

Error Bound Analysis for the Regularized Loss of Deep Linear Neural Networks

Po Chen*

Rujun Jiang[†]

Peng Wang[‡]

February 18, 2025

Abstract

The optimization foundations of deep linear networks have received significant attention lately. However, due to the non-convexity and hierarchical structure, analyzing the regularized loss of deep linear networks remains a challenging task. In this work, we study the local geometric landscape of the regularized squared loss of deep linear networks, providing a deeper understanding of its optimization properties. Specifically, we characterize the critical point set and establish an error-bound property for all critical points under mild conditions. Notably, we identify the sufficient and necessary conditions under which the error bound holds. To support our theoretical findings, we conduct numerical experiments demonstrating that gradient descent exhibits linear convergence when optimizing the regularized loss of deep linear networks.

Key words: Deep linear networks, critical points, error bounds, linear convergence

MSC numbers: 90C26, 68T07, 65K10

1 Introduction

Deep learning has been widely used in various fields, including computer vision [16], natural language processing [39], and healthcare [12], due to its exceptional empirical performance. Optimization is a key component of deep learning, playing a pivotal role in formulating learning objectives, training neural networks, and improving model generalization. In general, optimization problems arising in deep learning are highly non-convex and difficult to analyze due to the inherent nonlinearity and hierarchical structures of deep neural networks. Even in the context of linear neural networks, which are the most basic form of neural networks, our theoretical understanding remains far from complete and systematic, especially concerning the optimization properties. This motivates us to study the following problem based on deep linear networks:

$$\min_{\mathbf{W}} \|\mathbf{W}_L \dots \mathbf{W}_1 \mathbf{X} - \mathbf{Y}\|_F^2 + \sum_{l=1}^L \lambda_l \|\mathbf{W}_l\|_F^2, \quad (1)$$

*School of Data Science, Fudan University, Shanghai(chenp24@m.fudan.edu.cn).

[†]School of Data Science, Fudan University, Shanghai (rjjiang@fudan.edu.cn).

[‡]Corresponding author. Department of Electrical Engineering and Computer Science, University of Michigan, Ann Arbor (peng8wang@gmail.com).

where $(\mathbf{X}, \mathbf{Y}) \in \mathbb{R}^{d_0 \times N} \times \mathbb{R}^{d_L \times N}$ denotes the data input, $\mathbf{W}_l \in \mathbb{R}^{d_l \times d_{l-1}}$ denotes the l -th weight matrix for each $l = 1, \dots, L$, $\mathbf{W} = \{\mathbf{W}_l\}_{l=1}^L$ denotes the collection of all weight matrices, and $\lambda_l > 0$ for all $l \in [L]$ are regularization parameters. Notably, such a problem captures a wide range of deep learning problems arising in applications, including deep matrix factorization [2, 8, 38], neural collapse [14, 44, 50], and low-rank adaption (LoRA) of large language models [18, 45], to name just a few. Moreover, studying linear networks provides a valuable starting point for gaining insights into nonlinear networks, as they exhibit similar learning behaviors and phenomena to their nonlinear counterparts while maintaining a simpler structure [20, 35].

In practice, (stochastic) gradient descent (GD) and its variants are among the most widely used first-order methods for deep learning [25]. Over the past few years, substantial progress has been made in studying the convergence behavior of GD for solving Problem (1). In the literature, a large amount of work has been dedicated to studying an unregularized counterpart of Problem (1), i.e., $\min \|\mathbf{W}_L \dots \mathbf{W}_1 \mathbf{X} - \mathbf{Y}\|_F^2$, based on a gradient dynamic analysis. A noteworthy assumption in most of these studies is that the data is whitened, i.e., $\mathbf{X}\mathbf{X}^T = \mathbf{I}_d$, when analyzing the convergence rate of GD. In general, this assumption simplifies the analysis by ensuring that the input samples are uncorrelated and have unit variance. Now, suppose that the data is whitened. Bartlett et al. [5] showed that GD with the identity initialization converges to an ϵ -optimal solution within a polynomial number of iterations. Later, Arora et al. [3] further improved the convergence result, showing that GD converges linearly to a global optimum when $\min\{d_1, \dots, d_{L-1}\} \geq \min\{d_0, d_L\}$, the initial weight matrices are approximately balanced, and the initial loss is smaller than a threshold. Meanwhile, Wu et al. [43] showed that gradient descent with zero-asymmetric initialization avoids saddle points and converges to an ϵ -optimal solution in $O(L^3 \log(1/\epsilon))$ iterations. Other works also proved similar global convergence and convergence rate results of GD under different assumptions; see, e.g., [4, 6, 19, 32, 43, 49]. Despite these inspiring results, the existing analyses suffer from three notable limitations. First, they are highly specific to a particular problem, relying on the analysis of gradient dynamics under carefully designed weight initialization schemes. This raises questions about the generalizability of these analyses to GD or other first-order optimization methods with different initialization schemes. Second, regularization is commonly used during the training of neural networks to prevent overfitting, improve generalization, and accelerate convergence [26, 48]. However, the existing gradient dynamics analyses mainly focus on unregularized deep neural networks and cannot be directly applied to Problem (1), even when the data is whitened. Finally, the existing analyses only apply to analyze the convergence to global optimal solutions. However, Problem (1) and its regularized counterpart may have local minimizers, to which first-order methods, such as GD, are likely to converge. To the best of our knowledge, the convergence behavior of first-order methods when they approach a critical point—whether a global minimum, a local minimum, or even a saddle point—remains an open question in the literature. To sum up, it remains an unsolved challenge in the literature to develop a unified framework to analyze the convergence behavior of first-order methods to critical points when solving Problem (1).

To address the above challenge, a promising approach is to study the local geometric structure of Problem (1) associated with its objective function, such as the error bound condition [29, 34, 51], the Polyak-Łojasiewicz (PL) inequality [33], and quadratic growth [9]. When the

data is whitened, as is commonly assumed in the literature (see the review above), we let $\mathbf{Y} := \mathbf{Y}\mathbf{X}^T$ with a slight abuse of notation. Under this notation, Problem (1) reduces to

$$\min_{\mathbf{W}} F(\mathbf{W}) := \|\mathbf{W}_L \dots \mathbf{W}_1 - \mathbf{Y}\|_F^2 + \sum_{l=1}^L \lambda_l \|\mathbf{W}_l\|_F^2. \quad (2)$$

Notably, both deep matrix factorization [2, 8, 38] and the neural collapse problem [14, 50] are special instances of the above problem. In this work, we mainly focus on the *error bound* condition associated with Problem (2). Formally, let \mathcal{W} denote the set of critical points to Problem (2). We say that it possesses an error bound for \mathcal{W} if there exist constants $\epsilon, \kappa > 0$ such that for all \mathbf{W} satisfying $\text{dist}(\mathbf{W}, \mathcal{W}) \leq \epsilon$,

$$\text{dist}(\mathbf{W}, \mathcal{W}) \leq \kappa \|\nabla F(\mathbf{W})\|_F, \quad (3)$$

where $\text{dist}(\mathbf{W}, \mathcal{W}) := \min\{\|\mathbf{W} - \mathbf{X}\|_F : \mathbf{X} \in \mathcal{W}\}$ denotes the distance from \mathbf{W} to the critical point set \mathcal{W} . Intuitively, the error bound inequality (3) requires the distance from a point to the set of critical points to be bounded by its gradient norm. A unified framework leveraging this condition, along with some algorithm-dependent conditions, has been widely studied to analyze linear convergence of first-order methods [9, 51] or superlinear convergence of second-order methods [47] in convex optimization. Recently, Liao et al. [29] showed that even if the objective function is non-convex but smooth and satisfies the error bound, GD converges linearly to the solution set. Moreover, Liao et al. [29], Rebjock and Boumal [34] demonstrated that the error bound is equivalent to other regularity conditions, such as the PL inequality [33] and quadratic growth [9] under mild conditions. These regularity conditions are widely used in the literature to prove linear convergence of GD for optimizing non-convex problems [13]. Notably, this type of convergence analysis framework is not limited to studying convergence to global optimal solutions. Instead, it can also be applied to local minima or even saddle points, as long as the error bound holds for the targeted solution sets. This makes such a framework particularly effective for analyzing non-convex optimization problems, where the landscape often contains various types of critical points.

However, powerful as this approach may seem, a critical challenge is to prove the error bound inequality (3) for Problem (2). As far as we know, proving the error bound or other regularity conditions for deep neural networks remains relatively underexplored. Recently, Wang et al. [40] made progress in this direction by establishing the error bound for the set of global optimal solutions to a special instance of Problem (2): $\min \|\mathbf{W}_2 \mathbf{W}_1 - \mathbf{Y}\|_F^2 + \lambda_1 \|\mathbf{W}_1\|_F^2 + \lambda_2 \|\mathbf{W}_2\|_F^2$, where $\lambda_1, \lambda_2 > 0$ are regularization parameters and \mathbf{Y} is a membership matrix. However, this analysis is limited to global optimal solutions of 2-layer neural networks and a specific data matrix \mathbf{Y} . Extending these results to deeper networks, more general data, and broader sets of critical points remains an open problem.

In view of the above discussion, our goal in this paper is to establish an error bound for the set of critical points of Problem (2). Our main contributions are twofold. First, we explicitly characterize the critical point set of Problem (2) for general data \mathbf{Y} (see Theorem 1), despite the non-convexity and hierarchical structure of the problem. This characterization serves as a foundation for establishing the error bound for Problem (2). Second, leveraging this explicit

characterization, we show that all critical points of Problem (2) satisfy the error bound (see Theorem 2) under mild conditions. Notably, we identify the sufficient and necessary conditions on the relationship between the regularization parameters and the spectrum of the input data that ensure the error bound holds. Such a result is significant, as it expands the currently limited repertoire of non-convex problems for which the local loss geometry is well understood. Moreover, it is important to note that our work develops new techniques in the process of establishing the error bound to handle the repeated singular values in \mathbf{Y} and the complicated structure of the critical point set. These techniques could be of independent interest. The established error bound can be used to establish other regularity conditions, such as the PL inequality and quadratic growth. Building on the error bound of Problem (2), we demonstrate that first-order methods can achieve linear convergence to a critical point of Problem (2), provided that they satisfy certain algorithm-dependent properties. We conduct numerical experiments in deep linear networks and in more general settings and observe that gradient descent converges linearly to critical points. These results strongly support our theoretical findings.

Organization. The rest of this paper is organized as follows. In Section 2, we present the main results of this paper. We introduce key lemmas and propositions without proofs and prove the main results in Section 3, and provide the detailed proofs for these lemmas and propositions in Section 4. We report experimental results in Section 5 and conclude the paper in Section 6.

Notation. Given an integer n , we denote by $[n]$ the set $\{1, \dots, n\}$. Given a vector \mathbf{a} , let $\|\mathbf{a}\|$ denote the Euclidean norm of \mathbf{a} , a_i the i -th entry, and $\text{diag}(\mathbf{a})$ the diagonal matrix with \mathbf{a} as its diagonal. Unless specified otherwise, all vectors in this paper are column vectors. Given a matrix $\mathbf{A} \in \mathbb{R}^{N \times d}$, let $\|\mathbf{A}\|$ denote the spectral norm of \mathbf{A} , $\|\mathbf{A}\|_F$ denote the Frobenius norm, a_{ij} denote the (i, j) -th element, and $\sigma_i(\mathbf{A})$ denote the i -th largest singular value. We use $\mathbf{0}_{m \times n}$ to denote $m \times n$ all-zero matrix, $\mathbf{0}_m$ to denote $m \times m$ all-zero matrix, and simply write $\mathbf{0}$ when its dimension can be inferred from the context. We use $\mathcal{O}^{n \times d}$ to denote the set of all $n \times d$ matrices that have orthonormal columns (in particular, we use \mathcal{O}^n to denote the set of all $d \times d$ orthogonal matrices); \mathcal{P}^n to denote the set of all $n \times n$ permutation matrices; $\text{BlkDiag}(\mathbf{X}_1, \dots, \mathbf{X}_n)$ to denote the block diagonal matrix whose diagonal blocks are $\mathbf{X}_1, \dots, \mathbf{X}_n$. Given a matrix $\mathbf{X} \in \mathbb{R}^{m \times n}$ and a non-empty closed set $\mathcal{X} \subseteq \mathbb{R}^{m \times n}$, we use $\text{dist}(\mathbf{X}, \mathcal{X})$ to denote the Euclidean distance from \mathbf{X} to \mathcal{X} ; the distance between \mathcal{X} and another non-empty closed set \mathcal{Y} is defined as $\text{dist}(\mathcal{X}, \mathcal{Y}) = \min_{\mathbf{X} \in \mathcal{X}, \mathbf{Y} \in \mathcal{Y}} \|\mathbf{X} - \mathbf{Y}\|_F$. Given weight matrices $\mathbf{W}_1, \dots, \mathbf{W}_L$, let $\mathbf{W}_{i:1} := \mathbf{W}_i \mathbf{W}_{i-1} \dots \mathbf{W}_1$ for each $i = 2, \dots, L$ and $\mathbf{W}_{L:i} := \mathbf{W}_L \mathbf{W}_{L-1} \dots \mathbf{W}_i$ for each $i = 1, \dots, L-1$. In particular, we define $\mathbf{W}_{0:1} := \mathbf{I}$ and $\mathbf{W}_{L:L+1} := \mathbf{I}$. For all $i \geq j+1$, let $\mathbf{W}_{i:j} = \mathbf{W}_i \mathbf{W}_{i-1} \dots \mathbf{W}_j$ and $\mathbf{W}_{i:j} = \mathbf{W}_i$ when $i = j$. All other notation is standard.

2 Main Results

To begin, we compute the gradient of the objective function as follows:

$$\nabla_{\mathbf{W}_l} F(\mathbf{W}) = \mathbf{W}_{L:l+1}^T (\mathbf{W}_{L:1} - \mathbf{Y}) \mathbf{W}_{l-1:1}^T + \lambda_l \mathbf{W}_l, \quad \forall l \in [L].$$

Then, the set of critical points of Problem (2) is defined as

$$\mathcal{W} := \{\mathbf{W} = (\mathbf{W}_1, \dots, \mathbf{W}_L) : \nabla_{\mathbf{W}_l} F(\mathbf{W}) = \mathbf{0}, \quad \forall l \in [L]\}. \quad (4)$$

Before we proceed, it is important to highlight the challenges in characterizing the set of critical points and establishing the error bound inequality (3) for Problem (2). First, the objective function is non-convex, presenting a complex loss landscape that may include a variety of critical points such as global minimizers, local minimizers, and saddle points [31]. Second, the objective function admits a hierarchical structure, where the weight matrices are multiplied sequentially and the layers have varying sizes. This hierarchical composition introduces a rotational invariance to the solution set (4), i.e., $\mathbf{W}_{l+1}\mathbf{W}_l = (\mathbf{W}_{l+1}\mathbf{Q}^T)(\mathbf{Q}\mathbf{W}_l)$ for any $\mathbf{Q} \in \mathcal{O}^{d_l}$ while the objective remains unchanged. This invariance leads to equivalence classes of solutions, thereby significantly complicating the analysis.

2.1 Characterization of Critical Points

In this subsection, we characterize the critical point set of Problem (2). To proceed, we first introduce the following assumption on the width of network layers.

Assumption 1. *It holds that $\min\{d_1, \dots, d_{L-1}\} \geq \min\{d_0, d_L\}$.*

This assumption ensures that the width of each hidden layer is no less than that of the input or output layer. It is important to note that this assumption is not restrictive and aligns with common practices in deep learning, where hidden layers are often designed to have sufficient capacity to capture complex data representations [15]. In addition, this assumption is widely used in the literature to analyze the convergence behavior of GD for optimizing deep networks; see, e.g., [3, 10, 11, 27].

Throughout the rest of the paper, let $d_{\min} := \min\{d_0, d_1, \dots, d_L\}$. This, together with Assumption 1, implies $d_{\min} = \min\{d_0, d_L\}$. Let $r_Y := \text{rank}(\mathbf{Y})$ and

$$\mathbf{Y} = \mathbf{U}_Y \mathbf{\Sigma}_Y \mathbf{V}_Y^T \quad (5)$$

be a singular value decomposition (SVD) of $\mathbf{Y} \in \mathbb{R}^{d_L \times d_0}$, where $\mathbf{U}_Y \in \mathcal{O}^{d_L}$, $\mathbf{V}_Y \in \mathcal{O}^{d_0}$, and $\mathbf{\Sigma}_Y = \text{BlkDiag}(\text{diag}(y_1, y_2, \dots, y_{r_Y}), \mathbf{0}) \in \mathbb{R}^{d_L \times d_0}$ with $y_1 \geq y_2 \geq \dots \geq y_{r_Y} > 0$ being top r_Y positive singular values. In the literature, it is common to assume that the singular values of \mathbf{Y} are distinct to simplify the analysis; see, e.g., [1, 24]. However, this assumption is overly strict and does not hold in many practical scenarios. For example, when \mathbf{Y} is a membership matrix in K -classification problems, it typically has K repeated singular values. To address this challenge, we introduce the following key setup and notions to improve the analysis. Let p_Y denote the number of distinct positive singular values of \mathbf{Y} . In other words, there exist indices s_0, s_1, \dots, s_{p_Y} such that $0 = s_0 < s_1 < \dots < s_{p_Y} = r_Y$ and

$$y_{s_0+1} = \dots = y_{s_1} > y_{s_1+1} = \dots = y_{s_2} > \dots > y_{s_{p_Y-1}+1} = \dots = y_{s_{p_Y}} > 0, \quad (6a)$$

$$y_{s_{p_Y}+1} = \dots = y_{d_{\min}} = 0. \quad (6b)$$

Then, let $h_i := s_i - s_{i-1}$ denote the multiplicity of the i -th largest positive singular value for each $i \in [p_Y]$. Consequently, we have $\sum_{i=1}^{p_Y} h_i = r_Y$. Now, we are ready to characterize the set of critical points explicitly as follows:

Theorem 1. *Suppose that Assumption 1 holds. It holds that $\mathbf{W} \in \mathcal{W}$ if and only if*

$$\mathbf{W}_1 = \mathbf{Q}_2 \Sigma_1 \text{BlkDiag}(\mathbf{O}_1, \dots, \mathbf{O}_{p_Y}, \mathbf{O}_{p_Y+1}) \mathbf{V}_Y^T, \quad (7a)$$

$$\mathbf{W}_l = \mathbf{Q}_{l+1} \Sigma_l \mathbf{Q}_l^T, \quad l = 2, \dots, L-1, \quad (7b)$$

$$\mathbf{W}_L = \mathbf{U}_Y \text{BlkDiag}(\mathbf{O}_1^T, \dots, \mathbf{O}_{p_Y}^T, \hat{\mathbf{O}}_{p_Y+1}^T) \Sigma_L \mathbf{Q}_L^T, \quad (7c)$$

where $\mathbf{Q}_l \in \mathcal{O}^{d_l-1}$ for all $l = 2, \dots, L$, $\mathbf{O}_i \in \mathcal{O}^{h_i}$ for each $i \in [p_Y]$, $\mathbf{O}_{p_Y+1} \in \mathcal{O}^{d_0-r_Y}$, $\hat{\mathbf{O}}_{p_Y+1} \in \mathcal{O}^{d_L-r_Y}$, and $\Sigma_l = \text{BlkDiag}(\text{diag}(\boldsymbol{\sigma})/\sqrt{\lambda_l}, \mathbf{0}) \in \mathbb{R}^{d_l \times d_{l-1}}$ for each $l \in [L]$ with $\boldsymbol{\sigma} \in \mathbb{R}^{r_Y}$ satisfying

$$\sigma_i^{2L-1} - \left(\sqrt{\lambda_1 \dots \lambda_L} y_i \right) \sigma_i^{L-1} + (\lambda_1 \dots \lambda_L) \sigma_i = 0, \quad \sigma_i \geq 0, \quad \forall i \in [r_Y]. \quad (8)$$

Before we proceed, let us make some remarks about this theorem. First, despite the non-convexity and hierarchical structure of Problem (2), we still characterize the set of critical points of Problem (2) explicitly, where each weight matrix admits the SVD in (7). Note that all weight matrices share the same singular values defined in (8) up to scaling of the regularization parameters $\{\lambda_l\}_{l=1}^L$, and their left and right singular matrices are determined by the orthogonal matrices \mathbf{Q}_l and \mathbf{O}_i . Here, \mathbf{Q}_l for all l are introduced to handle the rotational invariance in the sequential matrix multiplication of the weight matrices, while \mathbf{O}_i for all i are used to address the repeated singular values in \mathbf{Y} .

Second, to the best of our knowledge, there is no complete characterization of the critical point set for regularized deep linear networks in the existing literature. Recently, Dang et al. [7] have studied the geometric properties of the global minimizers of the neural collapse problem with deep linear networks. Notably, this problem is a special instance of Problem (2), where \mathbf{Y} is a membership matrix. In contrast to their work, which focuses only on global minimizers, our approach provides an explicit characterization of all critical points for arbitrary \mathbf{Y} . In addition, considerable research has been conducted to study the critical points of unregularized deep linear networks. The most complete and recent result is [1], which provides both necessary and sufficient conditions for identifying first-order critical points. Nevertheless, their analysis is limited to unregularized settings and cannot be applied to our case due to regularization terms.

Third, the explicit characterization of the critical point set in (7) serves as a cornerstone for establishing the error bound for Problem (2). It is worth noting that each critical point is not an isolated point but a union of connected sets (see Proposition 2). This brings significant difficulty to compute the distance $\text{dist}(\mathbf{W}, \mathcal{W})$ from a point \mathbf{W} to the critical point set. Such an intricacy also underscores the importance of a thorough and precise characterization of the critical point set to facilitate rigorous error-bound analysis.

2.2 Error Bound for Deep Linear Networks

In this subsection, we present our main result, which establishes the error bound for Problem (2). Before we proceed, it is important to point out that there are some degenerate cases under which the error bound does not hold. For example, when $L = 2$, we consider

$$\min_{\mathbf{W}_1, \mathbf{W}_2 \in \mathbb{R}^{d \times d}} F(\mathbf{W}) = \frac{1}{2} \|\mathbf{W}_2 \mathbf{W}_1 - \lambda \mathbf{I}_d\|_F^2 + \frac{1}{2} \lambda (\|\mathbf{W}_1\|_F^2 + \|\mathbf{W}_2\|_F^2).$$

This, together with Theorem 1, yields that the set of critical points is $\mathcal{W} = \{(\mathbf{0}, \mathbf{0})\}$. Now, consider $\mathbf{W}_1 = \mathbf{W}_2 = x\mathbf{I}_d$. We compute

$$\begin{aligned}\nabla_{\mathbf{W}_1} F(\mathbf{W}) &= \mathbf{W}_2^T (\mathbf{W}_2 \mathbf{W}_1 - \lambda \mathbf{I}_d) + \lambda \mathbf{W}_1 = x^3 \mathbf{I}_d, \\ \nabla_{\mathbf{W}_2} F(\mathbf{W}) &= (\mathbf{W}_2 \mathbf{W}_1 - \lambda \mathbf{I}_d) \mathbf{W}_1^T + \lambda \mathbf{W}_2 = x^3 \mathbf{I}_d.\end{aligned}$$

Therefore, we have $\|\nabla F(\mathbf{W})\|_F^2 = 2dx^6$. Moreover, we compute $\text{dist}^2(\mathbf{W}, \mathcal{W}) = \|\mathbf{W}_1\|_F^2 + \|\mathbf{W}_2\|_F^2 = 2dx^2$. Consequently, we have $\|\nabla F(\mathbf{W})\|_F = \text{dist}^3(\mathbf{W}, \mathcal{W})/(2d)$, and thus the error bound does not hold as $x \rightarrow 0$. To avoid such degenerate cases, we impose the following conditions on the relationship between the regularization parameters and the data matrix \mathbf{Y} .

Assumption 2. *It holds for $L = 2$ that*

$$\lambda_1 \lambda_2 \neq y_i^2, \quad \forall i \in [r_Y]. \quad (9)$$

For all $L \geq 3$, it holds that

$$\lambda_1 \dots \lambda_L \neq y_i^{2(L-1)} \left(\left(\frac{L-2}{L} \right)^{\frac{L}{2(L-1)}} + \left(\frac{L}{L-2} \right)^{\frac{L-2}{2(L-1)}} \right)^{-2(L-1)}, \quad \forall i \in [r_Y]. \quad (10)$$

Notably, this assumption provides sufficient and necessary conditions under which the error bound holds for all critical points of Problem (2). Indeed, if Assumption 2 holds, the error bound can be rigorously established (see Theorem 2). Conversely, if Assumption 2 is violated, we can show that the error bound fails to hold (see Appendix B). This dual perspective highlights the critical role of Assumption 2 in establishing the error bound for Problem (2). Now, under Assumptions 1 and 2, we are ready to prove the error bound for Problem (2).

Theorem 2. *Suppose that Assumptions 1 and 2 hold. There exist constants $\epsilon_1, \kappa_1 > 0$ such that for all \mathbf{W} satisfying $\text{dist}(\mathbf{W}, \mathcal{W}) \leq \epsilon_1$,*

$$\text{dist}(\mathbf{W}, \mathcal{W}) \leq \kappa_1 \|\nabla F(\mathbf{W})\|_F. \quad (11)$$

This theorem is significant as it establishes the fact that the error bound holds at any critical point of Problem (2) under Assumption 2. Importantly, the constants κ_1, ϵ_1 can be explicitly derived through our proofs; see Section 3.3. Now, we discuss some implications and the related work on this result. First, the established error bound can be used to derive other regularity conditions for analyzing the convergence behavior of first-order methods. The PL inequality and quadratic growth of Problem (2) can be obtained in Corollary 1 when the error bound holds, whose proof is deferred to Appendix C.1.

Corollary 1. *Suppose that \mathbf{W}^* is a critical point such that (11) holds for all \mathbf{W} satisfying $\|\mathbf{W} - \mathbf{W}^*\|_F \leq \epsilon_1$. The following statements hold:*

(i) *There exists a constant $\mu_1 > 0$ such that*

$$\|\nabla F(\mathbf{W})\|_F^2 \geq \mu_1 (F(\mathbf{W}) - F(\mathbf{W}^*)). \quad (12)$$

(ii) *If \mathbf{W}^* is a local minimizer, there exists a constant $\mu_2 > 0$ such that*

$$\text{dist}^2(\mathbf{W}, \mathcal{W}) \leq \mu_2 (F(\mathbf{W}) - F(\mathbf{W}^*)).$$

Second, combining the error bound with other algorithm-dependent conditions, such as *sufficient decrease*, *cost-to-go estimate*, and *safeguard*, we obtain local linear convergence of first-order methods for solving Problem (2) (see Proposition 7). In contrast to algorithm-specific convergence rate analyses in [4, 6, 19, 32, 43, 49], which are tailored to the dynamics of individual algorithms, the error-bound-based framework provides a unified approach. It applies not only to gradient descent but also to a broad class of first-order methods capable of optimizing Problem (2). It is also worth mentioning that our experimental results in Section 5 demonstrate that GD achieves linear convergence to both optimal and non-optimal critical points across various network architectures, which further supports our theoretical findings.

Third, we emphasize our technical contributions to establishing the error bound. Unlike prior works [30, 21, 22, 42], where the critical point set \mathcal{W} had a relatively simple structure that allows for explicit computation of the distance $\text{dist}(\mathbf{W}, \mathcal{W})$, the presence of orthogonal permutations and hierarchical structures in \mathcal{W} (see Theorem 1) poses a significant challenge to our analysis. To address this, we carefully construct an intermediate point $\hat{\mathbf{W}}$ that leverages the singular vectors of \mathbf{W} and singular values of $\mathbf{W}^* \in \mathcal{W}$. Then, we bound $\text{dist}(\hat{\mathbf{W}}, \mathcal{W})$ and $\|\mathbf{W} - \hat{\mathbf{W}}\|_F$ by the gradient norm at \mathbf{W} , respectively. Combining these bounds with the triangle inequality yields (11). Notably, this construction technique provides a new approach to showing the error bound of non-convex problems with a complicated solution set.

3 Proofs of the Main Results

In this section, we prove our main theorems concerning the critical points (i.e., Theorem 1) and the error bound (i.e., Theorem 2) of Problem (2). To avoid interrupting the flow of the main proofs, we introduce key lemmas and propositions, with their detailed proofs deferred to Section 4. Based on these results, we present the complete proofs of Theorem 1 and Theorem 2.

3.1 Preliminary Results

To characterize the critical points and establish the error bound of Problem (2), we claim that it suffices to study the following problem:

$$\min_{\mathbf{W}} G(\mathbf{W}) := \left\| \mathbf{W}_L \dots \mathbf{W}_1 - \sqrt{\lambda} \mathbf{Y} \right\|_F^2 + \lambda \sum_{l=1}^L \|\mathbf{W}_l\|_F^2, \quad (13)$$

where $\lambda := \lambda_1 \dots \lambda_L$. We compute the gradient of $G(\mathbf{W})$, which will be frequently used, as follows:

$$\frac{1}{2} \nabla_{\mathbf{W}_l} G(\mathbf{W}) = \mathbf{W}_{L:l+1}^T \left(\mathbf{W}_{L:1} - \sqrt{\lambda} \mathbf{Y} \right) \mathbf{W}_{l-1:1}^T + \lambda \mathbf{W}_l, \quad \forall l \in [L]. \quad (14)$$

The critical point set of this problem is defined as

$$\mathcal{W}_G := \{ \mathbf{W} = (\mathbf{W}_1, \dots, \mathbf{W}_L) : \nabla_{\mathbf{W}_l} G(\mathbf{W}) = \mathbf{0}, \quad \forall l \in [L] \}. \quad (15)$$

Now, we present the following lemma to prove our claim.

Lemma 1. *Consider Problems (2) and (13). The following statements hold:*

(i) $(\mathbf{W}_1, \dots, \mathbf{W}_L) \in \mathcal{W}$ if and only if $(\sqrt{\lambda_1} \mathbf{W}_1, \dots, \sqrt{\lambda_L} \mathbf{W}_L) \in \mathcal{W}_G$.

(ii) Suppose that there exist constants $\epsilon, \kappa > 0$ such that for all \mathbf{Z} satisfying $\text{dist}(\mathbf{Z}, \mathcal{W}_G) \leq \epsilon$, it holds that

$$\text{dist}(\mathbf{Z}, \mathcal{W}_G) \leq \kappa \|\nabla G(\mathbf{Z})\|_F. \quad (16)$$

Then, for all \mathbf{W} satisfying $\text{dist}(\mathbf{W}, \mathcal{W}) \leq \epsilon/\sqrt{\lambda_{\max}}$, it holds that

$$\text{dist}(\mathbf{W}, \mathcal{W}) \leq \frac{\kappa\lambda}{\lambda_{\min}} \|\nabla F(\mathbf{W})\|_F.$$

(iii) Let (5) be an SVD of \mathbf{Y} . It holds that $(\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}_L)$ is a critical point of Problem (13) if and only if $(\mathbf{W}_1 \mathbf{V}_Y, \mathbf{W}_2, \dots, \mathbf{U}_Y^T \mathbf{W}_L)$ is a critical point of the following problem:

$$\min_{\mathbf{W}} \left\| \mathbf{W}_L \dots \mathbf{W}_1 - \sqrt{\lambda} \boldsymbol{\Sigma}_Y \right\|_F^2 + \lambda \sum_{l=1}^L \|\mathbf{W}_l\|_F^2.$$

Proof. (i) Let $\mathbf{W} := (\mathbf{W}_1, \dots, \mathbf{W}_L) \in \mathcal{W}$ be arbitrary. Using this and (15), we obtain

$$\nabla_{\mathbf{W}_l} F(\mathbf{W}) = 2\mathbf{W}_{L:l+1}^T (\mathbf{W}_{L:1} - \mathbf{Y}) \mathbf{W}_{l-1:1}^T + 2\lambda_l \mathbf{W}_l = \mathbf{0}, \quad \forall l \in [L]. \quad (17)$$

Now, let $\hat{\mathbf{W}} := (\sqrt{\lambda_1} \mathbf{W}_1, \dots, \sqrt{\lambda_L} \mathbf{W}_L)$. For each $l \in [L]$, we compute

$$\begin{aligned} \nabla_{\mathbf{W}_l} G(\hat{\mathbf{W}}) &\stackrel{(14)}{=} 2\sqrt{\lambda_L \dots \lambda_{l+1} \lambda_{l-1} \dots \lambda_1} \mathbf{W}_{L:l+1}^T \left(\sqrt{\lambda} \mathbf{W}_{L:1} - \sqrt{\lambda} \mathbf{Y} \right) \mathbf{W}_{l-1:1}^T + 2\lambda \sqrt{\lambda_l} \mathbf{W}_l \\ &= \frac{2\lambda}{\sqrt{\lambda_l}} \left(\mathbf{W}_{L:l+1}^T (\mathbf{W}_{L:1} - \mathbf{Y}) \mathbf{W}_{l-1:1}^T + \lambda_l \mathbf{W}_l \right) \stackrel{(17)}{=} \mathbf{0}. \end{aligned}$$

Therefore, we have $\hat{\mathbf{W}} \in \mathcal{W}_G$. Conversely, let $\hat{\mathbf{W}} \in \mathcal{W}_G$ be arbitrary. Using the same argument, we have $\mathbf{W} \in \mathcal{W}$.

(ii) Using (i), we express \mathcal{W}_G as

$$\mathcal{W}_G = \left\{ (\sqrt{\lambda_1} \mathbf{W}_1, \dots, \sqrt{\lambda_L} \mathbf{W}_L) : (\mathbf{W}_1, \dots, \mathbf{W}_L) \in \mathcal{W} \right\}.$$

Now, let $\mathbf{W} := (\mathbf{W}_1, \dots, \mathbf{W}_L)$ be arbitrary and $\mathbf{Z} := (\sqrt{\lambda_1} \mathbf{W}_1, \dots, \sqrt{\lambda_L} \mathbf{W}_L)$. Then, let $(\sqrt{\lambda_1} \hat{\mathbf{W}}_1^*, \dots, \sqrt{\lambda_L} \hat{\mathbf{W}}_L^*)$ with $\hat{\mathbf{W}}^* \in \mathcal{W}$ be such that $\text{dist}^2(\mathbf{Z}, \mathcal{W}_G) = \sum_{l=1}^L \lambda_l \|\mathbf{W}_l - \hat{\mathbf{W}}_l^*\|_F^2$ and $\mathbf{W}^* = (\mathbf{W}_1^*, \dots, \mathbf{W}_L^*) \in \mathcal{W}$ be such that $\text{dist}^2(\mathbf{W}, \mathcal{W}) = \sum_{l=1}^L \|\mathbf{W}_l - \mathbf{W}_l^*\|_F^2$. We obtain

$$\text{dist}^2(\mathbf{W}, \mathcal{W}) = \sum_{l=1}^L \|\mathbf{W}_l - \mathbf{W}_l^*\|_F^2 \geq \frac{\sum_{l=1}^L \|\sqrt{\lambda_l} \mathbf{W}_l - \sqrt{\lambda_l} \mathbf{W}_l^*\|_F^2}{\lambda_{\max}} \geq \frac{1}{\lambda_{\max}} \text{dist}^2(\mathbf{Z}, \mathcal{W}_G). \quad (18)$$

This, together with $\text{dist}(\mathbf{W}, \mathcal{W}) \leq \epsilon/\sqrt{\lambda_{\max}}$, implies $\text{dist}(\mathbf{Z}, \mathcal{W}_G) \leq \epsilon$. Using this and (16), we obtain

$$\text{dist}(\mathbf{Z}, \mathcal{W}_G) \leq \kappa \|\nabla G(\mathbf{Z})\|_F. \quad (19)$$

Next, we have

$$\text{dist}^2(\mathbf{Z}, \mathcal{W}_G) = \sum_{l=1}^L \lambda_l \|\mathbf{W}_l - \hat{\mathbf{W}}_l^*\|_F^2 \geq \lambda_{\min} \sum_{l=1}^L \|\mathbf{W}_l - \hat{\mathbf{W}}_l^*\|_F^2 \geq \lambda_{\min} \text{dist}^2(\mathbf{W}, \mathcal{W}). \quad (20)$$

Moreover, we have for each $l \in [L]$,

$$\begin{aligned} \nabla_{\mathbf{W}_l} G(\mathbf{Z}) &\stackrel{(17)}{=} 2\sqrt{\lambda_L \dots \lambda_{l+1} \lambda_{l-1} \dots \lambda_1} \mathbf{W}_{L:l+1}^T \left(\sqrt{\lambda} \mathbf{W}_{L:1} - \sqrt{\lambda} \mathbf{Y} \right) \mathbf{W}_{l-1:1}^T + 2\lambda \sqrt{\lambda_l} \mathbf{W}_l \\ &= \frac{\lambda}{\sqrt{\lambda_l}} \nabla_{\mathbf{W}_l} F(\mathbf{W}). \end{aligned} \quad (21)$$

This, together with (19) and (20), yields

$$\text{dist}(\mathbf{W}, \mathcal{W}) \leq \frac{1}{\sqrt{\lambda_{\min}}} \text{dist}(\mathbf{Z}, \mathcal{W}_G) \leq \frac{\kappa}{\sqrt{\lambda_{\min}}} \|\nabla G(\mathbf{W})\|_F \leq \frac{\kappa\lambda}{\lambda_{\min}} \|\nabla F(\mathbf{W})\|_F.$$

(iii) Obviously, each critical point of Problem (13) satisfies $\nabla_{\mathbf{W}_l} G(\mathbf{W}) = 0$. This, together with (5) and (14), directly implies the desired result. \square

Using (i) and (ii) of the above lemma, it suffices to characterize the critical points and establish the error bound for Problem (13) in the rest of this section. Moreover, using (iii) of the above lemma, we assume without loss of generality that

$$\mathbf{Y} = \text{BlkDiag} \left(\tilde{\Sigma}_Y, \mathbf{0}_{(d_L - d_{\min}) \times (d_0 - d_{\min})} \right) \in \mathbb{R}^{d_L \times d_0}, \quad (22)$$

where $\tilde{\Sigma}_Y = \text{diag}(y_1, \dots, y_{d_{\min}})$ with $y_1 \geq y_2 \geq \dots \geq y_{d_{\min}} \geq 0$ being singular values. According to (6), we write

$$\tilde{\Sigma}_Y = \text{BlkDiag} \left(y_{s_1} \mathbf{I}_{h_1}, \dots, y_{s_{p_Y}} \mathbf{I}_{h_{p_Y}}, \mathbf{0}_{d_{\min} - r_Y} \right) \in \mathbb{R}^{d_{\min} \times d_{\min}}. \quad (23)$$

Moreover, we define

$$\delta_y := \min \{ |y_{s_i} - y_{s_{i+1}}| : i \in [p_Y] \}. \quad (24)$$

Throughout this section, we will consistently use the above notation in all proofs.

3.2 Analysis of the Set of Critical Points

In this subsection, we focus on characterizing all critical points of Problem (2) explicitly by studying the critical points of Problem (13). To begin, we present a lemma to show that the weight matrices $\{\mathbf{W}_l\}_{l=1}^L$ at any critical point of Problem (13) are balanced.

Lemma 2. *Let $(\mathbf{W}_1, \dots, \mathbf{W}_L)$ be a critical point of Problem (13). The following statements hold:*

(i) *It holds that*

$$\mathbf{W}_l \mathbf{W}_l^T = \mathbf{W}_{l+1}^T \mathbf{W}_{l+1}, \quad \forall l \in [L-1]. \quad (25)$$

(ii) *It holds that*

$$(\mathbf{W}_l \mathbf{W}_l^T)^{L-1} \mathbf{W}_l - \sqrt{\lambda} \mathbf{W}_{L:l+1}^T \mathbf{Y} \mathbf{W}_{l-1:l}^T + \lambda \mathbf{W}_l = 0, \quad \forall l \in [L]. \quad (26)$$

The proof of this lemma is deferred to Section 4.1. Recall the notions in (6), (22), and (23). Leveraging Lemma 1 and Lemma 2, we are ready to characterize the set of critical points (15) as follows.

Proposition 1. *Suppose that Assumption 1 holds and $\mathbf{Y} \in \mathbb{R}^{d_L \times d_0}$ takes the form of (22). The critical point set (15) of Problem (13) can be expressed as*

$$\mathcal{W}_G = \left\{ \mathbf{W} : \begin{array}{l} \Sigma_l = \text{BlkDiag}(\text{diag}(\boldsymbol{\sigma}), \mathbf{0}) \in \mathbb{R}^{d_l \times d_{l-1}}, \quad \forall l \in [L], \quad (\boldsymbol{\sigma}, \boldsymbol{\Pi}) \in \mathcal{B}, \\ \mathbf{W}_1 = \mathbf{Q}_2 \Sigma_1 \text{BlkDiag}(\boldsymbol{\Pi}, \mathbf{I}) \text{BlkDiag}(\mathbf{O}_1, \dots, \mathbf{O}_{p_Y}, \mathbf{O}_{p_Y+1}), \\ \mathbf{W}_l = \mathbf{Q}_{l+1} \Sigma_l \mathbf{Q}_l^T, \quad l = 2, \dots, L-1, \quad \mathbf{Q}_l \in \mathcal{O}^{d_{l-1}}, \quad l = 2, \dots, L, \\ \mathbf{W}_L = \text{BlkDiag}(\mathbf{O}_1^T, \dots, \mathbf{O}_{p_Y}^T, \hat{\mathbf{O}}_{p_Y+1}^T) \text{BlkDiag}(\boldsymbol{\Pi}^T, \mathbf{I}) \Sigma_L \mathbf{Q}_L^T, \\ \mathbf{O}_i \in \mathcal{O}^{h_i}, \quad \forall i \in [p_Y], \quad \mathbf{O}_{p_Y+1} \in \mathcal{O}^{d_0 - r_Y}, \quad \hat{\mathbf{O}}_{p_Y+1} \in \mathcal{O}^{d_L - r_Y} \end{array} \right\}, \quad (27)$$

where

$$\mathcal{A} := \left\{ \boldsymbol{\sigma} \in \mathbb{R}^{d_{\min}} : \sigma_i^{2L-1} - \sqrt{\lambda} y_i \sigma_i^{L-1} + \lambda \sigma_i = 0, \sigma_i \geq 0, \forall i \in [d_{\min}] \right\}, \quad (28)$$

$$\mathcal{B} := \left\{ (\boldsymbol{\sigma}, \boldsymbol{\Pi}) \in \mathbb{R}^{d_{\min}} \times \mathcal{P}^{d_{\min}} : \boldsymbol{\sigma} = \boldsymbol{\Pi} \mathbf{a}, \mathbf{a} \in \mathcal{A}, \sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{d_{\min}} \right\}. \quad (29)$$

The proof of this proposition is deferred to Section 4.1. Notably, each \mathbf{W}_l in (27) is represented in an SVD form with singular values selected from \mathcal{A} up to a permutation determined by \mathcal{B} . Based on this proposition, we can further simplify the above structure of the critical point set by removing the permutation in (27).

Theorem 3. Suppose that Assumption 1 holds and $\mathbf{Y} \in \mathbb{R}^{d_L \times d_0}$ takes the form of (22). Then, the critical point set (15) of Problem (13) can be expressed as

$$\mathcal{W}_G = \left\{ \mathbf{W} : \begin{cases} \boldsymbol{\Sigma}_l = \text{BlkDiag}(\text{diag}(\boldsymbol{\sigma}), \mathbf{0}) \in \mathbb{R}^{d_l \times d_{l-1}}, \forall l \in [L], \boldsymbol{\sigma} \in \mathcal{A}, \\ \mathbf{W}_1 = \mathbf{Q}_2 \boldsymbol{\Sigma}_1 \text{BlkDiag}(\mathbf{O}_1, \dots, \mathbf{O}_{p_Y}, \mathbf{O}_{p_Y+1}), \\ \mathbf{W}_l = \mathbf{Q}_{l+1} \boldsymbol{\Sigma}_l \mathbf{Q}_l^T, l = 2, \dots, L-1, \mathbf{Q}_l \in \mathcal{O}^{d_{l-1}}, l = 2, \dots, L, \\ \mathbf{W}_L = \text{BlkDiag}(\mathbf{O}_1^T, \dots, \mathbf{O}_{p_Y}^T, \widehat{\mathbf{O}}_{p_Y+1}^T) \boldsymbol{\Sigma}_L \mathbf{Q}_L^T, \\ \mathbf{O}_i \in \mathcal{O}^{h_i}, \forall i \in [p_Y], \mathbf{O}_{p_Y+1} \in \mathcal{O}^{d_0 - r_Y}, \widehat{\mathbf{O}}_{p_Y+1} \in \mathcal{O}^{d_L - r_Y} \end{cases} \right\}, \quad (30)$$

where \mathcal{A} is defined in (28).

Proof. For ease of exposition, we denote the set on the right-hand side of (30) by \mathcal{X} . Let $(\mathbf{W}_1, \dots, \mathbf{W}_L) \in \mathcal{W}_G$ be arbitrary. It follows from Proposition 1 that there exists $(\boldsymbol{\sigma}, \boldsymbol{\Pi}) \in \mathcal{B}$ such that $(\mathbf{W}_1, \dots, \mathbf{W}_L)$ takes the form of (27). According to (28) and (29), there exists $\mathbf{a} \in \mathcal{A}$ such that $\boldsymbol{\Pi} \text{diag}(\mathbf{a}) \boldsymbol{\Pi}^T = \text{diag}(\boldsymbol{\sigma})$. Then, let

$$\begin{aligned} \boldsymbol{\Sigma}'_l &:= \text{BlkDiag}(\boldsymbol{\Pi}^T, \mathbf{I}_{d_l - d_{\min}}) \boldsymbol{\Sigma}_l \text{BlkDiag}(\boldsymbol{\Pi}, \mathbf{I}_{d_{l-1} - d_{\min}}) = \text{BlkDiag}(\text{diag}(\mathbf{a}), \mathbf{0}), \forall l \in [L], \\ \mathbf{Q}'_l &:= \mathbf{Q}_l \text{BlkDiag}(\boldsymbol{\Pi}, \mathbf{I}_{d_{l-1} - d_{\min}}) \in \mathcal{O}^{d_{l-1}}, \forall l \in [L]. \end{aligned}$$

Therefore, we have $(\mathbf{W}_1, \dots, \mathbf{W}_L) \in \mathcal{X}$.

Conversely, let $(\mathbf{W}_1, \dots, \mathbf{W}_L) \in \mathcal{X}$ be arbitrary. There exists $\boldsymbol{\sigma} \in \mathcal{A}$ such that $(\mathbf{W}_1, \dots, \mathbf{W}_L)$ takes the form of (30). Then, we choose a permutation matrix $\boldsymbol{\Pi} \in \mathcal{P}^{d_{\min}}$ such that $\boldsymbol{\sigma}' = \boldsymbol{\Pi} \boldsymbol{\sigma}$ satisfies $\sigma'_1 \geq \sigma'_2 \geq \dots \geq \sigma'_{d_{\min}}$. Let

$$\begin{aligned} \boldsymbol{\Sigma}'_l &:= \text{BlkDiag}(\boldsymbol{\Pi}, \mathbf{I}_{d_l - d_{\min}}) \boldsymbol{\Sigma}_l \text{BlkDiag}(\boldsymbol{\Pi}^T, \mathbf{I}_{d_{l-1} - d_{\min}}) = \text{BlkDiag}(\text{diag}(\boldsymbol{\sigma}'), \mathbf{0}), \forall l \in [L], \\ \mathbf{Q}'_l &:= \mathbf{Q}_l \text{BlkDiag}(\boldsymbol{\Pi}^T, \mathbf{I}_{d_{l-1} - d_{\min}}), \forall l \in [L]. \end{aligned}$$

This, together with (30), implies $(\mathbf{W}_1, \dots, \mathbf{W}_L) \in \mathcal{W}_G$. Then, we complete the proof. \square

Using Lemma 1 and Theorem 3, we can directly characterize the critical point set (4) of Problem (2), i.e., Theorem 1, as follows.

Proof of Theorem 1. Note that $y_i = 0$ for each $i = r_Y + 1, \dots, d_{\min}$. This, together with (28) and $\lambda > 0$, yields $\sigma_i = 0$ for each $i = r_Y + 1, \dots, d_{\min}$. Using this, (5), Lemma 1, and Theorem 3, we directly obtain Theorem 1. \square

Notably, Proposition 1 demonstrates that when \mathbf{Y} takes the form of (22), the critical point set (15) of Problem (13) can be expressed as

$$\mathcal{W}_G = \bigcup_{(\boldsymbol{\sigma}, \boldsymbol{\Pi}) \in \mathcal{B}} \mathcal{W}_{\boldsymbol{\sigma}, \boldsymbol{\Pi}},$$

where

$$\mathcal{W}_{\sigma, \Pi} := \left\{ \mathbf{W} : \begin{array}{l} \Sigma_l = \text{BlkDiag}(\text{diag}(\sigma), \mathbf{0}_{(d_l - d_{\min} \times (d_{l-1} - d_{\min}))}) \in \mathbb{R}^{d_l \times d_{l-1}}, \forall l \in [L], \\ \mathbf{W}_1 = \mathbf{Q}_2 \Sigma_1 \text{BlkDiag}(\Pi, \mathbf{I}) \text{BlkDiag}(\mathbf{O}_1, \dots, \mathbf{O}_{p_Y}, \mathbf{O}_{p_Y+1}), \\ \mathbf{W}_l = \mathbf{Q}_{l+1} \Sigma_l \mathbf{Q}_l^T, \quad l = 2, \dots, L-1, \quad \mathbf{Q}_l \in \mathcal{O}^{d_{l-1}}, \quad l = 2, \dots, L, \\ \mathbf{W}_L = \text{BlkDiag}(\mathbf{O}_1^T, \dots, \mathbf{O}_{p_Y}^T, \widehat{\mathbf{O}}_{p_Y+1}^T) \text{BlkDiag}(\Pi^T, \mathbf{I}) \Sigma_L \mathbf{Q}_L^T, \\ \mathbf{O}_i \in \mathcal{O}^{h_i}, \quad \forall i \in [p_Y], \quad \mathbf{O}_{p_Y+1} \in \mathcal{O}^{d_0 - r_Y}, \quad \widehat{\mathbf{O}}_{p_Y+1} \in \mathcal{O}^{d_L - r_Y} \end{array} \right\}. \quad (31)$$

Let

$$\mathcal{Y} := \bigcup_{i \in [d_{\min}]} \left\{ \sigma \geq 0 : \sigma^{2L-1} + \lambda \sigma - \sqrt{\lambda} y_i \sigma^{L-1} = 0 \right\}. \quad (32)$$

Now, we define

$$\delta_\sigma := \min \{ |x - y| : x \neq y \in \mathcal{Y} \}. \quad (33)$$

The following result elucidates the structure of the collection $\{\mathcal{W}_{\sigma, \Pi}\}_{(\sigma, \Pi) \in \mathcal{B}}$.

Proposition 2. *Let $(\sigma, \Pi) \in \mathcal{B}$ and $(\sigma', \Pi') \in \mathcal{B}$ be arbitrary. The following statements holds:*

- (i) *It holds that $\mathcal{W}_{\sigma, \Pi} = \mathcal{W}_{\sigma', \Pi'}$ if and only if $\sigma = \sigma'$.*
- (ii) *If $\sigma \neq \sigma'$, it holds that*

$$\text{dist}(\mathcal{W}_{\sigma, \Pi}, \mathcal{W}_{\sigma', \Pi'}) \geq \delta_\sigma. \quad (34)$$

The proof of this proposition is deferred to Section 4.1. This proposition demonstrates that for any pair $(\sigma, \Pi) \in \mathcal{B}$ and $(\sigma', \Pi') \in \mathcal{B}$, if $\sigma \neq \sigma'$, $\mathcal{W}_{\sigma, \Pi}$ is well separated from $\mathcal{W}_{\sigma', \Pi'}$, and otherwise they are identical. Therefore, for simplicity we write $\mathcal{W}_\sigma := \mathcal{W}_{\sigma, \Pi}$ for any Π satisfying $(\sigma, \Pi) \in \mathcal{B}$. Therefore, one can express the set of critical points (15) of Problem (13) as follows:

$$\mathcal{W}_G := \bigcup_{\sigma \in \mathcal{A}_{\text{sort}}} \mathcal{W}_\sigma, \quad \text{where } \mathcal{A}_{\text{sort}} := \left\{ \sigma \in \mathbb{R}^{d_{\min}} : (\sigma, \Pi) \in \mathcal{B} \right\}. \quad (35)$$

3.3 Analysis of the Error Bound

According to Proposition 2 and (35), for any \mathbf{W} satisfying $\text{dist}(\mathbf{W}, \mathcal{W}_G) \leq \delta_\sigma/2$, there exists a $\sigma^* \in \mathcal{A}_{\text{sort}}$ such that

$$\text{dist}(\mathbf{W}, \mathcal{W}_G) = \text{dist}(\mathbf{W}, \mathcal{W}_{\sigma^*}).$$

This observation simplifies our analysis, as it suffices to bound $\text{dist}(\mathbf{W}, \mathcal{W}_{\sigma^*})$ for each $\sigma^* \in \mathcal{A}_{\text{sort}}$. Note that $\mathbf{0} \in \mathcal{A}_{\text{sort}}$. For this case, we directly show the error bound as follows:

Proposition 3. *Suppose that Assumption 1 holds. Let $\mathcal{A}_{\text{sort}}$ be defined in (35) and consider $\sigma^* = \mathbf{0} \in \mathcal{A}_{\text{sort}}$. The following statements hold:*

- (i) *Suppose that $L = 2$ and Assumption 2 holds. For all \mathbf{W} satisfying*

$$\text{dist}(\mathbf{W}, \mathcal{W}_0) \leq \left(\frac{\sqrt{\lambda}}{2(\sqrt{\lambda} + y_1)} \min \left\{ \min_{i \in [s_{p_Y}]} |\lambda - y_i^2|, \lambda \right\} \right)^{1/2}, \quad (36)$$

it holds that

$$\text{dist}(\mathbf{W}, \mathcal{W}_0) \leq \frac{2(\sqrt{\lambda} + y_1)}{\sqrt{\lambda} \min \left\{ \min_{i \in [s_{p_Y}]} |\lambda - y_i^2|, \lambda \right\}} \|\nabla G(\mathbf{W})\|_F. \quad (37)$$

(ii) Suppose that $L \geq 3$. For all \mathbf{W} satisfying

$$\text{dist}(\mathbf{W}, \mathcal{W}_0) \leq \min \left\{ \left(\frac{\lambda}{3} \right)^{\frac{1}{2L-2}}, \left(\frac{\sqrt{\lambda}}{3y_1} \right)^{\frac{1}{L-2}} \right\}, \quad (38)$$

it holds that

$$\text{dist}(\mathbf{W}, \mathcal{W}_0) \leq \frac{3\sqrt{L}}{2\lambda} \|\nabla G(\mathbf{W})\|_F. \quad (39)$$

We defer the proof of this proposition to Appendix A. According to this proposition, we assume without loss of generality $\boldsymbol{\sigma}^* \neq \mathbf{0}$ from now on. Let $\mathbf{0} \neq \boldsymbol{\sigma}^* \in \mathcal{A}_{\text{sort}}$ be arbitrary and we define

$$\sigma_{\min}^* := \min \{ \sigma_i^* \neq 0 : i \in [d_{\min}] \}, \quad \sigma_{\max}^* := \max \{ \sigma_i^* : i \in [d_{\min}] \}, \quad r_{\sigma} := \|\boldsymbol{\sigma}^*\|_0, \quad (40)$$

where $\|\boldsymbol{\sigma}^*\|_0$ denotes the number of nonzero entries in $\boldsymbol{\sigma}$. From the definition of δ_{σ} in (33), we have $\delta_{\sigma} \leq \sigma_{\min}^*$. Now, we present a lemma and a corollary that establish some spectral properties of a point in the neighborhood of the set of critical points. Their proofs are provided in Section 4.2.

Lemma 3. Let $\mathbf{W} = (\mathbf{W}_1, \dots, \mathbf{W}_L)$ and $\mathbf{0} \neq \boldsymbol{\sigma}^* \in \mathcal{A}_{\text{sort}}$ be arbitrary such that

$$\text{dist}(\mathbf{W}, \mathcal{W}_{\boldsymbol{\sigma}^*}) < \frac{\sigma_{\min}^*}{2}. \quad (41)$$

The following statements hold:

(i) It holds that

$$\frac{\sigma_{\max}^*}{2} \leq \|\mathbf{W}_l\| \leq \frac{3\sigma_{\max}^*}{2}, \quad \forall l \in [L], \quad (42)$$

$$\sigma_i(\mathbf{W}_l) \geq \frac{\sigma_{\min}^*}{2}, \quad \forall l \in [L], \quad i \in [r_{\sigma}]. \quad (43)$$

(ii) It holds that

$$\|\mathbf{W}_{l+1}^T \mathbf{W}_{l+1} - \mathbf{W}_l \mathbf{W}_l^T\|_F \leq \frac{3\sqrt{2}\sigma_{\max}^*}{4\lambda} \|\nabla G(\mathbf{W})\|_F, \quad \forall l \in [L-1]. \quad (44)$$

(iii) It holds that

$$|\sigma_i(\mathbf{W}_l) - \sigma_i(\mathbf{W}_{l+1})| \leq \frac{3\sqrt{2}\sigma_{\max}^*}{4\lambda\sigma_{\min}^*} \|\nabla G(\mathbf{W})\|_F, \quad \forall l \in [L-1], \quad i \in [r_{\sigma}]. \quad (45)$$

(iv) It holds that

$$\|\mathbf{W}_{j:i}^T \mathbf{W}_{j:i} - (\mathbf{W}_i^T \mathbf{W}_i)^{j-i+1}\| \leq \frac{(j-i)(j-i+1)}{2\sqrt{2}\lambda} \left(\frac{3\sigma_{\max}^*}{2} \right)^{2j-2i+1} \|\nabla G(\mathbf{W})\|_F, \quad \forall i \leq j, \quad (46)$$

$$\|\mathbf{W}_{i:j} \mathbf{W}_{i:j}^T - (\mathbf{W}_i \mathbf{W}_i^T)^{i-j+1}\| \leq \frac{(i-j)(i-j+1)}{2\sqrt{2}\lambda} \left(\frac{3\sigma_{\max}^*}{2} \right)^{2i-2j+1} \|\nabla G(\mathbf{W})\|_F, \quad \forall i \geq j. \quad (47)$$

According to Lemma 2, each critical point satisfies (26). Leveraging this observation and Lemma 3, we can further show that any point in the neighborhood of the critical point set approximately satisfies (26) with the deviation bounded by its gradient norm.

Corollary 2. Consider the setting in Lemma 3. It holds for each $l \in [L]$ that

$$\left\| (\mathbf{W}_l \mathbf{W}_l^T)^{L-1} \mathbf{W}_l - \sqrt{\lambda} \mathbf{W}_{L:l+1}^T \mathbf{Y} \mathbf{W}_{l-1:1}^T + \lambda \mathbf{W}_l \right\|_F \leq c_1 \|\nabla G(\mathbf{W})\|_F, \quad (48)$$

where

$$c_1 := \max_{l \in [L]} \left\{ \left(\frac{3\sigma_{\max}^*}{2} \right)^{2L-2} \frac{(L-l)(L-l+1) + (l-1)l}{2\sqrt{2}\lambda} + \frac{1}{2} \right\}. \quad (49)$$

To handle the repeated singular values in $\boldsymbol{\sigma}^* \in \mathcal{A}_{\text{sort}}$, let $p \geq 1$ be the number of distinct positive singular values. Recall from (40) that r_σ denotes the number of positive singular values of $\boldsymbol{\sigma}^*$. Then there exist indices t_0, t_1, \dots, t_p such that $0 = t_0 < t_1 < \dots < t_p = r_\sigma$ and

$$\sigma_{t_0+1}^* = \dots = \sigma_{t_1}^* > \sigma_{t_1+1}^* = \dots = \sigma_{t_2}^* > \dots > \sigma_{t_{p-1}+1}^* = \dots = \sigma_{t_p}^* > 0. \quad (50)$$

Let $g_i := t_i - t_{i-1}$ be the multiplicity of the i -th largest positive value of $\boldsymbol{\sigma}^*$ for each $i \in [p]$ and $g_{\max} := \max\{g_i : i \in [p]\}$. Obviously, we have $r_\sigma = \sum_{i=1}^p g_i$. Moreover, let

$$\mathbf{W}_l = \mathbf{U}_l \boldsymbol{\Sigma}_l \mathbf{V}_l^T = \begin{bmatrix} \mathbf{U}_l^{(1)} & \dots & \mathbf{U}_l^{(p)} & \mathbf{U}_l^{(p+1)} \end{bmatrix} \begin{bmatrix} \boldsymbol{\Sigma}_l^{(1)} & \dots & \mathbf{0} & \mathbf{0} \\ \vdots & \ddots & \vdots & \vdots \\ \mathbf{0} & \dots & \boldsymbol{\Sigma}_l^{(p)} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \dots & \boldsymbol{\Sigma}_l^{(p+1)} \end{bmatrix} \begin{bmatrix} \mathbf{V}_l^{(1)T} \\ \vdots \\ \mathbf{V}_l^{(p)T} \\ \mathbf{V}_l^{(p+1)T} \end{bmatrix} \quad (51)$$

be an SVD of \mathbf{W}_l for each $l \in [L]$, where $\mathbf{U}_l \in \mathcal{O}^{d_l}$ with $\mathbf{U}_l^{(i)} \in \mathbb{R}^{d_l \times g_i}$ for each $i \in [p]$ and $\mathbf{U}_l^{(p+1)} \in \mathbb{R}^{d_l \times (d_l - r_\sigma)}$, $\boldsymbol{\Sigma}_l \in \mathbb{R}^{d_l \times d_{l-1}}$ with decreasing singular values, $\boldsymbol{\Sigma}_l^{(i)} \in \mathbb{R}^{g_i \times g_i}$ for each $i \in [p]$, and $\boldsymbol{\Sigma}_l^{(p+1)} \in \mathbb{R}^{(d_l - r_\sigma) \times (d_{l-1} - r_\sigma)}$, and $\mathbf{V}_l \in \mathcal{O}^{d_{l-1}}$ with $\mathbf{V}_l^{(i)} \in \mathbb{R}^{d_{l-1} \times g_i}$ for each $i \in [p]$ and $\mathbf{V}_l^{(p+1)} \in \mathbb{R}^{d_{l-1} \times (d_{l-1} - r_\sigma)}$.

With the above setup, we are ready to show that, for any point in the neighborhood of the set of critical points, the product of \mathbf{U}_{l-1} and \mathbf{V}_l is close to a block diagonal matrix, where the diagonal blocks are orthogonal matrices.

Proposition 4. Let $\mathbf{W} = (\mathbf{W}_1, \dots, \mathbf{W}_L)$ and $\mathbf{0} \neq \boldsymbol{\sigma}^* \in \mathcal{A}_{\text{sort}}$ be arbitrary such that

$$\text{dist}(\mathbf{W}, \mathcal{W}_{\boldsymbol{\sigma}^*}) \leq \frac{\delta_\sigma}{3}. \quad (52)$$

For each $l \in \{2, \dots, L\}$, there exist matrices $\mathbf{T}_l^{(i)} \in \mathcal{O}^{g_i}$ for all $i \in [p]$ and $\mathbf{T}_l^{(p+1)} \in \mathcal{O}^{d_{l-1} - r_\sigma}$ such that

$$\left\| \mathbf{U}_{l-1}^T \mathbf{V}_l - \text{BlkDiag} \left(\mathbf{T}_l^{(1)}, \dots, \mathbf{T}_l^{(p)}, \mathbf{T}_l^{(p+1)} \right) \right\|_F \leq \frac{9\sigma_{\max}^*}{4\delta_\sigma \lambda \sigma_{\min}^*} \|\nabla G(\mathbf{W})\|_F. \quad (53)$$

where \mathbf{U}_l and \mathbf{V}_l for each $l \in [L]$ are defined in (51).

The proof of this proposition is deferred to Section 4.2. Next, using the above proposition, we show that for weight matrices in the neighborhood of the set of critical points, the singular matrices $\boldsymbol{\Sigma}_1$ and $\boldsymbol{\Sigma}_L$ satisfy the following spectral inequalities.

Lemma 4. Consider the setting of Lemma 3. Suppose that there exist matrices $\mathbf{T}_l^{(i)} \in \mathcal{O}^{g_i}$ for all $i \in [p]$ and $\mathbf{T}_l^{(p+1)} \in \mathcal{O}^{d_{l-1}-r_\sigma}$ such that (53) holds for all $l = 2, \dots, L$. Then, we have

$$\left\| (\boldsymbol{\Sigma}_1 \boldsymbol{\Sigma}_1^T)^{L-1} \boldsymbol{\Sigma}_1 + \lambda \boldsymbol{\Sigma}_1 - \sqrt{\lambda} \text{BlkDiag}(\mathbf{A}_1, \dots, \mathbf{A}_p, \mathbf{A}_{p+1}) \mathbf{U}_L^T \mathbf{Y} \mathbf{V}_1 \right\|_F \leq c_2 \|\nabla G(\mathbf{W})\|_F, \quad (54)$$

$$\left\| (\boldsymbol{\Sigma}_L \boldsymbol{\Sigma}_L^T)^{L-1} \boldsymbol{\Sigma}_L + \lambda \boldsymbol{\Sigma}_L - \sqrt{\lambda} \mathbf{U}_L^T \mathbf{Y} \mathbf{V}_1 \text{BlkDiag}(\mathbf{B}_1, \dots, \mathbf{B}_p, \mathbf{B}_{p+1}) \right\|_F \leq c_2 \|\nabla G(\mathbf{W})\|_F, \quad (55)$$

where c_1 is defined in (49) and

$$\mathbf{A}_i := \left(\prod_{l=1}^{L-1} \mathbf{T}_{l+1}^{(i)} \right) (\boldsymbol{\Sigma}_1^{(i)})^{L-1}, \quad \forall i \in [p], \quad \mathbf{A}_{p+1} := \prod_{l=1}^{L-1} \mathbf{T}_{l+1}^{(p+1)} \boldsymbol{\Sigma}_{l+1}^{(p+1)T}, \quad (56)$$

$$\mathbf{B}_i := \left(\prod_{l=1}^{L-1} \mathbf{T}_{l+1}^{(i)} \right) (\boldsymbol{\Sigma}_L^{(i)})^{L-1}, \quad \forall i \in [p], \quad \mathbf{B}_{p+1} := \prod_{l=1}^{L-1} \boldsymbol{\Sigma}_l^{(p+1)T} \mathbf{T}_{l+1}^{(p+1)}, \quad (57)$$

$$\eta_1 := \frac{\sigma_{\max}^*}{\sigma_{\min}^*} \left(\frac{3\sqrt{2}\sigma_{\min}^*}{4\lambda} + \frac{81\sigma_{\max}^{*2}}{8\delta_\sigma \lambda} + \frac{9\sqrt{2g_{\max}}L\sigma_{\max}^*}{4\lambda} \right), \quad (58)$$

$$c_2 := \left(\frac{3}{2}\sigma_{\max}^* \right)^L \frac{3y_1 L}{2\sqrt{\lambda}\delta_\sigma \sigma_{\min}^*} + c_1 + y_1 p \sqrt{\lambda} \left(\frac{3\sigma_{\max}^*}{2} \right)^{L-2} \left(\frac{L^2 \eta_1}{2\sigma_{\min}^*} + \frac{3\sqrt{2g_{\max}}L^2\sigma_{\max}^*}{2\lambda\sigma_{\min}^*} \right). \quad (59)$$

The proof of this lemma is deferred to Section 4.2. Now, we bound the singular values and the associated singular vectors of weight matrices by the gradient norm when they lie in the neighborhood of the critical point set. Notably, according to $\lambda = \lambda_1 \cdots \lambda_L$ in (13), we obtain that when $L = 2$, (9) in Assumption 2 is equivalent to

$$\lambda \neq y_i^2, \quad \forall i \in [p_Y]. \quad (60)$$

When $L \geq 3$, (10) in Assumption 2 is equivalent to

$$\lambda \neq y_i^{2(L-1)} \left(\left(\frac{L-2}{L} \right)^{\frac{L}{2(L-1)}} + \left(\frac{L}{L-2} \right)^{\frac{L-2}{2(L-1)}} \right)^{-2(L-1)}, \quad \forall i \in [p_Y]. \quad (61)$$

Proposition 5. Let $\mathbf{0} \neq \boldsymbol{\sigma}^* \in \mathcal{A}_{\text{sort}}$ be arbitrary. The following statements hold:

(i) Suppose that $L = 2$ and (60) holds. Then for all \mathbf{W} satisfying

$$\text{dist}(\mathbf{W}, \mathcal{W}_{\boldsymbol{\sigma}^*}) \leq \min \left\{ \frac{\delta_\sigma}{3}, \frac{\sqrt{\lambda}}{\sqrt{3(\sqrt{\lambda} + y_1)}} \left(\min \left\{ \min_{i \in [s_{p_Y}]} |\sqrt{\lambda} - y_i|, \sqrt{\lambda} \right\} \right)^{\frac{1}{2}} \right\}, \quad (62a)$$

$$\|\nabla G(\mathbf{W})\|_F \leq \frac{\sqrt{2\lambda} \min \left\{ \min_{i \in [s_{p_Y}]} |\sqrt{\lambda} - y_i|, \sqrt{\lambda} \right\} \sigma_{\min}^*}{12c_2}, \quad (62b)$$

it holds for $l = 1, 2$ that

$$\sigma_i(\mathbf{W}_l) \leq c_3 \|\nabla G(\mathbf{W})\|_F, \quad \forall i = r_\sigma + 1, \dots, \min\{d_l, d_{l-1}\}, \quad (63)$$

where

$$c_3 := \frac{6c_2(y_1 + \sqrt{\lambda})}{\lambda \min \left\{ \min_{i \in [s_{p_Y}]} |\sqrt{\lambda} - y_i|, \sqrt{\lambda} \right\}}.$$

(ii) Suppose that $L = 3$ and (61) holds. Then for all \mathbf{W} satisfying

$$\text{dist}(\mathbf{W}, \mathcal{W}_{\sigma^*}) \leq \min \left\{ \frac{\delta_\sigma}{3}, \left(\frac{\sqrt{\lambda}}{2y_1} \right)^{1/(L-2)} \right\}, \quad (64)$$

it holds for all $l \in [L]$ that

$$\sigma_i(\mathbf{W}_l) \leq c_3 \|\nabla G(\mathbf{W})\|_F, \quad \forall i = r_\sigma + 1, \dots, \min\{d_l, d_{l-1}\}, \quad (65)$$

where

$$\eta_2 := c_1 + \left(\frac{3}{2} \sigma_{\max}^* \right)^L \frac{3(L-1)y_1}{2\delta_\sigma \sqrt{\lambda} \sigma_{\min}^*} \quad \text{and} \quad c_3 := \frac{2\eta_2}{\lambda}. \quad (66)$$

The proof of this proposition is deferred to Section 4.2. The above proposition bounds the smallest $\min\{d_l, d_{l-1}\} - r_\sigma$ singular values by the gradient norm. Next, we proceed to bound the leading r_σ singular values and all singular vectors by the gradient norm. For ease of exposition, we introduce an auxiliary function as follows:

$$\varphi(x) := \frac{x^{2L-1} + \lambda x}{\sqrt{\lambda} x^{L-1}}, \quad \forall x \neq 0. \quad (67)$$

Proposition 6. Let $\mathbf{0} \neq \sigma^* \in \mathcal{A}_{\text{sort}}$ be arbitrary and \mathbf{W} be arbitrary such that

$$\text{dist}(\mathbf{W}, \mathcal{W}_{\sigma^*}) \leq \min\{\delta_1, \delta_2\}, \quad (68a)$$

$$\|\nabla G(\mathbf{W})\|_F \leq \frac{\delta_y \sqrt{\lambda} (\sigma_{\min}^*)^{L-1}}{3 \cdot 2^{L-1} \sqrt{\eta_3^2 + \eta_4^2}}, \quad (68b)$$

where η_3, η_4 are respectively defined in (141) and (145), and δ_1, δ_2 are respectively described in (158) and (161). Suppose in addition Assumption 1 holds, (60) and (62) hold for $L = 2$, and (61) and (64) hold for $L \geq 3$. Then the following statements hold:

(i) There exist orthogonal matrices $\hat{\mathbf{U}}_L^{(i)} \in \mathcal{O}^{h_i}$ for each $i \in [p_Y]$, $\hat{\mathbf{U}}_L^{(p_Y+1)} \in \mathcal{O}^{d_L - r_Y}$, $\hat{\mathbf{V}}_1^{(p_Y+1)} \in \mathcal{O}^{d_0 - r_Y}$, $\mathbf{T}_l^{(i)} \in \mathcal{O}^{g_i}$ for each $i \in [p]$, $\mathbf{P} \in \mathcal{O}^{d_L - r_\sigma}$, $\mathbf{Q} \in \mathcal{O}^{d_0 - r_\sigma}$, and a permutation matrix $\mathbf{\Pi} \in \mathcal{P}^{d_{\min}}$ satisfying $(\sigma^*, \mathbf{\Pi}^T) \in \mathcal{B}$ such that

$$\|\tilde{\mathbf{U}}_L - \mathbf{U}_L\|_F \leq c_4 \|\nabla G(\mathbf{W})\|_F, \quad \|\tilde{\mathbf{V}}_1 - \mathbf{V}_1\|_F \leq c_4 \|\nabla G(\mathbf{W})\|_F, \quad (69)$$

where $c_4 > 0$ is a positive constant and

$$\tilde{\mathbf{U}}_L := \text{BlkDiag} \left(\hat{\mathbf{U}}_L^{(1)}, \dots, \hat{\mathbf{U}}_L^{(p_Y)}, \hat{\mathbf{U}}_L^{(p_Y+1)} \right) \text{BlkDiag}(\mathbf{\Pi}, \mathbf{I}_{d_L - d_{\min}}) \text{BlkDiag}(\mathbf{I}_{r_\sigma}, \mathbf{P}), \quad (70)$$

$$\tilde{\mathbf{V}}_1 := \text{BlkDiag} \left(\hat{\mathbf{U}}_L^{(1)}, \dots, \hat{\mathbf{U}}_L^{(p_Y)}, \hat{\mathbf{V}}_L^{(p_Y+1)} \right) \text{BlkDiag}(\mathbf{\Pi}, \mathbf{I}_{d_0 - d_{\min}}) \text{BlkDiag}(\mathbf{I}_{r_\sigma}, \mathbf{Q}) \hat{\mathbf{T}}, \quad (71)$$

$$\hat{\mathbf{T}} := \text{BlkDiag} \left(\prod_{l=2}^L \mathbf{T}_l^{(1)}, \dots, \prod_{l=2}^L \mathbf{T}_l^{(p)}, \mathbf{I}_{d_0 - r_\sigma} \right). \quad (72)$$

(ii) For each $l \in [L]$, there exists a constant $c_5 > 0$ such that

$$|\sigma_i(\mathbf{W}_l) - \sigma_i^*| \leq c_5 \|\nabla G(\mathbf{W})\|_F, \quad \forall i \in [r_\sigma]. \quad (73)$$

Equipped with all the above lemmas and propositions, we are ready to prove the error bound of Problem (13).

Theorem 4. Suppose that Assumptions 1 and (60) (resp., (61)) hold for $L = 2$ (resp., $L \geq 3$). Let \mathbf{W} be arbitrary such that (68) and (62) (resp., (64)) hold for $L = 2$ (resp., $L \geq 3$). It holds that

$$\text{dist}(\mathbf{W}, \mathcal{W}_G) \leq \sqrt{L} \left(\frac{9\sigma_{\max}^2}{4\delta_\sigma \lambda \sigma_{\min}^*} + c_3 \sqrt{d_{\max} - r_\sigma} + c_4 \sigma_{\max}^* + c_5 \sqrt{r_\sigma} \right) \|\nabla G(\mathbf{W})\|_F.$$

Proof. According to (62a) or (64), we have $\text{dist}(\mathbf{W}, \mathcal{W}) \leq \delta_\sigma/3$. This, together with Proposition 4, yields that there exist matrices $\mathbf{T}_l^{(i)} \in \mathcal{O}^{g_i}$ for each $i \in [p]$ and $\mathbf{T}_l^{(p+1)} \in \mathcal{O}^{d_{l-1}-r_\sigma}$ such that (53) holds. Moreover, it follows from Proposition 5, Proposition 6, we have (65), (69), (73) hold, where $\tilde{\mathbf{U}}_L$, $\tilde{\mathbf{V}}_1$, and $\hat{\mathbf{T}}$ are respectively defined in (70), (71), and (72). Since we have $\text{dist}(\mathbf{W}, \mathcal{W}_G) \leq \delta_\sigma/3$, it follows from Proposition 1 that there exists some $\boldsymbol{\sigma}^* \in \mathcal{A}_{\text{sort}}$ such that

$$\text{dist}(\mathbf{W}, \mathcal{W}_G) = \text{dist}(\mathbf{W}, \mathcal{W}_{\boldsymbol{\sigma}^*}).$$

Recall that $\mathbf{U}_l \in \mathcal{O}^{d_l}$ and $\mathbf{V}_l \in \mathcal{O}^{d_{l-1}}$ for all $l \in [L]$ are introduced in (51). For ease of exposition, we define $\boldsymbol{\Sigma}_l^* := \text{BlkDiag}(\text{diag}(\boldsymbol{\sigma}^*), \mathbf{0}) \in \mathbb{R}^{d_l \times d_{l-1}}$ and

$$\hat{\mathbf{W}}_1 := \mathbf{U}_1 \boldsymbol{\Sigma}_1^* \tilde{\mathbf{V}}_1^T, \quad \hat{\mathbf{W}}_l := \mathbf{U}_l \boldsymbol{\Sigma}_l^* \mathbf{V}_l^T, \quad l = 2, \dots, L-1, \quad \hat{\mathbf{W}}_L := \tilde{\mathbf{U}}_L \boldsymbol{\Sigma}_L^* \mathbf{V}_L^T. \quad (74)$$

Now, we compute

$$\begin{aligned} \|\hat{\mathbf{W}}_1 - \mathbf{W}_1\|_F &= \left\| \mathbf{U}_1 \boldsymbol{\Sigma}_1^* \tilde{\mathbf{V}}_1^T - \mathbf{U}_1 \boldsymbol{\Sigma}_1 \mathbf{V}_1^T \right\|_F \leq \left\| \mathbf{U}_1 \boldsymbol{\Sigma}_1^* \tilde{\mathbf{V}}_1^T - \mathbf{U}_1 \boldsymbol{\Sigma}_1^* \mathbf{V}_1^T \right\|_F + \\ &\quad \left\| \mathbf{U}_1 \boldsymbol{\Sigma}_1^* \mathbf{V}_1^T - \mathbf{U}_1 \boldsymbol{\Sigma}_1 \mathbf{V}_1^T \right\|_F \leq \sigma_{\max}^* \|\tilde{\mathbf{V}}_1 - \mathbf{V}_1\|_F + \|\boldsymbol{\Sigma}_1^* - \boldsymbol{\Sigma}_1\|_F \\ &\leq \left(c_4 \sigma_{\max}^* + c_3 \sqrt{d_{\max} - r_\sigma} + c_5 \sqrt{r_\sigma} \right) \|\nabla G(\mathbf{W})\|_F, \end{aligned} \quad (75)$$

where the last inequality uses (69), (65), and (73). Using the same argument, we obtain

$$\|\hat{\mathbf{W}}_L - \mathbf{W}_L\|_F \leq \left(c_4 \sigma_{\max}^* + c_3 \sqrt{d_{\max} - r_\sigma} + c_5 \sqrt{r_\sigma} \right) \|\nabla G(\mathbf{W})\|_F. \quad (76)$$

For all $l = 2, \dots, L-1$, we compute

$$\|\hat{\mathbf{W}}_l - \mathbf{W}_l\|_F = \|\boldsymbol{\Sigma}_l^* - \boldsymbol{\Sigma}_l\|_F \leq \left(c_3 \sqrt{d_{\max} - r_\sigma} + c_5 \sqrt{r_\sigma} \right) \|\nabla G(\mathbf{W})\|_F,$$

where the inequality uses (65) and (73). This, together with (75) and (76), yields

$$\|\mathbf{W} - \hat{\mathbf{W}}\|_F \leq \sqrt{L} \left(c_4 \sigma_{\max}^* + c_3 \sqrt{d_{\max} - r_\sigma} + c_5 \sqrt{r_\sigma} \right) \|\nabla G(\mathbf{W})\|_F. \quad (77)$$

Next, we further define

$$\mathbf{T}_l := \text{BlkDiag}(\mathbf{T}_l^{(1)}, \dots, \mathbf{T}_l^{(p)}, \mathbf{T}_l^{(p+1)}), \quad l = 2, \dots, L, \quad (78)$$

$$\mathbf{W}_1^* := \mathbf{Q}_2 \boldsymbol{\Sigma}_1^* \hat{\mathbf{T}}^T \tilde{\mathbf{V}}_1^T, \quad \mathbf{W}_l^* = \mathbf{Q}_{l+1} \boldsymbol{\Sigma}_l^* \mathbf{Q}_l^T, \quad l = 2, \dots, L-1, \quad \mathbf{W}_L^* = \tilde{\mathbf{U}}_L \boldsymbol{\Sigma}_L^* \mathbf{Q}_L^T, \quad (79)$$

where

$$\mathbf{Q}_l := \mathbf{V}_l \text{BlkDiag} \left(\left(\prod_{k=l+1}^L \mathbf{T}_k^{(1)} \right), \dots, \left(\prod_{k=l+1}^L \mathbf{T}_k^{(p)} \right), \mathbf{I}_{d_{l-1}-r_\sigma} \right), \quad l = 2, \dots, L-1, \quad \mathbf{Q}_L := \mathbf{V}_L.$$

Using Proposition 1, one can verify that

$$\mathbf{W}^* = (\mathbf{W}_1^*, \mathbf{W}_2^*, \dots, \mathbf{W}_L^*) \in \mathcal{W}_{\sigma^*}. \quad (80)$$

For each $l = 2, \dots, L - 2$, we compute

$$\begin{aligned} \|\hat{\mathbf{W}}_l - \mathbf{W}_l^*\|_F &= \|\mathbf{U}_l \Sigma_l^* \mathbf{V}_l^T - \mathbf{Q}_{l+1} \Sigma_l^* \mathbf{Q}_l^T\|_F = \|\mathbf{U}_l \Sigma_l^* - \mathbf{Q}_{l+1} \Sigma_l^* \mathbf{Q}_l^T \mathbf{V}_l\|_F \\ &= \left\| \mathbf{V}_{l+1}^T \mathbf{U}_l \Sigma_l^* - \text{BlkDiag} \left(\sigma_{t_1}^* (\mathbf{T}_{l+1}^{(1)})^T, \dots, \sigma_{t_p}^* (\mathbf{T}_{l+1}^{(p)})^T, \mathbf{0} \right) \right\|_F \\ &= \|(\mathbf{V}_{l+1}^T \mathbf{U}_l - \mathbf{T}_{l+1}^T) \Sigma_l^*\|_F \leq \frac{9\sigma_{\max}^{*2}}{4\delta_\sigma \lambda \sigma_{\min}^*} \|\nabla G(\mathbf{W})\|_F, \end{aligned} \quad (81)$$

where the third equality uses $\sigma^* = (\sigma_{t_1} \mathbf{I}_{g_1}, \dots, \sigma_{t_p} \mathbf{I}_{g_p}, \mathbf{0}_{d_{\min} - r_\sigma})$ due to (50), the last equality follows from the diagonal block forms of \mathbf{T}_l and Σ_l^* , and the inequality holds because of (53). Using (74) and (79), we compute

$$\begin{aligned} \|\hat{\mathbf{W}}_1 - \mathbf{W}_1^*\|_F &= \|\mathbf{U}_1 \Sigma_1^* \tilde{\mathbf{V}}_1^T - \mathbf{Q}_2 \Sigma_1^* \hat{\mathbf{T}}^T \tilde{\mathbf{V}}_1^T\|_F = \|\mathbf{V}_2^T \mathbf{U}_1 \Sigma_1^* - \mathbf{V}_2^T \mathbf{Q}_2 \Sigma_1^* \hat{\mathbf{T}}^T\|_F \\ &\leq \|\mathbf{V}_2^T \mathbf{U}_1 \Sigma_1^* - \mathbf{T}_2^T \Sigma_1^*\|_F + \|\mathbf{T}_2^T \Sigma_1^* - \mathbf{V}_2^T \mathbf{Q}_2 \Sigma_1^* \hat{\mathbf{T}}^T\|_F \\ &= \|(\mathbf{U}_1^T \mathbf{V}_2 - \mathbf{T}_2)^T \Sigma_1^*\|_F + \|\Sigma_1^* - \mathbf{T}_2 \mathbf{V}_2^T \mathbf{Q}_2 \Sigma_1^* \hat{\mathbf{T}}^T\|_F \\ &= \|(\mathbf{U}_1^T \mathbf{V}_2 - \mathbf{T}_2)^T \Sigma_1^*\|_F \leq \frac{9\sigma_{\max}^{*2}}{4\delta_\sigma \lambda \sigma_{\min}^*} \|\nabla G(\mathbf{W})\|_F, \end{aligned}$$

where the last equality follows from the forms of \mathbf{T}_l in (78), $\Sigma_l^* := \text{BlkDiag}(\text{diag}(\sigma^*), \mathbf{0})$ with $\sigma^* = (\sigma_{t_1} \mathbf{I}_{g_1}, \dots, \sigma_{t_p} \mathbf{I}_{g_p}, \mathbf{0}_{d_{\min} - r_\sigma})$, \mathbf{Q}_2 , and $\hat{\mathbf{T}}$, and the last inequality uses (53). This, together with (81) and the fact that $\mathbf{W}_L^* = \hat{\mathbf{W}}_L$ due to $\mathbf{Q}_L := \mathbf{V}_L$, yields

$$\|\hat{\mathbf{W}} - \mathbf{W}^*\|_F \leq \frac{9\sigma_{\max}^{*2}}{4\delta_\sigma \lambda \sigma_{\min}^*} \sqrt{L-1} \|\nabla G(\mathbf{W})\|_F. \quad (82)$$

Finally, using (80), we obtain

$$\begin{aligned} \text{dist}(\mathbf{W}, \mathcal{W}_{\sigma^*}) &\leq \|\mathbf{W} - \mathbf{W}^*\|_F \leq \|\mathbf{W} - \hat{\mathbf{W}}\|_F + \|\hat{\mathbf{W}} - \mathbf{W}^*\|_F \\ &\leq \sqrt{L} \left(\frac{9\sigma_{\max}^{*2}}{4\delta_\sigma \lambda \sigma_{\min}^*} + c_3 \sqrt{d_{\max} - r_\sigma} + c_4 \sigma_{\max}^* + c_5 \sqrt{r_\sigma} \right) \|\nabla G(\mathbf{W})\|_F, \end{aligned}$$

where the last inequality uses (77) and (82). Then, we complete the proof. \square

Remark 1. Now, we highlight our technical contributions for establishing the error bound. Prior works [30, 42] derived a closed-form expression of the critical point set, denoted by \mathcal{X} , of their considered non-convex problems. Leveraging the favorable structure of \mathcal{X} , they explicitly computed the distance $\text{dist}(\mathbf{X}, \mathcal{X})$ from a point \mathbf{X} to \mathcal{X} . In contrast, our analysis faces a significant challenge: the distance $\text{dist}(\mathbf{W}, \mathcal{W}_G)$ cannot be explicitly computed due to the complicated structure of \mathcal{W}_G in Theorem 3. To address this challenge, we construct an intermediate point $\hat{\mathbf{W}}$ in (74), which uses the singular vectors \mathbf{W} and singular values of $\mathbf{W}^* \in \mathcal{W}_G$. This enables us to bridge the gap between \mathbf{W} and \mathcal{W}_G . Then, we respectively bound $\text{dist}(\hat{\mathbf{W}}, \mathcal{W}_G)$ and $\|\mathbf{W} - \hat{\mathbf{W}}\|_F$ by the gradient norm at \mathbf{W} . Finally, we prove the error bound.

In Theorem 4, the local neighborhood condition involves both the gradient norm of \mathbf{W} (see (62b) when $L = 2$ and (68b)) and its distance to the critical point set. However, the local

neighborhood condition in Theorem 2 depends only on the distance to the critical point set. To prove Theorem 2, we leverage the continuity of the gradient of $G(\mathbf{W})$ to transform the condition involving the gradient norm into the distance between \mathbf{W} and the critical point set.

Proof of Theorem 2. According to (14) and Theorem 3, we obtain that the gradient norm is a continuous function with respect to \mathbf{W} and \mathcal{W}_G is a compact set. This implies that there exists a positive constant δ_g such that (62b) for $L = 2$ and (68b) hold whenever $\text{dist}(\mathbf{W}, \mathcal{W}_G) \leq \delta_g$. Since Assumption 2 holds, we have (60) (resp., (61)) hold for $L = 2$ (resp., $L \geq 3$). These, together with Theorem 4, implies that there exist constants $\epsilon, \kappa > 0$ such that for all \mathbf{W} satisfying $\text{dist}(\mathbf{W}, \mathcal{W}_G) \leq \epsilon$, we have

$$\text{dist}(\mathbf{W}, \mathcal{W}_G) \leq \kappa \|\nabla G(\mathbf{W})\|_F.$$

Using this and (ii) of Lemma 1, Theorem 2 holds with $\epsilon_1 = \epsilon/\sqrt{\lambda_{\max}}$ and $\kappa_1 = \kappa\lambda/\lambda_{\min}$. \square

4 Proofs of Technical Results

In this section, we present detailed proofs for the lemmas and propositions introduced in Section 3. Notably, the main tools employed in these proofs are basic and largely self-contained, relying minimally on advanced results from the existing literature.

4.1 Proofs of the Set of Critical Points

In this subsection, we provide detailed proofs of Lemma 2, Proposition 1, and Proposition 2 in Section 3.2. We begin with the proof of Lemma 2, which establishes the balanced structure of the weight matrices at each critical point.

Proof of Lemma 2. (i) According to $\nabla_{\mathbf{W}_l} G(\mathbf{W}) = \mathbf{0}$, we have $\nabla_{\mathbf{W}_l} G(\mathbf{W}) \mathbf{W}_l^T - \mathbf{W}_{l+1}^T \nabla_{\mathbf{W}_{l+1}} G(\mathbf{W}) = \mathbf{0}$ for all $l \in [L-1]$. This, together with (14), implies (25).

(ii) Recursively using (25), we have

$$\begin{aligned} \mathbf{W}_{L:l+1}^T \mathbf{W}_{L:l} \mathbf{W}_{l-1:l}^T &= \mathbf{W}_{l+1}^T \dots \mathbf{W}_{L-1}^T \mathbf{W}_L^T \mathbf{W}_L \mathbf{W}_{L-1} \dots \mathbf{W}_2 \mathbf{W}_1 \mathbf{W}_1^T \mathbf{W}_2^T \dots \mathbf{W}_{l-1}^T \\ &= \mathbf{W}_{l+1}^T \dots (\mathbf{W}_{L-1}^T \mathbf{W}_{L-1})^2 \dots (\mathbf{W}_2 \mathbf{W}_2^T)^2 \dots \mathbf{W}_{l-1}^T = (\mathbf{W}_l \mathbf{W}_l^T)^{L-1} \mathbf{W}_l. \end{aligned}$$

Substituting this into (14) yields (26). \square

Now, we present the detailed proof of Proposition 1. Although this proof is complicated and lengthy, the main idea is rather straightforward. It lies in studying the SVD of \mathbf{W}_l for each l and characterizing their singular values and singular matrices.

Proof of Proposition 1. Suppose that $(\mathbf{W}_1, \dots, \mathbf{W}_L)$ satisfies (27). Using the structures of \mathbf{W}_l for each $l \in [L]$ in (27), we directly verify that (25) holds. For any $\mathbf{a} \in \mathcal{A}$, it follows from $\tilde{\Sigma}_Y = \text{diag}(y_1, \dots, y_{d_{\min}})$ that $\text{diag}^{2L-1}(\mathbf{a}) - \sqrt{\lambda} \text{diag}^{L-1}(\mathbf{a}) \tilde{\Sigma}_Y + \lambda \text{diag}(\mathbf{a}) = \mathbf{0}$. This, together with $\boldsymbol{\sigma} = \mathbf{\Pi} \mathbf{a}$ due to $(\boldsymbol{\sigma}, \mathbf{\Pi}) \in \mathcal{B}$, implies $\mathbf{\Pi}^T \text{diag}(\boldsymbol{\sigma}) \mathbf{\Pi} = \text{diag}(\mathbf{a})$ and

$$\text{diag}^{2L-1}(\boldsymbol{\sigma}) - \sqrt{\lambda} \text{diag}^{L-1}(\boldsymbol{\sigma}) \mathbf{\Pi} \tilde{\Sigma}_Y \mathbf{\Pi}^T + \lambda \text{diag}(\boldsymbol{\sigma}) = \mathbf{0}. \quad (83)$$

For each $l = 2, \dots, L-1$, we compute

$$\begin{aligned}
\nabla_{\mathbf{W}_l} G(\mathbf{W}) &\stackrel{(14),(25)}{=} (\mathbf{W}_l^T \mathbf{W}_l)^{L-1} \mathbf{W}_l + \lambda \mathbf{W}_l - \sqrt{\lambda} \mathbf{W}_{L:l+1}^T \mathbf{Y} \mathbf{W}_{l-1:1}^T \\
&\stackrel{(27)}{=} \mathbf{Q}_{l+1} (\boldsymbol{\Sigma}_l^T \boldsymbol{\Sigma}_l)^{L-1} \boldsymbol{\Sigma}_l \mathbf{Q}_l^T + \lambda \mathbf{Q}_{l+1} \boldsymbol{\Sigma}_l \mathbf{Q}_l^T - \sqrt{\lambda} \mathbf{Q}_{l+1} \\
&\quad \left(\prod_{j=l+1}^L \boldsymbol{\Sigma}_j^T \right) \text{BlkDiag}(\boldsymbol{\Pi}, \mathbf{I}_{d_0-d_{\min}}) \mathbf{Y} \text{BlkDiag}(\boldsymbol{\Pi}^T, \mathbf{I}_{d_0-d_{\min}}) \left(\prod_{j=1}^{l-1} \boldsymbol{\Sigma}_j^T \right) \mathbf{Q}_l^T \\
&= \mathbf{Q}_{l+1} \text{BlkDiag} \left(\text{diag}^{2L-1}(\boldsymbol{\sigma}) + \lambda \text{diag}(\boldsymbol{\sigma}) - \sqrt{\lambda} \text{diag}^{L-1}(\boldsymbol{\sigma}) \boldsymbol{\Pi} \tilde{\boldsymbol{\Sigma}}_Y \boldsymbol{\Pi}^T, \mathbf{0} \right) \mathbf{Q}_l^T = \mathbf{0},
\end{aligned}$$

where the third equality is due to $\boldsymbol{\Sigma}_l = \text{BlkDiag}(\text{diag}(\boldsymbol{\sigma}), \mathbf{0})$ for each $l \in [L]$ and $\boldsymbol{\Pi} \tilde{\boldsymbol{\Sigma}}_Y \boldsymbol{\Pi}^T$ are all diagonal matrices, and the last equality follows from $(\boldsymbol{\sigma}, \boldsymbol{\Pi}) \in \mathcal{B}$ and (83). For the case $l = 1$, we compute

$$\begin{aligned}
\nabla_{\mathbf{W}_1} G(\mathbf{W}) &\stackrel{(14)}{=} \mathbf{W}_{L:2}^T \mathbf{W}_{L:1} - \sqrt{\lambda} \mathbf{W}_{L:2}^T \mathbf{Y} + \lambda \mathbf{W}_1 \stackrel{(25)}{=} (\mathbf{W}_1 \mathbf{W}_1^T)^{L-1} \mathbf{W}_1 + \lambda \mathbf{W}_1 - \sqrt{\lambda} \mathbf{W}_{L:2}^T \mathbf{Y} \\
&\stackrel{(27)}{=} \mathbf{Q}_2 ((\boldsymbol{\Sigma}_1 \boldsymbol{\Sigma}_1^T)^{L-1} \boldsymbol{\Sigma}_1 + \lambda \boldsymbol{\Sigma}_1) \text{BlkDiag}(\boldsymbol{\Pi}, \mathbf{I}_{d_0-d_{\min}}) \text{BlkDiag}(\mathbf{O}_1, \dots, \mathbf{O}_{p_Y}, \mathbf{O}_{p_Y+1}) \\
&\quad - \sqrt{\lambda} \mathbf{Q}_2 \left(\prod_{l=2}^L \boldsymbol{\Sigma}_l \right) \text{BlkDiag}(\boldsymbol{\Pi}, \mathbf{I}_{d_L-d_{\min}}) \text{BlkDiag}(\mathbf{O}_1, \dots, \mathbf{O}_{p_Y}, \hat{\mathbf{O}}_{p_Y+1}) \mathbf{Y} \\
&= \mathbf{Q}_2 \left((\boldsymbol{\Sigma}_1 \boldsymbol{\Sigma}_1^T)^{L-1} \boldsymbol{\Sigma}_1 + \lambda \boldsymbol{\Sigma}_1 - \sqrt{\lambda} \left(\prod_{l=2}^L \boldsymbol{\Sigma}_l \right) \text{BlkDiag}(\boldsymbol{\Pi}, \mathbf{I}) \mathbf{Y} \text{BlkDiag}(\boldsymbol{\Pi}^T, \mathbf{I}) \right) \\
&\quad \text{BlkDiag}(\boldsymbol{\Pi}, \mathbf{I}) \text{BlkDiag}(\mathbf{O}_1, \dots, \mathbf{O}_{p_Y}, \mathbf{O}_{p_Y+1}) = \mathbf{0},
\end{aligned}$$

where the fourth inequality uses the block structure of \mathbf{Y} and $\text{BlkDiag}(\mathbf{O}_1, \dots, \mathbf{O}_{p_Y}, \hat{\mathbf{O}}_{p_Y+1}) \mathbf{Y} = \mathbf{Y} \text{BlkDiag}(\mathbf{O}_1, \dots, \mathbf{O}_{p_Y}, \mathbf{O}_{p_Y+1})$, and the last equality follows from $\boldsymbol{\Sigma}_l = \text{BlkDiag}(\text{diag}(\boldsymbol{\sigma}), \mathbf{0})$ for each $l \in [L]$ and (83). Using a similar argument, we prove $\nabla_{\mathbf{W}_L} G(\mathbf{W}) = \mathbf{0}$. Consequently, we conclude that $(\mathbf{W}_1, \dots, \mathbf{W}_L)$ is a critical point.

Conversely, suppose that $(\mathbf{W}_1, \dots, \mathbf{W}_L)$ is a critical point. According to Lemma 2, we obtain that \mathbf{W}_l for all $l \in [L]$ share the same rank denoted by r , which satisfies $r \leq d_{\min}$. For each $l \in [L]$, let

$$\mathbf{W}_l = \mathbf{U}_l \boldsymbol{\Sigma}_l \mathbf{V}_l^T \quad (84)$$

be an SVD of \mathbf{W}_l , where $\mathbf{U}_l \in \mathcal{O}^{d_l}$, $\mathbf{V}_l \in \mathcal{O}^{d_{l-1}}$, and $\boldsymbol{\Sigma}_l = \text{BlkDiag}(\tilde{\boldsymbol{\Sigma}}_l, \mathbf{0}) \in \mathbb{R}^{d_l \times d_{l-1}}$ with $\tilde{\boldsymbol{\Sigma}}_l \in \mathbb{R}^{r \times r}$ being a diagonal matrix with positive diagonal entries. This, together with (25), yields

$$\mathbf{U}_l \boldsymbol{\Sigma}_l \boldsymbol{\Sigma}_l^T \mathbf{U}_l^T = \mathbf{V}_{l+1} \boldsymbol{\Sigma}_{l+1}^T \boldsymbol{\Sigma}_{l+1} \mathbf{V}_{l+1}^T, \quad \forall l \in [L-1]. \quad (85)$$

Since the above both sides are eigenvalue decompositions of the same matrix, with eigenvalues in decreasing order, we have

$$\boldsymbol{\Sigma}_l \boldsymbol{\Sigma}_l^T = \boldsymbol{\Sigma}_{l+1}^T \boldsymbol{\Sigma}_{l+1}, \quad \forall l \in [L-1].$$

This implies that $\mathbf{W}_1, \dots, \mathbf{W}_L$ have the same positive singular values. Next, let $\{\sigma_i\}_{i=1}^r$ denote the positive singular values of \mathbf{W}_l for each $l \in [L]$ and p denote the number of distinct elements of positive singular values. In other words, there exist indices $\hat{s}_0, \hat{s}_1, \dots, \hat{s}_p$ such that $0 = \hat{s}_0 < \hat{s}_1 < \dots < \hat{s}_p = r$ and

$$\sigma_{\hat{s}_0+1} = \dots = \sigma_{\hat{s}_1} > \sigma_{\hat{s}_1+1} = \dots = \sigma_{\hat{s}_2} > \dots > \sigma_{\hat{s}_{p-1}+1} = \dots = \sigma_{\hat{s}_p} > 0.$$

Let $\hat{h}_i := \hat{s}_i - \hat{s}_{i-1}$ be the multiplicity of the i -th largest positive value for each $i \in [p]$. Obviously, we have $\sum_{i=1}^p \hat{h}_i = r$ and

$$\tilde{\Sigma}_l = \tilde{\Sigma} := \text{BlkDiag} \left(\sigma_{\hat{s}_1} \mathbf{I}_{\hat{h}_1}, \dots, \sigma_{\hat{s}_p} \mathbf{I}_{\hat{h}_p} \right) \in \mathbb{R}^{r \times r}. \quad (86)$$

Based on the above block form, we write \mathbf{U}_l and \mathbf{V}_l in (84) for each $l \in [L]$ as

$$\mathbf{U}_l = \left[\mathbf{U}_l^{(1)}, \dots, \mathbf{U}_l^{(p)}, \mathbf{U}_l^{(p+1)} \right], \quad \mathbf{V}_l = \left[\mathbf{V}_l^{(1)}, \dots, \mathbf{V}_l^{(p)}, \mathbf{V}_l^{(p+1)} \right], \quad (87)$$

where $\mathbf{U}_l^{(i)} \in \mathcal{O}^{d_l \times \hat{h}_i}$ and $\mathbf{V}_l^{(i)} \in \mathcal{O}^{d_{l-1} \times \hat{h}_i}$ for all $i \in [p]$, $\mathbf{U}_l^{(p+1)} \in \mathcal{O}^{d_l \times (d_l - r)}$, and $\mathbf{V}_l^{(p+1)} \in \mathcal{O}^{d_{l-1} \times (d_{l-1} - r)}$. This, together with (85), (86), and [41, Lemma 8(i)], implies that there exists orthogonal matrix $\mathbf{Q}_l^{(i)} \in \mathcal{O}^{\hat{h}_i}$ such that

$$\mathbf{U}_l^{(i)} = \mathbf{V}_{l+1}^{(i)} \mathbf{Q}_l^{(i)}, \quad \forall l \in [L-1], \quad i \in [p+1]. \quad (88)$$

Using this, along with (84), (86), and the commutativity of orthogonal and identity matrices, we compute

$$\mathbf{W}_{L:1} = \mathbf{U}_L \tilde{\Sigma}_L \dots \tilde{\Sigma}_1 \mathbf{Q} \mathbf{V}_1^T = \mathbf{U}_L \text{BlkDiag} \left(\tilde{\Sigma}^L, \mathbf{0}_{(d_L - r) \times (d_0 - r)} \right) \mathbf{Q} \mathbf{V}_1^T, \quad (89)$$

where $\prod_{l=L-1}^1 \mathbf{Q}_l^{(j)} := \mathbf{Q}_{L-1}^{(j)} \dots \mathbf{Q}_1^{(j)}$ for each $j \in [p]$ and

$$\mathbf{Q} = \text{BlkDiag} \left(\tilde{\mathbf{Q}}, \mathbf{I}_{d_0 - r} \right) = \text{BlkDiag} \left(\prod_{l=L-1}^1 \mathbf{Q}_l^{(1)}, \dots, \prod_{l=L-1}^1 \mathbf{Q}_l^{(p)}, \mathbf{I}_{d_0 - r} \right) \in \mathcal{O}^{d_0}. \quad (90)$$

Right-multiplying (26) by \mathbf{W}_L^T when $l = L$ and left-multiplying (26) by \mathbf{W}_1^T when $l = 1$, we obtain

$$(\mathbf{W}_L \mathbf{W}_L^T)^L - \sqrt{\lambda} \mathbf{Y} \mathbf{W}_{L:1}^T + \lambda \mathbf{W}_L \mathbf{W}_L^T = \mathbf{0}, \quad \text{and} \quad (\mathbf{W}_1^T \mathbf{W}_1)^L - \sqrt{\lambda} \mathbf{W}_{L:1}^T \mathbf{Y} + \lambda \mathbf{W}_1^T \mathbf{W}_1 = \mathbf{0},$$

respectively. Substituting (84), (86), and (89) into the above equalities, together with $\mathbf{U}_L \in \mathcal{O}^{d_L}$ and $\mathbf{V}_1 \in \mathcal{O}^{d_0}$, yields

$$\mathbf{U}_L^T \mathbf{Y} \mathbf{V}_1 \mathbf{Q}^T \text{BlkDiag} \left(\tilde{\Sigma}^L, \mathbf{0}_{(d_0 - r) \times (d_L - r)} \right) = \frac{1}{\sqrt{\lambda}} \text{BlkDiag} \left(\tilde{\Sigma}^{2L} + \lambda \tilde{\Sigma}^2, \mathbf{0}_{d_L - r} \right), \quad (91)$$

$$\mathbf{Q}^T \text{BlkDiag} \left(\tilde{\Sigma}^L, \mathbf{0}_{(d_0 - r) \times (d_L - r)} \right) \mathbf{U}_L^T \mathbf{Y} \mathbf{V}_1 = \frac{1}{\sqrt{\lambda}} \text{BlkDiag} \left(\tilde{\Sigma}^{2L} + \lambda \tilde{\Sigma}^2, \mathbf{0}_{d_0 - r} \right). \quad (92)$$

According to (92) and $\mathbf{Q} \in \mathcal{O}^{d_0}$, we obtain

$$\begin{aligned} \text{BlkDiag} \left(\tilde{\Sigma}^L, \mathbf{0}_{(d_0 - r) \times (d_L - r)} \right) \mathbf{U}_L^T \mathbf{Y} \mathbf{V}_1 &= \frac{1}{\sqrt{\lambda}} \mathbf{Q} \text{BlkDiag} \left(\tilde{\Sigma}^{2L} + \lambda \tilde{\Sigma}^2, \mathbf{0}_{d_0 - r} \right) \\ &= \frac{1}{\sqrt{\lambda}} \text{BlkDiag} \left(\tilde{\Sigma}^{2L} + \lambda \tilde{\Sigma}^2, \mathbf{0}_{d_0 - r} \right) \mathbf{Q}, \end{aligned}$$

where the second equality follows from the block diagonal structures of $\tilde{\Sigma}$ in (86) and \mathbf{Q} in (90). Right-multiplying on both sides of the above equality by \mathbf{Q}^T yields

$$\text{BlkDiag} \left(\tilde{\Sigma}^L, \mathbf{0}_{(d_0 - r) \times (d_L - r)} \right) \mathbf{U}_L^T \mathbf{Y} \mathbf{V}_1 \mathbf{Q}^T = \frac{1}{\sqrt{\lambda}} \text{BlkDiag} \left(\tilde{\Sigma}^{2L} + \lambda \tilde{\Sigma}^2, \mathbf{0}_{d_0 - r} \right). \quad (93)$$

We now partition $\mathbf{C} := \mathbf{U}_L^T \mathbf{Y} \mathbf{V}_1 \mathbf{Q}^T \in \mathbb{R}^{d_L \times d_0}$ into the block form $\mathbf{C} = \begin{bmatrix} \mathbf{C}_1 & \mathbf{C}_2 \\ \mathbf{C}_3 & \mathbf{S} \end{bmatrix}$, where $\mathbf{C}_1 \in \mathbb{R}^{r \times r}$. This, together with (91) and (93), yields $\mathbf{C}_1 = (\tilde{\Sigma}^L + \lambda \tilde{\Sigma}^{2-L})/\sqrt{\lambda}$, $\mathbf{C}_2 = \mathbf{0}$, $\mathbf{C}_3 = \mathbf{0}$. Consequently, we obtain

$$\mathbf{U}_L^T \mathbf{Y} \mathbf{V}_1 \mathbf{Q}^T = \begin{bmatrix} \frac{1}{\sqrt{\lambda}} (\tilde{\Sigma}^L + \lambda \tilde{\Sigma}^{2-L}) & \mathbf{0} \\ \mathbf{0} & \mathbf{S} \end{bmatrix}. \quad (94)$$

Now, let $\mathbf{U}_S \Sigma_S \mathbf{V}_S = \mathbf{S}$ be an SVD of \mathbf{S} , where $\mathbf{U}_S \in \mathcal{O}^{d_L-r}$, $\mathbf{V}_S \in \mathcal{O}^{d_0-r}$, and $\Sigma_S \in \mathbb{R}^{(d_L-r) \times (d_0-r)}$. Substituting this into (94) and rearranging the terms yields

$$\begin{bmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{U}_S^T \end{bmatrix} \mathbf{U}_L^T \mathbf{Y} \mathbf{V}_1 \mathbf{Q}^T \begin{bmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{V}_S \end{bmatrix} = \begin{bmatrix} \frac{1}{\sqrt{\lambda}} (\tilde{\Sigma}^L + \lambda \tilde{\Sigma}^{2-L}) & \mathbf{0} \\ \mathbf{0} & \Sigma_S \end{bmatrix}. \quad (95)$$

Since \mathbf{Y} is a diagonal matrix, $\text{BlkDiag}(\mathbf{I}, \mathbf{U}_S^T) \mathbf{U}_L^T \in \mathcal{O}^{d_L}$, and $\mathbf{V}_1 \mathbf{Q}^T \text{BlkDiag}(\mathbf{I}, \mathbf{V}_S) \in \mathcal{O}^{d_0}$, the above left-hand side is an SVD of the right-hand diagonal matrix. Therefore, we obtain that the diagonal elements of $\text{BlkDiag}((\tilde{\Sigma}^L + \lambda \tilde{\Sigma}^{2-L})/\sqrt{\lambda}, \Sigma_S)$ is a permutation of those of \mathbf{Y} . This, together with $\mathbf{Y} \in \mathbb{R}^{d_L \times d_0}$ in (22), yields that there exists a permutation matrix $\Pi \in \mathcal{P}^{d_{\min}}$ such that

$$\begin{bmatrix} \frac{1}{\sqrt{\lambda}} (\tilde{\Sigma}^L + \lambda \tilde{\Sigma}^{2-L}) & \mathbf{0} \\ \mathbf{0} & \Sigma_S \end{bmatrix} = \begin{bmatrix} \Pi & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{d_L-d_{\min}} \end{bmatrix} \mathbf{Y} \begin{bmatrix} \Pi^T & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{d_0-d_{\min}} \end{bmatrix}.$$

Combining this with (95) yields

$$\mathbf{Y} = \left(\begin{bmatrix} \Pi^T & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{d_0-d_{\min}} \end{bmatrix} \begin{bmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{U}_S^T \end{bmatrix} \mathbf{U}_L^T \right) \mathbf{Y} \left(\mathbf{V}_1 \mathbf{Q}^T \begin{bmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{V}_S \end{bmatrix} \begin{bmatrix} \Pi & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{d_0-d_{\min}} \end{bmatrix} \right).$$

Since the right-hand side is an SVD of a diagonal matrix \mathbf{Y} in (22) and (23), there exist $\mathbf{O}_1 \in \mathcal{O}^{h_1}, \dots, \mathbf{O}_{p_Y} \in \mathcal{O}^{h_{p_Y}}, \mathbf{O}_{p_Y+1} \in \mathcal{O}^{d_0-r_Y}$ and $\hat{\mathbf{O}}_{p_Y+1} \in \mathcal{O}^{d_L-r_Y}$ such that

$$\begin{aligned} \begin{bmatrix} \Pi^T & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{d_0-d_{\min}} \end{bmatrix} \begin{bmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{V}_S^T \end{bmatrix} \mathbf{Q} \mathbf{V}_1^T &= \text{BlkDiag}(\mathbf{O}_1, \dots, \mathbf{O}_{p_Y}, \mathbf{O}_{p_Y+1}), \\ \begin{bmatrix} \Pi^T & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{d_L-d_{\min}} \end{bmatrix} \begin{bmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{U}_S^T \end{bmatrix} \mathbf{U}_L^T &= \text{BlkDiag}(\mathbf{O}_1, \dots, \mathbf{O}_{p_Y}, \hat{\mathbf{O}}_{p_Y+1}). \end{aligned}$$

This implies

$$\mathbf{V}_1 = \text{BlkDiag}(\mathbf{O}_1^T, \dots, \mathbf{O}_{p_Y}^T, \mathbf{O}_{p_Y+1}^T) \text{BlkDiag}(\Pi^T, \mathbf{I}_{d_0-d_{\min}}) \text{BlkDiag}(\mathbf{I}_r, \mathbf{V}_S^T) \mathbf{Q}, \quad (96)$$

$$\mathbf{U}_L = \text{BlkDiag}(\mathbf{O}_1^T, \dots, \mathbf{O}_{p_Y}^T, \hat{\mathbf{O}}_{p_Y+1}^T) \text{BlkDiag}(\Pi^T, \mathbf{I}_{d_L-d_{\min}}) \text{BlkDiag}(\mathbf{I}_r, \mathbf{U}_S^T). \quad (97)$$

Substituting (86) and (97) into (84) yields

$$\mathbf{W}_L = \text{BlkDiag}(\mathbf{O}_1^T, \dots, \mathbf{O}_{p_Y}^T, \hat{\mathbf{O}}_{p_Y+1}^T) \text{BlkDiag}(\Pi^T, \mathbf{I}_{d_L-d_{\min}}) \text{BlkDiag}(\tilde{\Sigma}, \mathbf{0}) \mathbf{V}_L^T. \quad (98)$$

Next, substituting (87) and (88) into (84) yields

$$\mathbf{W}_{L-1} = \mathbf{V}_L \text{BlkDiag}(\mathbf{Q}_{L-1}^{(1)}, \dots, \mathbf{Q}_{L-1}^{(p)}, \mathbf{I}) \text{BlkDiag}(\tilde{\Sigma}, \mathbf{0}) \mathbf{V}_{L-1}^T = \mathbf{V}_L \text{BlkDiag}(\tilde{\Sigma}, \mathbf{0}) \mathbf{P}_{L-1}^T,$$

where $\mathbf{P}_{L-1} := \mathbf{V}_{L-1} \text{BlkDiag} \left(\mathbf{Q}_{L-1}^{(1)T}, \dots, \mathbf{Q}_{L-1}^{(p)T}, \mathbf{I} \right) \in \mathcal{O}^{d_{L-2}}$. Using the same argument, we obtain

$$\mathbf{W}_l = \mathbf{P}_{l+1} \text{BlkDiag} \left(\tilde{\Sigma}, \mathbf{0} \right) \mathbf{P}_l^T, \quad l = 2, \dots, L-2, \quad (99)$$

where $\mathbf{P}_l := \mathbf{V}_l \text{BlkDiag} \left(\prod_{j=l}^{L-1} \mathbf{Q}_j^{(1)T}, \dots, \prod_{j=l}^{L-1} \mathbf{Q}_j^{(p)T}, \mathbf{I} \right) \in \mathcal{O}^{d_{l-1}}$ for all $l = 2, \dots, L-2$. Finally, using (90) and (96), we compute

$$\begin{aligned} \mathbf{W}_1 &= \mathbf{P}_2 \text{BlkDiag} \left(\tilde{\Sigma}, \mathbf{0} \right) \text{BlkDiag} \left(\prod_{j=1}^{L-1} \mathbf{Q}_j^{(1)}, \dots, \prod_{j=1}^{L-1} \mathbf{Q}_j^{(p)}, \mathbf{I} \right) \mathbf{V}_1^T \\ &= \mathbf{P}_2 \text{BlkDiag} \left(\tilde{\Sigma}, \mathbf{0} \right) \text{BlkDiag} \left(\mathbf{\Pi}, \mathbf{I}_{d_0-d_{\min}} \right) \text{BlkDiag} \left(\mathbf{O}_1, \dots, \mathbf{O}_{p_Y}, \mathbf{O}_{p_Y+1} \right). \end{aligned} \quad (100)$$

By letting $\boldsymbol{\sigma} := (\sigma_1, \sigma_2, \dots, \sigma_{d_{\min}}) \in \mathbb{R}^{d_{\min}}$, we write the singular value matrix Σ_l as

$$\Sigma_l = \text{BlkDiag} \left(\text{diag}(\boldsymbol{\sigma}), \mathbf{0}_{(d_l-d_{\min}) \times (d_{l-1}-d_{\min})} \right), \quad \forall l \in [L].$$

Next, it remains to show that $(\boldsymbol{\sigma}, \mathbf{\Pi}) \in \mathcal{B}$. Substituting (98), (99), and (100) into $\nabla_{\mathbf{W}_L} G(\mathbf{W}) = \mathbf{0}$ yields

$$\begin{aligned} \mathbf{0} &= \nabla_{\mathbf{W}_L} G(\mathbf{W}) \stackrel{(14)}{=} \mathbf{W}_{L:1} \mathbf{W}_{L-1:1}^T - \sqrt{\lambda} \mathbf{Y} \mathbf{W}_{L-1:1}^T + \lambda \mathbf{W}_L \\ &= \text{BlkDiag} \left(\mathbf{O}_1^T, \dots, \mathbf{O}_{p_Y}^T, \hat{\mathbf{O}}_{p_Y+1}^T \right) \text{BlkDiag} \left(\mathbf{\Pi}^T (\text{diag}^{2L-1}(\boldsymbol{\sigma}) + \lambda \text{diag}(\boldsymbol{\sigma})), \mathbf{0} \right) \mathbf{V}_L^T - \\ &\quad \sqrt{\lambda} \mathbf{Y} \text{BlkDiag} \left(\mathbf{O}_1^T, \dots, \mathbf{O}_{p_Y}^T, \mathbf{O}_{p_Y+1}^T \right) \text{BlkDiag} \left(\mathbf{\Pi}^T \text{diag}^{L-1}(\boldsymbol{\sigma}), \mathbf{0} \right) \mathbf{V}_L^T. \end{aligned}$$

This, together with $\mathbf{Y} \text{BlkDiag} \left(\mathbf{O}_1^T, \dots, \mathbf{O}_{p_Y}^T, \mathbf{O}_{p_Y+1}^T \right) = \text{BlkDiag} \left(\mathbf{O}_1^T, \dots, \mathbf{O}_{p_Y}^T, \hat{\mathbf{O}}_{p_Y+1}^T \right) \mathbf{Y}$ due to (23), implies

$$\text{BlkDiag} \left(\mathbf{\Pi}^T (\text{diag}^{2L-1}(\boldsymbol{\sigma}) + \lambda \text{diag}(\boldsymbol{\sigma})), \mathbf{0} \right) - \sqrt{\lambda} \mathbf{Y} \text{BlkDiag} \left(\mathbf{\Pi}^T \text{diag}^{L-1}(\boldsymbol{\sigma}), \mathbf{0} \right) = \mathbf{0}.$$

This directly implies each element of $\mathbf{\Pi}^T \boldsymbol{\sigma} \in \mathcal{A}$, and thus there exists $\mathbf{a} \in \mathcal{A}$ such that $\mathbf{\Pi} \mathbf{a} = \boldsymbol{\sigma}$. Consequently, $(\mathbf{a}, \mathbf{\Pi}) \in \mathcal{B}$. Then, we complete the proof. \square

Next, we proceed to the proof of Proposition 2. This proof mainly builds on the structure of the critical point set in Proposition 1.

Proof of Proposition 2. (i) Obviously, the “only if” direction is trivial since $\boldsymbol{\sigma}$ of $\mathcal{W}_{\boldsymbol{\sigma}, \mathbf{\Pi}}$ and $\boldsymbol{\sigma}'$ in $\mathcal{W}_{\boldsymbol{\sigma}', \mathbf{\Pi}'}$ share the same non-increasing singular values. Now, we are devoted to proving the “if” direction. Note that $(\boldsymbol{\sigma}, \mathbf{\Pi}) \in \mathcal{B}$ and $(\boldsymbol{\sigma}', \mathbf{\Pi}') \in \mathcal{B}$ satisfy $\boldsymbol{\sigma} = \boldsymbol{\sigma}'$. There exist $\mathbf{a}, \mathbf{a}' \in \mathcal{A}$ such that

$$\mathbf{a} = \mathbf{\Pi}^T \boldsymbol{\sigma}, \quad \mathbf{a}' = \mathbf{\Pi}'^T \boldsymbol{\sigma}'. \quad (101)$$

This, together with $\boldsymbol{\sigma} = \boldsymbol{\sigma}'$, implies that \mathbf{a}, \mathbf{a}' have the same positive elements but in a different order. Consider (28) in definition \mathcal{A} ,

$$x_i^{2L-1} - \sqrt{\lambda} y_i x_i^{L-1} + \lambda x_i = 0, \quad x_i \geq 0, \quad \forall i \in [d_{\min}]. \quad (102)$$

Note that if $x_i = x_j > 0$, we have $y_i = y_j$, which implies that if $y_i \neq y_j$, we have $x_i \neq x_j$ when $x_i, x_j > 0$. Since $\boldsymbol{\sigma} = \boldsymbol{\sigma}'$ and each element of \mathbf{a} and \mathbf{a}' is obtained by solving the equation (102), along with the fact that different y_i 's correspond to equations with no common positive root

and the partition of $(y_1, \dots, y_{d_{\min}})$ in (6), yields that the elements of \mathbf{a} and \mathbf{a}' in each partition of the form \mathbf{Y} differ only in their order. Therefore, we obtain the following conclusion: There exist $\mathbf{P}_i \in \mathcal{P}^{h_i}$ for all $i \in [p_Y]$ and $\mathbf{P}_{p_Y+1} \in \mathcal{P}^{d_{\min}-r_Y}$ such that

$$\text{diag}(\mathbf{a}) = \text{BlkDiag}(\mathbf{P}_1, \dots, \mathbf{P}_{p_Y}, \mathbf{P}_{p_Y+1}) \text{diag}(\mathbf{a}') \text{BlkDiag}(\mathbf{P}_1^T, \dots, \mathbf{P}_{p_Y}^T, \mathbf{P}_{p_Y+1}^T).$$

Substituting (101) into the above equality, together with $\boldsymbol{\sigma} = \boldsymbol{\sigma}'$, yields

$$\boldsymbol{\Pi}^T \text{diag}(\boldsymbol{\sigma}) \boldsymbol{\Pi} = \text{BlkDiag}(\mathbf{P}_1, \dots, \mathbf{P}_{p_Y}, \mathbf{P}_{p_Y+1}) \boldsymbol{\Pi}'^T \text{diag}(\boldsymbol{\sigma}) \boldsymbol{\Pi}' \text{BlkDiag}(\mathbf{P}_1^T, \dots, \mathbf{P}_{p_Y}^T, \mathbf{P}_{p_Y+1}^T). \quad (103)$$

Let $(\mathbf{W}_1, \dots, \mathbf{W}_L) \in \mathcal{W}_{\boldsymbol{\sigma}, \boldsymbol{\Pi}}$ be arbitrary. For ease of exposition, let

$$\tilde{\mathbf{Q}}_l := \mathbf{Q}_l \text{BlkDiag}(\boldsymbol{\Pi}, \mathbf{I}) \text{BlkDiag}(\mathbf{P}_1, \dots, \mathbf{P}_{p_Y}, \mathbf{P}_{p_Y+1}, \mathbf{I}) \text{BlkDiag}(\boldsymbol{\Pi}'^T, \mathbf{I}) \in \mathcal{O}^{d_l-1}, \quad \forall l \in [2, L].$$

Moreover, we have

$$\boldsymbol{\Sigma}_l \text{BlkDiag}(\boldsymbol{\Pi}, \mathbf{I}) = \text{BlkDiag}(\text{diag}(\boldsymbol{\sigma}) \boldsymbol{\Pi}, \mathbf{0}), \quad \forall l \in [L]. \quad (104)$$

Then, we have

$$\begin{aligned} \mathbf{W}_1 &= \mathbf{Q}_2 \boldsymbol{\Sigma}_1 \text{BlkDiag}(\boldsymbol{\Pi}, \mathbf{I}) \text{BlkDiag}(\mathbf{O}_1, \dots, \mathbf{O}_{p_Y}, \mathbf{O}_{p_Y+1}) \\ &= \tilde{\mathbf{Q}}_2 \boldsymbol{\Sigma}_1 \text{BlkDiag}(\boldsymbol{\Pi}', \mathbf{I}) \text{BlkDiag}(\mathbf{P}_1^T \mathbf{O}_1, \dots, \mathbf{P}_{p_Y}^T \mathbf{O}_{p_Y}, \text{BlkDiag}(\mathbf{P}_{p_Y+1}^T, \mathbf{I}_{d_0-d_{\min}}) \mathbf{O}_{p_Y+1}), \end{aligned}$$

where the second equality uses (103) and (104). Using the same argument, we have

$$\begin{aligned} \mathbf{W}_L &= \text{BlkDiag}(\mathbf{O}_1^T \mathbf{P}_1, \dots, \mathbf{O}_{p_Y}^T \mathbf{P}_Y, \hat{\mathbf{O}}_{p_Y+1}^T \text{BlkDiag}(\mathbf{P}_{p_Y+1}, \mathbf{I})) \text{BlkDiag}(\boldsymbol{\Pi}'^T, \mathbf{I}) \boldsymbol{\Sigma}_L \tilde{\mathbf{Q}}_L^T, \\ \mathbf{W}_l &= \tilde{\mathbf{Q}}_{l+1}^T \boldsymbol{\Sigma}_l \tilde{\mathbf{Q}}_l^T, \quad l = 2, \dots, L-1. \end{aligned}$$

Therefore, we obtain that $(\mathbf{W}_1, \dots, \mathbf{W}_L) \in \mathcal{W}_{\boldsymbol{\sigma}', \boldsymbol{\Pi}'}$ and thus $\mathcal{W}_{\boldsymbol{\sigma}, \boldsymbol{\Pi}} \subseteq \mathcal{W}_{\boldsymbol{\sigma}', \boldsymbol{\Pi}'}$. Applying the same argument, we also have $\mathcal{W}_{\boldsymbol{\sigma}', \boldsymbol{\Pi}'} \subseteq \mathcal{W}_{\boldsymbol{\sigma}, \boldsymbol{\Pi}}$. Therefore, we have $\mathcal{W}_{\boldsymbol{\sigma}, \boldsymbol{\Pi}} = \mathcal{W}_{\boldsymbol{\sigma}', \boldsymbol{\Pi}'}$.

(ii) Using Mirsky's inequality (see Lemma 11) and (31), we have for any $\mathbf{W} \in \mathcal{W}_{\boldsymbol{\sigma}, \boldsymbol{\Pi}}$ and $\mathbf{W}' \in \mathcal{W}_{\boldsymbol{\sigma}', \boldsymbol{\Pi}'}$,

$$\|\mathbf{W} - \mathbf{W}'\|_F \geq \|\boldsymbol{\sigma} - \boldsymbol{\sigma}'\|_2 \geq \delta_{\boldsymbol{\sigma}},$$

where the last inequality follows from (33). This implies (34). \square

4.2 Proofs of the Error Bound

In this subsection, we mainly present the detailed proofs of Lemma 3, Corollary 2, Proposition 4, Lemma 4, Proposition 5, and Proposition 6 in Section 3.3. To begin, we prove Lemma 3 and Corollary 2, which play a key role in the subsequent analysis.

Proof of Lemma 3. (i) Let $\mathbf{W}^* = (\mathbf{W}_1^*, \dots, \mathbf{W}_L^*) \in \mathcal{W}_{\boldsymbol{\sigma}^*}$ be such that $\text{dist}(\mathbf{W}, \mathcal{W}_{\boldsymbol{\sigma}^*}) = \|\mathbf{W} - \mathbf{W}^*\|_F$. Using the triangle inequality, we have for each $l \in [L]$,

$$\|\mathbf{W}_l\| \geq \|\mathbf{W}_l^*\| - \|\mathbf{W}_l - \mathbf{W}_l^*\| \geq \|\mathbf{W}_l^*\| - \|\mathbf{W}_l - \mathbf{W}_l^*\|_F \geq \sigma_{\max}^* - \frac{\sigma_{\min}^*}{2} \geq \frac{\sigma_{\max}^*}{2},$$

where the third inequality follows from (41) and $\mathbf{W}_l^* = \mathbf{Q}_{l+1}^* \boldsymbol{\Sigma}_l^* \mathbf{Q}_l^{*T}$ and $\boldsymbol{\Sigma}_l^* = \text{BlkDiag}(\boldsymbol{\sigma}^*, \mathbf{0})$ for each $\mathbf{Q}_l \in \mathcal{O}^{d_l-1}$ with $l = 2, \dots, L$ according to Theorem 3. Similarly, we have $\|\mathbf{W}_l\| \leq 3\sigma_{\max}^*/2$

for each $l \in [L]$. Therefore, we have (42). Using Weyl's inequality (see [17, Corollary 7.3.5]), we have for each $l \in [L]$ and $i \in [r_\sigma]$,

$$|\sigma_i(\mathbf{W}_l) - \sigma_i^*| \leq \|\mathbf{W}_l - \mathbf{W}_l^*\| \leq \frac{\sigma_{\min}^*}{2},$$

which implies

$$\sigma_i(\mathbf{W}_l) \geq \sigma_i^* - \frac{\sigma_{\min}^*}{2} \geq \frac{\sigma_{\min}^*}{2}.$$

(ii) According to (14), we compute

$$\begin{aligned} \frac{1}{2} \nabla_{\mathbf{W}_l} G(\mathbf{W}) \mathbf{W}_l^T &= \mathbf{W}_{L:l+1}^T \left(\mathbf{W}_{L:1} - \sqrt{\lambda} \mathbf{Y} \right) \mathbf{W}_{l:1}^T + \lambda \mathbf{W}_l \mathbf{W}_l^T, \\ \frac{1}{2} \mathbf{W}_{l+1}^T \nabla_{\mathbf{W}_{l+1}} G(\mathbf{W}) &= \mathbf{W}_{L:l+1}^T \left(\mathbf{W}_{L:1} - \sqrt{\lambda} \mathbf{Y} \right) \mathbf{W}_{l+1}^T + \lambda \mathbf{W}_{l+1}^T \mathbf{W}_{l+1}. \end{aligned}$$

Then, we have for each $l = 1, \dots, L-1$,

$$\frac{1}{2} (\nabla_{\mathbf{W}_l} G(\mathbf{W}) \mathbf{W}_l^T - \mathbf{W}_{l+1}^T \nabla_{\mathbf{W}_{l+1}} G(\mathbf{W})) = \lambda (\mathbf{W}_l \mathbf{W}_l^T - \mathbf{W}_{l+1}^T \mathbf{W}_{l+1}).$$

This implies

$$\begin{aligned} \|\mathbf{W}_l \mathbf{W}_l^T - \mathbf{W}_{l+1}^T \mathbf{W}_{l+1}\|_F &= \frac{1}{2\lambda} \|\nabla_{\mathbf{W}_l} G(\mathbf{W}) \mathbf{W}_l^T - \mathbf{W}_{l+1}^T \nabla_{\mathbf{W}_{l+1}} G(\mathbf{W})\|_F \\ &\leq \frac{3\sigma_{\max}^*}{4\lambda} (\|\nabla_{\mathbf{W}_l} G(\mathbf{W})\|_F + \|\nabla_{\mathbf{W}_{l+1}} G(\mathbf{W})\|_F) \leq \frac{3\sqrt{2}\sigma_{\max}^*}{4\lambda} \|\nabla G(\mathbf{W})\|_F, \end{aligned}$$

where the first inequality follows from the triangle inequality, $\|\mathbf{A}\mathbf{B}\|_F \leq \|\mathbf{A}\| \|\mathbf{B}\|_F$ for all \mathbf{A}, \mathbf{B} of the same size, and (42).

(iii) Using Weyl's inequality and (44), we obtain for each $l \in [L]$ and $i \in [r_\sigma]$,

$$|\sigma_i^2(\mathbf{W}_l) - \sigma_i^2(\mathbf{W}_{l+1})| \leq \|\mathbf{W}_l \mathbf{W}_l^T - \mathbf{W}_{l+1}^T \mathbf{W}_{l+1}\|_F \leq \frac{3\sqrt{2}\sigma_{\max}^*}{4\lambda} \|\nabla G(\mathbf{W})\|_F.$$

This implies

$$|\sigma_i(\mathbf{W}_l) - \sigma_i(\mathbf{W}_{l+1})| \leq \frac{3\sqrt{2}\sigma_{\max}^*}{4\lambda (\sigma_i(\mathbf{W}_l) + \sigma_i(\mathbf{W}_{l+1}))} \|\nabla G(\mathbf{W})\|_F \leq \frac{3\sqrt{2}\sigma_{\max}^*}{4\lambda \sigma_{\min}^*} \|\nabla G(\mathbf{W})\|_F,$$

where the last inequality follows from (43).

(iv) We first prove (46). Note that

$$\begin{aligned} \mathbf{W}_{j:i}^T \mathbf{W}_{j:i} - (\mathbf{W}_i^T \mathbf{W}_i)^{j-i+1} &= \mathbf{W}_{j-1:i}^T (\mathbf{W}_j^T \mathbf{W}_j - \mathbf{W}_{j-1}^T \mathbf{W}_{j-1}) \mathbf{W}_{j-1:i} \\ &\quad + \mathbf{W}_{j-2:i}^T ((\mathbf{W}_{j-1}^T \mathbf{W}_{j-1})^2 - (\mathbf{W}_{j-2}^T \mathbf{W}_{j-2})^2) \mathbf{W}_{j-2:i} + \dots \\ &\quad + \mathbf{W}_i^T ((\mathbf{W}_{i+1}^T \mathbf{W}_{i+1})^{j-i} - \mathbf{W}_i^T \mathbf{W}_i^{j-i}) \mathbf{W}_i \\ &= \sum_{k=1}^{j-i} \mathbf{W}_{j-k:i}^T \left((\mathbf{W}_{j-k+1}^T \mathbf{W}_{j-k+1})^k - (\mathbf{W}_{j-k}^T \mathbf{W}_{j-k})^k \right) \mathbf{W}_{j-k:i}. \end{aligned} \quad (105)$$

For each $k \in [j-i]$, we compute

$$\begin{aligned} &\left\| \mathbf{W}_{j-k:i}^T \left((\mathbf{W}_{j+1-k}^T \mathbf{W}_{j+1-k})^k - (\mathbf{W}_{j-k}^T \mathbf{W}_{j-k})^k \right) \mathbf{W}_{j-k:i} \right\|_F \\ &\leq \left\| (\mathbf{W}_{j+1-k}^T \mathbf{W}_{j+1-k})^k - (\mathbf{W}_{j-k}^T \mathbf{W}_{j-k})^k \right\|_F \prod_{l=i}^{j-k} \|\mathbf{W}_l\|^2 \\ &\leq k \left(\frac{3\sigma_{\max}^*}{2} \right)^{2(j-i)} \left\| \mathbf{W}_{j+1-k}^T \mathbf{W}_{j+1-k} - \mathbf{W}_{j-k}^T \mathbf{W}_{j-k} \right\|_F \\ &\leq k \left(\frac{3\sigma_{\max}^*}{2} \right)^{2(j-i)} \frac{3\sqrt{2}\sigma_{\max}^*}{4\lambda} \|\nabla G(\mathbf{W})\|_F = \frac{\sqrt{2}k}{2\lambda} \left(\frac{3\sigma_{\max}^*}{2} \right)^{2(j-i)+1} \|\nabla G(\mathbf{W})\|_F, \end{aligned}$$

where the last inequality uses (44) and the second inequality follows from

$$\begin{aligned}
& \left\| (\mathbf{W}_{j+1-k}^T \mathbf{W}_{j+1-k})^k - (\mathbf{W}_{j-k}^T \mathbf{W}_{j-k})^k \right\|_F \\
&= \left\| \sum_{l=1}^k (\mathbf{W}_{j+1-k}^T \mathbf{W}_{j+1-k})^{k-l} (\mathbf{W}_{j+1-k}^T \mathbf{W}_{j+1-k} - \mathbf{W}_{j-k}^T \mathbf{W}_{j-k}) (\mathbf{W}_{j-k}^T \mathbf{W}_{j-k})^{l-1} \right\|_F \\
&\leq \left\| \mathbf{W}_{j+1-k}^T \mathbf{W}_{j+1-k} - \mathbf{W}_{j-k}^T \mathbf{W}_{j-k} \right\|_F \sum_{l=1}^k \left\| \mathbf{W}_{j+1-k} \right\|^{2(k-l)} \left\| \mathbf{W}_{j-k} \right\|^{2(l-1)} \\
&\leq k \left(\frac{3\sigma_{\max}^*}{2} \right)^{2(k-1)} \left\| \mathbf{W}_{j+1-k}^T \mathbf{W}_{j+1-k} - \mathbf{W}_{j-k}^T \mathbf{W}_{j-k} \right\|_F.
\end{aligned}$$

This, together with (105), yields

$$\begin{aligned}
\left\| \mathbf{W}_{j:i}^T \mathbf{W}_{j:i} - (\mathbf{W}_i^T \mathbf{W}_i)^{j-i+1} \right\|_F &\leq \left(\sum_{k=1}^{j-i} k \right) \frac{\sqrt{2}}{2\lambda} \left(\frac{3\sigma_{\max}^*}{2} \right)^{2j-2i+1} \left\| \nabla G(\mathbf{W}) \right\|_F \\
&\leq \frac{(j-i)(j-i+1)}{2\sqrt{2}\lambda} \left(\frac{3\sigma_{\max}^*}{2} \right)^{2j-2i+1} \left\| \nabla G(\mathbf{W}) \right\|_F.
\end{aligned}$$

Using the same argument, we prove (47). Then, we complete the proof. \square

Proof of Corollary 2. For ease of exposition, let

$$\begin{aligned}
\mathbf{R}_1(l) &:= \mathbf{W}_{l-1:1} \mathbf{W}_{l-1:1}^T - (\mathbf{W}_l^T \mathbf{W}_l)^{l-1}, \quad l = 2, 3, \dots, L, \quad \mathbf{R}_1(1) := \mathbf{0}, \\
\mathbf{R}_2(l) &:= \mathbf{W}_{L:l+1}^T \mathbf{W}_{L:l+1} - (\mathbf{W}_l \mathbf{W}_l^T)^{L-l}, \quad l = 1, 2, \dots, L-1, \quad \mathbf{R}_2(L) := \mathbf{0}.
\end{aligned}$$

For each $l \in [L]$, we compute

$$\begin{aligned}
& \left\| (\mathbf{W}_{l-1} \mathbf{W}_{l-1}^T)^{l-1} - (\mathbf{W}_l^T \mathbf{W}_l)^{l-1} \right\|_F \leq \left\| (\mathbf{W}_{l-1} \mathbf{W}_{l-1}^T)^{l-2} (\mathbf{W}_{l-1} \mathbf{W}_{l-1}^T - \mathbf{W}_l^T \mathbf{W}_l) \right\|_F + \dots \\
&+ \left\| (\mathbf{W}_{l-1} \mathbf{W}_{l-1}^T - \mathbf{W}_l^T \mathbf{W}_l) (\mathbf{W}_l^T \mathbf{W}_l)^{l-2} \right\|_F \leq \frac{3\sqrt{2}(l-1)\sigma_{\max}^*}{4\lambda} \left(\frac{3\sigma_{\max}^*}{2} \right)^{2l-4} \left\| \nabla G(\mathbf{W}) \right\|_F,
\end{aligned}$$

where the last inequality uses (42) and (44) in Lemma 3. This, together with the triangular inequality and (47), yields for $l = 2, 3, \dots, L$,

$$\begin{aligned}
\left\| \mathbf{R}_1(l) \right\|_F &\leq \left\| \mathbf{W}_{l-1:1} \mathbf{W}_{l-1:1}^T - (\mathbf{W}_{l-1} \mathbf{W}_{l-1}^T)^{l-1} \right\|_F + \left\| (\mathbf{W}_{l-1} \mathbf{W}_{l-1}^T)^{l-1} - (\mathbf{W}_l^T \mathbf{W}_l)^{l-1} \right\|_F \\
&\leq \frac{l(l-1)}{2\sqrt{2}\lambda} \left(\frac{3\sigma_{\max}^*}{2} \right)^{2l-3} \left\| \nabla G(\mathbf{W}) \right\|_F.
\end{aligned} \tag{106}$$

Using the same argument, we have for $l = 1, 2, \dots, L-1$,

$$\left\| \mathbf{R}_2(l) \right\|_F \leq \frac{(L-l+1)(L-l)}{2\sqrt{2}\lambda} \left(\frac{3\sigma_{\max}^*}{2} \right)^{2(L-l)-1} \left\| \nabla G(\mathbf{W}) \right\|_F. \tag{107}$$

For each $l \in [L]$, we compute

$$\begin{aligned}
\frac{1}{2} \nabla_{\mathbf{W}_l} G(\mathbf{W}) &\stackrel{(13)}{=} \mathbf{W}_{L:l+1}^T (\mathbf{W}_L \cdots \mathbf{W}_1 - \sqrt{\lambda} \mathbf{Y}) \mathbf{W}_{l-1:1}^T + \lambda \mathbf{W}_l = (\mathbf{W}_l \mathbf{W}_l^T)^{L-1} \mathbf{W}_l - \\
&\quad \sqrt{\lambda} \mathbf{W}_{L:l+1}^T \mathbf{Y} \mathbf{W}_{l-1:1}^T + \lambda \mathbf{W}_l + \mathbf{R}_2(l) \mathbf{W}_l \mathbf{W}_{l-1:1} \mathbf{W}_{l-1:1}^T + (\mathbf{W}_l \mathbf{W}_l^T)^{L-l} \mathbf{W}_l \mathbf{R}_1(l).
\end{aligned}$$

This, together with the triangle inequality, (106), and (107), yields for each $l \in [L]$,

$$\begin{aligned}
& \left\| (\mathbf{W}_l \mathbf{W}_l^T)^{L-1} \mathbf{W}_l - \sqrt{\lambda} \mathbf{W}_{L:l+1}^T \mathbf{Y} \mathbf{W}_{l-1:1}^T + \lambda \mathbf{W}_l \right\|_F \\
& \leq \left\| \mathbf{R}_2(l) \mathbf{W}_l \mathbf{W}_{l-1:1} \mathbf{W}_{l-1:1}^T \right\|_F + \left\| (\mathbf{W}_l \mathbf{W}_l^T)^{L-l} \mathbf{W}_l \mathbf{R}_1(l) \right\|_F + \frac{1}{2} \|\nabla G(\mathbf{W})\|_F \\
& \leq \left(\left(\frac{3\sigma_{\max}^*}{2} \right)^{2L-2} \frac{(L-l)(L-l+1) + (l-1)l}{2\sqrt{2}\lambda} + \frac{1}{2} \right) \|\nabla G(\mathbf{W})\|_F,
\end{aligned} \tag{108}$$

which directly implies (48). \square

Next, using Lemma 3 and the SVD of \mathbf{W}_l in (51), we proceed to the proof of Proposition 4.

Proof of Proposition 4. Let $\mathbf{W}^* = (\mathbf{W}_1^*, \dots, \mathbf{W}_L^*) \in \mathcal{W}_{\sigma^*}$ be such that $\text{dist}(\mathbf{W}, \mathcal{W}_{\sigma^*}) = \|\mathbf{W} - \mathbf{W}^*\|_F$. For ease of exposition, let $\mathbf{H}_l := \mathbf{U}_{l-1}^T \mathbf{V}_l$ for each $l \in \{2, \dots, L\}$. According to (51), the (i, j) -th block of \mathbf{H}_l , denoted by $\mathbf{H}_l^{(i,j)}$, is

$$\mathbf{H}_l^{(i,j)} := \mathbf{U}_{l-1}^{(i)T} \mathbf{V}_l^{(j)}, \quad \forall (i, j) \in [p+1] \times [p+1]. \tag{109}$$

According to (52) with $\delta_\sigma \leq \sigma_{\min}^*$ and Lemma 3, we obtain (44). Substituting (51) and $\mathbf{H}_l = \mathbf{U}_{l-1}^T \mathbf{V}_l$ into (44) yields

$$\left\| \mathbf{H}_l \Sigma_l^T \Sigma_l - \Sigma_{l-1} \Sigma_{l-1}^T \mathbf{H}_l \right\|_F \leq \frac{3\sqrt{2}\sigma_{\max}^*}{4\lambda} \|\nabla G(\mathbf{W})\|_F.$$

Using this and the block structures of Σ_l in (51) and \mathbf{H}_l in (109), we obtain

$$\begin{aligned}
& \sum_{j=1}^{p+1} \sum_{i=1, i \neq j}^{p+1} \left\| \mathbf{H}_l^{(i,j)} \Sigma_l^{(j)T} \Sigma_l^{(j)} - \Sigma_{l-1}^{(i)} \Sigma_{l-1}^{(i)T} \mathbf{H}_l^{(i,j)} \right\|_F^2 \\
& \leq \sum_{j=1}^{p+1} \sum_{i=1}^{p+1} \left\| \mathbf{H}_l^{(i,j)} \Sigma_l^{(j)T} \Sigma_l^{(j)} - \Sigma_{l-1}^{(i)} \Sigma_{l-1}^{(i)T} \mathbf{H}_l^{(i,j)} \right\|_F^2 \leq \left(\frac{3\sqrt{2}\sigma_{\max}^*}{4\lambda} \right)^2 \|\nabla G(\mathbf{W})\|_F^2.
\end{aligned} \tag{110}$$

For each $l \in [L]$ and $i \in [p+1]$, we have $\|\Sigma_{l,i} - \sigma_{t_i}^* \mathbf{I}\| \leq \|\mathbf{W}_l - \mathbf{W}_l^*\| \leq \text{dist}(\mathbf{W}, \mathcal{W}_{\sigma^*}) \leq \delta_\sigma/3$, where the first inequality follows from Weyl's inequality and the last inequality uses (52). This implies

$$\sigma_{t_i}^* - \frac{\delta_\sigma}{3} \leq \lambda_{\min}(\Sigma_{l,i}) \leq \lambda_{\max}(\Sigma_{l,i}) \leq \sigma_{t_i}^* + \frac{\delta_\sigma}{3}. \tag{111}$$

For each $l \in [L]$ and $i, j \in [p+1]$ with $i > j$, we compute

$$\left\| \mathbf{H}_l^{(i,j)} \Sigma_l^{(j)T} \Sigma_l^{(j)} \right\|_F \geq \lambda_{\min}^2(\Sigma_l^{(j)}) \left\| \mathbf{H}_l^{(i,j)} \right\|_F \geq \left(\sigma_{t_j}^* - \frac{\delta_\sigma}{3} \right)^2 \left\| \mathbf{H}_l^{(i,j)} \right\|_F, \tag{112}$$

$$\left\| \Sigma_{l-1}^{(i)} \Sigma_{l-1}^{(i)T} \mathbf{H}_l^{(i,j)} \right\|_F^2 \leq \lambda_{\max}^2 \left(\Sigma_{l-1}^{(i)} \Sigma_{l-1}^{(i)T} \right) \left\| \mathbf{H}_l^{(i,j)} \right\|_F^2 \leq \left(\sigma_{t_i}^* + \frac{\delta_\sigma}{3} \right)^2 \left\| \mathbf{H}_l^{(i,j)} \right\|_F^2. \tag{113}$$

For each $l \in [L]$ and $i > j \in [p+1]$, we bound

$$\begin{aligned}
& \left\| \mathbf{H}_l^{(i,j)} \Sigma_l^{(j)T} \Sigma_l^{(j)} - \Sigma_{l-1}^{(i)} \Sigma_{l-1}^{(i)T} \mathbf{H}_l^{(i,j)} \right\|_F \geq \left\| \mathbf{H}_l^{(i,j)} \Sigma_l^{(j)T} \Sigma_l^{(j)} \right\|_F - \left\| \Sigma_{l-1}^{(i)} \Sigma_{l-1}^{(i)T} \mathbf{H}_l^{(i,j)} \right\|_F \\
& \geq \left(\left(\sigma_{t_j}^* - \frac{\delta_\sigma}{3} \right)^2 - \left(\sigma_{t_i}^* + \frac{\delta_\sigma}{3} \right)^2 \right) \left\| \mathbf{H}_l^{(i,j)} \right\|_F \geq \frac{2\delta_\sigma \sigma_{\min}^*}{3} \left\| \mathbf{H}_l^{(i,j)} \right\|_F,
\end{aligned}$$

where the second equality follows from (112) and (113), and the last inequality is due to (33). Applying the same argument to the case where $i < j$, we obtain the same result. These, together with (110), yield

$$\sum_{j=1}^{p+1} \sum_{i=1, i \neq j}^{p+1} \|\mathbf{H}_l^{(i,j)}\|_F^2 \leq \frac{81\sigma_{\max}^{*2}}{32\delta_\sigma^2 \lambda^2 \sigma_{\min}^{*2}} \|\nabla G(\mathbf{W})\|_F^2. \quad (114)$$

For each $l \in [L]$ and $i \in [p+1]$, let $\mathbf{H}_l^{(i,i)} = \mathbf{P}_l^{(i)} \mathbf{\Lambda}_l^{(i)} \mathbf{Q}_l^{(i)T}$ be a full SVD of $\mathbf{H}_l^{(i,i)}$, where $\mathbf{P}_l^{(i)}, \mathbf{Q}_l^{(i)}$ are square orthogonal matrix. Then, we compute

$$\begin{aligned} \|\mathbf{H}_l^{(i,i)} - \mathbf{P}_l^{(i)} \mathbf{Q}_l^{(i)T}\|_F^2 &= \sum_{k=1}^{g_i} \left(1 - \sigma_k(\mathbf{H}_l^{(i,i)})\right)^2 \leq \sum_{k=1}^{g_i} \left(1 - \sigma_k^2(\mathbf{H}_l^{(i,i)})\right)^2 \\ &= \left\| \mathbf{I} - \mathbf{H}_l^{(i,i)} \mathbf{H}_l^{(i,i)T} \right\|_F^2 = \sum_{j=1, j \neq i}^{p+1} \|\mathbf{H}_l^{(i,j)} \mathbf{H}_l^{(i,j)T}\|_F^2 \leq \sum_{j=1, j \neq i}^{p+1} \|\mathbf{H}_l^{(i,j)}\|_F^2, \end{aligned}$$

where the last equality is due to $[\mathbf{H}_l^{(i,1)}, \dots, \mathbf{H}_l^{(i,p+1)}] = \mathbf{U}_{l-1}^{(i)T} \mathbf{V}_l$ with $\mathbf{U}_{l-1}, \mathbf{V}_l \in \mathcal{O}^{d_{l-1}}$, and last inequality follows from $\|\mathbf{H}_l^{(i,j)}\| \leq 1$ for all $l \in [L]$ and $i, j \in [p+1]$. Therefore, we obtain

$$\sum_{i=1}^{p+1} \|\mathbf{H}_l^{(i,i)} - \mathbf{P}_l^{(i)} \mathbf{Q}_l^{(i)T}\|_F^2 \leq \sum_{i=1}^{p+1} \sum_{j=1, j \neq i}^{p+1} \|\mathbf{H}_l^{(i,j)}\|_F^2.$$

This, together with (114) and $\mathbf{T}_l^{(i)} = \mathbf{P}_l^{(i)} \mathbf{Q}_l^{(i)T}$, yields

$$\|\mathbf{H}_l - \text{BlkDiag}(\mathbf{T}_l^{(1)}, \dots, \mathbf{T}_l^{(p)}, \mathbf{T}_l^{(p+1)})\|_F^2 \leq \frac{81\sigma_{\max}^{*2}}{16\delta_\sigma^2 \lambda^2 \sigma_{\min}^{*2}} \|\nabla G(\mathbf{W})\|_F^2,$$

which directly implies (53). \square

Now, we prove two key inequalities in Lemma 4, which will be frequently used in the subsequent analysis.

Proof of Lemma 4. Let $\mathbf{W}^* = (\mathbf{W}_1^*, \dots, \mathbf{W}_L^*)$ be such that $\text{dist}(\mathbf{W}, \mathcal{W}_{\sigma^*}) = \|\mathbf{W} - \mathbf{W}^*\|_F$. It follows from Lemma 3 and Corollary 2 that (42)-(45) and (48) holds. Recall that (51) denotes an SVD of \mathbf{W}_l for each $l \in [L]$. For ease of exposition, let $\mathbf{H}_l := \mathbf{U}_{l-1}^T \mathbf{V}_l$ and $\mathbf{T}_l := \text{BlkDiag}(\mathbf{T}_l^{(1)}, \dots, \mathbf{T}_l^{(p)}, \mathbf{T}_l^{(p+1)})$ for each $l = 2, \dots, L$. This, together with (53), implies

$$\|\mathbf{H}_l - \mathbf{T}_l\|_F \leq \frac{9\sigma_{\max}^*}{4\delta_\sigma \lambda \sigma_{\min}^*} \|\nabla G(\mathbf{W})\|_F, \quad l = 2, \dots, L. \quad (115)$$

First, we are devoted to proving (55). It follows from (48) with $l = L$ that

$$\left\| (\mathbf{W}_L \mathbf{W}_L^T)^{L-1} \mathbf{W}_L - \sqrt{\lambda} \mathbf{Y} \mathbf{W}_{L-1:1}^T + \lambda \mathbf{W}_L \right\|_F \leq c_1 \|\nabla G(\mathbf{W})\|_F, \quad (116)$$

where c_1 is defined in (49). Using the SVD (51), we have

$$\begin{aligned}
& \left\| (\mathbf{W}_L \mathbf{W}_L^T)^{L-1} \mathbf{W}_L - \sqrt{\lambda} \mathbf{Y} \mathbf{W}_{L-1:1}^T + \lambda \mathbf{W}_L \right\|_F \\
&= \left\| (\boldsymbol{\Sigma}_L \boldsymbol{\Sigma}_L^T)^{L-1} \boldsymbol{\Sigma}_L - \sqrt{\lambda} \mathbf{U}_L^T \mathbf{Y} \mathbf{V}_1 \prod_{l=1}^{L-1} \boldsymbol{\Sigma}_l^T \mathbf{H}_{l+1} + \lambda \boldsymbol{\Sigma}_L \right\|_F \\
&\geq \left\| (\boldsymbol{\Sigma}_L \boldsymbol{\Sigma}_L^T)^{L-1} \boldsymbol{\Sigma}_L + \lambda \boldsymbol{\Sigma}_L - \sqrt{\lambda} \mathbf{U}_L^T \mathbf{Y} \mathbf{V}_1 \prod_{l=1}^{L-1} \boldsymbol{\Sigma}_l^T \mathbf{T}_{l+1} \right\|_F - \sqrt{\lambda} \left\| \mathbf{Y} \mathbf{V}_1 \boldsymbol{\Sigma}_1^T (\mathbf{H}_2 - \mathbf{T}_2) \prod_{l=2}^{L-1} \boldsymbol{\Sigma}_l^T \mathbf{H}_{l+1} \right\|_F \\
&\quad - \sqrt{\lambda} \left\| \mathbf{Y} \mathbf{V}_1 \boldsymbol{\Sigma}_1^T \mathbf{T}_2 \boldsymbol{\Sigma}_2^T (\mathbf{H}_3 - \mathbf{T}_3) \prod_{l=3}^{L-1} \boldsymbol{\Sigma}_l^T \mathbf{H}_{l+1} \right\|_F - \dots - \sqrt{\lambda} \left\| \mathbf{Y} \mathbf{V}_1 \left(\prod_{l=1}^{L-2} \boldsymbol{\Sigma}_l^T \mathbf{T}_{l+1} \right) \boldsymbol{\Sigma}_{L-1}^T (\mathbf{H}_L - \mathbf{T}_L) \right\|_F \\
&\geq \left\| (\boldsymbol{\Sigma}_L \boldsymbol{\Sigma}_L^T)^{L-1} \boldsymbol{\Sigma}_L + \lambda \boldsymbol{\Sigma}_L - \sqrt{\lambda} \mathbf{U}_L^T \mathbf{Y} \mathbf{V}_1 \prod_{l=1}^{L-1} \boldsymbol{\Sigma}_l^T \mathbf{T}_{l+1} \right\|_F - \left(\frac{3}{2} \sigma_{\max}^* \right)^L \frac{3y_1 L}{2\sqrt{\lambda} \delta_\sigma \sigma_{\min}^*} \|\nabla G(\mathbf{W})\|_F,
\end{aligned}$$

where the first inequality uses the triangle inequality and the last inequality follows from (42) in Lemma 3, $\|\mathbf{H}_l\| = 1$ and $\|\mathbf{T}_l\| = 1$ for all $l \in [L]$, and (115). This, together with (116), implies

$$\begin{aligned}
& \left\| (\boldsymbol{\Sigma}_L \boldsymbol{\Sigma}_L^T)^{L-1} \boldsymbol{\Sigma}_L + \lambda \boldsymbol{\Sigma}_L - \sqrt{\lambda} \mathbf{U}_L^T \mathbf{Y} \mathbf{V}_1 \prod_{l=1}^{L-1} \boldsymbol{\Sigma}_l^T \mathbf{T}_{l+1} \right\|_F \\
&\leq \left(\left(\frac{3}{2} \sigma_{\max}^* \right)^L \frac{3y_1 L}{2\sqrt{\lambda} \delta_\sigma \sigma_{\min}^*} + c_1 \right) \|\nabla G(\mathbf{W})\|_F.
\end{aligned} \tag{117}$$

According to the block structure of $\boldsymbol{\Sigma}_l$ in (51) and $\mathbf{T}_l = \text{BlkDiag}(\mathbf{T}_l^{(1)}, \dots, \mathbf{T}_l^{(p)}, \mathbf{T}_l^{(p+1)})$, we compute

$$\prod_{l=1}^{L-1} \boldsymbol{\Sigma}_l^T \mathbf{T}_{l+1} = \text{BlkDiag} \left(\prod_{l=1}^{L-1} \boldsymbol{\Sigma}_l^{(1)} \mathbf{T}_{l+1}^{(1)}, \dots, \prod_{l=1}^{L-1} \boldsymbol{\Sigma}_l^{(p)} \mathbf{T}_{l+1}^{(p)}, \prod_{l=1}^{L-1} \boldsymbol{\Sigma}_l^{(p+1)} \mathbf{T}_{l+1}^{(p+1)} \right).$$

Using (45) in Lemma 3, we have for each $l \in [L]$ and $i \in [p]$,

$$\begin{aligned}
\|\boldsymbol{\Sigma}_l^{(i)} - \boldsymbol{\Sigma}_L^{(i)}\|_F^2 &= \sum_{j=1}^{g_i} \left(\sigma_j(\boldsymbol{\Sigma}_l^{(i)}) - \sigma_j(\boldsymbol{\Sigma}_L^{(i)}) \right)^2 = \sum_{j=1}^{g_i} \left(\sum_{k=l}^{L-1} \left(\sigma_j(\boldsymbol{\Sigma}_k^{(i)}) - \sigma_j(\boldsymbol{\Sigma}_{k+1}^{(i)}) \right) \right)^2 \\
&\leq g_i (L-1)^2 \left(\frac{3\sqrt{2} \sigma_{\max}^*}{4\lambda \sigma_{\min}^*} \|\nabla G(\mathbf{W})\|_F \right)^2,
\end{aligned}$$

which implies

$$\|\boldsymbol{\Sigma}_l^{(i)} - \boldsymbol{\Sigma}_L^{(i)}\|_F \leq \frac{3\sqrt{2} g_{\max} L \sigma_{\max}^*}{4\lambda \sigma_{\min}^*} \|\nabla G(\mathbf{W})\|_F. \tag{118}$$

We claim that it holds for each $l \in [L]$ and $i \in [p]$ that

$$\|\mathbf{T}_l^{(i)} \boldsymbol{\Sigma}_L^{(i)} - \boldsymbol{\Sigma}_L^{(i)} \mathbf{T}_l^{(i)}\|_F \leq \frac{\eta_1}{\sigma_{\min}^*} \|\nabla G(\mathbf{W})\|_F, \tag{119}$$

where η_1 is defined in (58). Using the triangular inequality, we obtain

$$\begin{aligned} & \left\| \prod_{l=1}^{L-1} \Sigma_L^{(i)} \mathbf{T}_{l+1}^{(i)} - \left(\prod_{l=1}^{L-1} \mathbf{T}_{l+1}^{(i)} \right) (\Sigma_L^{(i)})^{L-1} \right\|_F \\ & \leq \sum_{j=2}^L \left\| \underbrace{\left(\prod_{l=1}^{L-j} \Sigma_L^{(i)} \mathbf{T}_{l+1}^{(i)} \right) \left(\Sigma_L^{(i)} \prod_{k=L-j+2}^L \mathbf{T}_k^{(i)} - \left(\prod_{k=L-j+2}^L \mathbf{T}_k^{(i)} \right) \Sigma_L^{(i)} \right) (\Sigma_L^{(i)})^{j-2}}_{\mathbf{R}_j} \right\|_F. \end{aligned} \quad (120)$$

For each $\|\mathbf{R}_j\|_F$, we bound

$$\begin{aligned} \|\mathbf{R}_j\|_F & \leq \left(\frac{3\sigma_{\max}^*}{2} \right)^{L-2} \left\| \Sigma_L^{(i)} \prod_{k=L-j+2}^L \mathbf{T}_k^{(i)} - \left(\prod_{k=L-j+2}^L \mathbf{T}_k^{(i)} \right) \Sigma_L^{(i)} \right\|_F \\ & \leq \left(\frac{3\sigma_{\max}^*}{2} \right)^{L-2} \left\| \left(\Sigma_L^{(i)} \mathbf{T}_{L-j+2}^{(i)} - \mathbf{T}_{L-j+2}^{(i)} \Sigma_L^{(i)} \right) \prod_{k=L-j+3}^L \mathbf{T}_k^{(i)} \right\|_F \\ & \quad + \left(\frac{3\sigma_{\max}^*}{2} \right)^{L-2} \left\| \mathbf{T}_{L-j+2}^{(i)} \left(\Sigma_L^{(i)} \mathbf{T}_{L-j+3}^{(i)} - \mathbf{T}_{L-j+3}^{(i)} \Sigma_L^{(i)} \right) \prod_{k=L-j+4}^L \mathbf{T}_k^{(i)} \right\|_F + \dots \\ & \quad + \left(\frac{3\sigma_{\max}^*}{2} \right)^{L-2} \left\| \prod_{k=L-j+2}^{L-1} \mathbf{T}_k^{(i)} (\Sigma_L^{(i)} \mathbf{T}_L^{(i)} - \mathbf{T}_L^{(i)} \Sigma_L^{(i)}) \right\|_F \\ & \leq (j-1) \left(\frac{3\sigma_{\max}^*}{2} \right)^{L-2} \frac{\eta_1}{\sigma_{\min}^*} \|\nabla G(\mathbf{W})\|_F, \end{aligned}$$

where the first inequality uses (i) of Lemma 3 and $\|\mathbf{T}_{l+1}^{(i)}\| = 1$ for all l and i , the second one is due the triangular inequality, and the last one follows from (119). This, together with (120), implies

$$\left\| \prod_{l=1}^{L-1} \Sigma_L^{(i)} \mathbf{T}_{l+1}^{(i)} - \left(\prod_{l=1}^{L-1} \mathbf{T}_{l+1}^{(i)} \right) (\Sigma_L^{(i)})^{L-1} \right\|_F \leq \left(\frac{3\sigma_{\max}^*}{2} \right)^{L-2} \frac{L^2 \eta_1}{2\sigma_{\min}^*} \|\nabla G(\mathbf{W})\|_F. \quad (121)$$

For each $i \in [p]$, we have

$$\begin{aligned} & \left\| \prod_{l=1}^{L-1} \Sigma_l^{(i)} \mathbf{T}_{l+1}^{(i)} - \left(\prod_{l=1}^{L-1} \mathbf{T}_{l+1}^{(i)} \right) (\Sigma_L^{(i)})^{L-1} \right\|_F \\ & \leq \left\| \prod_{l=1}^{L-1} \Sigma_L^{(i)} \mathbf{T}_{l+1}^{(i)} - \left(\prod_{l=1}^{L-1} \mathbf{T}_{l+1}^{(i)} \right) (\Sigma_L^{(i)})^{L-1} \right\|_F + \left\| (\Sigma_1^{(i)} - \Sigma_L^{(i)}) \mathbf{T}_2^{(i)} \prod_{l=2}^{L-1} \Sigma_l^{(i)} \mathbf{T}_{l+1}^{(i)} \right\|_F + \\ & \quad \left\| \Sigma_L^{(i)} \mathbf{T}_2^{(i)} (\Sigma_2^{(i)} - \Sigma_L^{(i)}) \mathbf{T}_3^{(i)} \prod_{l=3}^{L-1} \Sigma_l^{(i)} \mathbf{T}_{l+1}^{(i)} \right\|_F + \dots + \left\| \prod_{l=1}^{L-2} \Sigma_l^{(i)} \mathbf{T}_{l+1}^{(i)} (\Sigma_{L-1}^{(i)} - \Sigma_L^{(i)}) \mathbf{T}_L^{(i)} \right\|_F \\ & \leq \left(\frac{3\sigma_{\max}^*}{2} \right)^{L-2} \left(\frac{L^2 \eta_1}{2\sigma_{\min}^*} + \frac{3\sqrt{2g_{\max}} L^2 \sigma_{\max}^*}{4\lambda \sigma_{\min}^*} \right) \|\nabla G(\mathbf{W})\|_F, \end{aligned}$$

where the first inequality uses the triangular inequality and the last inequality follows from

(121), (42) in Lemma 3, $\|\mathbf{T}_l\| = 1$ for all $l \in [L]$, and (118). This, together with (57), yields

$$\begin{aligned}
& \left\| \prod_{l=1}^{L-1} \boldsymbol{\Sigma}_l^T \mathbf{T}_{l+1} - \text{BlkDiag}(\mathbf{B}_1, \dots, \mathbf{B}_p, \mathbf{B}_{p+1}) \right\|_F \\
& \leq \sum_{i=1}^p \left\| \prod_{l=1}^{L-1} \boldsymbol{\Sigma}_l^{(i)} \mathbf{T}_{l+1}^{(i)} - \left(\prod_{l=1}^{L-1} \mathbf{T}_{l+1}^{(i)} \right) \left(\boldsymbol{\Sigma}_L^{(i)} \right)^{L-1} \right\|_F \\
& \leq p \left(\frac{3\sigma_{\max}^*}{2} \right)^{L-2} \left(\frac{L^2 \eta_1}{2\sigma_{\min}^*} + \frac{3\sqrt{2g_{\max}} L^2 \sigma_{\max}^*}{4\lambda \sigma_{\min}^*} \right) \|\nabla G(\mathbf{W})\|_F. \tag{122}
\end{aligned}$$

Now, we are ready to prove (55). Specifically, we have

$$\begin{aligned}
& \left\| (\boldsymbol{\Sigma}_L \boldsymbol{\Sigma}_L^T)^{L-1} \boldsymbol{\Sigma}_L + \lambda \boldsymbol{\Sigma}_L - \sqrt{\lambda} \mathbf{U}_L^T \mathbf{Y} \mathbf{V}_1 \text{BlkDiag}(\mathbf{B}_1, \dots, \mathbf{B}_p, \mathbf{B}_{p+1}) \right\|_F \\
& \leq \left\| (\boldsymbol{\Sigma}_L \boldsymbol{\Sigma}_L^T)^{L-1} \boldsymbol{\Sigma}_L + \lambda \boldsymbol{\Sigma}_L - \sqrt{\lambda} \mathbf{U}_L^T \mathbf{Y} \mathbf{V}_1 \prod_{l=1}^{L-1} \boldsymbol{\Sigma}_l^T \mathbf{T}_{l+1} \right\|_F + \\
& \quad \sqrt{\lambda} y_1 \left\| \prod_{l=1}^{L-1} \boldsymbol{\Sigma}_l^T \mathbf{T}_{l+1} - \text{BlkDiag}(\mathbf{B}_1, \dots, \mathbf{B}_p, \mathbf{B}_{p+1}) \right\|_F \leq c_2 \|\nabla G(\mathbf{W})\|_F,
\end{aligned}$$

where c_2 is defined in (59) and the last inequality follows from (117) and (122). Applying the same argument to $\nabla_{\mathbf{W}_1} G(\mathbf{W})$, we obtain (54).

The rest of the proof is devoted to proving the claim (119). According to (51) and $\mathbf{H}_l = \mathbf{U}_{l-1}^T \mathbf{V}_l$ for each $l = 2, \dots, L$, we have for each $i \in [p]$,

$$\begin{aligned}
& \|\mathbf{W}_l^T \mathbf{W}_l - \mathbf{W}_{l-1} \mathbf{W}_{l-1}^T\|_F = \|\mathbf{H}_l \boldsymbol{\Sigma}_l^T \boldsymbol{\Sigma}_l \mathbf{H}_l^T - \boldsymbol{\Sigma}_{l-1} \boldsymbol{\Sigma}_{l-1}^T\|_F \\
& \geq \|\mathbf{T}_l \boldsymbol{\Sigma}_l^T \boldsymbol{\Sigma}_l \mathbf{T}_l^T - \boldsymbol{\Sigma}_{l-1} \boldsymbol{\Sigma}_{l-1}^T\|_F - 2\|\boldsymbol{\Sigma}_l\|^2 \|\mathbf{H}_l - \mathbf{T}_l\|_F \\
& \geq \|\mathbf{T}_l^{(i)} (\boldsymbol{\Sigma}_l^{(i)})^2 - (\boldsymbol{\Sigma}_{l-1}^{(i)})^2 \mathbf{T}_l^{(i)}\|_F - 2\|\boldsymbol{\Sigma}_l\|^2 \|\mathbf{H}_l - \mathbf{T}_l\|_F \\
& \geq \|\mathbf{T}_l^{(i)} (\boldsymbol{\Sigma}_L^{(i)})^2 - (\boldsymbol{\Sigma}_L^{(i)})^2 \mathbf{T}_l^{(i)}\|_F - 2\|\boldsymbol{\Sigma}_l\|^2 \|\mathbf{H}_l - \mathbf{T}_l\|_F - \|(\boldsymbol{\Sigma}_L^{(i)})^2 - (\boldsymbol{\Sigma}_l^{(i)})^2\|_F - \|(\boldsymbol{\Sigma}_L^{(i)})^2 - (\boldsymbol{\Sigma}_{l-1}^{(i)})^2\|_F \\
& \geq \|\mathbf{T}_l^{(i)} (\boldsymbol{\Sigma}_L^{(i)})^2 - (\boldsymbol{\Sigma}_L^{(i)})^2 \mathbf{T}_l^{(i)}\|_F - 2\|\boldsymbol{\Sigma}_l\|^2 \|\mathbf{H}_l - \mathbf{T}_l\|_F - 3\sigma_{\max}^* \left(\|\boldsymbol{\Sigma}_L^{(i)} - \boldsymbol{\Sigma}_l^{(i)}\|_F + \|\boldsymbol{\Sigma}_L^{(i)} - \boldsymbol{\Sigma}_{l-1}^{(i)}\|_F \right),
\end{aligned}$$

where the second inequality follows from $\mathbf{T}_l^{(i)} \in \mathcal{O}^{g_i}$ for each $i \in [p]$ and the last inequality follows from (42). This together with (44), (115), and (118), yields that for each $l = 2, \dots, L$ and $i \in [p]$, we have

$$\|\mathbf{T}_l^{(i)} (\boldsymbol{\Sigma}_L^{(i)})^2 - (\boldsymbol{\Sigma}_L^{(i)})^2 \mathbf{T}_l^{(i)}\|_F \leq \eta_1 \|\nabla G(\mathbf{W})\|_F. \tag{123}$$

For each $i \in [p]$, using Weyl's inequality yields for each $j \in [g_i]$,

$$\left| \sigma_j(\boldsymbol{\Sigma}_L^{(i)}) - \sigma_{t_i}^* \right| \leq \|\mathbf{W}_L - \mathbf{W}_L^*\| \stackrel{(41)}{\leq} \frac{\sigma_{\min}^*}{2},$$

which implies $\|\boldsymbol{\Sigma}_L^{(i)} - \sigma_{t_i}^* \mathbf{I}\| \leq \sigma_{\min}^*/2$. This, together with Lemma 5 and $\mathbf{T}_l^{(i)} \in \mathcal{O}^{g_i}$, yields

$$\|\mathbf{T}_l^{(i)} (\boldsymbol{\Sigma}_L^{(i)})^2 - (\boldsymbol{\Sigma}_L^{(i)})^2 \mathbf{T}_l^{(i)}\|_F \geq \sigma_{\min}^* \|\mathbf{T}_l^{(i)} \boldsymbol{\Sigma}_L^{(i)} - \boldsymbol{\Sigma}_L^{(i)} \mathbf{T}_l^{(i)}\|_F.$$

Using this and (123), we obtain (119). \square

Since the assumptions are different for the cases $L = 2$ and $L \geq 3$ to bound the singular values by the gradient norm, we prove Proposition 5 by addressing these two scenarios separately. Each case requires a distinct approach to account for the structural differences due to the varying depth of the networks.

Proof of Proposition 5. Let $\mathbf{W}^* = (\mathbf{W}_1^*, \dots, \mathbf{W}_L^*)$ be such that $\text{dist}(\mathbf{W}, \mathcal{W}_{\sigma^*}) = \|\mathbf{W} - \mathbf{W}^*\|_F$. According to (62a) or (64), we have $\text{dist}(\mathbf{W}, \mathcal{W}) \leq \delta_\sigma/3 < \sigma_{\min}^*/2$. This, together with Lemma 3 and Corollary 2, implies that (42)-(45) and (48) hold. According to $\text{dist}(\mathbf{W}, \mathcal{W}) \leq \delta_\sigma/3$, Proposition 4, and Lemma 4, there exist matrices $\mathbf{T}_l^{(i)} \in \mathcal{O}^{g_i}$ for all $i \in [p]$ and $\mathbf{T}_l^{(p+1)} \in \mathcal{O}^{d_{l-1}-r_\sigma}$ such that (53), (54), and (55) hold for all $l = 2, \dots, L$, where \mathbf{A}_i and \mathbf{B}_i for each $i \in [p+1]$ are respectively defined in (56) and (57). Recall that (51) denotes an SVD of \mathbf{W}_l and let $\mathbf{H}_l := \mathbf{U}_{l-1}^T \mathbf{V}_l$ and $\mathbf{T}_l := \text{BlkDiag}(\mathbf{T}_l^{(1)}, \dots, \mathbf{T}_l^{(p)}, \mathbf{T}_l^{(p+1)})$ for each $l = 2, \dots, L$. This, together with (53), implies

$$\|\mathbf{H}_l - \mathbf{T}_l\|_F \leq \frac{9\sigma_{\max}^*}{4\delta_\sigma \lambda \sigma_{\min}^*} \|\nabla G(\mathbf{W})\|_F, \quad l = 2, \dots, L. \quad (124)$$

(i) For ease of exposition, let $\Psi := \mathbf{U}_2^T \mathbf{Y} \mathbf{V}_1 \in \mathbb{R}^{d_2 \times d_0}$. Now, we partition Ψ into the following block form

$$\Psi = \begin{bmatrix} \Psi^{(1,1)} & \dots & \Psi^{(1,p+1)} \\ \vdots & \ddots & \vdots \\ \Psi^{(p+1,1)} & \dots & \Psi^{(p+1,p+1)} \end{bmatrix}, \quad (125)$$

where $\Psi^{(i,j)} \in \mathbb{R}^{g_i \times g_j}$ for each $i, j \in [p]$, and $\Psi^{(p+1,j)} \in \mathbb{R}^{(d_2-r_\sigma) \times g_j}$, $\Psi^{(i,p+1)} \in \mathbb{R}^{g_i \times (d_0-r_\sigma)}$ for each $i, j \in [p]$, and $\Psi^{(p+1,p+1)} \in \mathbb{R}^{(d_2-r_\sigma) \times (d_0-r_\sigma)}$. According to (54), we have

$$c_2^2 \|\nabla G(\mathbf{W})\|_F^2 \geq \left\| (\Sigma_1 \Sigma_1^T) \Sigma_1 + \lambda \Sigma_1 - \sqrt{\lambda} \text{BlkDiag}(\mathbf{A}_1, \dots, \mathbf{A}_p, \mathbf{A}_{p+1}) \Psi \right\|_F^2.$$

By dropping the diagonal blocks and the $(p+1)$ -th row blocks of the matrix on the right-hand side, we have

$$c_2^2 \|\nabla G(\mathbf{W})\|_F^2 \geq \lambda \sum_{i=1}^p \sum_{j=1, j \neq i}^{p+1} \|\mathbf{A}_i \Psi^{(i,j)}\|_F^2 \geq \frac{\lambda \sigma_{\min}^{*2}}{4} \sum_{i=1}^p \sum_{j=1, j \neq i}^{p+1} \|\Psi^{(i,j)}\|_F^2, \quad (126)$$

where the second inequality follows from (43) and (56). Applying the same argument to (55) yields

$$c_2^2 \|\nabla G(\mathbf{W})\|_F^2 \geq \frac{\lambda \sigma_{\min}^{*2}}{4} \sum_{j=1}^p \sum_{i=1, i \neq j}^{p+1} \|\Psi^{(i,j)}\|_F^2.$$

This, together with (126), yields

$$\frac{8c_2^2}{\lambda \sigma_{\min}^{*2}} \|\nabla G(\mathbf{W})\|_F^2 \geq \sum_{i=1}^{p+1} \sum_{j=1, j \neq i}^{p+1} \|\Psi^{(i,j)}\|_F^2. \quad (127)$$

Using (60), we define $\delta := \min \left\{ \min_{i \in [s_{p_Y}]} \left| \sqrt{\lambda} - y_i \right|, \sqrt{\lambda} \right\} > 0$. Substituting this into (62b) yields

$$\frac{\delta}{3} \geq \frac{2\sqrt{2}c_2}{\sqrt{\lambda} \sigma_{\min}^*} \|\nabla G(\mathbf{W})\|_F \geq \left\| \Psi - \text{BlkDiag}(\Psi^{(1,1)}, \dots, \Psi^{(p+1,p+1)}) \right\|_F, \quad (128)$$

where the second inequality uses (125) and (127). This, together with $\Psi = \mathbf{U}_2^T \mathbf{Y} \mathbf{V}_1$ and Weyl's inequality, yields for all $i \in [d_{\min}]$,

$$\begin{aligned} \frac{\delta}{3} &\geq \left| \sigma_i \left(\text{BlkDiag}(\Psi^{(1,1)}, \dots, \Psi^{(p+1,p+1)}) \right) - y_i \right| \\ &\geq \left| y_i - \sqrt{\lambda} \right| - \left| \sigma_i \left(\text{BlkDiag}(\Psi^{(1,1)}, \dots, \Psi^{(p+1,p+1)}) \right) - \sqrt{\lambda} \right| \\ &\geq \delta - \left| \sigma_i \left(\text{BlkDiag}(\Psi^{(1,1)}, \dots, \Psi^{(p+1,p+1)}) \right) - \sqrt{\lambda} \right|, \end{aligned}$$

which implies $\left| \sigma_i(\Psi^{(p+1,p+1)}) - \sqrt{\lambda} \right| \geq 2\delta/3$ for all $i \in [d_{\min} - r_\sigma]$. According to this, we obtain for all $i \in [d_{\min} - r_\sigma]$,

$$\left| \sigma_i^2(\Psi^{(p+1,p+1)}) - \lambda \right| = \left| \sigma_i(\Psi^{(p+1,p+1)}) - \sqrt{\lambda} \right| \left| \sigma_i(\Psi^{(p+1,p+1)}) + \sqrt{\lambda} \right| \geq \frac{2\delta\sqrt{\lambda}}{3}. \quad (129)$$

Using (54) and (55), we have

$$c_2 \|\nabla G(\mathbf{W})\|_F \geq \left\| (\Sigma_1^{(p+1)} \Sigma_1^{(p+1)T}) \Sigma_1^{(p+1)} + \lambda \Sigma_1^{(p+1)} - \sqrt{\lambda} \mathbf{T}_2^{(p+1)} \Sigma_2^{(p+1)T} \Psi^{(p+1,p+1)} \right\|_F, \quad (130)$$

$$c_2 \|\nabla G(\mathbf{W})\|_F \geq \left\| (\Sigma_2^{(p+1)} \Sigma_2^{(p+1)T}) \Sigma_2^{(p+1)} + \lambda \Sigma_2^{(p+1)} - \sqrt{\lambda} \Psi^{(p+1,p+1)} \Sigma_2^{(p+1)T} \mathbf{T}_2^{(p+1)} \right\|_F. \quad (131)$$

According to Mirsky's inequality (see Lemma 11), we have

$$\|\Sigma_1^{(p+1)}\|_F = \|\Sigma_1^{(p+1)} - \mathbf{0}\|_F \leq \|\Sigma_1 - \Sigma_1^*\|_F \leq \|\mathbf{W}_1 - \mathbf{W}_1^*\|_F \leq \text{dist}(\mathbf{W}, \mathcal{W}_{\sigma^*}). \quad (132)$$

Similarly, we have

$$\|\Sigma_2^{(p+1)}\|_F \leq \text{dist}(\mathbf{W}, \mathcal{W}_{\sigma^*}). \quad (133)$$

According to $\Psi = \mathbf{U}_2^T \mathbf{Y} \mathbf{V}_1$ and (125), we have

$$\|\Psi^{(p+1,p+1)}\| \leq \|\mathbf{Y}\| = y_1. \quad (134)$$

Using this and (130), we obtain

$$\begin{aligned} c_2 y_1 \|\nabla G(\mathbf{W})\|_F &\geq \left\| \left((\Sigma_1^{(p+1)} \Sigma_1^{(p+1)T}) \Sigma_1^{(p+1)} + \lambda \Sigma_1^{(p+1)} - \sqrt{\lambda} \mathbf{T}_2^{(p+1)} \Sigma_2^{(p+1)T} \Psi^{(p+1,p+1)} \right) \Psi^{(p+1,p+1)T} \right\|_F \\ &\geq \left\| \lambda \Sigma_1^{(p+1)} \Psi^{(p+1,p+1)T} - \sqrt{\lambda} \mathbf{T}_2^{(p+1)} \Sigma_2^{(p+1)T} \Psi^{(p+1,p+1)} \Psi^{(p+1,p+1)T} \right\|_F \\ &\quad - \left\| \Sigma_1^{(p+1)} \Sigma_1^{(p+1)T} \Sigma_1^{(p+1)} \Psi^{(p+1,p+1)T} \right\|_F \\ &\geq \left\| \lambda \mathbf{T}_2^{(p+1)T} \Sigma_1^{(p+1)T} \Psi^{(p+1,p+1)T} - \sqrt{\lambda} \Sigma_2^{(p+1)T} \Psi^{(p+1,p+1)} \Psi^{(p+1,p+1)T} \right\|_F - y_1 \|\Sigma_1^{(p+1)}\|_F^3, \end{aligned}$$

where the last inequality follows from $\mathbf{T}_2^{(p+1)} \in \mathcal{O}^{d_2 - r_\sigma}$ and (134). Using (131), we have

$$\begin{aligned} c_2 \sqrt{\lambda} \|\nabla G(\mathbf{W})\|_F &\geq \sqrt{\lambda} \left\| \lambda \Sigma_2^{(p+1)} - \sqrt{\lambda} \Psi^{(p+1,p+1)} \Sigma_1^{(p+1)T} \mathbf{T}_2^{(p+1)} \right\|_F - \sqrt{\lambda} \left\| \Sigma_2^{(p+1)} \Sigma_2^{(p+1)T} \Sigma_2^{(p+1)} \right\|_F \\ &\geq \left\| \lambda^{\frac{3}{2}} \Sigma_2^{(p+1)T} - \lambda \mathbf{T}_2^{(p+1)T} \Sigma_1^{(p+1)} \Psi^{(p+1,p+1)T} \right\|_F - \sqrt{\lambda} \left\| \Sigma_2^{(p+1)} \right\|_F^3. \end{aligned}$$

Summing up the above two inequalities yields

$$\begin{aligned} c_2 (y_1 + \sqrt{\lambda}) \|\nabla G(\mathbf{W})\|_F &\geq \left\| \lambda^{\frac{3}{2}} \Sigma_2^{(p+1)T} - \sqrt{\lambda} \Sigma_2^{(p+1)T} \Psi^{(p+1,p+1)} \Psi^{(p+1,p+1)T} \right\|_F \\ &\quad - \sqrt{\lambda} \left\| \Sigma_2^{(p+1)} \right\|_F^3 - y_1 \left\| \Sigma_1^{(p+1)} \right\|_F^3. \end{aligned} \quad (135)$$

Repeating the argument in (134)-(135) with multiplying $\sqrt{\lambda}$ on (130) and y_1 on (131), we obtain

$$\begin{aligned} c_2 (y_1 + \sqrt{\lambda}) \|\nabla G(\mathbf{W})\|_F &\geq \left\| \lambda^{\frac{3}{2}} \Sigma_1^{(p+1)T} - \sqrt{\lambda} \Psi^{(p+1,p+1)T} \Psi^{(p+1,p+1)} \Sigma_1^{(p+1)T} \right\|_F \\ &\quad - \sqrt{\lambda} \left\| \Sigma_1^{(p+1)} \right\|_F^3 - y_1 \left\| \Sigma_2^{(p+1)} \right\|_F^3. \end{aligned}$$

Summing up the above two inequalities yields

$$\begin{aligned}
& 2c_2(y_1 + \sqrt{\lambda})\|\nabla G(\mathbf{W})\|_F \\
& \geq \sqrt{\lambda}\sigma_{\min}\left(\Psi^{(p+1,p+1)}\Psi^{(p+1,p+1)T} - \lambda\mathbf{I}\right)\|\Sigma_2^{(p+1)}\|_F + \sqrt{\lambda}\sigma_{\min}\left(\Psi^{(p+1,p+1)T}\Psi^{(p+1,p+1)} - \lambda\mathbf{I}\right) \\
& \quad \|\Sigma_1^{(p+1)}\|_F - (\sqrt{\lambda} + y_1)(\|\Sigma_2^{(p+1)}\|_F^3 + \|\Sigma_1^{(p+1)}\|_F^3) \\
& \geq \left(\frac{2\lambda\delta}{3} - (\sqrt{\lambda} + y_1)\|\Sigma_1^{(p+1)}\|_F^2\right)\|\Sigma_1^{(p+1)}\|_F + \left(\frac{2\lambda\delta}{3} - (\sqrt{\lambda} + y_1)\|\Sigma_2^{(p+1)}\|_F^2\right)\|\Sigma_2^{(p+1)}\|_F \\
& \geq \frac{\lambda\delta}{3}(\|\Sigma_1^{(p+1)}\|_F + \|\Sigma_2^{(p+1)}\|_F),
\end{aligned}$$

where the second inequality follows from (129), and the last inequality follows from (62a), (132) and (133). This directly implies (63).

(ii) For ease of exposition, let $\Psi := \mathbf{U}_L^T \mathbf{Y} \mathbf{V}_1 \in \mathbb{R}^{d_L \times d_0}$. Let $l \in [L]$ be such that $l \in \arg \max\{\sigma_{r_{\sigma+1}}(\mathbf{W}_k) : k \in [L]\}$. Substituting (51) into (48) yields

$$\begin{aligned}
c_1\|\nabla G(\mathbf{W})\|_F & \geq \left\| (\Sigma_l \Sigma_l^T)^{L-1} \Sigma_l - \sqrt{\lambda} \left(\prod_{k=l+1}^L \mathbf{H}_k \Sigma_k^T \right) \Psi \left(\prod_{k=1}^{l-1} \Sigma_k^T \mathbf{H}_{k+1} \right) + \lambda \Sigma_l \right\|_F \\
& \geq \left\| (\Sigma_l \Sigma_l^T)^{L-1} \Sigma_l - \sqrt{\lambda} \left(\prod_{k=l+1}^L \mathbf{T}_k \Sigma_k^T \right) \Psi \left(\prod_{k=1}^{l-1} \Sigma_k^T \mathbf{T}_{k+1} \right) + \lambda \Sigma_l \right\|_F - \sqrt{\lambda} \|\Delta\|_F,
\end{aligned}$$

where the second inequality uses the triangular inequality and

$$\begin{aligned}
\Delta & := (\mathbf{H}_{l+1} - \mathbf{T}_{l+1}) \Sigma_{l+1}^T \left(\prod_{k=l+2}^L \mathbf{H}_k \Sigma_k^T \right) \Psi \left(\prod_{k=1}^{l-1} \Sigma_k^T \mathbf{H}_{k+1} \right) \\
& \quad + \mathbf{T}_{l+1} \Sigma_{l+1}^T (\mathbf{H}_{l+2} - \mathbf{T}_{l+2}) \Sigma_{l+2}^T \left(\prod_{k=l+3}^L \mathbf{H}_k \Sigma_k^T \right) \Psi \left(\prod_{k=1}^{l-1} \Sigma_k^T \mathbf{H}_{k+1} \right) + \dots \\
& \quad + \left(\prod_{k=l+1}^L \mathbf{T}_k \Sigma_k^T \right) \Psi \left(\prod_{k=1}^{l-2} \Sigma_k^T \mathbf{T}_{k+1} \right) \Sigma_{l-1}^T (\mathbf{H}_l - \mathbf{T}_l).
\end{aligned}$$

Using (42) and (124), we obtain

$$\|\Delta\| \leq \left(\frac{3}{2} \sigma_{\max}^* \right)^L \frac{3(L-1)y_1}{2\delta_{\sigma}\lambda\sigma_{\min}^*} \|\nabla G(\mathbf{W})\|_F.$$

This, together with the above inequality, yields

$$\left\| (\Sigma_l \Sigma_l^T)^{L-1} \Sigma_l - \sqrt{\lambda} \left(\prod_{k=l+1}^L \mathbf{T}_k \Sigma_k^T \right) \Psi \left(\prod_{k=1}^{l-1} \Sigma_k^T \mathbf{T}_{k+1} \right) + \lambda \Sigma_l \right\|_F \leq \eta_2 \|\nabla G(\mathbf{W})\|_F,$$

where η_2 is defined in (66). Using this inequality, the block structures of Σ_l and \mathbf{T}_l , and the fact that $l \in \arg \max\{\sigma_{r_{\sigma+1}}(\mathbf{W}_k) : k \in [L]\}$, we obtain

$$\begin{aligned}
\eta_2 \|\nabla G(\mathbf{W})\|_F & \geq \left\| \left(\Sigma_l^{(p+1)} \Sigma_l^{(p+1)T} \right)^{L-1} \Sigma_l^{(p+1)} - \sqrt{\lambda} \left(\prod_{k=l+1}^L \mathbf{T}_k^{(p+1)} \Sigma_k^{(p+1)T} \right) \right. \\
& \quad \left. \Psi_{r_{\sigma+1}:d_L, r_{\sigma+1}:d_0} \left(\prod_{k=1}^{l-1} \Sigma_k^{(p+1)T} \mathbf{T}_{k+1}^{(p+1)} \right) + \lambda \Sigma_l^{(p+1)} \right\| \\
& \geq \left\| \left(\Sigma_l^{(p+1)} \Sigma_l^{(p+1)T} \right)^{L-1} \Sigma_l^{(p+1)} + \lambda \Sigma_l^{(p+1)} \right\| - \sqrt{\lambda} \sigma_{r_{\sigma+1}}^{L-1}(\mathbf{W}_l) \|\Psi_{r_{\sigma+1}:d_L, r_{\sigma+1}:d_0}\| \\
& \geq \lambda \sigma_{r_{\sigma+1}}(\mathbf{W}_l) - \sqrt{\lambda} y_1 \sigma_{r_{\sigma+1}}^{L-1}(\mathbf{W}_l) \geq \sqrt{\lambda} \left(\sqrt{\lambda} - y_1 \sigma_{r_{\sigma+1}}^{L-2}(\mathbf{W}_l) \right) \sigma_{r_{\sigma+1}}(\mathbf{W}_l), \tag{136}
\end{aligned}$$

where $\Psi_{r_\sigma+1:d_L, r_\sigma+1:d_0}$ denotes the submatrix of Ψ with the rows being indexed from $r_\sigma + 1$ to d_L and the columns being indexed from $r_\sigma + 1$ to d_0 , and the third inequality is due to $\|\Psi_{r_\sigma+1:d_L, r_\sigma+1:d_0}\| \leq \|\mathbf{Y}\| = y_1$. Using Weyl's inequality and (64) with $L \geq 3$, we have

$$\sigma_{r_\sigma+1}(\mathbf{W}_l) \leq \text{dist}(\mathbf{W}, \mathcal{W}_{\sigma^*}) \leq \left(\frac{\sqrt{\lambda}}{2y_1} \right)^{1/(L-2)}.$$

This, together with (136), directly yields $\sigma_{r_\sigma+1}(\mathbf{W}_l) \leq 2\eta_2 \|\nabla G(\mathbf{W})\|_F / \lambda$, which further implies (65). \square

Equipped with Lemma 3, Corollary 2, Proposition 4, Lemma 4, Proposition 5, and Davis-Kahan Theorem (see Lemma 10), we are ready to prove Proposition 6.

Proof of Proposition 6. According to (62a) or (64), we have $\text{dist}(\mathbf{W}, \mathcal{W}) \leq \delta_\sigma / 3 < \sigma_{\min}^* / 2$. This, together with Lemma 3 and Corollary 2, implies that (42)-(45) and (48) hold. According to Proposition 4 and Lemma 4 with $\text{dist}(\mathbf{W}, \mathcal{W}) \leq \delta_\sigma / 3$, there exist matrices $\mathbf{T}_l^{(i)} \in \mathcal{O}^{g_i}$ for each $i \in [p]$ and $\mathbf{T}_l^{(p+1)} \in \mathcal{O}^{d_{l-1}-r_\sigma}$ such that (54) and (55) hold, where \mathbf{A}_i and \mathbf{B}_i for each $i \in [p+1]$ are defined in (56) and (57). According to (62) (resp., (64)) and Proposition 5, we have (63) (resp., (65)) hold. Recall that (51) denotes an SVD of \mathbf{W}_l for each $l \in [L]$. For ease of exposition, let

$$\Psi := \mathbf{U}_L^T \mathbf{Y} \mathbf{V}_1, \quad \hat{\Sigma}_l := \text{BlkDiag} \left(\Sigma_l^{(1)}, \dots, \Sigma_l^{(p)}, \mathbf{0}_{(d_l-r_\sigma) \times (d_{l-1}-r_\sigma)} \right), \quad \forall l \in [L]. \quad (137)$$

where $\Sigma_l^{(1)}, \dots, \Sigma_l^{(p)}$ are defined in (51) for each $l \in [L]$. Using the form of Σ_l in (51) and (63) (resp., (65)), we have for each $l \in [L]$,

$$\|\hat{\Sigma}_l - \Sigma_l\|_F = \|\Sigma_l^{(p+1)}\|_F \leq c_3 \sqrt{\min\{d_l, d_{l-1}\}} \|\nabla G(\mathbf{W})\|_F \leq c_3 \sqrt{d_{\max}} \|\nabla G(\mathbf{W})\|_F, \quad (138)$$

$$\begin{aligned} \left\| (\hat{\Sigma}_l \hat{\Sigma}_l^T)^{L-1} \hat{\Sigma}_l - (\Sigma_l \Sigma_l^T)^{L-1} \Sigma_l \right\|_F &= \left\| \left(\Sigma_l^{(p+1)} \Sigma_l^{(p+1)T} \right)^{L-1} \Sigma_l^{(p+1)} \right\|_F \\ &\leq c_3 \sqrt{d_{\max}} \left(\frac{3\sigma_{\max}^*}{2} \right)^{2(L-1)} \|\nabla G(\mathbf{W})\|_F, \end{aligned} \quad (139)$$

where the inequality uses (42) and (138).

(i) We first focus on $l = L$. Now, we bound

$$\begin{aligned} &\left\| (\hat{\Sigma}_L \hat{\Sigma}_L^T)^{L-1} \hat{\Sigma}_L + \lambda \hat{\Sigma}_L - \sqrt{\lambda} \Psi \text{BlkDiag}(\mathbf{B}_1, \dots, \mathbf{B}_p, \mathbf{0}) \right\|_F \\ &\leq \left\| (\Sigma_L \Sigma_L^T)^{L-1} \Sigma_L + \lambda \Sigma_L - \sqrt{\lambda} \Psi \text{BlkDiag}(\mathbf{B}_1, \dots, \mathbf{B}_p, \mathbf{B}_{p+1}) \right\|_F + \lambda \|\hat{\Sigma}_L - \Sigma_L\|_F \\ &\quad + \left\| (\hat{\Sigma}_L \hat{\Sigma}_L^T)^{L-1} \hat{\Sigma}_L - (\Sigma_L \Sigma_L^T)^{L-1} \Sigma_L \right\|_F + \sqrt{\lambda} \|\Psi\| \|\mathbf{B}_{p+1}\|_F \\ &\leq \left(c_2 + \left(\left(\frac{3\sigma_{\max}^*}{2} \right)^{2(L-1)} + \lambda \right) c_3 \sqrt{d_{\max}} \right) \|\nabla G(\mathbf{W})\|_F + \sqrt{\lambda} y_1 \|\mathbf{B}_{p+1}\|_F \leq \eta_3 \|\nabla G(\mathbf{W})\|_F, \end{aligned} \quad (140)$$

where the second inequality follows from (55), (57), (72), (138), and (139), and the last inequality uses (42), (57), (138), and

$$\eta_3 := c_2 + c_3 \sqrt{d_{\max}} \left(\left(\frac{3\sigma_{\max}^*}{2} \right)^{2(L-1)} + \sqrt{\lambda} y_1 \left(\frac{3\sigma_{\max}^*}{2} \right)^{L-2} + \lambda \right). \quad (141)$$

Substituting $\hat{\Sigma}_L$ in (137), (57), and (72) into (140) yields

$$\begin{aligned} & \sum_{i=1}^{r_\sigma} \left\| (\sigma_i^{2L-1}(\mathbf{W}_L) + \lambda \sigma_i(\mathbf{W}_L)) \mathbf{e}_i - \sqrt{\lambda} \sigma_i^{L-1}(\mathbf{W}_L) \Psi \hat{\mathbf{T}} \mathbf{e}_i \right\|^2 \\ & \leq \left\| (\hat{\Sigma}_L \hat{\Sigma}_L^T)^{L-1} \hat{\Sigma}_L + \lambda \hat{\Sigma}_L - \sqrt{\lambda} \Psi \hat{\mathbf{T}} \text{BlkDiag} \left((\Sigma_L^{(1)})^{L-1}, \dots, (\Sigma_L^{(p)})^{L-1}, \mathbf{0} \right) \right\|_F^2 \leq \eta_3^2 \|\nabla G(\mathbf{W})\|_F^2. \end{aligned} \quad (142)$$

Using this and the definition of $\varphi(\cdot)$ in (67), together with dividing the above inequality by $\sqrt{\lambda} \sigma_{r_\sigma}^{L-1}(\mathbf{W}_L)$ on both sides, yields

$$\sum_{i=1}^{r_\sigma} \left\| \varphi(\sigma_i(\mathbf{W}_L)) \mathbf{e}_i - \Psi \hat{\mathbf{T}} \mathbf{e}_i \right\|^2 \leq \frac{\eta_3^2}{\lambda \sigma_{r_\sigma}^{2(L-1)}(\mathbf{W}_L)} \|\nabla G(\mathbf{W})\|_F^2 \leq \frac{2^{2(L-1)} \eta_3^2}{\lambda (\sigma_{\min}^*)^{2L-2}} \|\nabla G(\mathbf{W})\|_F^2, \quad (143)$$

where the last inequality uses (43). Using the same argument of (140) for (54), we have

$$\left\| (\hat{\Sigma}_1 \hat{\Sigma}_1^T)^{L-1} \hat{\Sigma}_1 + \lambda \hat{\Sigma}_1 - \sqrt{\lambda} \text{BlkDiag}(\mathbf{A}_1, \dots, \mathbf{A}_p, \mathbf{0}) \Psi \right\|_F \leq \eta_3 \|\nabla G(\mathbf{W})\|_F. \quad (144)$$

Then we compute

$$\left\| \left(\hat{\Sigma}_1 \hat{\Sigma}_1^T \right)^{L-1} \hat{\Sigma}_1 + \lambda \hat{\Sigma}_1 - \sqrt{\lambda} \text{BlkDiag} \left((\Sigma_1^{(1)})^{L-1}, \dots, (\Sigma_1^{(p)})^{L-1}, \mathbf{0} \right) \Psi \hat{\mathbf{T}} \right\|_F \leq \eta_4 \|\nabla G(\mathbf{W})\|_F, \quad (145)$$

where

$$\eta_4 := \eta_3 + \frac{p \eta_1 (2L-1)L}{\sigma_{\min}^*} \left(\frac{3\sigma_{\max}^*}{2} \right)^{2L-2} + \frac{\lambda p \eta_1 L}{\sigma_{\min}^*}.$$

To maintain the flow of the main proof, we defer the proof of (137) to Appendix D.1. Applying a similar argument as in (143) to rows of the above inequality yields

$$\sum_{i=1}^{r_\sigma} \left\| \varphi(\sigma_i(\mathbf{W}_1)) \mathbf{e}_i^T - \mathbf{e}_i^T \Psi \hat{\mathbf{T}} \right\|^2 \leq \left(\frac{2^{L-1} \eta_4}{\sqrt{\lambda} \sigma_{\min}^{*L-1}} \right)^2 \|\nabla G(\mathbf{W})\|_F^2. \quad (146)$$

To proceed, we define $\mathbf{d} := (\varphi(\sigma_1(\mathbf{W}_L)), \dots, \varphi(\sigma_{r_\sigma}(\mathbf{W}_L))) \in \mathbb{R}^{r_\sigma}$ and

$$\mathbf{Z} := \text{BlkDiag} \left(\text{diag}(\mathbf{d}), (\Psi \hat{\mathbf{T}})_{r_\sigma+1:d_L, r_\sigma+1:d_0} \right),$$

where $(\Psi \hat{\mathbf{T}})_{r_\sigma+1:d_L, r_\sigma+1:d_0}$ denotes a submatrix of $\Psi \hat{\mathbf{T}}$ with rows indexed by $r_\sigma + 1 : d_L$ and columns indexed by $r_\sigma + 1 : d_0$. Now, we compute

$$\begin{aligned} \left\| \Psi \hat{\mathbf{T}} - \mathbf{Z} \right\|_F^2 &= \sum_{i=1}^{r_\sigma} \left\| (\Psi \hat{\mathbf{T}} - \mathbf{Z}) \mathbf{e}_i \right\|^2 + \sum_{i=r_\sigma+1}^{d_0} \left\| (\Psi \hat{\mathbf{T}} - \mathbf{Z}) \mathbf{e}_i \right\|^2 \\ &= \sum_{i=1}^{r_\sigma} \left\| \Psi \hat{\mathbf{T}} \mathbf{e}_i - \varphi(\sigma_i(\mathbf{W}_L)) \mathbf{e}_i \right\|^2 + \left\| (\Psi \hat{\mathbf{T}})_{1:r_\sigma, r_\sigma+1:d_0} \right\|_F^2 \\ &\leq \sum_{i=1}^{r_\sigma} \left\| \Psi \hat{\mathbf{T}} \mathbf{e}_i - \varphi(\sigma_i(\mathbf{W}_L)) \mathbf{e}_i \right\|^2 + \sum_{i=1}^{r_\sigma} \left\| \mathbf{e}_i^T \Psi \hat{\mathbf{T}} - \varphi(\sigma_i(\mathbf{W}_1)) \mathbf{e}_i^T \right\|^2 \\ &\leq \left(\frac{2^{L-1}}{\sqrt{\lambda} \sigma_{\min}^{*L-1}} \right)^2 (\eta_3^2 + \eta_4^2) \|\nabla G(\mathbf{W})\|_F^2, \end{aligned} \quad (147)$$

where the second equality uses the structure of \mathbf{Z} , the first inequality is due to

$$\left\| \mathbf{e}_i^T \Psi \hat{\mathbf{T}} - \varphi(\sigma_i(\mathbf{W}_1)) \mathbf{e}_i^T \right\|^2 \geq \left\| (\Psi \hat{\mathbf{T}})_{i, r_\sigma+1:d_0} \right\|^2, \quad \forall i \in [r_\sigma],$$

and the second inequality follows from (143) and (146). Let $(\Psi\hat{\mathbf{T}})_{r_\sigma+1:d_L, r_\sigma+1:d_0} = \mathbf{P}\mathbf{\Lambda}\mathbf{Q}^T$ be an SVD of $(\Psi\hat{\mathbf{T}})_{r_\sigma+1:d_L, r_\sigma+1:d_0}$, where $\mathbf{P} \in \mathcal{O}^{d_L-r_\sigma}$, $\mathbf{Q} \in \mathcal{O}^{d_0-r_\sigma}$, and $\mathbf{\Lambda} \in \mathbb{R}^{(d_L-r_\sigma) \times (d_0-r_\sigma)}$ with diagonal elements $\gamma_1 \geq \dots \geq \gamma_{d_{\min}-r_\sigma}$ being top $d_{\min}-r_\sigma$ singular values. Now, we define

$$\mathbf{c} := (\varphi(\sigma_1(\mathbf{W}_L)), \dots, \varphi(\sigma_{r_\sigma}(\mathbf{W}_L)), \gamma_1, \dots, \gamma_{d_{\min}-r_\sigma}) \in \mathbb{R}^{d_{\min}}. \quad (148)$$

Let $\mathbf{\Pi} \in \mathcal{P}^{d_{\min}}$ be a permutation matrix such that the entries $\mathbf{\Pi}\mathbf{c}$ are in a decreasing order and

$$\begin{aligned} \bar{\mathbf{U}}_L &:= \mathbf{U}_L \text{BlkDiag}(\mathbf{I}_{r_\sigma}, \mathbf{P}^T) \text{BlkDiag}(\mathbf{\Pi}^T, \mathbf{I}_{d_L-d_{\min}}), \\ \bar{\mathbf{V}}_1 &:= \mathbf{V}_1 \hat{\mathbf{T}} \text{BlkDiag}(\mathbf{I}_{r_\sigma}, \mathbf{Q}^T) \text{BlkDiag}(\mathbf{\Pi}^T, \mathbf{I}_{d_0-d_{\min}}), \\ \bar{\mathbf{Y}} &:= \text{BlkDiag}(\mathbf{\Pi} \text{diag}(\mathbf{c}) \mathbf{\Pi}^T, \mathbf{0}_{(d_L-d_{\min}) \times (d_0-d_{\min})}). \end{aligned} \quad (149)$$

Since both \mathbf{Y} and $\bar{\mathbf{Y}}$ are diagonal matrices with elements in decreasing order, by using Mirsky's inequality (see Lemma 11), we have

$$\|\mathbf{Y} - \bar{\mathbf{Y}}\|_F \leq \|\bar{\mathbf{U}}_L^T \mathbf{Y} \bar{\mathbf{V}}_1 - \bar{\mathbf{Y}}\|_F \leq \frac{2^{L-1}}{\sqrt{\lambda \sigma_{\min}^{*L-1}}} \sqrt{\eta_3^2 + \eta_4^2} \|\nabla G(\mathbf{W})\|_F, \quad (150)$$

where the last inequality uses (147). This, together with (68b), implies

$$\|\mathbf{Y} - \bar{\mathbf{Y}}\|_F \leq \frac{\delta_y}{3}, \quad \|\bar{\mathbf{Y}}\| \leq \|\bar{\mathbf{Y}} - \mathbf{Y}\| + \|\mathbf{Y}\| \leq y_1 + \frac{\delta_y}{3}. \quad (151)$$

According to the triangular inequality, we obtain

$$\begin{aligned} \|\bar{\mathbf{U}}_L^T \mathbf{Y} \mathbf{Y}^T \bar{\mathbf{U}}_L - \bar{\mathbf{Y}} \bar{\mathbf{Y}}^T\|_F &\leq \|(\bar{\mathbf{U}}_L^T \mathbf{Y} \bar{\mathbf{V}}_1 - \bar{\mathbf{Y}}) \bar{\mathbf{V}}_1^T \mathbf{Y}^T \bar{\mathbf{U}}_L\|_F + \|\bar{\mathbf{Y}} (\bar{\mathbf{U}}_L^T \mathbf{Y} \bar{\mathbf{V}}_1 - \bar{\mathbf{Y}})^T\|_F \\ &\leq \left(2y_1 + \frac{\delta_y}{3}\right) \frac{2^{L-1} \sqrt{\eta_3^2 + \eta_4^2}}{\sqrt{\lambda \sigma_{\min}^{*L-1}}} \|\nabla G(\mathbf{W})\|_F, \end{aligned} \quad (152)$$

where the second inequality uses (150) and (151). Next, we write $\bar{\mathbf{U}}_L$ and \mathbf{I}_{d_L} as follows:

$$\begin{aligned} \bar{\mathbf{U}}_L &= \begin{bmatrix} \bar{\mathbf{U}}_L^{(1)} \\ \vdots \\ \bar{\mathbf{U}}_L^{(p_Y)} \\ \bar{\mathbf{U}}_L^{(p_Y+1)} \end{bmatrix}, \quad \mathbf{I}_{d_L} = \begin{bmatrix} \mathbf{E}^{(1)} \\ \vdots \\ \mathbf{E}^{(p_Y)} \\ \mathbf{E}^{(p_Y+1)} \end{bmatrix}, \quad \mathbf{E}^{(1)} = [\mathbf{I}_{h_1}, \mathbf{0}_{h_1 \times (d_L-s_1)}], \\ \mathbf{E}^{(i)} &= [\mathbf{0}_{h_i \times s_{i-1}}, \mathbf{I}_{h_i}, \mathbf{0}_{h_i \times (d_L-s_i)}], \quad \mathbf{E}^{(p_Y+1)} = [\mathbf{0}_{(d_L-r_Y) \times r_Y}, \mathbf{I}_{d_L-r_Y}], \end{aligned}$$

where $\bar{\mathbf{U}}_L^{(i)}, \mathbf{E}^{(i)} \in \mathbb{R}^{h_i \times d_L}$ for all $i \in [p_Y]$, and $\bar{\mathbf{U}}_L^{(p_Y+1)}, \mathbf{E}^{(p_Y+1)} \in \mathbb{R}^{(d_L-r_Y) \times d_L}$. Applying the Davis-Kahan Theorem (see Lemma 10) to the matrices $\bar{\mathbf{U}}_L^T \mathbf{Y} \mathbf{Y}^T \bar{\mathbf{U}}_L$ (considered as \mathbf{A}), $\bar{\mathbf{Y}} \bar{\mathbf{Y}}^T$ (considered as $\hat{\mathbf{A}}$), $\bar{\mathbf{U}}_L^{(i)T}$ (considered as \mathbf{V}), and $\mathbf{E}^{(i)T}$ (considered as $\hat{\mathbf{V}}$) for each $i \in [p_Y+1]$, there exist orthogonal matrices $\hat{\mathbf{U}}_L^{(i)} \in \mathcal{O}^{h_i}$ for each $i \in [p_Y]$ and $\hat{\mathbf{U}}_L^{(p_Y+1)} \in \mathcal{O}^{d_L-r_Y}$ such that

$$\|\bar{\mathbf{U}}_L^{(i)} - \hat{\mathbf{U}}_L^{(i)} \mathbf{E}^{(i)}\|_F \leq \frac{4\|\bar{\mathbf{U}}_L^T \mathbf{Y} \mathbf{Y}^T \bar{\mathbf{U}}_L - \bar{\mathbf{Y}} \bar{\mathbf{Y}}^T\|_F}{\min\{\lambda_{i-1} - \lambda_i, \lambda_i - \lambda_{i+1}\}}, \quad \forall i \in [p_Y+1], \quad (153)$$

where λ_i is the i -th largest non-repeated eigenvalue of $\bar{\mathbf{U}}_L^T \mathbf{Y} \mathbf{Y}^T \bar{\mathbf{U}}_L$ and $\lambda_0 = \infty$ and $\lambda_{p_Y+2} = -\infty$ by convention. Here, according to the SVD of \mathbf{Y} in (22), we compute $\lambda_i = y_{s_i}^2$ for $i \in [p_Y]$ and $\lambda_{p_Y+1} = 0$ for $i \in [p_Y]$. Using (152) and the following inequalities

$$\begin{aligned} y_{s_j}^2 - y_{s_{j+1}}^2 &= (y_{s_j} - y_{s_{j+1}})(y_{s_j} + y_{s_{j+1}}) \geq \delta_y y_{s_{p_Y}}, \quad \forall j \in [p_Y], \\ \min\{y_{s_i}^2 - y_{s_{i+1}}^2, y_{s_{i+1}}^2 - y_{s_{i+1+1}}^2\} &\geq \delta_y y_{s_{p_Y}}, \quad \forall i \in [p_Y-1], \end{aligned}$$

(153) yields

$$\left\| \bar{\mathbf{U}}_L^{(i)} - \hat{\mathbf{U}}_L^{(i)} \mathbf{E}^{(i)} \right\|_F \leq \frac{2^{L+1}(6y_1 + \delta_y) \sqrt{\eta_3^2 + \eta_4^2}}{3\sqrt{\lambda} y_{s_{p_Y}} \delta_y \sigma_{\min}^{*L-1}} \|\nabla G(\mathbf{W})\|_F.$$

Therefore, we have

$$\left\| \bar{\mathbf{U}}_L - \text{BlkDiag} \left(\hat{\mathbf{U}}_L^{(1)}, \dots, \hat{\mathbf{U}}_L^{(p_Y+1)} \right) \right\|_F \leq \sum_{i=1}^{p_Y+1} \left\| \bar{\mathbf{U}}_L^{(i)} - \hat{\mathbf{U}}_L^{(i)} \mathbf{E}^{(i)} \right\|_F \leq \eta_5 \|\nabla G(\mathbf{W})\|_F, \quad (154)$$

where

$$\eta_5 := \frac{2^{L+1}(6y_1 + \delta_y)(p_Y + 1) \sqrt{\eta_3^2 + \eta_4^2}}{3\sqrt{\lambda} y_{s_{p_Y}} \delta_y \sigma_{\min}^{*L-1}}. \quad (155)$$

According to definition of $\tilde{\mathbf{U}}_L$ in (70), we have

$$\left\| \mathbf{U}_L - \tilde{\mathbf{U}}_L \right\|_F = \left\| \bar{\mathbf{U}}_L - \text{BlkDiag} \left(\hat{\mathbf{U}}_L^{(1)}, \dots, \hat{\mathbf{U}}_L^{(p_Y)}, \hat{\mathbf{U}}_L^{(p_Y+1)} \right) \right\|_F \leq \eta_5 \|\nabla G(\mathbf{W})\|_F,$$

where the second equality uses the definition of $\bar{\mathbf{U}}_L$ in (149) and the inequality is due to (154). Using the same argument in (152)-(154) to $\bar{\mathbf{V}}_1^T \mathbf{Y}^T \mathbf{Y} \bar{\mathbf{V}}_1$ and $\bar{\mathbf{Y}}^T \bar{\mathbf{Y}}$, we conclude that there exist orthogonal matrices $\hat{\mathbf{V}}_1^{(i)} \in \mathcal{O}^{h_i}$ for each $i \in [p]$ and $\hat{\mathbf{V}}_1^{(p_Y+1)} \in \mathcal{O}^{d_0-r_Y}$ such that

$$\left\| \bar{\mathbf{V}}_1 - \text{BlkDiag} \left(\hat{\mathbf{V}}_1^{(1)}, \dots, \hat{\mathbf{V}}_1^{(p_Y)}, \hat{\mathbf{V}}_1^{(p_Y+1)} \right) \right\|_F \leq \eta_5 \|\nabla G(\mathbf{W})\|_F. \quad (156)$$

Moreover, we have

$$\begin{aligned} & \left\| \text{BlkDiag} \left(\hat{\mathbf{U}}_L^{(1)}, \dots, \hat{\mathbf{U}}_L^{(p_Y)}, \hat{\mathbf{U}}_L^{(p_Y+1)} \right) \mathbf{Y} \text{BlkDiag} \left(\hat{\mathbf{V}}_1^{(1)}, \dots, \hat{\mathbf{V}}_1^{(p_Y)}, \hat{\mathbf{V}}_1^{(p_Y+1)} \right)^T - \mathbf{Y} \right\|_F \\ & \leq \left\| \left(\text{BlkDiag} \left(\hat{\mathbf{U}}_L^{(1)}, \dots, \hat{\mathbf{U}}_L^{(p_Y)}, \hat{\mathbf{U}}_L^{(p_Y+1)} \right) - \bar{\mathbf{U}}_L \right) \mathbf{Y} \text{BlkDiag} \left(\hat{\mathbf{V}}_1^{(1)}, \dots, \hat{\mathbf{V}}_1^{(p_Y)}, \hat{\mathbf{V}}_1^{(p_Y+1)} \right)^T \right\|_F + \\ & \quad \left\| \bar{\mathbf{U}}_L \mathbf{Y} \left(\text{BlkDiag} \left(\hat{\mathbf{V}}_1^{(1)}, \dots, \hat{\mathbf{V}}_1^{(p_Y)}, \hat{\mathbf{V}}_1^{(p_Y+1)} \right) - \bar{\mathbf{V}}_1 \right)^T \right\|_F + \|\bar{\mathbf{U}}_L \mathbf{Y} \bar{\mathbf{V}}_1^T - \bar{\mathbf{Y}}\|_F + \|\bar{\mathbf{Y}} - \mathbf{Y}\|_F \\ & \leq \left(\frac{2^L \sqrt{\eta_3^2 + \eta_4^2}}{\sqrt{\lambda} \sigma_{\min}^{*L-1}} + 2y_1 \eta_5 \right) \|\nabla G(\mathbf{W})\|_F, \end{aligned}$$

where the last inequality uses (150), (154), and (156). This, together with the block structure of \mathbf{Y} in (22) and (23), yields

$$\begin{aligned} & \left\| \text{BlkDiag} \left(\hat{\mathbf{U}}_L^{(1)}, \dots, \hat{\mathbf{U}}_L^{(p_Y)} \right) - \text{BlkDiag} \left(\hat{\mathbf{V}}_1^{(1)}, \dots, \hat{\mathbf{V}}_1^{(p_Y)} \right) \right\|_F \\ & \leq \frac{1}{y_{s_{p_Y}}} \left\| \text{BlkDiag} \left(y_{s_1} \hat{\mathbf{U}}_{L,1} \hat{\mathbf{V}}_{1,1}^T, \dots, y_{s_{p_Y}} \hat{\mathbf{U}}_{L,p_Y} \hat{\mathbf{V}}_{1,p_Y}^T \right) - \text{BlkDiag} \left(y_{s_1} \mathbf{I}_{h_1}, \dots, y_{s_{p_Y}} \mathbf{I}_{h_{p_Y}} \right) \right\|_F \\ & = \frac{1}{y_{s_{p_Y}}} \left\| \text{BlkDiag} \left(\hat{\mathbf{U}}_L^{(1)}, \dots, \hat{\mathbf{U}}_L^{(p_Y)}, \hat{\mathbf{U}}_L^{(p_Y+1)} \right) \mathbf{Y} \text{BlkDiag} \left(\hat{\mathbf{V}}_1^{(1)}, \dots, \hat{\mathbf{V}}_1^{(p_Y)}, \hat{\mathbf{V}}_1^{(p_Y+1)} \right)^T - \mathbf{Y} \right\|_F \\ & \leq \frac{1}{y_{s_{p_Y}}} \left(\frac{2^L \sqrt{\eta_3^2 + \eta_4^2}}{\sqrt{\lambda} \sigma_{\min}^{*L-1}} + 2y_1 \eta_5 \right) \|\nabla G(\mathbf{W})\|_F. \end{aligned}$$

Using the definition of $\tilde{\mathbf{V}}_1$, we compute

$$\begin{aligned}
\|\mathbf{V}_1 - \tilde{\mathbf{V}}_1\|_F &= \left\| \bar{\mathbf{V}}_1 - \text{BlkDiag}\left(\hat{\mathbf{U}}_L^{(1)}, \dots, \hat{\mathbf{U}}_L^{(p_Y)}, \hat{\mathbf{V}}_1^{(p_Y+1)}\right) \right\|_F \\
&\leq \left\| \bar{\mathbf{V}}_1 - \text{BlkDiag}\left(\hat{\mathbf{V}}_1^{(1)}, \dots, \hat{\mathbf{V}}_1^{(p_Y)}, \hat{\mathbf{V}}_1^{(p_Y+1)}\right) \right\|_F + \\
&\quad \left\| \text{BlkDiag}\left(\hat{\mathbf{V}}_1^{(1)}, \dots, \hat{\mathbf{V}}_1^{(p_Y)}\right) - \text{BlkDiag}\left(\hat{\mathbf{U}}_L^{(1)}, \dots, \hat{\mathbf{U}}_L^{(p_Y)}\right) \right\|_F \\
&\leq \left(\eta_5 + \frac{1}{y_{s_{p_Y}}} \left(\frac{2^L \sqrt{\eta_3^2 + \eta_4^2}}{\sqrt{\lambda} \sigma_{\min}^{*L-1}} + 2y_1 \eta_5 \right) \right) \|\nabla G(\mathbf{W})\|_F,
\end{aligned}$$

where the last inequality follows from (156) and the above inequality.

The rest of the proof is devoted to showing that $(\boldsymbol{\sigma}^*, \boldsymbol{\Pi}^T) \in \mathcal{B}$. According to (151) and the definition of $\bar{\mathbf{Y}}$ in (149), we have $\|\mathbf{c} - \boldsymbol{\Pi}^T \mathbf{y}\| \leq \delta_y/3$, where $\mathbf{y} := (y_1, y_2, \dots, y_{d_{\min}})$. This implies that there exists a permutation $\pi : [d_{\min}] \rightarrow [d_{\min}]$ such that¹

$$|c_i - y_{\pi^{-1}(i)}| \leq \frac{\delta_y}{3}, \quad \forall i \in [d_{\min}]. \quad (157)$$

Noting that $\varphi(x)$ is continuous and differentiable at $x \neq 0$ and $\delta^* > 0$, there exists a positive constant $\delta_1 > 0$ such that for all $|x - \sigma_i^*| \leq \min\{\delta_1, \frac{\delta_\sigma}{3}\}$ and all $i \in [r_\sigma]$,

$$|\varphi(x) - \varphi(\sigma_i^*)| \leq \frac{\delta_y}{3}. \quad (158)$$

Using Weyl's inequality, (68a), and (62a) (resp. 64), we have

$$|\sigma_i(\mathbf{W}_L) - \sigma_i^*| \leq \text{dist}(\mathbf{W}, \mathcal{W}_{\boldsymbol{\sigma}^*}) \leq \min\{\delta_1, \frac{\delta_\sigma}{3}\},$$

This, together with (148) and (158), implies

$$|c_i - \varphi(\sigma_i^*)| = |\varphi(\sigma_i(\mathbf{W}_L)) - \varphi(\sigma_i^*)| \leq \frac{\delta_y}{3}, \quad \forall i \in [r_\sigma]. \quad (159)$$

According to (28), there exists $k_i \in [d_{\min}]$ such that $\varphi(\sigma_i^*) = y_{k_i}$ for each $i \in [r_\sigma]$. For each $i \in [r_\sigma]$, we compute

$$|y_{k_i} - y_{\pi^{-1}(i)}| = |\varphi(\sigma_i^*) - y_{\pi^{-1}(i)}| \leq |\varphi(\sigma_i^*) - c_i| + |c_i - y_{\pi^{-1}(i)}| \leq \frac{2\delta_y}{3},$$

where the last inequality uses (157) and (158). This, together with the definition of δ_y in (24), yields $\varphi(\sigma_i^*) = y_{k_i} = y_{\pi^{-1}(i)}$ for each $i \in [r_\sigma]$. Using this and (67), we have

$$(\sigma_i^*)^{2L-1} + \lambda \sigma_i^* - \sqrt{\lambda} y_{\pi^{-1}(i)} (\sigma_i^*)^{L-1} = 0, \quad \forall i \in [r_\sigma].$$

For each $i \in [r_\sigma + 1, d_{\min}]$, we note that $\sigma_i^* = 0$. It is trivial to see that the above equation holds. These, together with the definition of \mathcal{B} in (29) and Lemma 12 in the appendix, yields that $(\boldsymbol{\sigma}^*, \boldsymbol{\Pi}^T) \in \mathcal{B}$. Then, we complete the proof.

¹Let $\pi : [d_{\min}] \rightarrow [d_{\min}]$ denote a permutation of the elements in $[d_{\min}]$. Note that there is an one-to-one correspondence between permutation matrix $\boldsymbol{\Pi} \in \mathbb{R}^{d_{\min} \times d_{\min}}$ and a permutation π , i.e., $\Pi_{ij} = 1$ if $j = \pi(i)$ and $\Pi_{ij} = 0$ otherwise for each $i \in [d_{\min}]$. Now, suppose that $\boldsymbol{\Pi}$ corresponds to π and $\boldsymbol{\Pi}^T$ corresponds to π^{-1} .

(ii) Using (148) and the fact that $y_{\pi^{-1}(i)} = \varphi(\sigma_i^*)$ for each $i \in [r_\sigma]$, we have

$$\begin{aligned} |\varphi(\sigma_i(\mathbf{W}_L)) - \varphi(\sigma_i^*)| &= |c_i - y_{\pi^{-1}(i)}| \leq \|\mathbf{c} - \mathbf{\Pi}^T \mathbf{y}\| = \|\mathbf{\Pi} \mathbf{c} - \mathbf{y}\| \\ &= \|\bar{\mathbf{Y}} - \mathbf{Y}\|_F \leq \frac{2^{L-1}}{\sqrt{\lambda} \sigma_{\min}^{*L-1}} \sqrt{\eta_3^2 + \eta_4^2} \|\nabla G(\mathbf{W})\|_F, \end{aligned} \quad (160)$$

where the last inequality uses (150). By (61) and the continuity of $\varphi'(x)$, Lemma 7 holds, and thus we conclude that there exists a positive constant $\delta_2 > 0$ such that for all $|x - \sigma_i^*| \leq \delta_2$ and all $i \in [r_\sigma]$, we have for all $L \geq 3$,

$$|\varphi'(x) - \varphi'(\sigma_i^*)| \leq \frac{|\varphi'(\sigma_i^*)|}{2}. \quad (161)$$

Note that there exists $\delta_2 > 0$ such that (161) holds when $L = 2$ (see Remark 2). Using (68a) and Weyl's inequality, it holds that

$$|\sigma_i(\mathbf{W}_L) - \sigma_i^*| \leq \text{dist}(\mathbf{W}, \mathcal{W}_{\sigma^*}) \leq \delta_2. \quad (162)$$

Applying the mean-value theorem to $\varphi(\cdot)$, there exists x between $\sigma_i(\mathbf{W}_L)$ and σ_i^* such that

$$|\varphi(\sigma_i(\mathbf{W}_L)) - \varphi(\sigma_i^*)| = |\varphi'(x)| |\sigma_i(\mathbf{W}_L) - \sigma_i^*| \geq \frac{|\varphi'(\sigma_i^*)|}{2} |\sigma_i(\mathbf{W}_L) - \sigma_i^*|, \quad (163)$$

where the inequality follows from $|x - \sigma_i^*| \leq |\sigma_i(\mathbf{W}_L) - \sigma_i^*| \leq \delta_2$ due to (162) and $|\varphi'(x)| \geq |\varphi'(\sigma_i^*)|/2$ due to (161). This, together with (160), yields for all $i \in [r_\sigma]$,

$$|\sigma_i(\mathbf{W}_L) - \sigma_i^*| \leq \frac{2^L \sqrt{\eta_3^2 + \eta_4^2}}{\sqrt{\lambda} \sigma_{\min}^{*L-1} \min_{i \in [r_\sigma]} |\varphi'(\sigma_i^*)|} \|\nabla G(\mathbf{W})\|_F.$$

Using this and (45) in Lemma 3, we obtain for each $l \in [L]$,

$$|\sigma_i(\mathbf{W}_l) - \sigma_i^*| \leq \left(\frac{2^L \sqrt{\eta_3^2 + \eta_4^2}}{\sqrt{\lambda} \sigma_{\min}^{*L-1} \min_{i \in [r_\sigma]} |\varphi'(\sigma_i^*)|} + \frac{3\sqrt{2}(L-1)\sigma_{\max}^*}{4\lambda\sigma_{\min}^*} \right) \|\nabla G(\mathbf{W})\|_F.$$

Then, we complete the proof. \square

5 Experimental Results

In this section, we conduct experiments under different settings to validate our theoretical results. Specifically, we employ GD to solve Problem (1) using the `PyTorch` library. Our codes are implemented in Python on a workstation equipped with 24 GB of RAM and an AMD Ryzen-7 8845H processor with integrated Radeon 780M Graphics operating at 3.80 GHz. We terminate the algorithm when the squared gradient norm satisfies $\|\nabla f(\mathbf{W}^k)\|_F^2 \leq 10^{-6}$ and the function value change satisfies $|f(\mathbf{W}^k) - f(\mathbf{W}^{k-1})| \leq 10^{-7}$ for all $k = 1, 2, \dots$, where \mathbf{W}^k is the k -th iterate.

5.1 Linear Convergence to Critical Points

In this subsection, we investigate the convergence behavior of GD to different critical points of Problem (2). In our experiments, we set $d_L = 20$, $d_0 = 10$, and $d_l = 32$ for each $l = 2, \dots, L-1$. Then, we i.i.d. sample each entry of $\mathbf{Y} \in \mathbb{R}^{d_L \times d_0}$ from the standard Gaussian distribution, i.e., $y_{ij} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$. Moreover, we set the regularization parameters $\lambda_l = 10^{-4}$ for all $l \in [L]$ and the learning rate 4.5×10^{-4} . To ensure convergence to different critical points, we initialize the weight matrices in the neighborhood of two different critical points of Problem (2). Now, we specify how to construct these critical points as follows. We apply the SVD to \mathbf{Y} and solve the equation (8) to get its roots. Using these results and Theorem 1, we can respectively construct an optimal solution, denoted by $\mathbf{W}_{\text{opt}}^*$, and a non-global critical point, denoted by $\mathbf{W}_{\text{crit}}^*$, of Problem (2). Then, we set $\mathbf{W}^0 = \mathbf{W}^* + 0.01\mathbf{\Delta}$, where \mathbf{W}^* is $\mathbf{W}_{\text{opt}}^*$ or $\mathbf{W}_{\text{crit}}^*$ and each entry of $\mathbf{\Delta}$ is i.i.d. sampled from the standard Gaussian distribution. For each initialization, we run GD for solving Problem (2) with different depths $L \in \{2, 4, 6\}$.

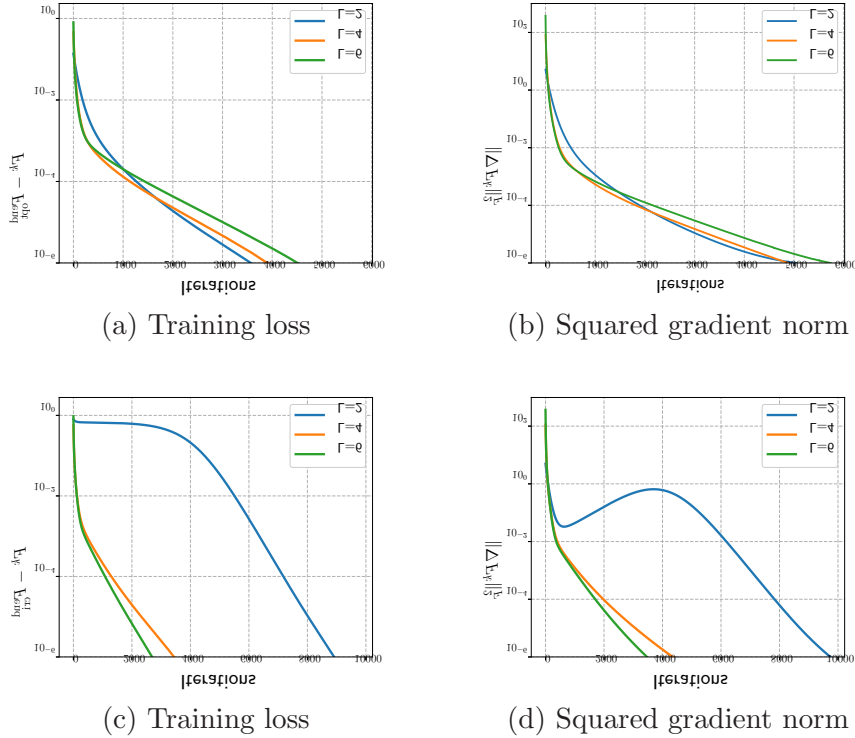


Figure 1: Linear convergence of GD to a critical point of Problem (2). In (a) and (b), $\{\mathbf{W}^k\}$ converges to an optimal solution $\mathbf{W}_{\text{opt}}^*$; in (c) and (d), $\{\mathbf{W}^k\}$ converges to a non-optimal critical point $\mathbf{W}_{\text{crit}}^*$. F^{end} denotes the training loss at the final step when the stop condition is triggered.

We respectively plot the function value gap $F(\mathbf{W}^k) - F(\mathbf{W}^{\text{end}})$ and the squared gradient norm $\|\nabla F(\mathbf{W}^k)\|_F^2$ against the iteration number in Figure 1(a) and (b) (resp., Figure 1(c) and (d)) for $\mathbf{W}_{\text{opt}}^{\text{end}}$ (resp., $\mathbf{W}_{\text{crit}}^{\text{end}}$). We also report the function values of the final iterate, denoted by \mathbf{W}^{end} , in different settings in Table 1. As observed from Figure 1, and Table 1, GD converges to an optimal solution at a linear rate for solving Problem (2) with different network depths, and similarly, it converges to a non-optimal critical point at a linear rate. This aligns with our linear convergence analysis for all critical points of Problem (2) in Proposition 7, which leverages the

error bound of Problem (2) and supports Theorem 2.

	$L = 2$	$L = 4$	$L = 6$
$F(\mathbf{W}_{\text{opt}}^*)$	9.2900×10^{-3}	8.1723×10^{-3}	9.5264×10^{-3}
$F(\mathbf{W}_{\text{opt}}^{\text{end}})$	9.2985×10^{-3}	8.2020×10^{-3}	9.5741×10^{-3}
$F(\mathbf{W}_{\text{cri}}^*)$	6.8038×10^{-1}	6.7906×10^{-1}	6.8022×10^{-1}
$F(\mathbf{W}_{\text{cri}}^{\text{end}})$	9.2980×10^{-3}	6.7909×10^{-1}	6.8027×10^{-1}

Table 1: The function values at the final iterate \mathbf{W}^{end} and the critical point \mathbf{W}^* .

From Figure 1(c) and (d) and Table 1, we observe that the convergence behavior of GD to a critical point when $L = 2$ differs significantly from those when $L = 4$ and $L = 6$. This difference arises from the benign global loss landscape of Problem (2) when $L = 2$ as shown in [50]. That is when $L = 2$, each critical point of Problem (2) is either a global minimizer or a strict saddle point. This, together with the result in [28], implies that GD almost surely avoids strict saddle points and converges to an optimal solution. Consequently, even when the starting point is initialized near a non-optimal critical point, GD will escape from it and eventually converge to a global optimum, although it may require more iterations to do so. In contrast, when $L = 4$ and $L = 6$, we conjecture that the loss landscape of Problem (2) is not benign, which contains non-optimal local minimizers, to which GD may converge. In the future, we will study the global optimization loss landscape of Problem (2) for different L .

5.2 Convergence Behavior in General Settings

In this subsection, we investigate the convergence behavior of GD in more general setups extending beyond linear networks. Specifically, the network depth is fixed as $L = 4$ and network widths at different layers are set as follows: $d_L = 16$, $d_0 = 10$, and $d_l = 32$ for each $l = 2, \dots, L - 1$. The regularization parameter is set as $\lambda_l = 5 \times 10^{-5}$ for all $l \in [L]$ and the learning rate of GD is set as 10^{-3} . We use the default initialization scheme in `PyTorch` to initialize the weights for GD. The data matrix \mathbf{X} is generated according to a uniform distribution using the `xavier_uniform` function in `PyTorch`, while the target matrix \mathbf{Y} is generated using the same approach as described in Section 5.1. Below, we outline the different setups used in our experiments.

General data input. In this experiment, we consider general data inputs $\mathbf{X}_1, \mathbf{X}_2, \mathbf{X}_3$ instead of orthogonal inputs \mathbf{X} . Then, we apply GD to optimize Problem (1) for each of these different data matrices.

Linear networks with bias. Next, we study the regularized loss of deep linear networks with bias terms as follows:

$$\min_{\{\mathbf{W}_l, \mathbf{b}_l\}} \sum_{i=1}^N \|\mathbf{W}_L(\mathbf{W}_{L-1} \dots (\mathbf{W}_1 \mathbf{x}_i + \mathbf{b}_1) + \mathbf{b}_{L-1}) + \mathbf{b}_L - \mathbf{y}_i\|^2 + \sum_{l=1}^L \lambda_l \left(\|\mathbf{W}_l\|_F^2 + \|\mathbf{b}_l\|^2 \right),$$

where \mathbf{x}_i and \mathbf{y}_i respectively denote the i -th column of \mathbf{X} and \mathbf{Y} . We apply GD to solve the above problem when the input is either an identity matrix or general input data.

Deep nonlinear networks. Finally, we study the performance of GD for solving the regularized loss of deep nonlinear networks with different activation functions:

$$\min_{\{\mathbf{W}_l, \mathbf{b}_l\}} \sum_{i=1}^N \|\mathbf{W}_L \sigma(\mathbf{W}_{L-1} \dots \sigma(\mathbf{W}_1 \mathbf{x}_i + \mathbf{b}_1) + \mathbf{b}_{L-1}) + \mathbf{b}_L - \mathbf{y}_i\|_F^2 + \sum_{l=1}^L \lambda_l \left(\|\mathbf{W}_l\|_F^2 + \|\mathbf{b}_l\|_F^2 \right),$$

where $\sigma(\cdot)$ denotes an activation function. We set \mathbf{X} as the identity matrix and use GD for solving the above problem when the activation functions are chosen as ReLU, Leaky ReLU, and tanh, respectively.

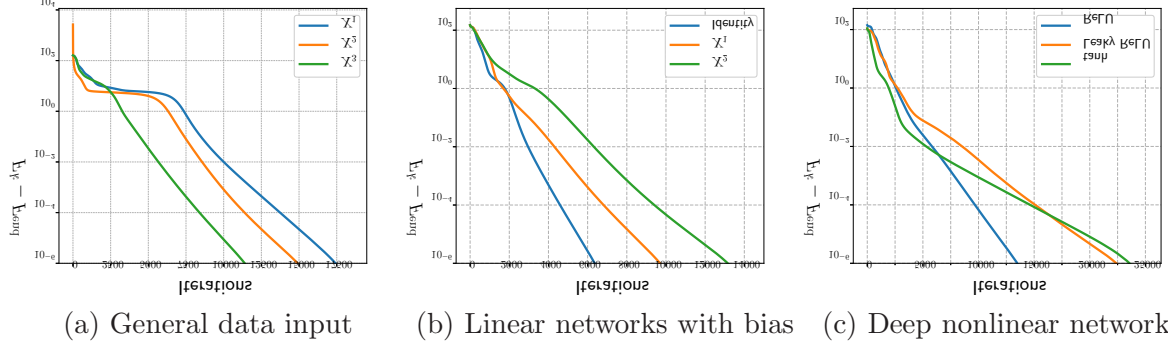


Figure 2: Linear convergence of GD under different settings: (a) Three different matrices \mathbf{X}_1 , \mathbf{X}_2 , and \mathbf{X}_3 , with condition numbers of 43.43, 36.04, and 16.36, respectively. (b) Identity matrix, \mathbf{X}_1 , and \mathbf{X}_2 , with condition numbers of 1, 4.82, and 14.62, respectively. (c) Identity matrix as data input with ReLU, Leaky ReLU, and tanh as activation functions. F^{end} denotes the training loss at the final step when the stop condition is triggered.

For the above three different settings, we plot the function value gap $F(\mathbf{W}^k) - F(\mathbf{W}^*)$ against the iteration number in Figure 2. It is observed from Figure 2 that GD converges to a solution at a linear rate across these settings. This consistent behavior leads us to conjecture that the error-bound condition may hold for deep networks in more general scenarios. Additionally, we observe that the number of iterations required to meet the stopping criterion generally increases as the condition number of the input data becomes larger. Exploring this phenomenon presents an interesting direction for future research.

6 Conclusions

In this paper, we studied the regularized squared loss of deep linear networks (i.e., Problem (2)) and proved its error bound, a critical regularity condition that characterizes local geometry around the critical point set. This result is not only theoretically significant but also lays the foundation for establishing strong convergence guarantees for various methods for solving Problem (2). To establish the error bound, we explicitly characterized the critical point set of (2) and developed new analytic techniques to show the error bound, which may be of independent interest. Our numerical results across different settings provide strong support for our theoretical findings. One future direction is to extend our analysis to Problem (1) with more general data input \mathbf{X} and loss functions. Another interesting direction is to investigate the regularized loss for deep nonlinear networks.

References

- [1] E. M. Achour, F. Malgouyres, and S. Gerchinovitz. The loss landscape of deep linear neural networks: a second-order analysis. *Journal of Machine Learning Research*, 25(242):1–76, 2024.
- [2] S. Arora, N. Cohen, W. Hu, and Y. Luo. Implicit regularization in deep matrix factorization. *Advances in Neural Information Processing Systems*, 32, 2019.
- [3] S. Arora, N. Golowich, N. Cohen, and W. Hu. A convergence analysis of gradient descent for deep linear neural networks. In *International Conference on Learning Representations*, 2019.
- [4] B. Bah, H. Rauhut, U. Terstiege, and M. Westdickenberg. Learning deep linear neural networks: Riemannian gradient flows and convergence to global minimizers. *Information and Inference: A Journal of the IMA*, 11(1):307–353, 2022.
- [5] P. Bartlett, D. Helmbold, and P. Long. Gradient descent with identity initialization efficiently learns positive definite linear transformations by deep residual networks. In *International Conference on Machine Learning*, pages 521–530. PMLR, 2018.
- [6] Y. Chitour, Z. Liao, and R. Couillet. A geometric approach of gradient descent algorithms in linear neural networks. *Mathematical Control and Related Fields*, 13(3):918–945, 2023.
- [7] H. Dang, T. M. Nguyen, T. Tran, H. T. Tran, H. Tran, and N. Ho. Neural collapse in deep linear networks: From balanced to imbalanced data. In *International Conference on Machine Learning*, 2023.
- [8] P. De Handschutter, N. Gillis, and X. Siebert. A survey on deep matrix factorizations. *Computer Science Review*, 42:100423, 2021.
- [9] D. Drusvyatskiy and A. S. Lewis. Error bounds, quadratic growth, and linear convergence of proximal methods. *Mathematics of Operations Research*, 43(3):919–948, 2018.
- [10] S. Du and W. Hu. Width provably matters in optimization for deep linear neural networks. In *International Conference on Machine Learning*, pages 1655–1664. PMLR, 2019.
- [11] S. S. Du, X. Zhai, B. Póczos, and A. Singh. Gradient descent provably optimizes over-parameterized neural networks. In *International Conference on Learning Representations*, 2019.
- [12] A. Esteva, A. Robicquet, B. Ramsundar, V. Kuleshov, M. DePristo, K. Chou, C. Cui, G. Corrado, S. Thrun, and J. Dean. A guide to deep learning in healthcare. *Nature Medicine*, 25(1):24–29, 2019.
- [13] S. Frei and Q. Gu. Proxy convexity: A unified framework for the analysis of neural networks trained by gradient descent. *Advances in Neural Information Processing Systems*, 34:7937–7949, 2021.

- [14] X. Han, V. Pappayan, and D. L. Donoho. Neural collapse under MSE loss: Proximity to and dynamics on the central path. In *International Conference on Learning Representations*, 2021.
- [15] B. Hanin. Universal function approximation by deep neural nets with bounded width and relu activations. *Mathematics*, 7(10):992, 2019.
- [16] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016.
- [17] R. Horn and C. Johnson. *Matrix Analysis*. Cambridge University Press, 2012. ISBN 9781139788885.
- [18] E. J. Hu, yelong shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022.
- [19] W. Hu, L. Xiao, and J. Pennington. Provable benefit of orthogonal initialization in optimizing deep linear networks. In *International Conference on Learning Representations*, 2020.
- [20] M. Huh, H. Mobahi, R. Zhang, B. Cheung, P. Agrawal, and P. Isola. The low-rank simplicity bias in deep networks. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL <https://openreview.net/forum?id=bCiNWDm1Y2>.
- [21] R. Jiang and D. Li. Novel reformulations and efficient algorithms for the generalized trust region subproblem. *SIAM Journal on Optimization*, 29(2):1603–1633, 2019.
- [22] R. Jiang and X. Li. Hölderian error bounds and Kurdyka-Łojasiewicz inequality for the trust region subproblem. *Mathematics of Operations Research*, 47(4):3025–3050, 2022.
- [23] H. Karimi, J. Nutini, and M. Schmidt. Linear convergence of gradient and proximal-gradient methods under the Polyak-Łojasiewicz condition. In *Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2016, Riva del Garda, Italy, September 19-23, 2016, Proceedings, Part I 16*, pages 795–811. Springer, 2016.
- [24] K. Kawaguchi. Deep learning without poor local minima. *Advances in Neural Information Processing Systems*, 29, 2016.
- [25] D. P. Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [26] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. *Advances in Neural Information Processing Systems*, 25, 2012.
- [27] T. Laurent and J. Brecht. Deep linear networks with arbitrary loss: All local minima are global. In *International Conference on Machine Learning*, pages 2902–2907. PMLR, 2018.

- [28] J. D. Lee, I. Panageas, G. Piliouras, M. Simchowitz, M. I. Jordan, and B. Recht. First-order methods almost always avoid strict saddle points. *Mathematical Programming*, 176: 311–337, 2019.
- [29] F.-Y. Liao, L. Ding, and Y. Zheng. Error bounds, PL condition, and quadratic growth for weakly convex functions, and linear convergences of proximal point methods. In *6th Annual Learning for Dynamics & Control Conference*, pages 993–1005. PMLR, 2024.
- [30] H. Liu, A. M.-C. So, and W. Wu. Quadratic optimization with orthogonality constraint: explicit Lojasiewicz exponent and linear convergence of retraction-based line-search and stochastic variance-reduced gradient methods. *Mathematical Programming*, 178(1):215–262, 2019.
- [31] D. Mehta, T. Chen, T. Tang, and J. D. Hauenstein. The loss surface of deep linear networks viewed through the algebraic geometry lens. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(9):5664–5680, 2021.
- [32] G. M. Nguegnang, H. Rauhut, and U. Terstiege. Convergence of gradient descent for learning linear neural networks. *Advances in Continuous and Discrete Models*, 2024(1):23, 2024.
- [33] B. T. Polyak. Gradient methods for the minimisation of functionals. *USSR Computational Mathematics and Mathematical Physics*, 3(4):864–878, 1963.
- [34] Q. Rebjock and N. Boumal. Fast convergence to non-isolated minima: four equivalent conditions for C^2 functions. *Mathematical Programming*, pages 1–49, 2024.
- [35] A. M. Saxe, J. L. McClelland, and S. Ganguli. A mathematical theory of semantic development in deep neural networks. *Proceedings of the National Academy of Sciences*, 116(23):11537–11546, 2019.
- [36] R. Schneider and A. Uschmajew. Convergence results for projected line-search methods on varieties of low-rank matrices via Lojasiewicz inequality. *SIAM Journal on Optimization*, 25(1):622–646, 2015.
- [37] G. W. Stewart and J.-g. Sun. Matrix perturbation theory. 1990.
- [38] G. Trigeorgis, K. Bousmalis, S. Zafeiriou, and B. W. Schuller. A deep matrix factorization method for learning attribute representations. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(3):417–429, 2016.
- [39] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, 30, 2017.
- [40] P. Wang, H. Liu, C. Yaras, L. Balzano, and Q. Qu. Linear convergence analysis of neural collapse with unconstrained features. In *OPT 2022: Optimization for Machine Learning (NeurIPS 2022 Workshop)*, 2022.

- [41] P. Wang, X. Li, C. Yaras, Z. Zhu, L. Balzano, W. Hu, and Q. Qu. Understanding deep representation learning via layerwise feature compression and discrimination. *arXiv preprint arXiv:2311.02960*, 2023.
- [42] P. Wang, H. Liu, and A. M.-C. So. Linear convergence of a proximal alternating minimization method with extrapolation for L1-norm principal component analysis. *SIAM Journal on Optimization*, 33(2):684–712, 2023.
- [43] L. Wu, Q. Wang, and C. Ma. Global convergence of gradient descent for deep linear residual networks. *Advances in Neural Information Processing Systems*, 32, 2019.
- [44] C. Yaras, P. Wang, Z. Zhu, L. Balzano, and Q. Qu. Neural collapse with normalized features: A geometric analysis over the riemannian manifold. In *Advances in Neural Information Processing Systems*, 2022.
- [45] C. Yaras, P. Wang, L. Balzano, and Q. Qu. Compressible dynamics in deep overparameterized low-rank learning & adaptation. *Forty-first International Conference on Machine Learning*, 2024.
- [46] Y. Yu, T. Wang, and R. J. Samworth. A useful variant of the Davis–Kahan theorem for statisticians. *Biometrika*, 102(2):315–323, 2015.
- [47] M.-C. Yue, Z. Zhou, and A. M.-C. So. A family of inexact SQA methods for non-smooth convex minimization with provable convergence guarantees based on the luo–tseng error bound property. *Mathematical Programming*, 174(1):327–358, 2019.
- [48] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115, 2021.
- [49] H. Zhao and J. Xu. Convergence analysis and trajectory comparison of gradient descent for overparameterized deep linear networks. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856.
- [50] J. Zhou, X. Li, T. Ding, C. You, Q. Qu, and Z. Zhu. On the optimization landscape of neural collapse under MSE loss: Global optimality with unconstrained features. In *International Conference on Machine Learning*, pages 27179–27202. PMLR, 2022.
- [51] Z. Zhou and A. M.-C. So. A unified approach to error bounds for structured convex optimization problems. *Mathematical Programming*, 165:689–728, 2017.

Supplementary Material

The appendix is organized as follows. In Appendix A, we provide supplementary results and proofs for Section 3. In Appendix B, we show that Assumption 2 is a necessary condition for the error bound to hold for Problem (2). In Appendix C, we provide a unified framework to establish linear convergence under the error bound. Finally, we present auxiliary results to prove our main results for completeness.

A Supplementary Results and Proofs for Section 3

Lemma 5. Let $a > 0$ be a constant, $\Sigma \in \mathbb{R}^{n \times n}$ be a diagonal matrix, and $Q \in \mathcal{O}^n$ be an orthogonal matrix. Then, if $\|\Sigma - aI\| \leq a/2$, it holds that

$$\|Q\Sigma^2 - \Sigma^2Q\|_F \geq a\|Q\Sigma - \Sigma Q\|_F.$$

Proof. For ease of exposition, let $\Delta := \Sigma - aI$. We compute

$$\begin{aligned} \|Q\Sigma^2 - \Sigma^2Q\|_F &= \|Q(\Delta + aI)^2 - (\Delta + aI)^2Q\|_F = \|Q\Delta^2 - \Delta^2Q + 2a(Q\Delta - \Delta Q)\|_F \\ &\geq 2a\|Q\Delta - \Delta Q\|_F - \|Q\Delta^2 - \Delta^2Q\|_F \\ &\geq 2(a - \|\Delta\|)\|Q\Sigma - \Sigma Q\|_F \geq a\|Q\Sigma - \Sigma Q\|_F, \end{aligned}$$

where the second inequality follows from $\|Q\Delta^2 - \Delta^2Q\|_F = \|\Delta(\Delta Q - Q\Delta) + (\Delta Q - Q\Delta)\Delta\|_F \leq 2\|\Delta\|\|Q\Sigma - \Sigma Q\|_F$ and $Q\Delta - \Delta Q = Q\Sigma - \Sigma Q$, and the last inequality uses $\|\Delta\| \leq a/2$. \square

Proof of Proposition 3. (i) According to $L = 2$, (31), and (35), we have $\mathcal{W}_0 = \{(\mathbf{0}, \mathbf{0})\}$. Then, we have

$$\text{dist}^2(\mathbf{W}, \mathcal{W}_0) = \|\mathbf{W}_1\|_F^2 + \|\mathbf{W}_2\|_F^2. \quad (164)$$

Without loss of generality, assume that $\|\mathbf{W}_1\|_F \geq \|\mathbf{W}_2\|_F$. Now, we compute

$$\frac{1}{2}\nabla_{\mathbf{W}_1}G(\mathbf{W}) = \mathbf{W}_2^T (\mathbf{W}_2\mathbf{W}_1 - \sqrt{\lambda}\mathbf{Y}) + \lambda\mathbf{W}_1, \quad (165)$$

$$\frac{1}{2}\nabla_{\mathbf{W}_2}G(\mathbf{W}) = (\mathbf{W}_2\mathbf{W}_1 - \sqrt{\lambda}\mathbf{Y})\mathbf{W}_1^T + \lambda\mathbf{W}_2. \quad (166)$$

Using (165), we have

$$\begin{aligned} \frac{\sqrt{\lambda}}{2}\|\nabla_{\mathbf{W}_1}G(\mathbf{W})\|_F &\geq \left\|\lambda\mathbf{W}_2^T\mathbf{Y} - \lambda^{\frac{3}{2}}\mathbf{W}_1\right\|_F - \sqrt{\lambda}\|\mathbf{W}_2^T\mathbf{W}_2\mathbf{W}_1\|_F \\ &\geq \left\|\lambda\mathbf{W}_2^T\mathbf{Y} - \lambda^{\frac{3}{2}}\mathbf{W}_1\right\|_F - \sqrt{\lambda}\|\mathbf{W}_2\|_F^2\|\mathbf{W}_1\|_F. \end{aligned} \quad (167)$$

Moreover, multiplying \mathbf{Y}^T on the both sides of (166), together with the triangular inequality, yields

$$\frac{1}{2}\|\mathbf{Y}^T\nabla_{\mathbf{W}_2}G(\mathbf{W})\|_F \geq \left\|\sqrt{\lambda}\mathbf{Y}^T\mathbf{Y}\mathbf{W}_1^T - \lambda\mathbf{Y}^T\mathbf{W}_2\right\|_F - \|\mathbf{Y}^T\mathbf{W}_2\mathbf{W}_1\mathbf{W}_1^T\|_F,$$

which implies

$$\frac{y_1}{2}\|\nabla_{\mathbf{W}_2}G(\mathbf{W})\|_F \geq \left\|\sqrt{\lambda}\mathbf{Y}^T\mathbf{Y}\mathbf{W}_1^T - \lambda\mathbf{Y}^T\mathbf{W}_2\right\|_F - y_1\|\mathbf{W}_2\|_F\|\mathbf{W}_1\|_F^2.$$

Summing up this inequality with (167) and using the triangular inequality, we obtain

$$\begin{aligned} &\frac{\sqrt{\lambda}}{2}\|\nabla_{\mathbf{W}_1}G(\mathbf{W})\|_F + \frac{y_1}{2}\|\nabla_{\mathbf{W}_2}G(\mathbf{W})\|_F \\ &\geq \sqrt{\lambda}\|\mathbf{W}_1(\lambda I - \mathbf{Y}^T\mathbf{Y})\|_F - \|\mathbf{W}_1\|_F\|\mathbf{W}_2\|_F(y_1\|\mathbf{W}_1\|_F + \sqrt{\lambda}\|\mathbf{W}_2\|_F) \\ &\geq \left(\sqrt{\lambda}\min\left\{\min_{i \in [s_{p_Y}]}|\lambda - y_i^2|, \lambda\right\}\right)\|\mathbf{W}_1\|_F - (\sqrt{\lambda} + y_1)\|\mathbf{W}_1\|_F^3 \\ &\geq \frac{\sqrt{\lambda}}{2}\min\left\{\min_{i \in [s_{p_Y}]}|\lambda - y_i^2|, \lambda\right\}\|\mathbf{W}_1\|_F, \end{aligned}$$

where the second inequality follows from (22), $\lambda \neq y_i^2$ for all $i \in [p_Y]$ due to (60), and $\|\mathbf{A}\mathbf{B}\|_F \geq \sigma_{\min}(\mathbf{A})\|\mathbf{B}\|_F$ when \mathbf{A} is non-degenerate, and the last inequality uses (36). This, together with (164), directly implies (37).

(ii) According to (31) and (35), we have $\mathcal{W}_0 = \{(\mathbf{0}, \mathbf{0}, \dots, \mathbf{0})\}$. Then, we have

$$\text{dist}^2(\mathbf{W}, \mathcal{W}_0) = \sum_{l=1}^L \|\mathbf{W}_l\|_F^2. \quad (168)$$

Let $k \in [L]$ be such that $k \in \arg \max_{l \in [L]} \|\mathbf{W}_l\|_F$. Then, we compute

$$\begin{aligned} \frac{1}{2} \|\nabla_{\mathbf{W}_k} G(\mathbf{W})\|_F &\stackrel{(14)}{=} \left\| \mathbf{W}_{L:k+1}^T \left(\mathbf{W}_{L:1} - \sqrt{\lambda} \mathbf{Y} \right) \mathbf{W}_{k-1:1}^T + \lambda \mathbf{W}_k \right\|_F \\ &\geq \lambda \|\mathbf{W}_k\|_F - \left(\|\mathbf{W}_k\|_F^{2L-1} + \sqrt{\lambda} y_1 \|\mathbf{W}_k\|_F^{L-1} \right) \\ &= \left(\lambda - \|\mathbf{W}_k\|_F^{2L-2} - \sqrt{\lambda} y_1 \|\mathbf{W}_k\|_F^{L-2} \right) \|\mathbf{W}_k\|_F \geq \frac{\lambda}{3} \|\mathbf{W}_k\|_F, \end{aligned}$$

where the last inequality uses (38). This, together with (168), yields

$$\text{dist}^2(\mathbf{W}, \mathcal{W}_0) \leq L \|\mathbf{W}_k\|_F^2 \leq \frac{9L}{4\lambda^2} \|\nabla_{\mathbf{W}_k} G(\mathbf{W})\|_F^2,$$

which directly implies (39). \square

B Showing the Necessity of Assumption 2

Before we proceed, we claim that the error bound of F (see Problem (2)) holds if and only if the error bound of G (see Problem (13)) holds. Indeed, it follows from (ii) in Lemma 1 that the “if” direction holds. Now, it remains to show the “only if” direction. Suppose that the error bound of Problem (2) holds (see (11)) for all critical points. Using the proof setup in Lemma 1, it follows (18) that

$$\frac{\lambda}{\lambda_{\max}} \text{dist}(\mathbf{Z}, \mathcal{W}_G) \leq \frac{\lambda}{\sqrt{\lambda_{\max}}} \text{dist}(\mathbf{W}, \mathcal{W}) \leq \frac{\lambda \kappa_1}{\sqrt{\lambda_{\max}}} \|\nabla F(\mathbf{W})\|_F \leq \kappa_1 \|\nabla G(\mathbf{Z})\|_F,$$

where the second inequality follows from (11) and the last inequality uses (21). Then, we prove the claim. Based on the error-bound equivalence between F and G , establishing the necessity of Assumption 2 for F is equivalent to doing so for G .

B.1 The Case $L = 2$

For ease of exposition, we denote $\widehat{\mathcal{W}}_{\sigma^*} := \mathcal{W}_{\text{sort}(\sigma^*)}$ for each $\sigma^* \in \mathcal{A}$, where $\text{sort}(\cdot)$ is a sorting function that arranges the elements of a vector in decreasing order.

Lemma 6. *Suppose that $L = 2$ and (60) does not hold. For any $\sigma^* \in \mathcal{A}$, the error bound for the critical point set $\widehat{\mathcal{W}}_{\sigma^*}$ of Problem (13) fails to hold, i.e., for any $\kappa, \delta > 0$, there exists \mathbf{W} satisfying $\text{dist}(\mathbf{W}, \widehat{\mathcal{W}}_{\sigma^*}) \leq \delta$ such that $\|\nabla G(\mathbf{W})\|_F \leq \kappa \cdot \text{dist}(\mathbf{W}, \mathcal{W}_G)$.*

Proof. Recall that $\lambda = \lambda_1 \lambda_2$ when $L = 2$. Since (60) does not hold and , there exists $i \in [r_Y]$ such that $y_i = \sqrt{\lambda}$. Now, we define $\mathbf{W}(t) := (\mathbf{W}_1(t), \mathbf{W}_2(t))$ as follows:

$$\begin{cases} \mathbf{W}_1(t) = \mathbf{Q}_2 \boldsymbol{\Sigma}_1(t) \text{BlkDiag}(\mathbf{O}_1, \dots, \mathbf{O}_{p_Y}, \mathbf{O}_{p_Y+1}), \\ \mathbf{W}_2(t) = \text{BlkDiag}(\mathbf{O}_1^T, \dots, \mathbf{O}_{p_Y}^T, \widehat{\mathbf{O}}_{p_Y+1}^T) \boldsymbol{\Sigma}_2(t) \mathbf{Q}_2^T, \\ \boldsymbol{\Sigma}_l(t) = \text{BlkDiag}(\text{diag}(\boldsymbol{\sigma}^* + t\mathbf{e}_i), \mathbf{0}) \in \mathbb{R}^{d_l \times d_l-1}, \forall l = 1, 2, \\ \mathbf{O}_i \in \mathcal{O}^{h_i}, \forall i \in [p_Y], \mathbf{O}_{p_Y+1} \in \mathcal{O}^{d_0-r_Y}, \widehat{\mathbf{O}}_{p_Y+1} \in \mathcal{O}^{d_2-r_Y}, \end{cases}$$

where $\mathbf{e}_i \in \mathbb{R}^{d_{\min}}$ is a standard basis vector with the i -th entry being 1 and all other entries being 0. According to Theorem 3, one can verify that $\mathbf{W}(0) \in \widehat{\mathcal{W}}_{\boldsymbol{\sigma}^*}$. Using $y_i = \sqrt{\lambda}$, we obtain $\sigma_i^* \in \{x : x^3 - \sqrt{\lambda}y_i x + \lambda x = 0, x \geq 0\} = \{0\}$. Then, we compute

$$\frac{1}{2} \|\nabla_{\mathbf{W}_l} G(\mathbf{W}(t))\|_F \stackrel{(14)}{=} |t^3 - \sqrt{\lambda}y_i t + \lambda t| = |t^3|, \quad l = 1, 2, \quad (169)$$

where the second equality follows from $y_i = \sqrt{\lambda}$. Moreover, we have

$$\begin{aligned} \text{dist}^2(\mathbf{W}(t), \widehat{\mathcal{W}}_{\boldsymbol{\sigma}^*}) &= \min_{(\mathbf{W}_1, \mathbf{W}_2) \in \widehat{\mathcal{W}}_{\boldsymbol{\sigma}^*}} \|\mathbf{W}_1(t) - \mathbf{W}_1\|_F^2 + \|\mathbf{W}_2(t) - \mathbf{W}_2\|_F^2 \\ &\geq \|\boldsymbol{\Sigma}_1(t) - \boldsymbol{\Sigma}_1(0)\|_F^2 + \|\boldsymbol{\Sigma}_2(t) - \boldsymbol{\Sigma}_2(0)\|_F^2 = 2t^2, \end{aligned} \quad (170)$$

where the first inequality follows from Mirsky's inequality (see Lemma 11). When $t \leq \delta_\sigma/3$, we have

$$\text{dist}(\mathbf{W}(t), \widehat{\mathcal{W}}_{\boldsymbol{\sigma}^*}) \leq \|\mathbf{W}(t) - \mathbf{W}(0)\|_F \leq \frac{\sqrt{2}\delta_\sigma}{3}.$$

Using this inequality and Proposition 2(ii), we have for each $\boldsymbol{\sigma}^* \in \mathcal{A}$ satisfying $\widehat{\mathcal{W}}_{\boldsymbol{\sigma}^*} \neq \widehat{\mathcal{W}}_{\bar{\boldsymbol{\sigma}}^*}$,

$$\text{dist}(\mathbf{W}(t), \widehat{\mathcal{W}}_{\bar{\boldsymbol{\sigma}}^*}) \geq \text{dist}(\widehat{\mathcal{W}}_{\boldsymbol{\sigma}^*}, \widehat{\mathcal{W}}_{\bar{\boldsymbol{\sigma}}^*}) - \text{dist}(\mathbf{W}(t), \widehat{\mathcal{W}}_{\boldsymbol{\sigma}^*}) \geq \left(1 - \frac{\sqrt{2}}{3}\right) \delta_\sigma > \text{dist}(\mathbf{W}(t), \widehat{\mathcal{W}}_{\boldsymbol{\sigma}^*}).$$

Consequently, we obtain $\text{dist}(\mathbf{W}(t), \mathcal{W}_G) = \text{dist}(\mathbf{W}(t), \widehat{\mathcal{W}}_{\boldsymbol{\sigma}^*})$. This, together with (169) and (170), implies

$$\|\nabla G(\mathbf{W}(t))\|_F = 2\sqrt{2}|t|^3 \leq 2^{3/4} \text{dist}^{\frac{3}{2}}(\mathbf{W}(t), \widehat{\mathcal{W}}_{\boldsymbol{\sigma}^*}) = 2^{3/4} \text{dist}^{\frac{3}{2}}(\mathbf{W}(t), \mathcal{W}_G).$$

Obviously, the error bound does not hold when $t \rightarrow 0$. □

B.2 The Case $L \geq 3$

To begin, we show the following lemma to derive the equivalent condition of (61):

Lemma 7. *When $L \geq 3$, (61) holds if and only if for all $0 \neq \sigma^* \in \mathcal{Y}$, we have $\varphi'(\sigma^*) \neq 0$, where \mathcal{Y} and $\varphi(\cdot)$ are defined in (32) and (67), respectively.*

Proof. Note that $\varphi'(\sigma^*) = 0$ is equivalent to

$$\varphi'(\sigma^*) = \frac{L}{\sqrt{\lambda}} \sigma^{*(L-1)} + \sqrt{\lambda}(2-L) \sigma^{*(1-L)} = 0. \quad (171)$$

According to $\sigma^* \in \mathcal{Y}$, there exists $j \in [d_{\min}]$ such that

$$\varphi(\sigma^*) = \frac{\sigma^{*(2L-1)} + \lambda\sigma^*}{\sqrt{\lambda}\sigma^{*(L-1)}} = y_j. \quad (172)$$

Combining (171) and (172), we obtain

$$y_j = \left[\left(\frac{L-2}{L} \right)^{\frac{L}{2(L-1)}} + \left(\frac{L}{L-2} \right)^{\frac{L-2}{2(L-1)}} \right] \lambda^{\frac{1}{2(L-1)}}.$$

This implies that (61) does not hold.

Conversely, suppose that (61) does not hold. This implies that there exists $j \in [d_{\min}]$ such that $y_j = \left[\left(\frac{L-2}{L} \right)^{\frac{L}{2(L-1)}} + \left(\frac{L}{L-2} \right)^{\frac{L-2}{2(L-1)}} \right] \lambda^{\frac{1}{2(L-1)}}$. Using this, one can verify that

$$x^{2L-1} - \sqrt{\lambda}y_jx^{L-1} + \lambda x = 0 \quad (173)$$

have a positive root such that

$$x^* = \left(\frac{\lambda(L-1)}{L} \right)^{\frac{1}{2(L-1)}} \text{ and } \varphi'(x^*) = 0.$$

By the definition of \mathcal{Y} , we have $x^* \in \mathcal{Y}$. Then, we complete the proof. \square

Remark 2. When $L = 2$, we compute

$$\varphi'(\sigma^*) = \frac{2}{\sqrt{\lambda}}\sigma^* > 0, \quad \forall 0 \neq \sigma^* \in \mathcal{Y}.$$

Therefore, $\varphi'(\sigma^*) \neq 0$ holds trivially.

Lemma 8. Suppose that $L \geq 3$ and (61) does not hold. There exists $\sigma^* \in \mathcal{A}$ such that the error bound for the critical point set $\widehat{\mathcal{W}}_{\sigma^*}$ of Problem (13) fails to hold, i.e., for any $\kappa, \delta > 0$, there exist \mathbf{W} satisfying $\text{dist}(\mathbf{W}, \widehat{\mathcal{W}}_{\sigma^*}) \leq \delta$ such that $\|\nabla G(\mathbf{W})\|_F \leq \kappa \cdot \text{dist}(\mathbf{W}, \mathcal{W}_G)$.

Proof. According to Lemma 7 and the fact that (61) does not hold, there exists $\sigma^* \in \mathcal{A}$ such that $\varphi'(\sigma_i^*) = 0$ for some $i \in [d_{\min}]$. This implies

$$f(\sigma_i^*) := (\sigma_i^*)^{2L-1} - \sqrt{\lambda}y_i(\sigma_i^*)^{L-1} + \lambda\sigma_i^* = 0, \quad \varphi(\sigma_i^*) = y_i.$$

Note that $f(x) = \sqrt{\lambda}x^{L-1}\varphi(x) - \sqrt{\lambda}y_ix^{L-1}$ and we compute

$$f'(\sigma_i^*) = \sqrt{\lambda}(\sigma_i^*)^{L-1}\varphi'(\sigma_i^*) = 0.$$

Now, we define $\mathbf{W}(t) = (\mathbf{W}_1(t), \dots, \mathbf{W}_L(t))$ as follow:

$$\left\{ \begin{array}{l} \mathbf{W}_1(t) = \mathbf{Q}_2 \Sigma_1(t) \text{BlkDiag}(\mathbf{O}_1, \dots, \mathbf{O}_{p_Y}, \mathbf{O}_{p_Y+1}), \\ \mathbf{W}_l(t) = \mathbf{Q}_{l+1} \Sigma_l(t) \mathbf{Q}_l^T, \quad l = 2, \dots, L-1, \mathbf{Q}_l \in \mathcal{O}^{d_{l-1}}, \quad l = 2, \dots, L, \\ \mathbf{W}_L(t) = \text{BlkDiag}(\mathbf{O}_1^T, \dots, \mathbf{O}_{p_Y}^T, \widehat{\mathbf{O}}_{p_Y+1}^T) \Sigma_L(t) \mathbf{Q}_L^T, \\ \Sigma_l(t) = \text{BlkDiag}(\text{diag}(\sigma^* + te_i), \mathbf{0}) \in \mathbb{R}^{d_l \times d_{l-1}}, \quad \forall l \in [L], \\ \mathbf{O}_i \in \mathcal{O}^{h_i}, \forall i \in [p_Y], \quad \mathbf{O}_{p_Y+1} \in \mathcal{O}^{d_0 - r_Y}, \widehat{\mathbf{O}}_{p_Y+1} \in \mathcal{O}^{d_L - r_Y}. \end{array} \right.$$

It follows from Theorem 3 that $\mathbf{W}(0) \in \mathcal{W}_G$. Therefore, we obtain for all $l \in [L]$,

$$\frac{1}{2} \|\nabla_{\mathbf{W}_l} G(\mathbf{W}(t))\|_F \stackrel{(14)}{=} \left| (\sigma_i^* + t)^{2L-1} - \sqrt{\lambda} y_i (\sigma_i^* + t)^{L-1} + \lambda (\sigma_i^* + t) \right| = |f(\sigma_i^* + t)|.$$

Applying the Taylor expansion to $f(\sigma_i^* + t)$ at σ_i^* , together with $f(\sigma_i^*) = 0$ and $f'(\sigma_i^*) = 0$, yields that when $t \rightarrow 0$, $\|\nabla G(\mathbf{W}(t))\|_F = O(t^2)$. We also note that

$$\text{dist}^2(\mathbf{W}(t), \widehat{\mathcal{W}}_{\sigma^*}) = \|\mathbf{W}_1(t) - \mathbf{W}_1^*\|_F^2 + \dots + \|\mathbf{W}_L(t) - \mathbf{W}_L^*\|_F^2 \geq Lt^2,$$

where the inequality follows from Weyl's inequality. Using the same argument in Lemma 6, we conclude that $\text{dist}(\mathbf{W}(t), \widehat{\mathcal{W}}_{\sigma^*}) = \text{dist}(\mathbf{W}(t), \mathcal{W}_G)$ when t is sufficient small. Then we have $\|\nabla G(\mathbf{W})\|_F = O(\text{dist}^2(\mathbf{W}(t), \mathcal{W}_G))$, which implies that the error bound fails to hold. \square

Remark 3. Under Assumption 1, it follows from Theorem 2 that Assumption 2 is a sufficient condition for the error bound to hold for the critical point set \mathcal{W}_G . According to Lemma 6 and Lemma 8, we conclude that Assumption 2 is also a necessary condition. Therefore, we establish the sufficient and necessary condition for the error bound of the critical point set \mathcal{W}_G .

C Linear Convergence under the Error Bound

Proposition 7 (Linear Convergence Analysis [23, 36, 51]). Suppose that Theorem 2 holds and the sequence $\{\mathbf{W}^k\}_{k \geq k_1}$ for some index $k_1 \geq 0$ satisfies the following conditions:

(i) (Sufficient Decrease) There exists a constant $\kappa_1 > 0$ such that

$$F(\mathbf{W}^{k+1}) - F(\mathbf{W}^k) \leq -\kappa_1 \|\mathbf{W}^{k+1} - \mathbf{W}^k\|_F^2.$$

(ii) (Cost-to-Go Estimate) There exists a constant $\kappa_2 > 0$ such that

$$F(\mathbf{W}^{k+1}) - F(\mathbf{W}^*) \leq \kappa_2 \left(\text{dist}^2(\mathbf{W}^k, \mathcal{W}) + \|\mathbf{W}^{k+1} - \mathbf{W}^k\|_F^2 \right).$$

(iii) (Safeguard) There exists a constant $\kappa_3 > 0$ such that

$$\|\nabla F(\mathbf{W}^k)\|_F \leq \kappa_3 \|\mathbf{W}^{k+1} - \mathbf{W}^k\|_F.$$

Then, the sequence $\{F(\mathbf{W}^k)\}_{k \geq k_1}$ converges Q -linearly to $F(\mathbf{W}^*)$ and $\{\mathbf{W}^k\}_{k \geq k_1}$ converges R -linearly to some $\mathbf{W}^* \in \mathcal{W}$.

We should point out that first-order methods, such as GD and proximal alternating minimization method, for solving Problem (2) satisfy the above conditions provided that an appropriate step size is chosen.

C.1 Proof of the PL Inequality and Quadratic Growth

Lemma 9 (Local Lipschitz Property). Let \mathbf{W} and $\bar{\mathbf{W}}$ be arbitrary such that

$$\text{dist}(\mathbf{W}, \mathcal{W}) \leq \delta_3 \text{ and } \text{dist}(\bar{\mathbf{W}}, \mathcal{W}) \leq \delta_3, \tag{174}$$

where $\delta_3 > 0$ is a constant and \mathcal{W} denotes the critical point set. Then, it holds that

$$\|\nabla F(\mathbf{W}) - \nabla F(\bar{\mathbf{W}})\|_F \leq L_F \|\mathbf{W} - \bar{\mathbf{W}}\|_F,$$

where $L_F > 0$ is a constant.

Proof. According to Theorem 1, we note that \mathcal{W} is a compact set. It implies that there exists a positive constant M such that $\|\mathbf{W}_l^*\| \leq M$ for each $l \in [L]$. This, together with (174) and the triangular inequality yields that there exists a constant M such that

$$\|\mathbf{W}_l\| \leq M + \delta_3, \quad \|\bar{\mathbf{W}}_l\| \leq M + \delta_3, \quad \forall l \in [L].$$

Then, we have

$$\frac{1}{2} \nabla_{\mathbf{W}_l} F(\mathbf{W}) - \frac{1}{2} \nabla_{\bar{\mathbf{W}}_l} F(\bar{\mathbf{W}}) = \lambda_l (\mathbf{W}_l - \bar{\mathbf{W}}_l) + \mathbf{R}_1 + \mathbf{R}_2,$$

where

$$\begin{aligned} \mathbf{R}_1 &:= (\mathbf{W}_{l+1} - \bar{\mathbf{W}}_{l+1})^T \mathbf{W}_{L:l+2}^T \mathbf{W}_{L:l} \mathbf{W}_{l-1:1}^T + \cdots + \bar{\mathbf{W}}_{L:l+1}^T \bar{\mathbf{W}}_{L:l} \bar{\mathbf{W}}_{l-2:1}^T (\mathbf{W}_{l-1} - \bar{\mathbf{W}}_{l-1})^T, \\ \mathbf{R}_2 &:= (\mathbf{W}_{l+1} - \bar{\mathbf{W}}_{l+1})^T \mathbf{W}_{L:l+2}^T \mathbf{Y} \mathbf{W}_{l-1:1}^T + \cdots + \bar{\mathbf{W}}_{L:l+1}^T \mathbf{Y} \bar{\mathbf{W}}_{l-2:1}^T (\mathbf{W}_{l-1} - \bar{\mathbf{W}}_{l-1})^T. \end{aligned}$$

Then we compute

$$\|\mathbf{R}_1\|_F + \|\mathbf{R}_2\|_F \leq (\|\mathbf{Y}\|(L-1)(M + \delta_3)^{L-1} + (2L-1)(M + \delta_3)^{2L-1}) \|\mathbf{W}_l - \bar{\mathbf{W}}_l\|_F.$$

It implies that

$$\|\nabla F(\mathbf{W}) - \nabla F(\bar{\mathbf{W}})\|_F \leq 2L(\lambda_{\max} + L(M + \delta_3)^{L-1} \|\mathbf{Y}\| + 2L(M + \delta_3)^{2L-1}) \|\mathbf{W} - \bar{\mathbf{W}}\|_F.$$

□

Proof of Corollary 1. (i) Let \mathcal{U} be a neighborhood in which both Theorem 2 and Lemma 9 hold and $F(\mathbf{W}) \leq F(\mathbf{W}^*) + \nu$, where $\nu > 0$ is a constant. For any $\mathbf{W} \in \mathcal{U}$, let $\mathbf{W}^* \in \mathcal{W}$ be such that

$$\text{dist}(\mathbf{W}, \mathcal{W}) = \|\mathbf{W} - \mathbf{W}^*\|_F. \quad (175)$$

According to Lemma 9, it follows that

$$\begin{aligned} F(\mathbf{W}) &\leq F(\mathbf{W}^*) + \langle \nabla F(\mathbf{W}^*), \mathbf{W} - \mathbf{W}^* \rangle + \frac{L_F}{2} \|\mathbf{W} - \mathbf{W}^*\|_F^2 \\ &= F(\mathbf{W}^*) + \frac{L_F}{2} \|\mathbf{W} - \mathbf{W}^*\|_F^2 \leq F(\mathbf{W}^*) + \frac{\kappa_1 L_F}{2} \|\nabla F(\mathbf{W})\|_F^2, \end{aligned}$$

where the equality is because \mathbf{W}^* is a critical point, and the last inequality uses (11) and (175).

(ii) This result has been established in [34, Proposition 2.2]. □

D Auxiliary Results

To bound the spectral gap between eigenvectors associated with repeated eigenvalues between two symmetric matrices, we introduce the Davis-Kahan theorem [46].

Lemma 10. *Let $\mathbf{A}, \hat{\mathbf{A}} \in \mathbb{R}^{p \times p}$ be symmetric with eigenvalues $\lambda_1 \geq \cdots \geq \lambda_p$ and $\hat{\lambda}_1 \geq \cdots \geq \hat{\lambda}_p$ respectively. Fix $1 \leq r \leq s \leq p$ and assume that $\min\{\lambda_{r-1} - \lambda_r, \lambda_s - \lambda_{s+1}\} > 0$, where $\lambda_0 := \infty$ and $\lambda_{p+1} := -\infty$. Let $d := s - r + 1$ and $\mathbf{V} = (\mathbf{v}_r, \mathbf{v}_{r+1}, \dots, \mathbf{v}_s) \in \mathbb{R}^{p \times d}$ and $\hat{\mathbf{V}} = (\hat{\mathbf{v}}_r, \hat{\mathbf{v}}_{r+1}, \dots, \hat{\mathbf{v}}_s) \in \mathbb{R}^{p \times d}$ have orthonormal columns satisfying $\mathbf{A}\mathbf{v}_j = \lambda_j \mathbf{v}_j$ and $\hat{\mathbf{A}}\hat{\mathbf{v}}_j = \hat{\lambda}_j \hat{\mathbf{v}}_j$ for all $j = r, r+1, \dots, s$. Then, there exists an orthogonal matrix $\mathbf{O} \in \mathcal{O}^d$ such that*

$$\|\hat{\mathbf{V}}\mathbf{O} - \mathbf{V}\|_F \leq \frac{4\|\hat{\mathbf{A}} - \mathbf{A}\|_F}{\min\{\lambda_{r-1} - \lambda_r, \lambda_s - \lambda_{s+1}\}}.$$

Lemma 11 (Mirsky Inequality [37]). *For any matrices $\mathbf{X}, \tilde{\mathbf{X}} \in \mathbb{R}^{m \times n}$ with singular values*

$$\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_l \quad \text{and} \quad \tilde{\sigma}_1 \geq \tilde{\sigma}_2 \geq \cdots \geq \tilde{\sigma}_l,$$

where $l = \min\{m, n\}$, then for any unitarily invariant norm (e.g., $\|\cdot\|_F$), we have

$$\|\text{diag}(\tilde{\sigma}_1 - \sigma_1, \dots, \tilde{\sigma}_l - \sigma_l)\| \leq \|\tilde{\mathbf{X}} - \mathbf{X}\|.$$

Lemma 12. *Suppose that $(\boldsymbol{\sigma}, \boldsymbol{\Pi}) \in \mathbb{R}^{d_{\min}} \times \mathcal{P}^{d_{\min}}$ satisfies the following conditions:*

$$\begin{cases} \sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_{d_{\min}} \geq 0, \\ \sigma_i^{2L-1} - \sqrt{\lambda} y_{\pi(i)} \sigma_i^{L-1} + \lambda \sigma_i = 0, \quad \forall i \in [d_{\min}], \end{cases} \quad (176)$$

where $\pi : [d_{\min}] \rightarrow [d_{\min}]$ is the permutation corresponding to $\boldsymbol{\Pi}$. Then, the pair $(\boldsymbol{\sigma}, \boldsymbol{\Pi})$ belongs to the set \mathcal{B} .

Proof. According to (176), we conclude

$$\sigma_{\pi^{-1}(i)}^{2L-1} - \sqrt{\lambda} y_i \sigma_{\pi^{-1}(i)}^{L-1} + \lambda \sigma_{\pi^{-1}(i)} = 0, \quad \sigma_{\pi^{-1}(i)} \geq 0, \quad \forall i \in [d_{\min}].$$

Combining this with the definition of \mathcal{A} in (28), we have $\mathbf{a} := \boldsymbol{\Pi}^T \boldsymbol{\sigma} = (\sigma_{\pi^{-1}(1)}, \dots, \sigma_{\pi^{-1}(d_{\min})})^T \in \mathcal{A}$. Since $\boldsymbol{\Pi} \mathbf{a} = \boldsymbol{\sigma}$ and by the definition of \mathcal{B} in (29), it follows that the pair $(\boldsymbol{\sigma}, \boldsymbol{\Pi}) \in \mathcal{B}$. \square

D.1 Proof of (145)

To simplify the notation, for each $i \in [p]$, we introduce the following definition: $\hat{\mathbf{T}}_i := \prod_{l=2}^L \mathbf{T}_l^{(i)}$. Then we have

$$\begin{aligned} & \left\| \left(\hat{\boldsymbol{\Sigma}}_1 \hat{\boldsymbol{\Sigma}}_1^T \right)^{L-1} \hat{\boldsymbol{\Sigma}}_1 + \lambda \hat{\boldsymbol{\Sigma}}_1 - \sqrt{\lambda} \text{BlkDiag} \left((\boldsymbol{\Sigma}_1^{(1)})^{L-1}, \dots, (\boldsymbol{\Sigma}_1^{(p)})^{L-1}, \mathbf{0} \right) \boldsymbol{\Psi} \hat{\mathbf{T}} \right\|_F \\ &= \left\| \text{BlkDiag} \left(\hat{\mathbf{T}}_1, \dots, \hat{\mathbf{T}}_p, \mathbf{I}_{d_1 - r_\sigma} \right) \left(\left(\hat{\boldsymbol{\Sigma}}_1 \hat{\boldsymbol{\Sigma}}_1^T \right)^{L-1} \hat{\boldsymbol{\Sigma}}_1 + \lambda \hat{\boldsymbol{\Sigma}}_1 \right) \hat{\mathbf{T}}^T \right. \\ & \quad \left. - \sqrt{\lambda} \text{BlkDiag} (\mathbf{A}_1, \dots, \mathbf{A}_p, \mathbf{0}) \boldsymbol{\Psi} \right\|_F \\ &\leq \left\| \left(\hat{\boldsymbol{\Sigma}}_1 \hat{\boldsymbol{\Sigma}}_1^T \right)^{L-1} \hat{\boldsymbol{\Sigma}}_1 + \lambda \hat{\boldsymbol{\Sigma}}_1 - \sqrt{\lambda} \text{BlkDiag} (\mathbf{A}_1, \dots, \mathbf{A}_p, \mathbf{0}) \boldsymbol{\Psi} \right\|_F \\ & \quad + \sum_{i=1}^p \left\| \hat{\mathbf{T}}_i (\boldsymbol{\Sigma}_1^{(i)})^{2L-1} \hat{\mathbf{T}}_i^T - (\boldsymbol{\Sigma}_1^{(i)})^{2L-1} \right\|_F + \lambda \sum_{i=1}^p \left\| \hat{\mathbf{T}}_i \boldsymbol{\Sigma}_1^{(i)} \hat{\mathbf{T}}_i^T - \boldsymbol{\Sigma}_1^{(i)} \right\|_F \\ &\leq \eta_3 \|\nabla G(\mathbf{W})\|_F + \sum_{i=1}^p \left\| \hat{\mathbf{T}}_i (\boldsymbol{\Sigma}_1^{(i)})^{2L-1} - (\boldsymbol{\Sigma}_1^{(i)})^{2L-1} \hat{\mathbf{T}}_i \right\|_F + \lambda \sum_{i=1}^p \left\| \hat{\mathbf{T}}_i \boldsymbol{\Sigma}_1^{(i)} - \boldsymbol{\Sigma}_1^{(i)} \hat{\mathbf{T}}_i \right\|_F \\ &\leq \left(\eta_3 + \frac{p\eta_1(2L-1)L}{\sigma_{\min}^*} \left(\frac{3\sigma_{\max}^*}{2} \right)^{2L-2} + \frac{\lambda p\eta_1 L}{\sigma_{\min}^*} \right) \|\nabla G(\mathbf{W})\|_F, \end{aligned}$$

where the first equality follows from the definition of \mathbf{A}_i for each $i \in [p]$ in (56); the second inequality follows from (144); and the last inequality follows from (119) and the following fact (see

(177) and (178)). For each i in $[p]$, we have:

$$\begin{aligned}
& \left\| \hat{\mathbf{T}}_i (\boldsymbol{\Sigma}_1^{(i)})^{2L-1} - (\boldsymbol{\Sigma}_1^{(i)})^{2L-1} \hat{\mathbf{T}}_i \right\|_F = \left\| \left(\prod_{l=2}^L \mathbf{T}_l^{(i)} \right) (\boldsymbol{\Sigma}_1^{(i)})^{2L-1} - (\boldsymbol{\Sigma}_1^{(i)})^{2L-1} \prod_{l=2}^L \mathbf{T}_l^{(i)} \right\|_F \\
& \leq \left\| \left(\prod_{l=2}^{L-1} \mathbf{T}_l^{(i)} \right) \left(\mathbf{T}_L^{(i)} (\boldsymbol{\Sigma}_1^{(i)})^{2L-1} - (\boldsymbol{\Sigma}_1^{(i)})^{2L-1} \mathbf{T}_L^{(i)} \right) \right\|_F + \\
& \quad \left\| \left(\prod_{l=2}^{L-2} \mathbf{T}_l^{(i)} \right) \left(\mathbf{T}_{L-1}^{(i)} (\boldsymbol{\Sigma}_1^{(i)})^{2L-1} - (\boldsymbol{\Sigma}_1^{(i)})^{2L-1} \mathbf{T}_{L-1}^{(i)} \right) \mathbf{T}_L^{(i)} \right\|_F + \dots \\
& \quad + \left\| \left(\mathbf{T}_2^{(i)} (\boldsymbol{\Sigma}_1^{(i)})^{2L-1} - (\boldsymbol{\Sigma}_1^{(i)})^{2L-1} \mathbf{T}_2^{(i)} \right) \left(\prod_{l=3}^L \mathbf{T}_l^{(i)} \right) \right\|_F \\
& \leq \sum_{l=2}^L \sum_{j=1}^{2L-1} \left\| (\boldsymbol{\Sigma}_1^{(i)})^{j-1} \left(\mathbf{T}_l^{(i)} \boldsymbol{\Sigma}_1^{(i)} - \boldsymbol{\Sigma}_1^{(i)} \mathbf{T}_l^{(i)} \right) (\boldsymbol{\Sigma}_1^{(i)})^{2L-j-1} \right\|_F \\
& \leq \frac{\eta_1 (2L-1)L}{\sigma_{\min}^*} \left(\frac{3\sigma_{\max}^*}{2} \right)^{2L-2} \|\nabla G(\mathbf{W})\|_F,
\end{aligned} \tag{177}$$

where the last inequality follows from (119) and (42). Similarly, we have

$$\left\| \hat{\mathbf{T}}_i \boldsymbol{\Sigma}_1^{(i)} - \boldsymbol{\Sigma}_1^{(i)} \hat{\mathbf{T}}_i \right\|_F \leq \frac{\eta_1 L}{\sigma_{\min}^*} \left(\frac{3\sigma_{\max}^*}{2} \right)^{2L-2} \|\nabla G(\mathbf{W})\|_F. \tag{178}$$