




# Convergence Analysis for Learning Orthonormal Deep Linear Neural Networks

Zhen Qin , Xuwei Tan , and Zhihui Zhu , *Member, IEEE*

**Abstract**—Enforcing orthonormal or isometric property for the weight matrices has been shown to enhance the training of deep neural networks by mitigating gradient exploding/vanishing and increasing the robustness of the learned networks. However, despite its practical performance, the theoretical analysis of orthonormality in neural networks is still lacking; for example, how orthonormality affects the convergence of the training process. In this letter, we aim to bridge this gap by providing convergence analysis for training orthonormal deep linear neural networks. Specifically, we show that Riemannian gradient descent with an appropriate initialization converges at a linear rate for training orthonormal deep linear neural networks with a class of loss functions. Unlike existing works that enforce orthonormal weight matrices for all the layers, our approach excludes this requirement for one layer, which is crucial to establish the convergence guarantee. Our results shed light on how increasing the number of hidden layers can impact the convergence speed. Experimental results validate our theoretical analysis.

**Index Terms**—Convergence analysis, deep neural networks, orthonormal structure, Riemannian optimization.

## I. INTRODUCTION

**E**NCORING orthonormal or isometric properties of the weight matrices has numerous advantages for the practice of deep learning: (i) it provides a better initialization [1], [2], (ii) it mitigates the problem of exploding/vanishing gradients during training [3], [4], [5], [6], [7], (iii) the resulting *orthonormal neural networks* [2], [8], [9], [10], [11], [12], [13], [14], [15], [16] exhibit improved robustness [17] and reduced overfitting issues [18].

Various approaches have been proposed for training neural networks, mainly falling into two categories: soft orthonormality and hard orthonormality. The first category of methods, such as those in [10], [12], [17], [19], adds an additional orthonormality regularization term to the training loss, resulting in weight matrices that are approximately orthonormal. In contrast, the other methods, as found in [9], [11], [20], learn weight matrices that are exactly orthonormal through the use of Riemannian optimization algorithms on the Stiefel manifold.

While orthonormal neural networks demonstrate strong practical performance, there remains a gap in the theoretical analysis

of orthonormality in neural networks. For example, convergence analysis for training neural networks has been extensively studied [21], [22], [23], [24], [25], [26], [27], [28]. However, all these results focus on standard training without orthonormal constraints, making them inapplicable to the training of orthonormal neural networks. To the best of our knowledge, there is a lack of rigorous convergence analysis even for orthonormal deep linear neural networks (ODLNNs). Despite its linear structure, a deep linear neural network still presents a non-convex training problem and has served as a testbed for understanding deep neural networks [21], [22], [23], [24]. In this letter, we aim to understand the effect of the orthonormal structure on the training process by studying ODLNN.

**Our contribution:** Specifically, we provide a local convergence rate of Riemannian gradient descent (RGD) for training the ODLNN. To achieve this, unlike existing works [9], [11], [17] that impose orthonormal constraints on all the weight matrices, we exclude such a constraint for one layer (say the weight matrix in the first hidden layer). The exclusion of a specific layer plays a crucial role in analyzing the convergence rate. Our findings demonstrate that within a specific class of loss functions, adhering to the restricted correlated gradient condition [29], the RGD algorithm exhibits linear convergence speed when appropriately initialized. Notably, our results also indicate that as the number of layers in the network increases, the rate of convergence only experiences a polynomial decrease. The validity of our theoretical analysis has been confirmed by experiments.

**Notation:** We use bold capital letters (e.g.,  $\mathbf{A}$ ) to denote matrices, bold lowercase letters (e.g.,  $\mathbf{a}$ ) to denote vectors, and italic letters (e.g.,  $a$ ) to denote scalar quantities. The superscript  $(\cdot)^\top$  denotes the transpose.  $\|\mathbf{A}\|$  and  $\|\mathbf{A}\|_F$  respectively represent the spectral norm and Frobenius norm of  $\mathbf{A}$ .  $\sigma_{\min}(\mathbf{A})$  is the smallest singular value of  $\mathbf{A}$ . The condition number of  $\mathbf{A}$  is defined as  $\kappa(\mathbf{A}) = \frac{\|\mathbf{A}\|}{\sigma_{\min}(\mathbf{A})}$ .  $\|\mathbf{a}\|_2$  is the  $l_2$  norm of  $\mathbf{a}$ . For a positive integer  $K$ ,  $[K]$  denotes the set  $\{1, \dots, K\}$ .  $b = \Omega(a)$  represents  $b \geq ca$  for some universal constant  $c$ .

## II. RIEMANNIAN GRADIENT DESCENT FOR ORTHONORMAL DEEP LINEAR NEURAL NETWORKS

**Problem statement:** Given a training set  $\{(\mathbf{x}_i, \mathbf{y}_i^*)\}_{i=1}^n \in \mathbb{R}^{d_x} \times \mathbb{R}^{d_y}$ , our goal is to estimate a hypothesis (predictor) from a parametric family  $\mathcal{H} := \{h_\theta : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_y} | \theta \in \Theta\}$  by minimizing the following empirical risk:

$$\min_{\theta \in \Theta} g(\theta) = \frac{1}{n} \sum_{i=1}^n l(h_\theta(\mathbf{x}_i); \mathbf{y}_i^*), \quad (1)$$

where  $l(h_\theta(\mathbf{x}_i); \mathbf{y}_i^*)$  is a suitable loss that captures the difference between the network prediction  $h_\theta(\mathbf{x}_i)$  and the label  $\mathbf{y}_i^*$ . For

Manuscript received 24 November 2023; revised 23 February 2024; accepted 28 February 2024. Date of publication 6 March 2024; date of current version 21 March 2024. This work was supported by NSF under Grant CCF-2240708 and Grant CCF-2241298. The associate editor coordinating the review of this manuscript and approving it for publication was Prof. Jianxin Li. (*Corresponding author: Zhihui Zhu.*)

The authors are with the Department of Computer Science and Engineering, Ohio State University, Columbus, OH 43210 USA (e-mail: qin.660@osu.edu; tan.1206@osu.edu; zhu.3440@osu.edu).

Digital Object Identifier 10.1109/LSP.2024.3374085

convenience, we stack all the training samples together as  $\mathbf{X} = [\mathbf{x}_1 \cdots \mathbf{x}_n]$  and  $\mathbf{Y}^* = [\mathbf{y}_1^* \cdots \mathbf{y}_n^*]$ .

Our main focus is on orthonormal deep linear neural networks (ODLNNs), which are fully-connected neural networks of form  $h_\theta(\mathbf{x}_i) = \mathbf{W}_N \cdots \mathbf{W}_1 \mathbf{x}_i$  with  $\mathbf{W}_i \in \mathbb{R}^{d_i \times d_{i-1}}$  for  $i \in [N]$ , where  $d_0 = d_x$  and  $d_N = d_y$ . In ODLNNs, we further assume that the weight matrices to be row orthogonal or column orthogonal depending on the dimension. Without loss of generality, we assume that all the matrices  $\{\mathbf{W}_i\}_{i \geq 2}$  are column orthonormal, except for  $\mathbf{W}_1$ . This is different to the previous works [9], [11], [17] which impose orthonormal constraints on all the weight matrices. Allowing  $\mathbf{W}_1$  to be unstructured offers more flexibility, as otherwise  $\mathbf{W}_N \cdots \mathbf{W}_1$  can only represent an orthonormal matrix, which restricts the output the same energy as input (i.e.,  $\|\mathbf{y}_i^*\|_2 = \|\mathbf{x}_i\|_2$ ). The choice of free weight matrix can vary, and the following analysis would still hold. Now the training loss can be written as

$$\min_{\substack{\mathbf{W}_i \in \mathbb{R}^{d_i \times d_{i-1}} \\ i \in [N]}} g(\mathbf{W}_N, \dots, \mathbf{W}_1) = L(\mathbf{W}_N \cdots \mathbf{W}_1 \mathbf{X}; \mathbf{Y}^*),$$

$$\text{s. t. } \mathbf{W}_i^\top \mathbf{W}_i = \mathbf{I}_{d_{i-1}}, \quad i = 2, \dots, N, \quad (2)$$

where  $L$  denotes a loss function encompassing all samples.

**Definition 1 (Data model):** Following the previous work on deep linear neural networks [21], [22], we assume that the dataset  $\mathbf{X}$  is whitened, i.e., its empirical covariance matrix is an identity matrix as  $\mathbf{X}\mathbf{X}^\top = \mathbf{I}_{d_x}$ . Also assume that the output is generated by a teacher ODLNN model,<sup>1</sup> i.e.,  $\mathbf{Y}^* = \mathbf{W}_N^* \cdots \mathbf{W}_1^* \mathbf{X}$ , where  $\mathbf{W}_i^* \in \mathbb{R}^{d_i \times d_{i-1}}$  and  $\{\mathbf{W}_i^*\}_{i \geq 2}$  are column orthonormal matrices.

**Stiefel manifold:** The Stiefel manifold  $\text{St}(m, n) = \{\mathbf{C} \in \mathbb{R}^{m \times n} : \mathbf{C}^\top \mathbf{C} = \mathbf{I}_n\}$  is a Riemannian manifold that is composed of all  $m \times n$  orthonormal matrices. We can regard  $\text{St}(m, n)$  as an embedded submanifold of a Euclidean space and further define  $\text{T}_{\mathbf{C}}\text{St} := \{\mathbf{A} \in \mathbb{R}^{m \times n} : \mathbf{A}^\top \mathbf{C} + \mathbf{C}^\top \mathbf{A} = \mathbf{0}\}$  as its tangent space at the point  $\mathbf{C} \in \text{St}(m, n)$ . For any  $\mathbf{B} \in \mathbb{R}^{m \times n}$ , the projection of  $\mathbf{B}$  onto  $\text{T}_{\mathbf{C}}\text{St}$  is given by [30]

$$\mathcal{P}_{\text{T}_{\mathbf{C}}\text{St}}(\mathbf{B}) = \mathbf{B} - \frac{1}{2} \mathbf{C}(\mathbf{B}^\top \mathbf{C} + \mathbf{C}^\top \mathbf{B}), \quad (3)$$

and its orthogonal complement is  $\mathcal{P}_{\text{T}_{\mathbf{C}}\text{St}}^\perp(\mathbf{B}) = \mathbf{B} - \mathcal{P}_{\text{T}_{\mathbf{C}}\text{St}}(\mathbf{B}) = \frac{1}{2} \mathbf{C}(\mathbf{B}^\top \mathbf{C} + \mathbf{C}^\top \mathbf{B})$ . When we have a gradient  $\mathbf{B}$  defined in the Hilbert space, we can use the projection operator (3) to compute the Riemannian gradient  $\mathcal{P}_{\text{T}_{\mathbf{C}}\text{St}}(\mathbf{B})$  on the tangent space of the Stiefel manifold. To project  $\hat{\mathbf{C}} = \mathbf{C} - c\mathcal{P}_{\text{T}_{\mathbf{C}}\text{St}}(\mathbf{B})$  with any positive constant  $c$  back onto the Stiefel manifold, we can utilize the polar decomposition-based retraction, i.e.,

$$\text{Retr}_{\mathbf{C}}(\hat{\mathbf{C}}) = \hat{\mathbf{C}}(\hat{\mathbf{C}}^\top \hat{\mathbf{C}})^{-\frac{1}{2}}. \quad (4)$$

**Riemannian gradient descent (RGD):** Given the gradient  $\nabla_{\mathbf{W}_i} g(\mathbf{W}_N^{(t)}, \dots, \mathbf{W}_1^{(t)})$ , we can compute the Riemannian gradient  $\mathcal{P}_{\text{T}_{\mathbf{W}_i}\text{St}}(\nabla_{\mathbf{W}_i} g(\mathbf{W}_N^{(t)}, \dots, \mathbf{W}_1^{(t)}))$  on the Stiefel manifold via (3). To streamline the notation, let us represent  $\nabla_{\mathbf{W}_i} g(\mathbf{W}_N^{(t)}, \dots, \mathbf{W}_1^{(t)})$  as  $\nabla_{\mathbf{W}_i^{(t)}} g$ . Now the weight matrices can be updated via the following RGD:

$$\mathbf{W}_1^{(t+1)} = \mathbf{W}_1^{(t)} - \mu \gamma \nabla_{\mathbf{W}_1^{(t)}} g,$$

$$\mathbf{W}_i^{(t+1)} = \text{Retr}_{\mathbf{W}_i}(\mathbf{W}_i^{(t)} - \mu \mathcal{P}_{\text{T}_{\mathbf{W}_i}\text{St}}(\nabla_{\mathbf{W}_i^{(t)}} g)), i \geq 2, \quad (5)$$

<sup>1</sup>Here, for the sake of simplifying the subsequent analysis, we designate the teacher model as the ODLNN, which can encompass any linear model.

where  $\mu > 0$  is the learning rate for  $\{\mathbf{W}_i\}$  and  $\gamma > 0$  controls the ratio between the learning rates for  $\mathbf{W}_1$  and  $\{\mathbf{W}_i\}_{i \geq 2}$ . The discrepant learning rates in (5) are used to accelerate the convergence rate of  $\mathbf{W}_1$  since the energy of  $\mathbf{W}_1^*$  and  $\{\mathbf{W}_i^*\}_{i \geq 2}$  are unbalanced, i.e.,  $\|\mathbf{W}_1^*\|^2 = \|\mathbf{Y}^*\|^2$  and  $\|\mathbf{W}_i^*\|^2 = 1$  when the dataset  $\mathbf{X}$  is whitened, i.e.  $\mathbf{X}\mathbf{X}^\top = \mathbf{I}_{d_x}$ .

### III. CONVERGENCE ANALYSIS

In this section, we will delve into the convergence rate analysis of RGD for training ODLNNs. Towards that goal, we will use the teacher model introduced in Definition 1 that the training samples  $(\mathbf{X}, \mathbf{Y}^*)$  are generated according to  $\mathbf{Y}^* = \mathbf{W}_N^* \cdots \mathbf{W}_1^* \mathbf{X}$ . Given the nonlinear nature of the retraction operation in the RGD, we will study the convergence in terms of the weight matrices  $\mathbf{W} = \{\mathbf{W}_i\}$  and  $\mathbf{W}^* = \{\mathbf{W}_i^*\}$ . But we will show that the convergence can be equivalently established in terms of the outputs.

**Distance measure:** Consider that the factors in  $\mathbf{W}^*$  are identifiable up to orthonormal transforms since  $\mathbf{W}_N^* \cdots \mathbf{W}_1^* \mathbf{X} = \mathbf{W}_N^* \mathbf{R}_{N-1} \mathbf{R}_{N-1}^\top \mathbf{W}_{N-1}^* \mathbf{R}_{N-2} \cdots \mathbf{W}_1^* \mathbf{X}$  for any orthonormal matrices  $\mathbf{R}_i \in \mathbb{O}^{d_i \times d_i}$ ,  $i \in [N-1]$ . Also,  $\mathbf{W}_1^*$  and  $\mathbf{W}_i^*$  could be imbalanced as  $\mathbf{W}_i^*$  are orthonormal for  $i \geq 2$ . This discrepancy can be quantified by observing that  $\|\mathbf{W}_1^*\| = \|\mathbf{Y}^*\|$  (since the input matrix  $\mathbf{X}$  is whitened) and  $\|\mathbf{W}_i^*\| = 1$  for all  $i \geq 2$ . Thus, we propose the following measure to capture the distance between two sets of factors:

$$\text{dist}^2(\mathbf{W}, \mathbf{W}^*) = \min_{\substack{\mathbf{R}_i \in \mathbb{O}^{d_i \times d_i} \\ i \in [N-1]}} \sum_{i=2}^N \|\mathbf{Y}^*\|^2 \|\mathbf{W}_i - \mathbf{R}_i^\top \mathbf{W}_i^* \mathbf{R}_{i-1}\|_F^2 + \|\mathbf{W}_1 - \mathbf{R}_1^\top \mathbf{W}_1^*\|_F^2. \quad (6)$$

Here the coefficient  $\|\mathbf{Y}^*\|^2$  is to harmonize the energy levels between  $\mathbf{W} = \{\mathbf{W}_i\}_{i \geq 2}$  and  $\mathbf{W}_1$ . The following result elucidates the connection between  $\text{dist}^2(\mathbf{W}, \mathbf{W}^*)$  and  $\|\mathbf{Y} - \mathbf{Y}^*\|_F^2$ , guaranteeing the convergence of  $\mathbf{Y}$  as  $\mathbf{W}$  approaches the global minima.

**Lemma 1:** Assume a whitened input  $\mathbf{X} \in \mathbb{R}^{d_x \times n}$ , i.e.  $\mathbf{X}\mathbf{X}^\top = \mathbf{I}_{d_x}$ . Let  $\mathbf{Y} = \mathbf{W}_N \cdots \mathbf{W}_1 \mathbf{X}$  and  $\mathbf{Y}^* = \mathbf{W}_N^* \cdots \mathbf{W}_1^* \mathbf{X}$ , where  $\mathbf{W}_i, \mathbf{W}_i^* \in \mathbb{R}^{d_i \times d_{i-1}}$  are orthonormal for  $i = 2, \dots, N$ . Given that  $\|\mathbf{W}_1\|^2 \leq \frac{9\|\mathbf{Y}^*\|^2}{4}$  and  $\|\mathbf{W}_1^*\|^2 = \|\mathbf{Y}^*\|^2$ , we can get

$$\|\mathbf{Y} - \mathbf{Y}^*\|_F^2 \geq \frac{1}{(16N-8)\kappa^2(\mathbf{Y}^*)} \text{dist}^2(\mathbf{W}, \mathbf{W}^*), \quad (7)$$

$$\|\mathbf{Y} - \mathbf{Y}^*\|_F^2 \leq \frac{9N}{4} \text{dist}^2(\mathbf{W}, \mathbf{W}^*). \quad (8)$$

**Proof:** Using the result [29, eq. (E.4)], for any  $j \geq 2$ , we have that  $\|\mathbf{W}_N \cdots \mathbf{W}_j - \mathbf{W}_N^* \cdots \mathbf{W}_j^* \mathbf{R}_{j-1}\|_F^2 \leq \frac{4\|\mathbf{Y} - \mathbf{Y}^*\|_F^2}{\sigma_{\min}^2(\mathbf{Y}^*)}$  for any  $\mathbf{R}_i \in \mathbb{O}^{d_i \times d_i}$ . It follows that

$$\begin{aligned} & \|\mathbf{W}_{N-1} - \mathbf{R}_{N-1}^\top \mathbf{W}_{N-1}^* \mathbf{R}_{N-2}\|_F^2 \\ &= \|\mathbf{W}_N^* \mathbf{R}_{N-1} \mathbf{W}_{N-1} - \mathbf{W}_N^* \mathbf{W}_{N-1}^* \mathbf{R}_{N-2}\|_F^2 \\ &\leq 2\|\mathbf{W}_N^* \mathbf{R}_{N-1} - \mathbf{W}_N\|_F^2 \|\mathbf{W}_{N-1}\|_F^2 \\ &\quad + 2\|\mathbf{W}_N \mathbf{W}_{N-1} - \mathbf{W}_N^* \mathbf{W}_{N-1}^* \mathbf{R}_{N-2}\|_F^2 \\ &\leq \frac{16\|\mathbf{Y} - \mathbf{Y}^*\|_F^2}{\sigma_{\min}^2(\mathbf{Y}^*)}. \end{aligned} \quad (9)$$

Similarly, we get  $\|\mathbf{W}_i - \mathbf{R}_i^\top \mathbf{W}_i^* \mathbf{R}_{i-1}\|_F^2 \leq \frac{16\|\mathbf{Y} - \mathbf{Y}^*\|_F^2}{\sigma_{\min}^2(\mathbf{Y}^*)}$  for  $i = 2, \dots, N-2$ . We now bound  $\|\mathbf{W}_1 - \mathbf{R}_1^\top \mathbf{W}_1^*\|_F^2$  by

$$\begin{aligned} & \|\mathbf{W}_1 - \mathbf{R}_1^\top \mathbf{W}_1^*\|_F^2 \\ &= \|\mathbf{W}_N^* \cdots \mathbf{W}_2^* \mathbf{R}_1 \mathbf{W}_1 - \mathbf{W}_N^* \cdots \mathbf{W}_1^*\|_F^2 \\ &\leq 2\|\mathbf{W}_1\|^2 \|\mathbf{W}_N \cdots \mathbf{W}_2 - \mathbf{W}_N^* \cdots \mathbf{W}_2^* \mathbf{R}_1\|_F^2 \\ &\quad + 2\|\mathbf{W}_N \cdots \mathbf{W}_1 - \mathbf{W}_N^* \cdots \mathbf{W}_1^*\|_F^2 \\ &\leq \frac{20\|\mathbf{Y}^*\|^2 \|\mathbf{Y} - \mathbf{Y}^*\|_F^2}{\sigma_{\min}^2(\mathbf{Y}^*)}, \end{aligned} \quad (10)$$

where the second inequality uses the fact that  $\mathbf{Y} = \mathbf{W}_N \cdots \mathbf{W}_1 \mathbf{X}$  and  $\|\mathbf{A}\mathbf{X}\|_F = \|\mathbf{A}\|_F$  for any  $\mathbf{A}$  since  $\mathbf{X}$  is whitened. Based on the preceding discussion and the definition of  $\text{dist}^2(\mathbf{W}, \mathbf{W}^*)$ , we can conclude (7).

Finally, we can prove the other direction by

$$\begin{aligned} & \|\mathbf{Y} - \mathbf{Y}^*\|_F^2 \\ &= \left\| \sum_{i=1}^N \mathbf{W}_N^* \cdots \mathbf{W}_{i+1}^* \mathbf{R}_i (\mathbf{W}_i - \mathbf{R}_i^\top \mathbf{W}_i^* \mathbf{R}_{i-1}) \mathbf{W}_{i-1} \cdots \mathbf{X} \right\|_F^2 \\ &\leq N \left( \sum_{i=2}^N \frac{9\|\mathbf{Y}^*\|^2}{4} \|\mathbf{W}_i - \mathbf{R}_i^\top \mathbf{W}_i^* \mathbf{R}_{i-1}\|_F^2 \right. \\ &\quad \left. + \|\mathbf{W}_1 - \mathbf{R}_1^\top \mathbf{W}_1^*\|_F^2 \right) \leq \frac{9N}{4} \text{dist}^2(\mathbf{W}, \mathbf{W}^*). \end{aligned} \quad (11)$$

**Main results:** To establish the convergence rate of RGD, we require the loss function  $L$  to satisfy a certain property. Given that our primary focus is the analysis of the local convergence, we will assume that the loss function behaves well only in a local region. Specifically, we will consider a category of loss functions that satisfies the so-called restricted correlated gradient (RCG) condition [29]:

**Definition 2:** We say the loss function  $L(\cdot; \mathbf{Y}^*)$  satisfies RCG( $\alpha, \beta, \mathcal{C}$ ) condition for  $\alpha, \beta > 0$  and the set  $\mathcal{C}$  if

$$\begin{aligned} & \langle \nabla L(\mathbf{Y}_1; \mathbf{Y}^*) - \nabla L(\mathbf{Y}_2; \mathbf{Y}^*), \mathbf{Y}_1 - \mathbf{Y}_2 \rangle \\ &\geq \alpha \|\mathbf{Y}_1 - \mathbf{Y}_2\|_F^2 + \beta \|\nabla L(\mathbf{Y}_1; \mathbf{Y}^*) - \nabla L(\mathbf{Y}_2; \mathbf{Y}^*)\|_F^2 \end{aligned} \quad (12)$$

for any  $\mathbf{Y}_1, \mathbf{Y}_2 \in \mathcal{C}$ .

The RCG condition is a generalization of the strong convexity. When  $L$  represents the MSE loss, i.e.,  $L(\mathbf{Y}, \mathbf{Y}^*) = \|\mathbf{Y} - \mathbf{Y}^*\|_F^2$ , which is commonly used in the convergence analysis of training deep linear networks [22], [23], [31], it satisfies the RCG condition with  $\alpha = \beta = 1$  and  $\mathcal{C} = \mathbb{R}^{d_y \times n}$ . The RCG condition may also accommodate other loss functions such as the cross entropy (CE) loss.

Based on Definition 2, we can initially deduce the Riemannian regularity condition as an extension of the regularity condition found in matrix factorization [32], [33], ensuring that gradients remain well-behaved within a defined region. Specifically, we have

**Lemma 2:** (Riemannian regularity condition) Suppose the training data  $(\mathbf{X}, \mathbf{Y}^*)$  is generated according to the data model in Definition 1. Also assume that the loss function  $L$  in (2) adheres to the RCG( $\alpha, \beta, \mathcal{C}$ ) condition where  $\mathcal{C} \triangleq \{\mathbf{Y} : \|\mathbf{Y} - \mathbf{Y}^*\|_F^2 \leq \frac{\alpha\beta\sigma_{\min}^2(\mathbf{Y}^*)}{72(2N-1)^2(N^2-1)\kappa^2(\mathbf{Y}^*)}\}$ . Under this assumption, for

any  $\mathbf{W} \in \{\mathbf{W} : \text{dist}^2(\mathbf{W}, \mathbf{W}^*) \leq \frac{\alpha\beta\sigma_{\min}^2(\mathbf{Y}^*)}{9(2N-1)(N^2-1)}\}$ , the function  $g$  in (2) satisfies the Riemannian regularity condition as following:

$$\begin{aligned} & \sum_{i=2}^N \langle \mathbf{W}_i - \mathbf{R}_i^\top \mathbf{W}_i^* \mathbf{R}_{i-1}, \mathcal{P}_{\mathbf{T}_{\mathbf{W}_i} \text{St}}(\nabla \mathbf{W}_i g) \rangle \\ &\quad + \langle \mathbf{W}_1 - \mathbf{R}_1 \mathbf{W}_1^*, \nabla \mathbf{W}_1 g \rangle \\ &\geq \frac{\alpha}{16(2N-1)\kappa^2(\mathbf{Y}^*)} \text{dist}^2(\mathbf{W}, \mathbf{W}^*) + \frac{\beta}{(9N-5)\|\mathbf{Y}^*\|^2} \\ &\quad \times \left( \sum_{i=2}^N \|\mathcal{P}_{\mathbf{T}_{\mathbf{W}_i} \text{St}}(\nabla \mathbf{W}_i g)\|_F^2 + \|\mathbf{Y}^*\|^2 \|\nabla \mathbf{W}_1 g\|_F^2 \right). \end{aligned} \quad (13)$$

*Proof:* To begin with, we can derive  $\|\mathbf{W}_1\|^2 \leq 2\|\mathbf{W}_1^*\|^2 + 2\|\mathbf{W}_1 - \mathbf{R}_1^\top \mathbf{W}_1^*\|^2 \leq 2\|\mathbf{Y}^*\|^2 + 2\text{dist}^2(\mathbf{W}, \mathbf{W}^*) \leq 2\|\mathbf{Y}^*\|^2 + \frac{2\alpha\beta\sigma_{\min}^2(\mathbf{Y}^*)}{9(2N-1)(N^2-1)} \leq \frac{9\|\mathbf{Y}^*\|^2}{4}$  where  $\alpha\beta \leq \frac{1}{4}$  [29] is used.

Next, through the gradients  $\nabla \mathbf{W}_i g = \mathbf{W}_{i+1}^\top \cdots \mathbf{W}_N^\top \nabla L(\mathbf{Y}; \mathbf{Y}^*) \mathbf{X}^\top \cdots \mathbf{W}_{i-1}^\top$ ,  $i \in [N]$ , we need to derive

$$\|\nabla \mathbf{W}_1 g\|_F^2 \leq \|\nabla L(\mathbf{Y}; \mathbf{Y}^*) - \nabla L(\mathbf{Y}^*; \mathbf{Y}^*)\|_F^2, \quad (14)$$

$$\begin{aligned} \|\nabla \mathbf{W}_i g\|_F^2 &\leq \|\mathbf{W}_1\|^2 \|\nabla L(\mathbf{Y}; \mathbf{Y}^*)\|_F^2 \\ &\leq \frac{9\|\mathbf{Y}^*\|^2}{4} \|\nabla L(\mathbf{Y}; \mathbf{Y}^*) - \nabla L(\mathbf{Y}^*; \mathbf{Y}^*)\|_F^2, \end{aligned} \quad (15)$$

where we employ  $\nabla L(\mathbf{Y}^*; \mathbf{Y}^*) = 0$ . Combining (14) and (15), we can obtain

$$\begin{aligned} & \sum_{i=2}^N \|\mathcal{P}_{\mathbf{T}_{\mathbf{W}_i} \text{St}}(\nabla \mathbf{W}_i g)\|_F^2 + \|\mathbf{Y}^*\|^2 \|\nabla \mathbf{W}_1 g\|_F^2 \\ &\leq \sum_{i=2}^N \|\nabla \mathbf{W}_i g\|_F^2 + \|\mathbf{Y}^*\|^2 \|\nabla \mathbf{W}_1 g\|_F^2 \\ &\leq \frac{9N-5}{4} \|\mathbf{Y}^*\|^2 \|\nabla L(\mathbf{Y}; \mathbf{Y}^*) - \nabla L(\mathbf{Y}^*; \mathbf{Y}^*)\|_F^2, \end{aligned} \quad (16)$$

in which the first inequality follows from the fact that for any matrix  $\mathbf{B} = \mathcal{P}_{\mathbf{T}_{L(\mathbf{X}_i)} \text{St}}(\mathbf{B}) + \mathcal{P}_{\mathbf{T}_{L(\mathbf{X}_i)}^\perp \text{St}}(\mathbf{B})$  where  $\mathcal{P}_{\mathbf{T}_{L(\mathbf{X}_i)} \text{St}}(\mathbf{B})$  and  $\mathcal{P}_{\mathbf{T}_{L(\mathbf{X}_i)}^\perp \text{St}}(\mathbf{B})$  are orthogonal, we have  $\|\mathcal{P}_{\mathbf{T}_{L(\mathbf{X}_i)} \text{St}}(\mathbf{B})\|_F^2 \leq \|\mathbf{B}\|_F^2$ .

Before analyzing the lower bound of cross term in (12), we need to establish the upper bound for the inner product between the orthogonal complement and the gradient as following:

$$\begin{aligned} T &= \sum_{i=2}^N \langle \mathcal{P}_{\mathbf{T}_{\mathbf{W}_i} \text{St}}^\perp(\mathbf{W}_i - \mathbf{R}_i^\top \mathbf{W}_i^* \mathbf{R}_{i-1}), \nabla \mathbf{W}_i g \rangle \\ &\leq \sum_{i=2}^N \frac{1}{2} \|\mathbf{W}_i\| \|\mathbf{W}_i - \mathbf{R}_i^\top \mathbf{W}_i^* \mathbf{R}_{i-1}\|_F^2 \\ &\quad \times \|\nabla L(\mathbf{Y}; \mathbf{Y}^*) - \nabla L(\mathbf{Y}^*; \mathbf{Y}^*)\|_F \|\mathbf{W}_1\| \\ &\leq \frac{\beta}{4} \|\nabla L(\mathbf{Y}; \mathbf{Y}^*) - \nabla L(\mathbf{Y}^*; \mathbf{Y}^*)\|_F^2 \\ &\quad + \frac{9(N-1)\|\mathbf{Y}^*\|^2}{16\beta} \sum_{i=2}^N \|\mathbf{W}_i - \mathbf{R}_i^\top \mathbf{W}_i^* \mathbf{R}_{i-1}\|_F^4 \\ &\leq \frac{\beta}{4} \|\nabla L(\mathbf{Y}; \mathbf{Y}^*) - \nabla L(\mathbf{Y}^*; \mathbf{Y}^*)\|_F^2 \\ &\quad + \frac{9(N-1)}{16\beta\|\mathbf{Y}^*\|^2} \text{dist}^4(\mathbf{W}, \mathbf{W}^*), \end{aligned} \quad (17)$$



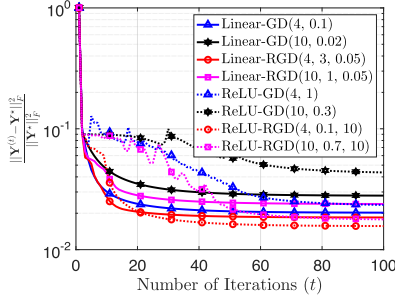


Fig. 1. Convergence analysis for GD( $N, \mu$ ) and RGD( $N, \mu, \gamma$ ) with different activation functions and  $N$ .

where  $\mathcal{P}_{\mathbf{W}_i}^\perp \text{St}(\mathbf{W}_i - \mathbf{R}_i^\top \mathbf{W}_i^* \mathbf{R}_{i-1}) = \frac{1}{2} \mathbf{W}_i (\mathbf{W}_i - \mathbf{R}_i^\top \mathbf{W}_i^* \mathbf{R}_{i-1})^\top (\mathbf{W}_i - \mathbf{R}_i^\top \mathbf{W}_i^* \mathbf{R}_{i-1})$  [30] and  $\nabla L(\mathbf{Y}^*; \mathbf{Y}^*) = 0$ .

Let us now introduce the notation  $\mathbf{H} = \mathbf{Y}^* - \mathbf{W}_N \cdots \mathbf{W}_2 \mathbf{R}_1^\top \mathbf{W}_1^* \mathbf{X} + \sum_{i=1}^N \mathbf{W}_N \cdots \mathbf{W}_{i+1} (\mathbf{W}_i - \mathbf{R}_i^\top \mathbf{W}_i^* \mathbf{R}_{i-1}) \mathbf{W}_{i-1} \cdots \mathbf{W}_1 \mathbf{X}$ , enabling us to simplify the expression of the cross term within (13). Then (13) can be rewritten as

$$\begin{aligned} & \sum_{i=2}^N \left\langle \mathbf{W}_i - \mathbf{R}_i^\top \mathbf{W}_i^* \mathbf{R}_{i-1}, \mathcal{P}_{\mathbf{W}_i}^\perp \text{St}(\nabla_{\mathbf{W}_i} g) \right\rangle \\ & + \langle \mathbf{W}_1 - \mathbf{R}_1^\top \mathbf{W}_1^*, \nabla_{\mathbf{W}_1} g \rangle \\ & = \langle \nabla L(\mathbf{Y}; \mathbf{Y}^*) - \nabla L(\mathbf{Y}^*; \mathbf{Y}^*), \mathbf{Y} - \mathbf{Y}^* + \mathbf{H} \rangle - T \\ & \geq \alpha \|\mathbf{Y} - \mathbf{Y}^*\|_F^2 + \frac{\beta}{4} \|\nabla L(\mathbf{Y}; \mathbf{Y}^*) - \nabla L(\mathbf{Y}^*; \mathbf{Y}^*)\|_F^2 \\ & \quad - \frac{9(N^2 - 1)}{16 \|\mathbf{Y}^*\|^2} \text{dist}^4(\mathbf{W}, \mathbf{W}^*) \\ & \geq \frac{\alpha}{16(2N - 1)\kappa^2(\mathbf{Y}^*)} \text{dist}^2(\mathbf{W}, \mathbf{W}^*) + \frac{\beta}{(9N - 5)\|\mathbf{Y}^*\|^2} \\ & \quad \times \left( \sum_{i=2}^N \|\mathcal{P}_{\mathbf{W}_i}^\perp \text{St}(\nabla_{\mathbf{W}_i} g)\|_F^2 + \|\mathbf{Y}^*\|^2 \|\nabla_{\mathbf{W}_1} g\|_F^2 \right), \quad (18) \end{aligned}$$

where the first inequality follows the RCG condition, (17) and  $\|\mathbf{H}\|_F^2 \leq \frac{9N(N-1)}{8\|\mathbf{Y}^*\|^2} \text{dist}^4(\mathbf{W}, \mathbf{W}^*)$  which is established through a mathematical transformation and the use of norm inequalities. The detailed proof for the upper bound of  $\|\mathbf{H}\|_F^2$  has been omitted here due to space limitations. In the last line, we leverage (16), Lemma 1 and  $\text{dist}^2(\mathbf{W}, \mathbf{W}^*) \leq \frac{\alpha\beta\sigma_{\min}^2(\mathbf{Y}^*)}{9(2N-1)(N^2-1)}$ . ■

We note that according to (7) in Lemma 1, the set  $\{\mathbf{W} : \text{dist}^2(\mathbf{W}, \mathbf{W}^*) \leq \frac{\alpha\beta\sigma_{\min}^2(\mathbf{Y}^*)}{9(2N-1)(N^2-1)}\}$  implies the region  $\mathcal{C}$ . By leveraging Riemannian regularity condition in Lemma 2 and the nonexpansiveness property of the polar decomposition-based retraction in [30, Lemma 1], we ultimately reach the following conclusion:

**Theorem 1:** In accordance with the identical conditions outlined in Lemma 2, we assume the initialization satisfies  $\text{dist}^2(\mathbf{W}^{(0)}, \mathbf{W}^*) \leq \frac{\alpha\beta\sigma_{\min}^2(\mathbf{Y}^*)}{9(2N-1)(N^2-1)}$ . When employing the learning rate  $\mu \leq \frac{2\beta}{(9N-5)\|\mathbf{Y}^*\|^2}$  and  $\gamma = \|\mathbf{Y}^*\|^2$  in RGD (5), we have

$$\text{dist}^2(\mathbf{W}^{(t+1)}, \mathbf{W}^*) \leq \left(1 - \frac{\alpha\sigma_{\min}^2(\mathbf{Y}^*)\mu}{8(2N-1)}\right) \text{dist}^2(\mathbf{W}^{(t)}, \mathbf{W}^*).$$

Our results reveal that the RGD demonstrates a linear convergence rate with polynomial decay concerning  $N$ . In

addition, by Lemma 1, we can easily obtain  $\|\mathbf{Y}^{(t)} - \mathbf{Y}^*\|_F^2 \leq (1 - \frac{\alpha\sigma_{\min}^2(\mathbf{Y}^*)\mu}{8(2N-1)})^t \frac{\alpha\beta\sigma_{\min}^2(\mathbf{Y}^*)}{8N^2-12N+4}$ . It is worth noting that, through analogous analysis, Theorem 1 can be extended to a broader scenario wherein  $\mathbf{W}_j^*$  takes on an arbitrary matrix, and matrices  $\mathbf{W}_i^*$  exhibit row orthogonality for  $i < j$  and column orthogonality for  $i > j$ . Moreover, we emphasize that our focus is primarily on the local convergence property of the RGD, and does not cover initialization methods extensively investigated in prior research, such as those discussed in [2], [21], [34], [35].

The research most closely related to our work is the convergence analysis of gradient descent in deep linear neural models across multiple layers, as demonstrated in [22], wherein the MSE loss function is taken into account. It is established that  $\|\mathbf{Y}^{(t)} - \mathbf{Y}^*\|_F^2 \leq \epsilon$  can be deduced when  $t \geq \Omega(\frac{N^3\|\mathbf{Y}^*\|_F^6}{c^2} \log(\frac{1}{\epsilon}))$ , provided that  $c \leq \sigma_{\min}(\mathbf{Y}^*)$ . Applying a similar derivation as presented in [22, Theorem 1], utilizing both Theorem 1 and Lemma 1, we can infer that  $t \geq \Omega(N^2\kappa^2(\mathbf{Y}^*) \log(\frac{1}{\epsilon}))$  in the RGD is sufficient to meet the same requirement. This further highlights the convergence advantage of the RGD.

#### IV. EXPERIMENTS

In this section, we conduct experiments to compare the performance of the RGD with gradient descent (GD). Specifically, we concentrate on the multi-class classification task using the MNIST dataset, where the input feature dimension is 784, while the output feature is 10, represented as  $\mathbf{y}_i \in \mathbb{R}^{10}$ . In this representation, each  $\mathbf{y}_i$  is designed such that it holds a value of 1 solely at the position that aligns with its categorical label, leaving the other positions assigned to 0. For our model architecture, we can deploy a multi-layer perceptron (MLP) with  $N$  layers where  $\mathbf{W}_1 \in \mathbb{R}^{100 \times 784}$ ,  $\{\mathbf{W}_i\}_{i=2}^{N-2} \in \mathbb{R}^{100 \times 100}$ ,  $\mathbf{W}_{N-1} \in \mathbb{R}^{50 \times 100}$ , and  $\mathbf{W}_N \in \mathbb{R}^{10 \times 50}$ . Each layer in the MLP is connected by a linear or rectified linear unit (Relu) activation function. While our theoretical result is only established for linear networks, we will also test the performance on a nonlinear MLP without bias terms.

We apply the GD and RGD for the CE loss function [36] in combination with the softmax function to train MLP models. The weight matrices are initialized using orthogonal initialization [2], except for  $\mathbf{W}_1$  in the RGD, which follows a uniform distribution within the range of  $(-\frac{1}{\sqrt{784}}, \frac{1}{\sqrt{784}})$  [34]. In addition, we perform a grid search to fine-tune the hyperparameters ( $\mu$  and  $\gamma$ ).

In Fig. 1, it is evident that RGD achieves a quicker convergence compared to GD. Notably, with an increase in the number of layers, the convergence rate decreases in alignment with our theoretical analysis. Moreover, due to the impact of nonlinear activation functions, algorithms that employ ReLU demonstrate a comparatively slower convergence rate compared to those utilizing linear activation functions. However, despite this, the error of the RGD with Relu outperforms that of the RGD using a linear activation function.

#### V. CONCLUSION

In this letter, we have provided a convergence analysis of the Riemannian gradient descent for a specific class of loss functions within orthonormal deep linear neural networks. Remarkably, our analysis guarantees a linear convergence rate, provided appropriate initialization. This will serve as a stepping stone for future explorations of training nonlinear orthonormal deep neural networks with adaptive learning rates.

## REFERENCES

- [1] D. Mishkin and J. Matas, "All you need is a good init," 2015, *arXiv:1511.06422*.
- [2] A. M. Saxe, J. L. McClelland, and S. Ganguli, "Exact solutions to the nonlinear dynamics of learning in deep linear neural networks," 2013, *arXiv:1312.6120*.
- [3] Y. Bengio, P. Simard, and P. Frasconi, "Learning long-term dependencies with gradient descent is difficult," *IEEE Trans. Neural Netw.*, vol. 5, no. 2, pp. 157–166, Mar. 1994.
- [4] R. Pascanu, T. Mikolov, and Y. Bengio, "On the difficulty of training recurrent neural networks," in *Proc. Int. Conf. Mach. Learn.*, 2013, pp. 1310–1318.
- [5] Q. V. Le, N. Jaitly, and G. E. Hinton, "A simple way to initialize recurrent networks of rectified linear units," 2015, *arXiv:1504.00941*.
- [6] M. Arjovsky, A. Shah, and Y. Bengio, "Unitary evolution recurrent neural networks," in *Proc. Int. Conf. Mach. Learn.*, 2016, pp. 1120–1128.
- [7] B. Hanin, "Which neural net architectures give rise to exploding and vanishing gradients?," in *Proc. Adv. Neural Inf. Process. Syst.*, 2018, vol. 31, pp. 580–589.
- [8] M. Harandi and B. Fernando, "Generalized backpropagation Etude De Cas: Orthogonality," 2016, *arXiv:1611.05927*.
- [9] S. Li, K. Jia, Y. Wen, T. Liu, and D. Tao, "Orthogonal deep neural networks," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 43, no. 4, pp. 1352–1368, Apr. 2021.
- [10] L. Huang et al., "Controllable orthogonalization in training DNNs," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 6429–6438.
- [11] J. Li, F. Li, and S. Todorovic, "Efficient Riemannian optimization on the Stiefel manifold via the cayley transform," in *Proc. Int. Conf. Learn. Representations*, 2020, pp. 1–16.
- [12] J. Wang, Y. Chen, R. Chakraborty, and S. X. Yu, "Orthogonal convolutional neural networks," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 11505–11515.
- [13] F. Malgouyres and F. Mamalet, "Existence, stability and scalability of orthogonal convolutional neural networks," *J. Mach. Learn. Res.*, vol. 23, pp. 15743–15798, 2022.
- [14] V. Dorobantu, P. A. Stromhaug, and J. Renteria, "DizzyRNN: Reparameterizing recurrent neural networks for norm-preserving backpropagation," 2016, *arXiv:1612.04035*.
- [15] Z. Mhammedi, A. Hellicar, A. Rahman, and J. Bailey, "Efficient orthogonal parametrisation of recurrent neural networks using householder reflections," in *Proc. Int. Conf. Mach. Learn.*, 2017, pp. 2401–2409.
- [16] E. Vorontsov, C. Trabelsi, S. Kadoury, and C. Pal, "On orthogonality and learning recurrent networks with long term dependencies," in *Proc. Int. Conf. Mach. Learn.*, 2017, pp. 3570–3578.
- [17] M. Cisse, P. Bojanowski, E. Grave, Y. Dauphin, and N. Usunier, "Parseval networks: Improving robustness to adversarial examples," in *Proc. Int. Conf. Mach. Learn.*, 2017, pp. 854–863.
- [18] M. Cogswell, F. Ahmed, R. Girshick, L. Zitnick, and D. Batra, "Reducing overfitting in deep networks by decorrelating representations," 2015, *arXiv:1511.06068*.
- [19] N. Bansal, X. Chen, and Z. Wang, "Can we gain more from orthogonality regularizations in training deep networks?," in *Proc. Adv. Neural Inf. Process. Syst.*, 2018, vol. 31, pp. 4266–4276.
- [20] L. Huang, X. Liu, B. Lang, A. Yu, Y. Wang, and B. Li, "Orthogonal weight normalization: Solution to optimization over multiple dependent stiefel manifolds in deep neural networks," in *Proc. AAAI Conf. Artif. Intell.*, 2018, pp. 3271–3278.
- [21] P. Bartlett, D. Helmbold, and P. Long, "Gradient descent with identity initialization efficiently learns positive definite linear transformations by deep residual networks," in *Proc. Int. Conf. Mach. Learn.*, 2018, pp. 521–530.
- [22] S. Arora, N. Cohen, N. Golowich, and W. Hu, "A convergence analysis of gradient descent for deep linear neural networks," 2018, *arXiv:1810.02281*.
- [23] D. Zou, P. M. Long, and Q. Gu, "On the global convergence of training deep linear ResNets," 2020, *arXiv:2003.01094*.
- [24] O. Shamir, "Exponential convergence time of gradient descent for one-dimensional deep linear neural networks," in *Proc. Conf. Learn. Theory*, 2019, pp. 2691–2713.
- [25] Z. Allen-Zhu, Y. Li, and Z. Song, "On the convergence rate of training recurrent neural networks," in *Proc. Adv. Neural Inf. Process. Syst.*, 2019, vol. 32, pp. 6676–6688.
- [26] S. Chatterjee, "Convergence of gradient descent for deep neural networks," 2022, *arXiv:2203.16462*.
- [27] M. Zhou, R. Ge, and C. Jin, "A local convergence theory for mildly over-parameterized two-layer neural network," in *Proc. Conf. Learn. Theory*, 2021, pp. 4577–4632.
- [28] X. Zhang, Y. Yu, L. Wang, and Q. Gu, "Learning one-hidden-layer relu networks via gradient descent," in *Proc. 22nd Int. Conf. Artif. Intell. Statist.*, 2019, pp. 1524–1534.
- [29] R. Han, R. Willett, and A. R. Zhang, "An optimal statistical and computational framework for generalized tensor estimation," 2020, *arXiv:2002.11255*.
- [30] X. Li, S. Chen, Z. Deng, Q. Qu, Z. Zhu, and A. Man-Cho So, "Weakly convex optimization over Stiefel manifold using Riemannian subgradient-type methods," *SIAM J. Optim.*, vol. 31, no. 3, pp. 1605–1634, 2021.
- [31] Z. Zhu, D. Soudry, Y. C. Eldar, and M. B. Wakin, "The global optimization geometry of shallow linear neural networks," *J. Math. Imag. Vis.*, vol. 62, pp. 279–292, 2020.
- [32] S. Tu, R. Boczar, M. Simchowitz, M. Soltanolkotabi, and B. Recht, "Low-rank solutions of linear matrix equations via procrustes flow," in *Proc. Int. Conf. Mach. Learn.*, 2016, pp. 964–973.
- [33] Z. Zhu, Q. Li, G. Tang, and M. B. Wakin, "The global optimization geometry of low-rank matrix optimization," *IEEE Trans. Inf. Theory*, vol. 67, no. 2, pp. 1308–1331, Feb. 2021.
- [34] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2015, pp. 1026–1034.
- [35] L. Xiao, Y. Bahri, J. Sohl-Dickstein, S. Schoenholz, and J. Pennington, "Dynamical isometry and a mean field theory of CNNs: How to train 10,000-layer vanilla convolutional neural networks," in *Proc. Int. Conf. Mach. Learn.*, 2018, pp. 5393–5402.
- [36] J. Zhou et al., "Are all losses created equal: A neural collapse perspective," in *Proc. Adv. Neural Inf. Process. Syst.*, 2022, vol. 35, pp. 31697–31710.