\{h(\beta, \eta\lambda_{\min}(H)), h(\beta, \eta\lambda_{\max}(H))\}}}$ , where  $h(\beta, \cdot)$  is the constant defined in Theorem 5 and  $H = H_0 := \frac{1}{m^{L-1}d_y} \sum_{l=1}^L [(W_0^{(l-1:1)} X)^\top (W_0^{(l-1:1)} X) \otimes W_0^{(L:l+1)} (W_0^{(L:l+1)})^\top] \in \mathbb{R}^{d_y n \times d_y n}$ , as defined in Lemma 4. We also denote  $\beta = (1 - \frac{1}{2}\sqrt{\eta\lambda})^2 := \beta_*^2$ . As mentioned in the main text, following Du & Hu (2019); Hu et al. (2020b), we will further assume that (A1) there exists a  $W^*$  such that  $Y = W^* X$ ,  $X \in \mathbb{R}^{d \times \bar{r}}$ , and  $\bar{r} = \text{rank}(X)$ , which is actually without loss of generality (see e.g. the discussion in Appendix B of Du & Hu (2019)). We summarize the notions in Table 2.

Notation	definition (or value)	meaning
$\mathcal{N}_W^{L\text{-linear}}(x)$	$\mathcal{N}_W^{L\text{-linear}}(x) := \frac{1}{\sqrt{m^{L-1}d_y}} W^{(L)} W^{(L-1)} \dots W^{(1)} x$ ,	output of the deep linear network
$H_0$	$H_0 := \frac{1}{m^{L-1}d_y} \sum_{l=1}^L [(W_0^{(l-1:1)} X)^\top (W_0^{(l-1:1)} X) \otimes W_0^{(L:l+1)} (W_0^{(L:l+1)})^\top] \in \mathbb{R}^{d_y n \times d_y n}$	$H$ in (8) is $H = H_0$ (Lemma 4)
$\lambda_{\max}(H_0)$	$\lambda_{\max}(H_0) \leq L\sigma_{\max}^2(X)/d_y$ (Lemma 13)	the largest eigenvalue of $H_0$
$\lambda_{\min}(H_0)$	$\lambda_{\min}(H_0) \geq L\sigma_{\min}^2(X)/d_y$ (Lemma 13)	the least eigenvalue of $H_0$
$\lambda$	$\lambda := L\sigma_{\min}^2(X)/d_y$	(a lower bound of) the least eigenvalue of $H_0$
$\kappa$	$\kappa := \frac{\lambda_1(X^\top X)}{\lambda_{\bar{r}}(X^\top X)} = \frac{\sigma_{\max}^2(X)}{\sigma_{\min}^2(X)}$ (A1)	the condition number of the data matrix $X$
$\hat{\kappa}$	$\hat{\kappa} := \frac{\lambda_{\max}(H_0)}{\lambda_{\min}(H_0)} \leq \frac{\sigma_{\max}^2(X)}{\sigma_{\min}^2(X)} = \kappa$ (Lemma 13)	the condition number of $H_0$
$\eta$	$\eta = \frac{d_y}{L\sigma_{\max}^2(X)}$	step size
$\beta$	$\beta = (1 - \frac{1}{2}\sqrt{\eta\lambda})^2 = (1 - \frac{1}{2\sqrt{\kappa}})^2 := \beta_*^2$	momentum parameter
$\beta_*$	$\beta_* = \sqrt{\beta} = 1 - \frac{1}{2}\sqrt{\eta\lambda}$	squared root of $\beta$
$\theta$	$\theta = \beta_* + \frac{1}{4}\sqrt{\eta\lambda} = 1 - \frac{1}{4}\sqrt{\eta\lambda} = 1 - \frac{1}{4\sqrt{\kappa}}$	the convergence rate
$C_0$	$C_0 := \frac{\sqrt{2(\beta+1)}}{\sqrt{\min\{h(\beta, \eta\lambda_{\min}(H_0)), h(\beta, \eta\lambda_{\max}(H_0))\}}}$	the constant used in Theorem 5

Table 2. Summary of the notations for proving Theorem 10. We will simply use  $\kappa$  to represent the condition number of the matrix  $H_0$  in the analysis since we have  $\hat{\kappa} \leq \kappa$ .

**Lemma 13.** [Lemma 4.2 in (Hu et al., 2020b)] *By the orthogonal initialization, we have*

$$\lambda_{\min}(H_0) \geq L\sigma_{\min}^2(X)/d_y, \quad \lambda_{\max}(H_0) \leq L\sigma_{\max}^2(X)/d_y.$$

$$\sigma_{\max}(W_0^{(j:i)}) = m^{\frac{j-i+1}{2}}, \quad \sigma_{\min}(W_0^{(j:i)}) = m^{\frac{j-i+1}{2}}$$

Furthermore, with probability  $1 - \delta$ ,

$$\ell(W_0) \leq B_0^2 = O\left(1 + \frac{\log(\bar{r}/\delta)}{d_y} + \|W^*\|_2^2\right),$$

for some constant  $B_0 > 0$ .

We remark that Lemma 13 implies that the condition number of  $H_0$  satisfies

$$\hat{\kappa} := \frac{\lambda_{\max}(H_0)}{\lambda_{\min}(H_0)} \leq \frac{\sigma_{\max}^2(X)}{\sigma_{\min}^2(X)} = \kappa. \quad (83)$$

**Lemma 14.** *Following the setting as Theorem 10, denote  $\theta := \beta_* + \frac{1}{4}\sqrt{\eta\lambda} = 1 - \frac{1}{4}\sqrt{\eta\lambda}$ . If we have (I) for any  $s \leq t$ , the residual dynamics satisfies  $\left\| \begin{bmatrix} \xi_s \\ \xi_{s-1} \end{bmatrix} \right\| \leq \theta^s \cdot \nu C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|$ , for some constant  $\nu > 0$ , and (II) for all  $l \in [L]$  and for any  $s \leq t$ ,  $\|W_s^{(l)} - W_0^{(l)}\|_F \leq R^{L\text{-linear}} := \frac{64\|X\|_2\sqrt{d_y}}{L\sigma_{\min}^2(X)} \nu C_0 B_0$ , then*

$$\|\phi_t\| \leq \frac{43\sqrt{d_y}}{\sqrt{m}\|X\|_2} \theta^{2t} \nu^2 C_0^2 \left(\frac{\|\xi_0\|}{1-\theta}\right)^2, \quad \|\psi_t\| \leq \frac{43\sqrt{d_y}}{\sqrt{m}\|X\|_2} \theta^{2(t-1)} \nu^2 C_0^2 \left(\frac{\|\xi_0\|}{1-\theta}\right)^2, \quad \|t_t\| \leq \frac{\eta\lambda}{80} \theta^t \nu C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|.$$

Consequently,  $\varphi_t$  in Lemma 4 satisfies

$$\|\varphi_t\| \leq \frac{1920\sqrt{d_y}}{\sqrt{m}\|X\|_2} \frac{1}{\eta\lambda} \theta^{2t} \nu^2 C_0^2 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|^2 + \frac{\eta\lambda}{80} \theta^t \nu C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|.$$

*Proof.* By Lemma 4,  $\varphi_t = \phi_t + \psi_t + \iota_t \in \mathbb{R}^{d_y n}$ , we have

$$\phi_t := \frac{1}{\sqrt{m^{L-1}d_y}} \text{vec}(\Phi_t X), \text{ with } \Phi_t := \Pi_l \left( W_t^{(l)} - \eta M_{t,l} \right) - W_t^{(L:1)} + \eta \sum_{l=1}^L W_t^{(L:l+1)} M_{t,l} W_t^{(l-1:1)}, \quad (84)$$

and

$$\psi_t := \frac{1}{\sqrt{m^{L-1}d_y}} \text{vec} \left( (L-1)\beta W_t^{(L:1)} X + \beta W_{t-1}^{(L:1)} X - \beta \sum_{l=1}^L W_t^{(L:l+1)} W_{t-1}^{(l)} W_t^{(l-1:1)} X \right). \quad (85)$$

and

$$\iota_t := \eta(H_0 - H_t)\xi_t. \quad (86)$$

So if we can bound  $\|\phi_t\|$ ,  $\|\psi_t\|$ , and  $\|\iota_t\|$  respectively, then we can bound  $\|\varphi_t\|$  by the triangle inequality.

$$\|\varphi_t\| \leq \|\phi_t\| + \|\psi_t\| + \|\iota_t\|. \quad (87)$$

Let us first upper-bound  $\|\phi_t\|$ . Note that  $\Phi_t$  is the sum of all the high-order (of  $\eta$ 's) term in the product,

$$W_{t+1}^{(L:1)} = \Pi_l \left( W_t^{(l)} - \eta M_{t,l} \right) = W_t^{(L:1)} - \eta \sum_{l=1}^L W_t^{(L:l+1)} M_{t,l} W_t^{(l-1:1)} + \Phi_t. \quad (88)$$

By induction, we can bound the gradient norm of each layer as

$$\begin{aligned} \left\| \frac{\partial \ell(W_s^{(L:1)})}{\partial W_s^{(l)}} \right\|_F &\leq \frac{1}{\sqrt{m^{L-1}d_y}} \|W_s^{(L:l+1)}\|_2 \|U_s - Y\|_F \|W_s^{(l-1:1)}\|_2 \|X\|_2 \\ &\leq \frac{1}{\sqrt{m^{L-1}d_y}} 1.1m^{\frac{L-l}{2}} \theta^s \nu C_0 2\sqrt{2} \|U_0 - Y\|_F 1.1m^{\frac{l-1}{2}} \|X\|_2 \\ &\leq \frac{4\|X\|_2}{\sqrt{d_y}} \theta^s \nu C_0 \|U_0 - Y\|_F, \end{aligned} \quad (89)$$

where the second inequality we use Lemma 16 and that  $\left\| \begin{bmatrix} \xi_s \\ \xi_{s-1} \end{bmatrix} \right\| \leq \theta^s \nu C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|$  and  $\|\xi_s\| = \|U_s - Y\|_F$ .

So the momentum term of each layer can be bounded as

$$\begin{aligned} \|M_{t,l}\|_F &= \left\| \sum_{s=0}^t \beta^{t-s} \frac{\partial \ell(W_s^{(L:1)})}{\partial W_s^{(l)}} \right\|_F \leq \sum_{s=0}^t \beta^{t-s} \left\| \frac{\partial \ell(W_s^{(L:1)})}{\partial W_s^{(l)}} \right\|_F \\ &\leq \frac{4\|X\|_2}{\sqrt{d_y}} \sum_{s=0}^t \beta^{t-s} \theta^s \nu C_0 \|U_0 - Y\|_F \\ &\leq \frac{4\|X\|_2}{\sqrt{d_y}} \sum_{s=0}^t \theta^{2(t-s)} \theta^s \nu C_0 \|U_0 - Y\|_F \\ &\leq \frac{4\|X\|_2}{\sqrt{d_y}} \frac{\theta^t}{1-\theta} \nu C_0 \|U_0 - Y\|_F, \end{aligned} \quad (90)$$

where in the second to last inequality we use  $\beta = \beta_*^2 \leq \theta^2$ .

Combining all the pieces together, we can bound  $\|\frac{1}{\sqrt{m^{L-1}d_y}}\Phi_t X\|_F$  as

$$\begin{aligned}
 & \|\frac{1}{\sqrt{m^{L-1}d_y}}\Phi_t X\|_F \\
 & \stackrel{(a)}{\leq} \frac{1}{\sqrt{m^{L-1}d_y}} \sum_{j=2}^L \binom{L}{j} \left( \eta \frac{4\|X\|_2}{\sqrt{d_y}} \frac{\theta^t}{1-\theta} \nu C_0 \|U_0 - Y\|_F \right)^j (1.1)^{j+1} m^{\frac{L-j}{2}} \|X\|_2 \\
 & \stackrel{(b)}{\leq} 1.1 \frac{1}{\sqrt{m^{L-1}d_y}} \sum_{j=2}^L L^j \left( \eta \frac{4.4\|X\|_2}{\sqrt{d_y}} \frac{\theta^t}{1-\theta} \nu C_0 \|U_0 - Y\|_F \right)^j m^{\frac{L-j}{2}} \|X\|_2 \\
 & \leq 1.1 \sqrt{\frac{m}{d_y}} \|X\|_2 \sum_{j=2}^L \left( \eta \frac{4.4L\|X\|_2}{\sqrt{md_y}} \frac{\theta^t}{1-\theta} \nu C_0 \|U_0 - Y\|_F \right)^j, \tag{91}
 \end{aligned}$$

where (a) uses (90) and Lemma 16 for bounding a  $j \geq 2$  higher-order terms like  $\frac{1}{\sqrt{m^{L-1}d_y}} \beta W_t^{(L:k_j+1)} \cdot (-\eta M_{t,k_j}) W_t^{(k_j-1:k_{j-1}+1)} \cdot (-\eta M_{t,k_{j-1}}) \cdots (-\eta M_{t,k_1}) \cdot W_t^{(k_1-1:1)}$ , where  $1 \leq k_1 < \cdots < k_j \leq L$  and (b) uses that  $\binom{L}{j} \leq \frac{L^j}{j!}$

To proceed, let us bound  $\eta \frac{4.4L\|X\|_2}{\sqrt{md_y}} \frac{\theta^t}{1-\theta} \nu C_0 \|U_0 - Y\|_F$  in the sum above. We have

$$\begin{aligned}
 \eta \frac{4.4L\|X\|_2}{\sqrt{md_y}} \frac{\theta^t}{1-\theta} \nu C_0 \|U_0 - Y\|_F & \leq 4.4 \sqrt{\frac{d_y}{m}} \frac{1}{\|X\|_2} \frac{\theta^t}{1-\theta} \nu C_0 \|U_0 - Y\|_F \\
 & \leq 0.5, \tag{92}
 \end{aligned}$$

where the last inequality uses that  $\tilde{C}_1 \frac{d_y B_0^2 C_0^2}{\|X\|_2^2} \frac{1}{(1-\theta)^2} \leq \tilde{C}_1 \frac{d_y B_0^2 C_0^2}{\|X\|_2^2} \frac{1}{\eta \lambda} \leq \tilde{C}_2 \frac{d_y B_0^2 \kappa^2}{\|X\|_2^2} \leq m$ , for some sufficiently large constant  $\tilde{C}_1, \tilde{C}_2 > 0$ . Combining the above results, we have

$$\begin{aligned}
 \|\phi_t\| & = \|\frac{1}{\sqrt{m^{L-1}d_y}}\Phi_t X\|_F \\
 & \leq 1.1 \sqrt{\frac{m}{d_y}} \|X\|_2 \left( \eta \frac{4.4L\|X\|_2}{\sqrt{md_y}} \frac{\theta^t}{1-\theta} \nu C_0 \|U_0 - Y\|_F \right)^2 \sum_{j=2}^{L-2} (0.5)^{j-2} \\
 & \leq 2.2 \sqrt{\frac{m}{d_y}} \|X\|_2 \left( \eta \frac{4.4L\|X\|_2}{\sqrt{md_y}} \frac{\theta^t}{1-\theta} \nu C_0 \|U_0 - Y\|_F \right)^2 \\
 & \leq \frac{43\sqrt{d_y}}{\sqrt{m}\|X\|_2} \left( \frac{\theta^t}{1-\theta} \nu C_0 \|U_0 - Y\|_F \right)^2. \tag{93}
 \end{aligned}$$

Now let us switch to upper-bound  $\|\psi_t\|$ . It is equivalent to upper-bounding the Frobenius norm of  $\frac{1}{\sqrt{m^{L-1}d_y}} \beta (L-1) W_t^{(L:1)} X + \frac{1}{\sqrt{m^{L-1}d_y}} \beta W_{t-1}^{(L:1)} X - \frac{1}{\sqrt{m^{L-1}d_y}} \beta \sum_{l=1}^L W_t^{(L:l+1)} W_{t-1}^{(l)} W_t^{(l-1:1)} X$ , which can be rewritten as

$$\begin{aligned}
 & \underbrace{\frac{1}{\sqrt{m^{L-1}d_y}} \beta (L-1) \cdot \prod_{l=1}^L (W_{t-1}^{(l)} - \eta M_{t-1,l}) X}_{\text{first term}} + \underbrace{\frac{1}{\sqrt{m^{L-1}d_y}} \beta W_{t-1}^{(L:1)} X}_{\text{second term}} \\
 & - \underbrace{\frac{1}{\sqrt{m^{L-1}d_y}} \beta \sum_{l=1}^L \prod_{i=l+1}^L (W_{t-1}^{(i)} - \eta M_{t-1,i}) W_{t-1}^{(l)} \prod_{j=1}^{l-1} (W_{t-1}^{(j)} - \eta M_{t-1,j}) X}_{\text{third term}}. \tag{94}
 \end{aligned}$$

The above can be written as  $B_0 + \eta B_1 + \eta^2 B_2 + \dots + \eta^L B_L$  for some matrices  $B_0, \dots, B_L \in \mathbb{R}^{d_y \times n}$ . Specifically, we have

$$\begin{aligned}
 B_0 &= \underbrace{\frac{1}{\sqrt{m^{L-1}d_y}}(L-1)\beta W_{t-1}^{(L:1)}X}_{\text{due to the first term}} + \underbrace{\frac{1}{\sqrt{m^{L-1}d_y}}\beta W_{t-1}^{(L:1)}X}_{\text{due to the second term}} - \underbrace{\frac{1}{\sqrt{m^{L-1}d_y}}\beta L W_{t-1}^{(L:1)}X}_{\text{due to the third term}} = 0 \\
 B_1 &= -\underbrace{\frac{1}{\sqrt{m^{L-1}d_y}}(L-1)\beta \sum_{l=1}^L W_{t-1}^{(L:l+1)}M_{t-1,l}W_{t-1}^{(l-1:1)}}_{\text{due to the first term}} + \underbrace{\frac{1}{\sqrt{m^{L-1}d_y}}\beta \sum_{l=1}^L \sum_{k \neq l} W_{t-1}^{(L:k+1)}M_{t-1,k}W_{t-1}^{(k-1:1)}}_{\text{due to the third term}} = 0.
 \end{aligned} \tag{95}$$

So what remains on (94) are all the higher-order terms (in terms of the power of  $\eta$ ), i.e. those with  $\eta M_{t-1,i}$  and  $\eta M_{t-1,j}$ ,  $\forall i \neq j$  or higher.

To continue, observe that for a fixed  $(i, j)$ ,  $i < j$ , the second-order term that involves  $\eta M_{t-1,i}$  and  $\eta M_{t-1,j}$  on (94) is with coefficient  $\frac{1}{\sqrt{m^{L-1}d_y}}\beta$ , because the first term on (94) contributes to  $\frac{1}{\sqrt{m^{L-1}d_y}}(L-1)\beta$ , while the third term on (94) contributes to  $-\frac{1}{\sqrt{m^{L-1}d_y}}(L-2)\beta$ . Furthermore, for a fixed  $(i, j, k)$ ,  $i < j < k$ , the third-order term that involves  $\eta M_{t-1,i}$ ,  $\eta M_{t-1,j}$ , and  $\eta M_{t-1,k}$  on (94) is with coefficient  $-2\frac{1}{\sqrt{m^{L-1}d_y}}\beta$ , as the first term on (94) contributes to  $-\frac{1}{\sqrt{m^{L-1}d_y}}(L-1)\beta$ , while the third term on (94) contributes to  $\frac{1}{\sqrt{m^{L-1}d_y}}(L-3)\beta$ . Similarly, for a  $p$ -order term  $\underbrace{\eta M_{t-1,*}, \dots, \eta M_{t-1,**}}_{p \text{ terms}}$ , the coefficient is  $(p-1)\frac{1}{\sqrt{m^{L-1}d_y}}\beta(-1)^p$ .

By induction (see (90)), we can bound the norm of the momentum at layer  $l$  as

$$\|M_{t-1,l}\|_F \leq \frac{4\|X\|_2}{\sqrt{d_y}} \frac{\theta^{t-1}}{1-\theta} \nu C_0 \|U_0 - Y\|_F. \tag{96}$$

Combining all the pieces together, we have

$$\begin{aligned}
 &\left\| \frac{1}{\sqrt{m^{L-1}d_y}}\beta(L-1)W_t^{(L:1)}X + \frac{1}{\sqrt{m^{L-1}d_y}}\beta W_{t-1}^{(L:1)}X - \frac{1}{\sqrt{m^{L-1}d_y}}\beta \sum_{l=1}^L W_t^{(L:l+1)}W_{t-1}^{(l)}W_t^{(l-1:1)}X \right\|_F \\
 &\stackrel{(a)}{\leq} \frac{\beta}{\sqrt{m^{L-1}d_y}} \sum_{j=2}^L (j-1) \binom{L}{j} \left( \eta \frac{4\|X\|_2}{\sqrt{d_y}} \frac{\theta^{t-1}}{1-\theta} \nu C_0 \|U_0 - Y\|_F \right)^j (1.1)^{j+1} m^{\frac{L-j}{2}} \|X\|_2 \\
 &\stackrel{(b)}{\leq} 1.1 \frac{\beta}{\sqrt{m^{L-1}d_y}} \sum_{j=2}^L L^j \left( \eta \frac{4.4\|X\|_2}{\sqrt{d_y}} \frac{\theta^{t-1}}{1-\theta} \nu C_0 \|U_0 - Y\|_F \right)^j m^{\frac{L-j}{2}} \|X\|_2 \\
 &\leq 1.1\beta \sqrt{\frac{m}{d_y}} \|X\|_2 \sum_{j=2}^L \left( \eta \frac{4.4L\|X\|_2}{\sqrt{md_y}} \frac{\theta^{t-1}}{1-\theta} \nu C_0 \|U_0 - Y\|_F \right)^j,
 \end{aligned} \tag{97}$$

where (a) uses (96), the above analysis of the coefficients of the higher-order terms and Lemma 16 for bounding a  $j \geq 2$  higher-order terms like  $\frac{1}{\sqrt{m^{L-1}d_y}}\beta(j-1)(-1)^j W_{t-1}^{(L:k_j+1)} \cdot (-\eta M_{t-1,k_j}) W_{t-1}^{(k_j-1:k_{j-1}+1)} \cdot (-\eta M_{t-1,k_{j-1}}) \dots (-\eta M_{t-1,k_1}) \cdot W_{t-1}^{(k_1-1:1)}$ , where  $1 \leq k_1 < \dots < k_j \leq L$  and (b) uses that  $\binom{L}{j} \leq \frac{L^j}{j!}$

Let us bound  $\eta \frac{4.4L\|X\|_2}{\sqrt{md_y}} \frac{\theta^{t-1}}{1-\theta} \nu C_0 \|U_0 - Y\|_F$  in the sum above. We have

$$\begin{aligned}
 \eta \frac{4.4L\|X\|_2}{\sqrt{md_y}} \frac{\theta^{t-1}}{1-\theta} \nu C_0 \|U_0 - Y\|_F &\leq 4.4 \sqrt{\frac{d_y}{m}} \frac{1}{\|X\|_2} \frac{\theta^{t-1}}{1-\theta} \nu C_0 \|U_0 - Y\|_F \\
 &\leq 0.5,
 \end{aligned} \tag{98}$$



where the last inequality uses that  $\tilde{C}_1 \frac{d_y B_0^2 C_0^2}{\|X\|_2^2} \frac{1}{(1-\theta)^2} \leq \tilde{C}_1 \frac{d_y B_0^2 C_0^2}{\|X\|_2^2} \frac{1}{\eta\lambda} \leq \tilde{C}_2 \frac{d_y B_0^2 \kappa^2}{\|X\|_2^2} \leq m$ , for some sufficiently large constant  $\tilde{C}_1, \tilde{C}_2 > 0$ . Combining the above results, i.e. (97) and (98), we have

$$\begin{aligned}
 \|\psi_t\| &\leq \left\| \frac{1}{\sqrt{m^{L-1}d_y}} \beta(L-1)W_t^{(L:1)}X + \frac{1}{\sqrt{m^{L-1}d_y}} \beta W_{t-1}^{(L:1)}X - \frac{1}{\sqrt{m^{L-1}d_y}} \beta \sum_{l=1}^L W_t^{(L:l+1)}W_{t-1}^{(l)}W_t^{(l-1:1)}X \right\|_F \\
 &\leq 1.1\beta \sqrt{\frac{m}{d_y}} \|X\|_2 \left( \eta \frac{4.4L\|X\|_2}{\sqrt{md_y}} \frac{\theta^{t-1}}{1-\theta} \nu C_0 \|U_0 - Y\|_F \right)^2 \sum_{j=2}^{L-2} (0.5)^{j-2} \\
 &\leq 2.2\beta \sqrt{\frac{m}{d_y}} \|X\|_2 \left( \eta \frac{4.4L\|X\|_2}{\sqrt{md_y}} \frac{\theta^{t-1}}{1-\theta} \nu C_0 \|U_0 - Y\|_F \right)^2 \\
 &\leq \frac{43\sqrt{d_y}}{\sqrt{m}\|X\|_2} \left( \frac{\theta^{t-1}}{1-\theta} \nu C_0 \|U_0 - Y\|_F \right)^2,
 \end{aligned} \tag{99}$$

where the last inequality uses  $\eta \leq \frac{d_y}{L\|X\|_2^2}$ .

Now let us switch to bound  $\|\iota_t\|$ . We have

$$\begin{aligned}
 \|\iota_t\| &= \|\eta(H_t - H_0)\xi_t\| \\
 &= \frac{\eta}{m^{L-1}d_y} \left\| \sum_{l=1}^L W_t^{(L:l+1)}(W_t^{(L:l+1)})^\top (U_t - Y)(W_t^{(l-1:1)}X)^\top W_t^{(l-1:1)}X \right. \\
 &\quad \left. - \sum_{l=1}^L W_0^{(L:l+1)}(W_0^{(L:l+1)})^\top (U_t - Y)(W_0^{(l-1:1)}X)^\top W_0^{(l-1:1)}X \right\|_F \\
 &\leq \frac{\eta}{m^{L-1}d_y} \sum_{l=1}^L \left\| W_t^{(L:l+1)}(W_t^{(L:l+1)})^\top (U_t - Y)(W_t^{(l-1:1)}X)^\top W_t^{(l-1:1)}X \right. \\
 &\quad \left. - W_0^{(L:l+1)}(W_0^{(L:l+1)})^\top (U_t - Y)(W_0^{(l-1:1)}X)^\top W_0^{(l-1:1)}X \right\|_F \\
 &\leq \frac{\eta}{m^{L-1}d_y} \sum_{l=1}^L \underbrace{\left\| \left( W_t^{(L:l+1)}(W_t^{(L:l+1)})^\top - W_0^{(L:l+1)}(W_0^{(L:l+1)})^\top \right) (U_t - Y)(W_t^{(l-1:1)}X)^\top W_t^{(l-1:1)}X \right\|_F}_{\text{first term}} \\
 &\quad + \underbrace{\left\| W_0^{(L:l+1)}(W_0^{(L:l+1)})^\top (U_t - Y) \left( W_t^{(l-1:1)}X \right)^\top W_t^{(l-1:1)}X - (W_0^{(l-1:1)}X)^\top W_0^{(l-1:1)}X \right\|_F}_{\text{second term}}.
 \end{aligned} \tag{100}$$

Now let us bound the first term. We have

$$\begin{aligned}
 &\underbrace{\left\| \left( W_t^{(L:l+1)}(W_t^{(L:l+1)})^\top - W_0^{(L:l+1)}(W_0^{(L:l+1)})^\top \right) (U_t - Y)(W_t^{(l-1:1)}X)^\top W_t^{(l-1:1)}X \right\|_F}_{\text{first term}} \\
 &\leq \|W_t^{(L:l+1)}(W_t^{(L:l+1)})^\top - W_0^{(L:l+1)}(W_0^{(L:l+1)})^\top\|_2 \|U_t - Y\|_F \|(W_t^{(l-1:1)}X)^\top W_t^{(l-1:1)}X\|_2.
 \end{aligned} \tag{101}$$

For  $\|(W_t^{(l-1:1)}X)^\top W_t^{(l-1:1)}X\|_2$ , by using Lemma 15 and Lemma 16, we have

$$\|(W_t^{(l-1:1)}X)^\top W_t^{(l-1:1)}X\|_2 \leq \left( \sigma_{\max}(W_t^{(l-1:1)}X) \right)^2 \leq \left( 1.1m^{\frac{l-1}{2}} \sigma_{\max}(X) \right)^2. \tag{102}$$

For  $\|W_t^{(L:l+1)}(W_t^{(L:l+1)})^\top - W_0^{(L:l+1)}(W_0^{(L:l+1)})^\top\|_2$ , denote  $W_t^{(L:l+1)} = W_0^{(L:l+1)} + \Delta_t^{(L:l+1)}$ , we have

$$\begin{aligned}
 & \|W_t^{(L:l+1)}(W_t^{(L:l+1)})^\top - W_0^{(L:l+1)}(W_0^{(L:l+1)})^\top\|_2 \\
 & \leq \|\Delta_t^{(L:l+1)}(W_t^{(L:l+1)})^\top + W_t^{(L:l+1)}(\Delta_t^{(L:l+1)})^\top + \Delta_t^{(L:l+1)}(\Delta_t^{(L:l+1)})^\top\|_2 \\
 & \leq 2\|\Delta_t^{(L:l+1)}\|_2 \cdot \sigma_{\max}(W_t^{(L:l+1)}) + \|\Delta_t^{(L:l+1)}\|_2^2 \\
 & \leq 2\|\Delta_t^{(L:l+1)}\|_2 \cdot \left(1.1m^{\frac{L-l}{2}}\right) + \|\Delta_t^{(L:l+1)}\|_2^2.
 \end{aligned} \tag{103}$$

Therefore, we have to bound  $\|\Delta_t^{(L:l+1)}\|_2$ . We have for any  $1 \leq i \leq j \leq L$ .

$$W_t^{(j:i)} = \left(W_0^{(j)} + \Delta_j\right) \cdots \left(W_0^{(i)} + \Delta_i\right), \tag{104}$$

where  $\|\Delta_i\|_2 \leq \|W_t^{(i)} - W_0^{(i)}\|_F \leq D := \frac{64\|X\|_2\sqrt{d_y}}{L\sigma_{\min}^2(X)}\nu C_0 B_0$  by Lemma 15. The product (104) above minus  $W_0^{(j:i)}$  can be written as a finite sum of some terms of the form

$$W_0^{(j:k_l+1)} \Delta_{k_l} W_0^{(k_l-1:k_{l-1}+1)} \Delta_{k_{l-1}} \cdots \Delta_{k_1} W_0^{(k_1-1:i)}, \tag{105}$$

where  $i \leq k_1 < \cdots < k_l \leq j$ . Recall that  $\|W_0^{(j':i')}\|_2 = m^{\frac{j'-i'+1}{2}}$  by Lemma 13. Thus, we can bound

$$\begin{aligned}
 \|\Delta_t^{(j:i)}\|_2 & \leq \|W_t^{(j:i)} - W_0^{(j:i)}\|_F \leq \sum_{l=1}^{j-i+1} \binom{j-i+1}{l} (D)^l m^{\frac{j-i+1-l}{2}} = (\sqrt{m} + D)^{j-i+1} - (\sqrt{m})^{j-i+1} \\
 & = (\sqrt{m})^{j-i+1} \left( (1 + D/\sqrt{m})^{j-i+1} - 1 \right) \leq (\sqrt{m})^{j-i+1} \left( (1 + D/\sqrt{m})^L - 1 \right) \\
 & \stackrel{(a)}{=} \left( 1 + \frac{1}{\sqrt{C'L\kappa}} \right)^L - 1 \left( \sqrt{m} \right)^{j-i+1} \stackrel{(b)}{\leq} \left( \exp\left(\frac{1}{\sqrt{C'L\kappa}}\right) - 1 \right) (\sqrt{m})^{j-i+1} \\
 & \stackrel{(c)}{\leq} \left( 1 + (e-1)\frac{1}{\sqrt{C'L\kappa}} - 1 \right) (\sqrt{m})^{j-i+1} \stackrel{(d)}{\leq} \frac{1}{480\kappa} (\sqrt{m})^{j-i+1},
 \end{aligned} \tag{106}$$

where (a) uses  $\frac{D}{\sqrt{m}} \leq \frac{1}{\sqrt{C'L\kappa}}$ , for some constant  $C' > 0$ , since  $C' \frac{d_y C_0^2 B_0^2 \kappa^4}{\|X\|_2^2} \leq C' \frac{d_y B_0^2 \kappa^5}{\|X\|_2^2} \leq m$ , (b) follows by the inequality  $(1 + x/n)^n \leq e^x, \forall x \geq 0, n > 0$ , (c) from Bernoulli's inequality  $e^r \leq 1 + (e-1)r, \forall 0 \leq r \leq 1$ , and (d) by choosing any sufficiently larger  $C'$ .

From (106), we have  $\|\Delta_t^{(L:l+1)}\|_2 \leq \frac{1}{480\kappa} (\sqrt{m})^{L-l}$ . Combining this with (101), (102), and (103), we have

$$\begin{aligned}
 & \underbrace{\| \left( W_t^{(L:l+1)}(W_t^{(L:l+1)})^\top - W_0^{(L:l+1)}(W_0^{(L:l+1)})^\top \right) (U_t - Y)(W_t^{(l-1:1)}X)^\top W_t^{(l-1:1)}X \|_F}_{\text{first term}} \\
 & \leq \left( 2\|\Delta_t^{(L:l+1)}\|_2 \cdot \left(1.1m^{\frac{L-l}{2}}\right) + \|\Delta_t^{(L:l+1)}\|_2^2 \right) \left(1.1m^{\frac{l-1}{2}} \sigma_{\max}(X)\right)^2 \|U_t - Y\|_F \\
 & \leq \left( 2\frac{1}{480\kappa} (\sqrt{m})^{L-l} \cdot \left(1.1m^{\frac{L-l}{2}}\right) + \left(\frac{1}{480\kappa} (\sqrt{m})^{L-l}\right)^2 \right) \left(1.1m^{\frac{l-1}{2}} \sigma_{\max}(X)\right)^2 \|U_t - Y\|_F \\
 & \leq \frac{\sigma_{\min}^2(X)}{160} m^{L-1} \|U_t - Y\|_F,
 \end{aligned} \tag{107}$$

where in the last inequality we use  $\kappa := \frac{\sigma_{\max}^2(X)}{\sigma_{\min}^2(X)}$ .

Now let us switch to bound the second term, we have

$$\begin{aligned}
 & \underbrace{\| (W_0^{(L:l+1)}(W_0^{(L:l+1)})^\top)^\top (U_t - Y) \left( W_t^{(l-1:1)}X \right)^\top W_t^{(l-1:1)}X - (W_0^{(l-1:1)}X)^\top W_0^{(l-1:1)}X \|_F}_{\text{second term}} \\
 & \leq \| (W_0^{(L:l+1)}(W_0^{(L:l+1)})^\top)^\top \|_2 \|U_t - Y\|_F \| (W_t^{(l-1:1)}X)^\top W_t^{(l-1:1)}X - (W_0^{(l-1:1)}X)^\top W_0^{(l-1:1)}X \|_2.
 \end{aligned} \tag{108}$$

For  $\|W_0^{(L:l+1)}(W_0^{(L:l+1)})^\top\|_2$ , based on Lemma 13, we have

$$\|W_0^{(L:l+1)}(W_0^{(L:l+1)})^\top\|_2 \leq m^{L-l}. \quad (109)$$

To bound  $\|(W_t^{(l-1:1)}X)^\top W_t^{(l-1:1)}X - (W_0^{(l-1:1)}X)^\top W_0^{(l-1:1)}X\|_2$ , we proceed as follows. Denote  $W_t^{(l-1:1)} = W_0^{(l-1:1)} + \Delta_t^{(l-1:1)}$ , we have

$$\begin{aligned} & \|(W_t^{(l-1:1)}X)^\top W_t^{(l-1:1)}X - (W_0^{(l-1:1)}X)^\top W_0^{(l-1:1)}X\|_2 \\ & \leq 2\|(\Delta_t^{(l-1:1)}X)^\top W_t^{(l-1:1)}X\|_2 + \|\Delta_t^{(l-1:1)}X\|_2^2 \\ & \leq \left(2\|\Delta_t^{(l-1:1)}\| \|W_t^{(l-1:1)}\|_2 + \|\Delta_t^{(l-1:1)}\|_2^2\right) \|X\|_2^2 \\ & \leq \left(2\frac{1}{480\kappa}m^{\frac{l-1}{2}}1.1m^{\frac{l-1}{2}} + \left(\frac{1}{480\kappa}m^{\frac{l-1}{2}}\right)^2\right) \|X\|_2^2 \\ & \leq \frac{\sigma_{\min}^2(X)}{160}m^{l-1}, \end{aligned} \quad (110)$$

where the second to last inequality uses (106), Lemma 15, and Lemma 16, while the last inequality uses  $\kappa := \frac{\sigma_{\max}^2(X)}{\sigma_{\min}^2(X)}$ . Combining (108), (109), (110), we have

$$\begin{aligned} & \underbrace{\|(W_0^{(L:l+1)}(W_0^{(L:l+1)})^\top(U_t - Y) (W_t^{(l-1:1)}X)^\top W_t^{(l-1:1)}X - (W_0^{(l-1:1)}X)^\top W_0^{(l-1:1)}X)\|_F}_{\text{second term}} \\ & \leq \frac{\sigma_{\min}^2(X)}{160}m^{L-1}\|U_t - Y\|_F. \end{aligned} \quad (111)$$

Now combing (100), (107), and (111), we have

$$\|t_t\| \leq \frac{\eta}{m^{L-1}d_y}L\frac{\sigma_{\min}^2(X)}{80}m^{L-1}\|U_t - Y\|_F = \frac{\eta\lambda}{80}\|\xi_t\|, \quad (112)$$

where we use  $\lambda := \frac{L\sigma_{\min}^2(X)}{d_y}$ .

Now we have (93), (99), and (112), which leads to

$$\begin{aligned} \|\varphi_t\| & \leq \|\phi_t\| + \|\psi_t\| + \|t_t\| \\ & \leq \frac{43\sqrt{d_y}}{\sqrt{m}\|X\|_2}(\theta^{2t} + \theta^{2(t-1)})\nu^2C_0^2\left(\frac{\|\xi_0\|}{1-\theta}\right)^2 + \frac{\eta\lambda}{80}\nu C_0\left\|\begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix}\right\| \\ & \leq \frac{1920\sqrt{d_y}}{\sqrt{m}\|X\|_2}\frac{1}{\eta\lambda}\theta^{2t}\nu^2C_0^2\left\|\begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix}\right\|^2 + \frac{\eta\lambda}{80}\nu C_0\left\|\begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix}\right\|. \end{aligned} \quad (113)$$

where the last inequality uses that  $1 \leq \frac{16}{9}\theta^2$  as  $\eta\lambda \leq 1$  so that  $\theta \geq \frac{3}{4}$ .  $\square$

**Lemma 15.** *Following the setting as Theorem 10, denote  $\theta := \beta_* + \frac{1}{4}\sqrt{\eta\lambda} = 1 - \frac{1}{4}\sqrt{\eta\lambda}$ . If for any  $s \leq t$ , the residual dynamics satisfies  $\left\|\begin{bmatrix} \xi_s \\ \xi_{s-1} \end{bmatrix}\right\| \leq \theta^s \cdot \nu C_0 \left\|\begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix}\right\|$ , for some constant  $\nu > 0$ , then*

$$\|W_t^{(l)} - W_0^{(l)}\|_F \leq R^{L\text{-linear}} := \frac{64\|X\|_2\sqrt{d_y}}{L\sigma_{\min}^2(X)}\nu C_0 B_0.$$

*Proof.* We have

$$\begin{aligned}
 \|W_{t+1}^{(l)} - W_0^{(l)}\|_F &\stackrel{(a)}{\leq} \eta \sum_{s=0}^t \|M_{s,l}\|_F \stackrel{(b)}{=} \eta \sum_{s=0}^t \left\| \sum_{\tau=0}^s \beta^{s-\tau} \frac{\partial \ell(W_\tau^{(L:1)})}{\partial W_\tau^{(L)}} \right\|_F \leq \eta \sum_{s=0}^t \sum_{\tau=0}^s \beta^{s-\tau} \left\| \frac{\partial \ell(W_\tau^{(L:1)})}{\partial W_\tau^{(L)}} \right\|_F \\
 &\stackrel{(c)}{\leq} \eta \sum_{s=0}^t \sum_{\tau=0}^s \beta_*^{2(s-\tau)} \frac{4\|X\|_2}{\sqrt{d_y}} \theta^\tau \nu C_0 \|U_0 - Y\|_F. \\
 &\stackrel{(d)}{\leq} \eta \sum_{s=0}^t \frac{\theta^s}{1-\theta} \frac{4\|X\|_2}{\sqrt{d_y}} \nu C_0 \|U_0 - Y\|_F. \\
 &\leq \frac{4\eta\|X\|_2}{\sqrt{d_y}} \frac{1}{(1-\theta)(1-\theta)} \nu C_0 \|U_0 - Y\|_F \\
 &\stackrel{(e)}{\leq} \frac{64\|X\|_2}{\lambda\sqrt{d_y}} \nu C_0 \|U_0 - Y\|_F \\
 &\stackrel{(f)}{\leq} \frac{64\|X\|_2 \sqrt{d_y}}{L\sigma_{\min}^2(X)} \nu C_0 B_0, \tag{114}
 \end{aligned}$$

where (a), (b) is by the update rule of momentum, which is  $W_{t+1}^{(l)} - W_t^{(l)} = -\eta M_{t,l}$ , where  $M_{t,l} := \sum_{s=0}^t \beta^{t-s} \frac{\partial \ell(W_s^{(L:1)})}{\partial W_s^{(L)}}$ , (c) is because  $\left\| \frac{\partial \ell(W_s^{(L:1)})}{\partial W_s^{(L)}} \right\|_F = \frac{4\|X\|_2}{\sqrt{d_y}} \theta^s \nu C_0 \|U_0 - Y\|_F$  (see (89)), (d) is because that  $\beta = \beta_*^2 \leq \theta^2$ , (e) is because that  $\frac{1}{(1-\theta)^2} = \frac{16}{\eta\lambda}$ , and (f) uses the upper-bound  $B_0 \geq \|U_0 - Y\|$  defined in Lemma 13 and  $\lambda := \frac{L\sigma_{\min}^2(X)}{d_y}$ . The proof is completed.  $\square$

**Lemma 16.** (Hu et al., 2020b) Let  $R^{L\text{-linear}}$  be an upper bound that satisfies  $\|W_t^{(l)} - W_t^{(l)}\|_F \leq R^{L\text{-linear}}$  for all  $l$  and  $t$ . Suppose the width  $m$  satisfies  $m > C(LR^{L\text{-linear}})^2$ , where  $C$  is any sufficiently large constant. Then,

$$\sigma_{\max}(W_t^{(j:i)}) \leq 1.1m^{\frac{j-i+1}{2}}, \quad \sigma_{\min}(W_t^{(j:i)}) \geq 0.9m^{\frac{j-i+1}{2}}.$$

*Proof.* The lemma has been proved in proof of Claim 4.4 and Claim 4.5 in (Hu et al., 2020b). For completeness, let us replicate the proof here.

We have for any  $1 \leq i \leq j \leq L$ .

$$W_t^{(j:i)} = \left(W_0^{(j)} + \Delta_j\right) \cdots \left(W_0^{(i)} + \Delta_i\right), \tag{115}$$

where  $\Delta_i = W_t^{(i)} - W_0^{(i)}$ . The product above minus  $W_0^{(j:i)}$  can be written as a finite sum of some terms of the form

$$W_0^{(j:k_l+1)} \Delta_{k_l} W_0^{(k_l-1:k_{l-1}+1)} \Delta_{k_{l-1}} \cdots \Delta_{k_1} W_0^{(k_1-1:i)}, \tag{116}$$

where  $i \leq k_1 < \cdots < k_l \leq j$ . Recall that  $\|W_0^{(j':i')}\|_2 = m^{\frac{j'-i'+1}{2}}$ . Thus, we can bound

$$\begin{aligned}
 \|W_t^{(j:i)} - W_0^{(j:i)}\|_F &\leq \sum_{l=1}^{j-i+1} \binom{j-i+1}{l} (R^{L\text{-linear}})^l m^{\frac{j-i+1-l}{2}} = (\sqrt{m} + R^{L\text{-linear}})^{j-i+1} - (\sqrt{m})^{j-i+1} \\
 &= (\sqrt{m})^{j-i+1} \left( (1 + R^{L\text{-linear}}/\sqrt{m})^{j-i+1} - 1 \right) \leq (\sqrt{m})^{j-i+1} \left( (1 + R^{L\text{-linear}}/\sqrt{m})^L - 1 \right) \\
 &\leq 0.1(\sqrt{m})^{j-i+1}, \tag{117}
 \end{aligned}$$

where the last step uses  $m > C(LR^{L\text{-linear}})^2$ . By combining this with Lemma 13, one can obtain the result.  $\square$

**Remark:** In the proof of Lemma 14, we obtain a tighter bound of the distance  $\|W_t^{(j:i)} - W_0^{(j:i)}\|_F \leq O\left(\frac{1}{\kappa}(\sqrt{m})^{j-i+1}\right)$ . However, to get the upper-bound  $\sigma_{\max}(W_t^{(j:i)})$  shown in Lemma 16, (117) is sufficient for the purpose.

**G.1. Proof of Theorem 10**

*Proof.* (of Theorem 10) Denote  $\lambda := L\sigma_{\min}^2(X)/d_y$ . By Lemma 13,  $\lambda_{\min}(H) \geq \lambda$ . Also, denote  $\beta_* := 1 - \frac{1}{2}\sqrt{\eta\lambda}$  and  $\theta := \beta_* + \frac{1}{4}\sqrt{\eta\lambda} = 1 - \frac{1}{4}\sqrt{\eta\lambda}$ . Let  $\nu = 2$  in Lemma 14, 15, and let  $C_1 = C_3 = C_0$  and  $C_2 = \frac{1}{4}\sqrt{\eta\lambda}$  in Theorem 6. The goal is to show that  $\left\| \begin{bmatrix} \xi_t \\ \xi_{t-1} \end{bmatrix} \right\| \leq \theta^t 2C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|$  for all  $t$  by induction. To achieve this, we will also use induction to show that for all iterations  $s$ ,

$$\forall l \in [L], \|W_t^{(l)} - W_0^{(l)}\| \leq R^{L\text{-linear}} := \frac{64\|X\|_2\sqrt{d_y}}{L\sigma_{\min}^2(X)} C_0 B_0, \quad (118)$$

which is clearly true in the base case  $s = 0$ .

By Lemma 4, 13, 14, 15, Theorem 6 and Corollary 1, it suffices to show that  $\left\| \begin{bmatrix} \xi_s \\ \xi_{s-1} \end{bmatrix} \right\| \leq \theta^s \cdot 2C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|$  and  $\forall l \in [L], \|W_s^{(l)} - W_0^{(l)}\| \leq R^{L\text{-linear}}$  hold at  $s = 0, 1, \dots, t-1$ , one has

$$\left\| \sum_{s=0}^{t-1} A^{t-s-1} \begin{bmatrix} \varphi_s \\ 0 \end{bmatrix} \right\| \leq \theta^t C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|, \quad (119)$$

$$\forall l \in [L], \|W_t^{(l)} - W_0^{(l)}\| \leq R^{L\text{-linear}} := \frac{64\|X\|_2\sqrt{d_y}}{L\sigma_{\min}^2(X)} C_0 B_0, \quad (120)$$

where the matrix  $A$  and the vector  $\varphi_t$  are defined in Lemma 4, and  $B_0$  is a constant such that  $B_0 \geq \|Y - U_0\|_F$  with probability  $1 - \delta$  by Lemma 13. The inequality (119) is the required condition for using the result of Theorem 6, while the inequality (120) helps us to show (119) through invoking Lemma 14 to bound the terms  $\{\varphi_s\}$  as shown in the following.

Let us show (119) first. We have

$$\begin{aligned} \left\| \sum_{s=0}^{t-1} A^{t-1-s} \begin{bmatrix} \varphi_s \\ 0 \end{bmatrix} \right\| &\stackrel{(a)}{\leq} \sum_{s=0}^{t-1} \beta_*^{t-1-s} C_0 \|\varphi_s\| \\ &\stackrel{(b)}{\leq} \frac{1920\sqrt{d_y}}{\sqrt{m}\|X\|_2} \frac{1}{\eta\lambda} \sum_{s=0}^{t-1} \beta_*^{t-1-s} \theta^{2s} 4C_0^3 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|^2 + \sum_{s=0}^{t-1} \beta_*^{t-1-s} \frac{\eta\lambda}{80} \theta^{2s} 2C_0^2 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\| \\ &\stackrel{(c)}{\leq} \frac{1920\sqrt{d_y}}{\sqrt{m}\|X\|_2} \frac{1}{\eta\lambda} \sum_{s=0}^{t-1} \beta_*^{t-1-s} \theta^{2s} 4C_0^3 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|^2 + \frac{2\sqrt{\eta\lambda}}{15} \theta^t C_0^2 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\| \\ &\stackrel{(d)}{\leq} \frac{1920\sqrt{d_y}}{\sqrt{m}\|X\|_2} \frac{16}{3(\eta\lambda)^{3/2}} \theta^t 4C_0^3 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|^2 + \frac{2\sqrt{\eta\lambda}}{15} \theta^t C_0^2 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\| \\ &\stackrel{(e)}{\leq} \frac{1}{3} \theta^t C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\| + \frac{2\sqrt{\eta\lambda}}{15} \theta^t C_0^2 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\| \\ &\stackrel{(f)}{\leq} \theta^t C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|, \end{aligned} \quad (121)$$

where (a) uses Theorem 5 with  $\beta = \beta_*^2$ , (b) is by Lemma 14, (c) uses  $\sum_{s=0}^{t-1} \beta_*^{t-1-s} \theta^{2s} = \theta^{t-1} \sum_{s=0}^{t-1} \left(\frac{\beta_*}{\theta}\right)^{t-1-s} \leq \theta^{t-1} \sum_{s=0}^{t-1} \theta^{t-1-s} \leq \theta^{t-1} \frac{4}{\sqrt{\eta\lambda}} \leq \theta^t \frac{16}{3\sqrt{\eta\lambda}}$ ,  $\beta_* = 1 - \frac{1}{2}\sqrt{\eta\lambda} \geq \frac{1}{2}$ , and  $\theta = 1 - \frac{1}{4}\sqrt{\eta\lambda} \geq \frac{3}{4}$ , (d) uses  $\sum_{s=0}^{t-1} \beta_*^{t-1-s} \theta^{2s} \leq \sum_{s=0}^{t-1} \theta^{t-1+s} \leq \frac{\theta^{t-1}}{1-\theta} \leq \theta^t \frac{16}{3\sqrt{\eta\lambda}}$ , (e) is because  $C' \frac{d_y C_0^4 B_0^2}{\|X\|_2^2} \frac{1}{(\eta\lambda)^3} \leq C \frac{d_y B_0^2}{\|X\|_2^2} \kappa^5 \leq m$  for some sufficiently large constants  $C', C > 0$ , and (f) uses that  $\eta\lambda = \frac{1}{\kappa}$  and  $C_0 \leq 4\sqrt{\kappa}$  by Corollary 1. Hence, we have shown (119). Therefore, by Theorem 6, we have  $\left\| \begin{bmatrix} \xi_t \\ \xi_{t-1} \end{bmatrix} \right\| \leq \theta^t 2C_0 \left\| \begin{bmatrix} \xi_0 \\ \xi_{-1} \end{bmatrix} \right\|$ .

By Lemma 15, we have (120). Thus, we have completed the proof.  $\square$

## H. Experiment

### H.1. ReLU network

We report a proof-of-concept experiment for training the ReLU network. We sample  $n = 5$  points from the normal distribution, and then scale the size to the unit norm. We generate the labels uniformly random from  $\{1, -1\}$ . We let  $m = 1000$  and  $d = 10$ . We compare vanilla GD and gradient descent with Polyak’s momentum. Denote  $\hat{\lambda}_{\max} := \lambda_{\max}(H_0)$ ,  $\hat{\lambda}_{\min} := \lambda_{\min}(H_0)$ , and  $\hat{\kappa} := \hat{\lambda}_{\max}/\hat{\lambda}_{\min}$ . Then, for gradient descent with Polyak’s momentum, we set the step size  $\eta = 1/\left(\hat{\lambda}_{\max}\right)$  and set the momentum parameter  $\beta = \left(1 - \frac{1}{2} \frac{1}{\sqrt{\hat{\kappa}}}\right)^2$ . For gradient descent, we set the same step size. The result is shown on Figure 1.

We also report the percentiles of pattern changes over iterations. Specifically, we report the quantity

$$\frac{\sum_{i=1}^n \sum_{r=1}^m \mathbb{1}\{\text{sign}(x_i^\top w_t^{(r)}) \neq \text{sign}(x_i^\top w_0^{(r)})\}}{mn},$$

as there are  $mn$  patterns. For gradient descent with Polyak’s momentum, the percentiles of pattern changes is approximately 0.76%; while for vanilla gradient descent, the percentiles of pattern changes is 0.55%.

### H.2. Deep linear network

We let the input and output dimension  $d = d_y = 20$ , the width of the intermediate layers  $m = 50$ , the depth  $L = 100$ . We sampled a  $X \in \mathbb{R}^{20 \times 5}$  from the normal distribution. We let  $W^* = I_{20} + 0.1\bar{W}$ , where  $\bar{W} \in \mathbb{R}^{20 \times 20}$  is sampled from the normal distribution. Then, we have  $Y = W^*X$ ,  $\eta = \frac{d_y}{L\sigma_{\max}^2(X)}$  and  $\beta = \left(1 - \frac{1}{2}\sqrt{\eta\lambda}\right)^2$ , where  $\lambda = \frac{L\sigma_{\min}^2(X)}{d_y}$ . Vanilla GD also uses the same step size. The network is initialized by the orthogonal initialization and both algorithms start from the same initialization. The result is shown on Figure 2.

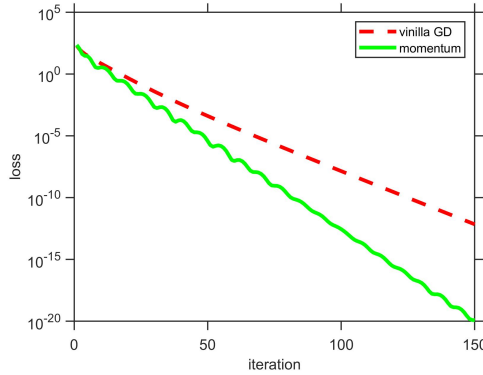


Figure 2. Training a 100-layer deep linear network. Here “momentum” stands for gradient descent with Polyak’s momentum.