

LEARNING DEEP MODELS: CRITICAL POINTS AND LOCAL OPENNESS

Anonymous authors

Paper under double-blind review

ABSTRACT

With the increasing interest in deeper understanding of the loss surface of many non-convex *deep models*, this paper presents a unifying framework to study the local/global optima equivalence of the optimization problems arising from training of such non-convex models. Using the *local openness* property of the underlying training models, we provide simple sufficient conditions under which any local optimum of the resulting optimization problem is globally optimal. We first *completely characterize the local openness of matrix multiplication mapping in its range*. Then we use our characterization to: 1) show that every local optimum of two layer linear networks is globally optimal. Unlike many existing results in the literature, our result requires no assumption on the target data matrix \mathbf{Y} , and input data matrix \mathbf{X} . 2) develop *almost complete* characterization of the local/global optima equivalence of multi-layer linear neural networks. We provide various counterexamples to show the necessity of each of our assumptions. 3) show global/local optima equivalence of non-linear deep models having certain pyramidal structure. Unlike some existing works, our result requires no assumption on the differentiability of the activation functions and can go beyond “full-rank” cases.

1 INTRODUCTION

Deep learning models have recently led to significant practical successes in various fields ranging from computer vision to natural language processing. Despite these significant empirical successes, the theoretical understanding of the behavior of these models is still very limited. While some recent works have tried to explain these successes through the lens of *expressivity* by showing the power of these models in learning large class of mappings, other works find the root of the success in the *generalizability* of these models from learning perspective.

From optimization perspective, training deep models require solving non-convex optimization problems, where non-convexity arises from the “deep” structure of the model. In fact, it has been shown by Blum & Rivest (1989) that training neural networks to global optimality is NP-complete in the worst case even for the simple case of three node networks. Despite this worst case barrier, the practical success of deep learning may suggest that most of the local optimal points of these models are close to the global optimal points. In particular, Choromanska et al. (2015) uses spin glass theory and empirical experiments to show that the local optima of deep neural network optimization problem are close to the global optima.

In an effort to better understand the landscape of training deep neural networks, Kawaguchi (2016); Lu & Kawaguchi (2017); Yun et al. (2017); Hardt & Ma (2016) studied the linear neural networks and provided sufficient conditions under which critical points (or local optimal points) of the training optimization problems are globally optimal. For non-linear neural networks, multiple works have shown that when the number of parameters of the model is larger than the data dimension, local optima of the resulting optimization problems can be easily found using local search procedures; see, e.g., Soltanolkotabi et al. (2017); Soudry & Carmon (2016); Nguyen & Hein (2017); Xie et al. (2017).

Despite the growing interest in studying the landscape of deep optimization problems, many of the results and mathematical analyses are problem specific and cannot be generalized to other problems and network structures easily. As a first step toward reaching a unifying theory for these results, we propose the use of open mappings for characterizing the properties of the local optima of these “deep” optimization problems.

To study the landscape of shallow/deep models, we study the general optimization problem

$$\underset{\mathbf{w} \in \mathcal{W}}{\text{minimize}} \ell(\mathcal{F}(\mathbf{w})), \quad (1)$$

where $\ell(\cdot)$ is the loss function and $\mathcal{F}(\cdot)$ represents a statistical model with parameter \mathbf{w} that needs to be learned by solving the above optimization problem. A simple example is the popular linear regression problem

$$\underset{\mathbf{w}}{\text{minimize}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2,$$

where \mathbf{y} is a given constant response vector and \mathbf{X} is a given constant feature matrix. In this example, the loss function is the ℓ_2 loss, i.e., $\ell(\mathbf{z}) = \|\mathbf{z} - \mathbf{y}\|_2^2$, and the fitted model \mathcal{F} is a linear model, i.e., $\mathcal{F}(\mathbf{w}) = \mathbf{X}\mathbf{w}$. While this linear regression problem is convex and easy, fitting many practical models, such as deep neural networks, requires solving non-trivial non-convex optimization problems.

In this paper, we use the local openness of the mapping \mathcal{F} to provide sufficient conditions under which every local optimum of (1) is in fact global optimum. To proceed, let us define our notations that will be used throughout the paper. We use $\mathbf{A}_{l,\cdot}$ and $\mathbf{A}_{\cdot,l}$ to denote the l^{th} row and column of the matrix \mathbf{A} , respectively. The notation $\mathbf{I}_d \in \mathbb{R}^{d \times d}$ is used to denote the $d \times d$ -dimensional identity matrix. Let $\|\mathbf{A}\|$, $\mathcal{N}(\mathbf{A})$, $\mathcal{C}(\mathbf{A})$, $\text{rank}(\mathbf{A})$ be respectively the Frobenius norm, null-space, column-space, and the rank of the matrix \mathbf{A} . Given subspaces \mathbf{U} and \mathbf{V} , we say $\mathbf{U} \perp \mathbf{V}$ if \mathbf{U} is orthogonal to \mathbf{V} , and $\mathbf{U} = \mathbf{V}^\perp$ if \mathbf{U} is the orthogonal complement of \mathbf{V} . We say matrix $\mathbf{A} \in \mathbb{R}^{d_1 \times d_0}$ is rank deficient if $\text{rank}(\mathbf{A}) < \min\{d_1, d_0\}$, and full rank if $\text{rank}(\mathbf{A}) = \min\{d_1, d_0\}$. We call a point $\mathbf{W} = (\mathbf{W}_h, \dots, \mathbf{W}_1)$, with $\mathbf{W}_i \in \mathbb{R}^{d_i \times d_{i-1}}$, non-degenerate if $\text{rank}(\mathbf{W}_h \cdots \mathbf{W}_1) = \min_{0 \leq i \leq h} d_i$, and degenerate if $\text{rank}(\mathbf{W}_h \cdots \mathbf{W}_1) < \min_{0 \leq i \leq h} d_i$. We also say a point $\bar{\mathbf{W}}$ is a *second order saddle point* of an unconstrained optimization problem if the gradient of the objective function is zero at $\bar{\mathbf{W}}$ and the hessian of the objective function at $\bar{\mathbf{W}}$ has a negative eigenvalue. Let us start by briefly explaining the training problem of feedforward neural networks which will also be used as a motivation for our analysis:

Example: Training Feedforward Neural Networks. Consider the following multiple layer feedforward neural network optimization problem:

$$\underset{\mathbf{W}}{\text{minimize}} \frac{1}{2} \|\mathcal{F}_h(\mathbf{W}) - \mathbf{Y}\|^2$$

where \mathcal{F}_h is defined in a recursive manner:

$$\mathcal{F}_k(\mathbf{W}) \triangleq \sigma_k(\mathbf{W}_k \mathcal{F}_{k-1}(\mathbf{W})), \text{ for } k \in \{2, \dots, h\},$$

with

$$\mathcal{F}_1(\mathbf{W}) \triangleq \sigma_1(\mathbf{W}_1 \mathbf{X}).$$

Here h is the number of hidden units in our network; $\sigma_k(\cdot)$ denotes the activation function of layer k ; the matrix $\mathbf{W}_k \in \mathbb{R}^{d_k \times d_{k-1}}$ is the weight of layer k with $\mathbf{W} \triangleq (\mathbf{W}_i)_{i=1}^h$ being our optimization variable. The matrix $\mathbf{X} \in \mathbb{R}^{d_0 \times n}$ is the input training data; and $\mathbf{Y} \in \mathbb{R}^{d_h \times n}$ is the target training data where n is the number of samples; see, e.g. Goodfellow & Courville (2016). Notice that this problem is a special case of the optimization problem in (1) which can be obtained simply by setting our loss function to the ℓ_2 loss, and setting $\mathcal{F} = \mathcal{F}_h$.

A special instance of this optimization problem was studied in Nguyen & Hein (2017) which considers the non-linear neural network with pyramidal structure (i.e. $d_i \leq d_{i-1} \forall i = 1, \dots, h$ and $d_0 \geq n$). Note that this special network structure does not allow wide intermediate layers. (Nguyen & Hein, 2017, Theorem 3.8) shows that under some conditions, among which are the differentiability of the loss function $\ell(\cdot)$ and the activation function $\sigma(\cdot)$, if \mathbf{W} is a critical point with \mathbf{W}_i 's being full row rank then it is a global minimum. In this paper, we will relax the differentiability assumption on both $\ell(\cdot)$ and $\sigma(\cdot)$; and we will show any local optimum is a global optimum of the objective function. Another special case is the linear feedforward network where the mapping $\sigma_k(\cdot)$ is the identity map in all layers, which leads to the optimization problem:

$$\underset{\mathbf{W}}{\text{minimize}} \frac{1}{2} \|\mathbf{W}_h \cdots \mathbf{W}_1 \mathbf{X} - \mathbf{Y}\|^2. \quad (2)$$

For this optimization problem, Lu & Kawaguchi (2017) showed that every local optimum of the objective function is globally optimal under some assumptions. More precisely, by using perturbation analysis, (Lu & Kawaguchi, 2017, Theorem 2.2) prove that when \mathbf{X} and \mathbf{Y} are full row rank, every local optimum in problem (2) is a local optimum of the following problem:

$$\begin{aligned} & \underset{\mathbf{Z} \in \mathbb{R}^{d_h \times d_0}}{\text{minimum}} \frac{1}{2} \|\mathbf{Z}\mathbf{X} - \mathbf{Y}\|^2 \\ & \text{subject to} \quad \text{rank}(\mathbf{Z}) \leq d_p \triangleq \min_{0 \leq i \leq h} d_i. \end{aligned} \quad (3)$$

Moreover, they show that when \mathbf{X} is full row rank, every local optimum of problem (3) is a global optimum. Thus, with the sufficient condition that \mathbf{X} and \mathbf{Y} are both full row rank, every local optimum of problem (2) is a global

optimum. Another recent work Yun et al. (2017) shows the same result under similar set of assumptions. It is in fact not hard to see that one cannot relax the full rankness assumption of \mathbf{Y} due to the following simple counterexample:

$$\mathbf{X} = \mathbf{I} \quad \mathbf{W}_3 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad \mathbf{W}_2 = [0], \quad \mathbf{W}_1 = [1 \quad 0], \quad \mathbf{Y} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}.$$

It is not hard to check that the point $\mathbf{W} = (\mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3)$ is a local optimum of a 3-layer deep linear model (problem (2) with $h = 3$) that is not a global optimum. However, we will show that if a given local optimum is non-degenerate (which is a simple checkable condition), the full rankness of \mathbf{Y} can be relaxed. Moreover, for degenerate local optima, we show that if there exist $1 \leq p_1 < p_2 \leq h - 1$ with $d_h > d_{p_2}$ and $d_0 > d_{p_1}$, we can find \mathbf{Y} and \mathbf{X} such that problem (2) has a local minimum that is not global. Otherwise, given any \mathbf{X} and \mathbf{Y} , we present a method for constructing a descent direction from any given degenerate critical point that is not a global optimum; thus we show every degenerate local minimum is global.

Other examples: Matrix Factorization and Matrix Completion. In addition to the training of deep neural networks, the matrix completion problem also lies in the category of non-convex problems in (1). For the matrix completion problem, Park et al. (2016) shows that the non-convex matrix factorization formulation of the non-square matrix sensing problem has no spurious local optimum under restricted isometry property (RIP) conditions. Similar results were obtained for the symmetric matrix multiplication problem by Ge et al. (2016), and the non-convex factorized low-rank matrix recovery problem by Bhojanapalli et al. (2016). Like the analysis in Ge et al. (2016), we start with the fully observed matrix completion scenario:

$$\underset{\mathbf{W}_1 \in \mathbb{R}^{d_1 \times d_0}, \mathbf{W}_2 \in \mathbb{R}^{d_2 \times d_1}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{W}_2 \mathbf{W}_1 - \mathbf{Y}\|^2. \quad (4)$$

This problem, which is also referred to as the low rank matrix estimation problem in Srebro & Jaakkola (2003), can also be viewed as a 2-layer linear neural network optimization problem with the input data matrix $\mathbf{X} = \mathbf{I}$. Clearly, this problem is much simpler than the general matrix completion problem and we only study it as a first step. Moreover, this optimization problem is a special case of (1) with the loss function being the ℓ_2 loss, and the mapping \mathcal{F} being defined as $\mathcal{F}(\mathbf{W}_1, \mathbf{W}_2) = \mathbf{W}_2 \mathbf{W}_1$. In this paper, using our framework, we show that every critical point of (4) is either a global minimum or a second-order saddle point. This result can be generalized to general loss function $\ell(\cdot)$ for degenerate critical points.

In addition to these results, one of our main contributions is the complete characterization of the local openness of the matrix multiplication mapping in its range. These results could be used in many other optimization problems for characterizing the local/global equivalence.

2 MATHEMATICAL FRAMEWORK

As discussed in the previous section, we are interested in solving

$$\underset{\mathbf{w} \in \mathcal{W}}{\text{minimize}} \quad \ell(\mathcal{F}(\mathbf{w})), \quad (5)$$

where $\mathcal{F} : \mathcal{W} \mapsto \mathcal{S}$ is a mapping and $\ell : \mathcal{S} \mapsto \mathbb{R}$ is a loss function. Here we assume that the set \mathcal{W} is closed and the mapping \mathcal{F} is continuous. In non-convex scenarios, this optimization problem can only be solved up to ‘‘local optimality’’ by local search procedures; see Lee et al. (2016) for an example. To proceed, let us define the auxiliary optimization problem

$$\underset{s \in \mathcal{S}}{\text{minimize}} \quad \ell(s), \quad (6)$$

where \mathcal{S} is the range of the mapping \mathcal{F} . Since problem (6) minimizes the function $\ell(\cdot)$ over the range of the mapping \mathcal{F} , the global optimal objective values for problems (5) and (6) are the same. Moreover, there is a clear relation between the global optima of the two optimization problem through the mapping \mathcal{F} . However, the connection between the local optima of the two optimization problems is not clear. This connection, in particular, is important when the local optima of (6) are ‘‘nice’’ (e.g. globally optimal or close to optimal). In what follows, we establish the connection between the local optima of the optimization problems (5) and (6) under some simple sufficient conditions. This connection is then used to study the relation between local and global optima of (5) and (6) for various deep learning models. Let us first define the following concepts, which will help us state our simple sufficient condition.

- **Open mapping:** A mapping $\mathcal{F} : \mathcal{W} \rightarrow \mathcal{S}$ is said to be open, if for every open set $U \in \mathcal{W}$, $\mathcal{F}(U)$ is (relatively) open in \mathcal{S} .
- **Locally open mapping:** A mapping $\mathcal{F}(\cdot)$ is said to be locally open at w if for every $\epsilon > 0$, there exists $\delta > 0$ such that $\mathcal{B}_\delta(\mathcal{F}(w)) \subseteq \mathcal{F}(\mathcal{B}_\epsilon(w))$. Here $\mathcal{B}_\delta(w) \subseteq \mathcal{W}$ is an open ball with radius δ centered at w , and $\mathcal{B}_\epsilon(\mathcal{F}(w)) \subseteq \mathcal{S}$ is the ball of radius ϵ centered at $\mathcal{F}(w)$.

By definition, openness of a mapping is stronger than local openness. Furthermore, it is not hard to see that a mapping is locally open everywhere if and only if it is open. A useful property of (locally) open mappings is that *the composition of two (locally) open maps is (locally) open*.

The following simple intuitive observation, which establishes a connection between the local optima of (5) and (6), is a major building block of our analyses.

Observation 1. *Suppose $\mathcal{F}(\cdot)$ is locally open at \bar{w} . If \bar{w} is a local minimum of problem (5), then $\bar{s} = \mathcal{F}(\bar{w})$ is a local minimum of problem (6).*

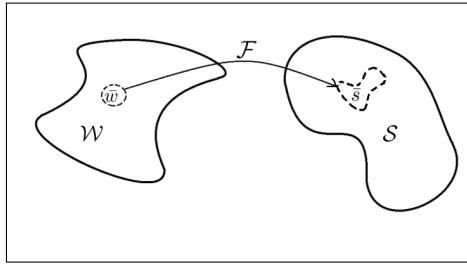


Figure 1: Sketch of the Proof of Observation 1.

Proof. Let \bar{w} be a local minimum of problem (5). Then there exists an $\epsilon > 0$ such that $\ell(\mathcal{F}(\bar{w})) \leq \ell(\mathcal{F}(w))$, $\forall w \in \mathcal{B}_\epsilon(\bar{w})$. By the definition of local openness,

$$\exists \delta > 0 \text{ such that } \mathcal{B}_\delta(\bar{s}) \subset \mathcal{F}(\mathcal{B}_\epsilon(\bar{w})).$$

where $\bar{s} = \mathcal{F}(\bar{w})$. Therefore, $\ell(\bar{s}) \leq \ell(s)$, $\forall s \in \mathcal{B}_\delta(\bar{s})$, which implies \bar{s} is a local minimum of problem (6). \square

The above observation can be used to map multiple local optima of the original problem (5) to one local optimum of the auxiliary problem (6); and potentially make the problem easier to understand. This mapping is particularly interesting in neural networks since permuting the weights in each layer does not change the objective function. Hence, by nature, the optimization problem has multiple (disconnected) global optima; and hence it is non-convex. However, collapsing these multiple local optima to one could potentially simplify the problem. In other words, instead of understanding the problem in the original variables, we can analyze it in the space of the resulted mapping. Let us clarify this point through the following simple examples:

Example 2. *Consider the optimization problem*

$$\underset{w \in \mathbb{R}}{\text{minimize}} (w^2 - 1)^2, \tag{7}$$

and its corresponding auxiliary problem

$$\underset{z \geq 0}{\text{minimize}} (z - 1)^2. \tag{8}$$

Plots of these two problems can be found in Figure 2a and Figure 2b. Since $\mathcal{F}(x) \triangleq w^2$ is an open mapping in its range, it follows from Observation 1 that every local minimum in problem (7) is a local minimum of problem (8). Thus the two local minima $w = -1$ and $w = +1$ in (7) are mapped to a single local minimum $z = 1$ of problem (8). Moreover, since the optimization problem (8) is convex, the local minimum is global; and hence the original local optima $w = -1$ and $w = +1$ should be both global despite non-convexity of (7).

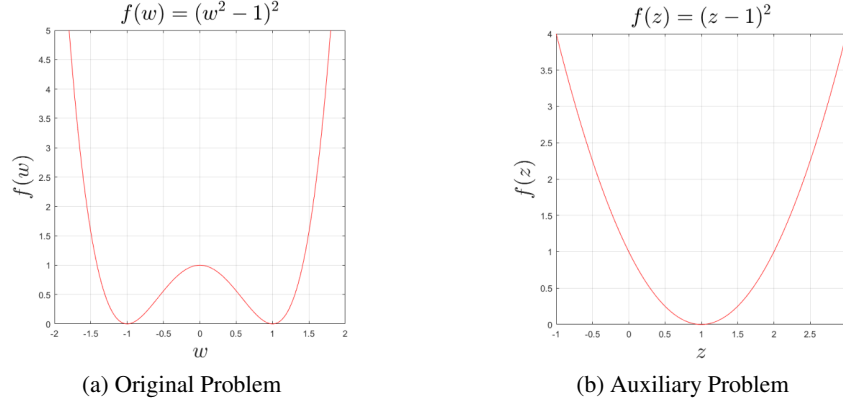


Figure 2: Two local minima $w = -1$ and $w = +1$ in (a) are mapped to a single local minimum $z = 1$ in (b).

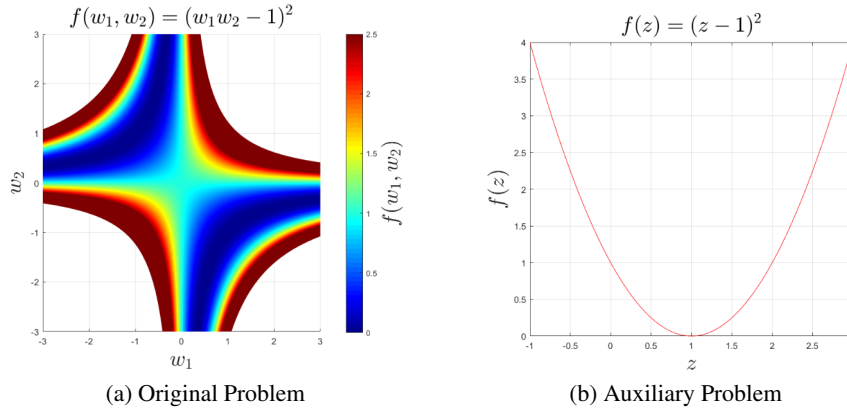


Figure 3: All the points in the set $\{(w_1, w_2) \mid w_1 w_2 = 1\}$ are local minima in (a) and are mapped to a single local minimum $z = 1$ in (b).

Example 3. Another example is related to the widely used matrix multiplication mapping $\mathbf{W}_1 \mathbf{W}_2$. Let $(\bar{\mathbf{W}}_1, \bar{\mathbf{W}}_2)$ be a local minimum of the optimization problem

$$\underset{\mathbf{W}_1 \in \mathbb{R}^{m \times k}, \mathbf{W}_2 \in \mathbb{R}^{k \times n}}{\text{minimize}} \quad \ell(\mathbf{W}_1 \mathbf{W}_2).$$

Then, any point in the set $\mathcal{S} \triangleq \{(\bar{\mathbf{W}}_1 \mathbf{Q}_1, \mathbf{Q}_2 \bar{\mathbf{W}}_2) \text{ with } \mathbf{Q}_1 \mathbf{Q}_2 = \mathbf{I}\}$ is also a local minimum. If the matrix product $\mathbf{W}_1 \mathbf{W}_2$ is locally open at the point $(\bar{\mathbf{W}}_1, \bar{\mathbf{W}}_2)$, then all points in \mathcal{S} are mapped to a single local minimum $\mathbf{Z} = \bar{\mathbf{W}}_1 \bar{\mathbf{W}}_2$ in the corresponding auxiliary problem. A simple one dimensional example is plotted in Figure 3a and Figure 3b.

This motivates us to study the local openness of the matrix multiplication mapping defined as

$$\mathcal{M} : \mathbb{R}^{m \times k} \times \mathbb{R}^{k \times n} \mapsto \mathcal{R}_{\mathcal{M}} \quad \text{with} \quad \mathcal{M}(\mathbf{W}_1, \mathbf{W}_2) \triangleq \mathbf{W}_1 \mathbf{W}_2, \quad (9)$$

where $\mathcal{R}_{\mathcal{M}} \triangleq \{\mathbf{Z} \in \mathbb{R}^{m \times n} \mid \text{rank}(\mathbf{Z}) \leq \min\{m, n, k\}\}$ is the range of the mapping \mathcal{M} .

Although matrix multiplication mappings $\mathcal{M}(\mathbf{W}_1, \mathbf{W}_2)$ naturally appears in deep models and is widely used as a non-convex factorization for rank constrained problems, see Wang et al. (2016); Bhojanapalli et al. (2016); Ge et al. (2016); Srebro & Jaakkola (2003); Sun (2015), to our knowledge, the complete characterization of the local openness of this mapping has not been studied in the optimization literature before.

While the classical open mapping theorem in Rudin (1973) states that surjective continuous linear operators are open, this is not true in general for bilinear mappings such as matrix product. In fact, by providing a simple counterexample

of a bilinear mapping that is not open, Horowitz (1975) shows that the linear case cannot be generally extended to multilinear maps. Several papers, see Balcerzak et al. (2013; 2005); Behrends (2011), investigate this bilinear mapping and provide a characterization of the points where this mapping is open. Moreover, Behrends (2017) studies the matrix multiplication mapping \mathcal{M} which is a special example of bilinear mappings and provides an almost complete characterization of the points where the mapping is locally open. However, the openness is studied in $\mathbb{R}^{m \times n}$; while the range of the mapping is $\mathcal{R}_{\mathcal{M}}$; and the (relative) local openness should be studied with respect to this range in our framework. This, in particular causes trouble when $\mathbb{R}^{m \times n} \neq \mathcal{R}_{\mathcal{M}}$, i.e., when $k < \min\{m, n\}$.

For the above reason, we study the local openness of the mapping \mathcal{M} in its range $\mathcal{R}_{\mathcal{M}}$ and characterize it completely. An intuitive (and unofficial) definition of local openness of $\mathcal{M}(\cdot)$ at $(\bar{\mathbf{W}}_1, \bar{\mathbf{W}}_2)$ in $\mathcal{R}_{\mathcal{M}}$ is as follows. We say the multiplication mapping is locally open at $(\bar{\mathbf{W}}_1, \bar{\mathbf{W}}_2)$ if for any small perturbation $\tilde{\mathbf{Z}} \in \mathcal{R}_{\mathcal{M}}$ of $\mathbf{Z} = \bar{\mathbf{W}}_1 \bar{\mathbf{W}}_2$, there exists a pair $(\tilde{\mathbf{W}}_1, \tilde{\mathbf{W}}_2)$, a small perturbation of $(\bar{\mathbf{W}}_1, \bar{\mathbf{W}}_2)$, such that $\tilde{\mathbf{Z}} = \tilde{\mathbf{W}}_1 \tilde{\mathbf{W}}_2$.

Notice that when $k \geq \min\{m, n\}$, we get $\mathcal{R}_{\mathcal{M}} = \mathbb{R}^{m \times n}$. However, in the case where $k < \min\{m, n\}$ the mapping is definitely not locally open in $\mathbb{R}^{m \times n}$, but can still be locally open in $\mathcal{R}_{\mathcal{M}}$. As a simple example, consider $\bar{\mathbf{W}}_1 = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$ and $\bar{\mathbf{W}}_2 = \begin{bmatrix} 1 & 1 \end{bmatrix}$. In this example there does not exist $\tilde{\mathbf{W}}_1, \tilde{\mathbf{W}}_2$ perturbations of $\bar{\mathbf{W}}_1$ and $\bar{\mathbf{W}}_2$ respectively such that $\tilde{\mathbf{W}}_1 \tilde{\mathbf{W}}_2 = \tilde{\mathbf{Z}}$ when $\tilde{\mathbf{Z}}$ is a full rank perturbation of $\mathbf{Z} = \bar{\mathbf{W}}_1 \bar{\mathbf{W}}_2$; however, for any rank 1 perturbation $\tilde{\mathbf{Z}}$, we can find a perturbed pair $(\tilde{\mathbf{W}}_1, \tilde{\mathbf{W}}_2)$ such that $\tilde{\mathbf{Z}} = \tilde{\mathbf{W}}_1 \tilde{\mathbf{W}}_2$. Motivated by Observation 1, we study in the next section the local openness/openness of the mapping \mathcal{M} . We later use these results to analyze the behavior of local optima of deep neural networks.

3 LOCAL OPENNESS OF THE MATRIX MULTIPLICATION MAPPING

When $\mathbf{W}_1 \in \mathbb{R}^{m \times k}$ and $\mathbf{W}_2 \in \mathbb{R}^{k \times n}$ with $k \geq \min\{m, n\}$, the range of the mapping $\mathcal{M}(\mathbf{W}_1, \mathbf{W}_2) = \mathbf{W}_1 \mathbf{W}_2$ is the entire space $\mathbb{R}^{m \times n}$. In this case, which we refer to as the full rank case, (Behrends, 2017, Theorem 2.5) provides a complete characterization of the pairs $(\mathbf{W}_1, \mathbf{W}_2)$ for which the mapping is locally open. However, when $k < \min\{m, n\}$, which we refer to as the rank-deficient case, the characterization of the set of points for which the mapping is locally open has not been resolved before. We settled this question in Theorem 5 by providing a complete characterization of points $(\mathbf{W}_1, \mathbf{W}_2)$ for which the mapping \mathcal{M} is locally open when $k < \min\{m, n\}$. We start by restating the main result in Behrends (2017):

Proposition 4. (Behrends, 2017, Theorem 2.5 Rephrased) *Let $\mathcal{M}(\mathbf{W}_1, \mathbf{W}_2) = \mathbf{W}_1 \mathbf{W}_2$ denote the matrix multiplication mapping with $\mathbf{W}_1 \in \mathbb{R}^{m \times k}$ and $\mathbf{W}_2 \in \mathbb{R}^{k \times n}$. Assume $k \geq \min\{m, n\}$. Then the the following statements are equivalent:*

1. $\mathcal{M}(\cdot, \cdot)$ is locally open at $(\bar{\mathbf{W}}_1, \bar{\mathbf{W}}_2)$.
2.
$$\begin{cases} \exists \tilde{\mathbf{W}}_1 \in \mathbb{R}^{m \times k} \text{ such that } \tilde{\mathbf{W}}_1 \bar{\mathbf{W}}_2 = \mathbf{0} \text{ and } \bar{\mathbf{W}}_1 + \tilde{\mathbf{W}}_1 \text{ is full row rank.} \\ \text{or} \\ \exists \tilde{\mathbf{W}}_2 \in \mathbb{R}^{k \times n} \text{ such that } \bar{\mathbf{W}}_1 \tilde{\mathbf{W}}_2 = \mathbf{0} \text{ and } \bar{\mathbf{W}}_2 + \tilde{\mathbf{W}}_2 \text{ is full column rank.} \end{cases}$$
3. $\dim(\mathcal{N}(\bar{\mathbf{W}}_1) \cap \mathcal{C}(\bar{\mathbf{W}}_2)) \leq k - m$ or $n - (\text{rank}(\bar{\mathbf{W}}_2) - \dim(\mathcal{N}(\bar{\mathbf{W}}_1) \cap \mathcal{C}(\bar{\mathbf{W}}_2))) \leq k - \text{rank}(\bar{\mathbf{W}}_1)$.

The above proposition provides a checkable condition which completely characterizes the local openness of the mapping \mathcal{M} at different points when the range of the mapping is the entire space. Now, let us state our result that characterizes the local openness of the mapping \mathcal{M} in its range when $k < \min\{m, n\}$.

Theorem 5. *Let $\mathcal{M}(\mathbf{W}_1, \mathbf{W}_2) = \mathbf{W}_1 \mathbf{W}_2$ denote the matrix multiplication mapping with $\mathbf{W}_1 \in \mathbb{R}^{m \times k}$ and $\mathbf{W}_2 \in \mathbb{R}^{k \times n}$. Assume $k < \min\{m, n\}$. Then if $\text{rank}(\bar{\mathbf{W}}_1) \neq \text{rank}(\bar{\mathbf{W}}_2)$, $\mathcal{M}(\cdot, \cdot)$ is not locally open at $(\bar{\mathbf{W}}_1, \bar{\mathbf{W}}_2)$. Else, if $\text{rank}(\bar{\mathbf{W}}_1) = \text{rank}(\bar{\mathbf{W}}_2)$, then the following statements are equivalent:*

- i) $\exists \tilde{\mathbf{W}}_1 \in \mathbb{R}^{m \times k}$ such that $\tilde{\mathbf{W}}_1 \bar{\mathbf{W}}_2 = \mathbf{0}$ and $\bar{\mathbf{W}}_1 + \tilde{\mathbf{W}}_1$ is full column rank.

ii) $\exists \widetilde{\mathbf{W}}_2 \in \mathbb{R}^{k \times n}$ such that $\widetilde{\mathbf{W}}_1 \widetilde{\mathbf{W}}_2 = \mathbf{0}$ and $\widetilde{\mathbf{W}}_2 + \widetilde{\mathbf{W}}_2$ is full row rank.

iii) $\dim(\mathcal{N}(\widetilde{\mathbf{W}}_1) \cap \mathcal{C}(\widetilde{\mathbf{W}}_2)) = 0$.

iv) $\dim(\mathcal{N}(\widetilde{\mathbf{W}}_2^T) \cap \mathcal{C}(\widetilde{\mathbf{W}}_1^T)) = 0$.

v) $\mathcal{M}(\cdot, \cdot)$ is locally open at $(\widetilde{\mathbf{W}}_1, \widetilde{\mathbf{W}}_2)$ in its range $\mathcal{R}_{\mathcal{M}}$.

Note that the proof of Theorem 5, which can be found in the appendix section, is different than the proof of Proposition 4, as in the former we need to work with the set of low rank matrices. Besides, the conditions in Theorem 5 are different than the ones in Proposition 4. For example, while conditions i) and ii) are equivalent in the rank-deficient case, they are not equivalent in the full-rank case. Moreover, unlike the full-rank case, the condition $\text{rank}(\widetilde{\mathbf{W}}_1) = \text{rank}(\widetilde{\mathbf{W}}_2)$ is necessary for local openness in the low rank case.

How much perturbation is needed? As previously mentioned, local openness can be described in terms of perturbation analysis. For example, $\mathcal{M}(\cdot, \cdot)$ is locally open at $(\mathbf{W}_1, \mathbf{W}_2)$ if for a given $\epsilon > 0$, there exists $\delta > 0$ such that for any $\widetilde{\mathbf{Z}} = \mathbf{Z} + \mathbf{R}_\delta \in \mathcal{R}_{\mathcal{M}}$ with $\|\mathbf{R}_\delta\| \leq \delta$, there exists $\widetilde{\mathbf{W}}_1, \widetilde{\mathbf{W}}_2$ with $\|\widetilde{\mathbf{W}}_1\| \leq \epsilon, \|\widetilde{\mathbf{W}}_2\| \leq \epsilon$, such that $\widetilde{\mathbf{Z}} = (\mathbf{W}_1 + \widetilde{\mathbf{W}}_1)(\mathbf{W}_2 + \widetilde{\mathbf{W}}_2)$. As a perturbation bound on δ , we show that for any locally open pair $(\mathbf{W}_1, \mathbf{W}_2)$, given an $\epsilon > 0$, the chosen δ is of order ϵ , i.e., $\delta = \mathcal{O}(\epsilon)$. The details of our analysis can be found in the proof of Theorem 5 in Appendix B.

Remark 1 It follows from Theorem 5 that when \mathbf{W}_1 is full column rank, and \mathbf{W}_2 is full row rank, the mapping $\mathcal{M}(\cdot, \cdot)$ is locally open at $(\mathbf{W}_1, \mathbf{W}_2)$. This result was observed in other works; see, e.g., (Sun, 2015, Proposition 4.2). Also when $k < \min\{m, n\}$ if only one of the two matrices is full rank, then the mapping is not locally open. We have showed this result in the proof of Theorem 5, and below is a simple example for this phenomenon:

Let

$$\mathbf{W}_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad \mathbf{W}_2 = [0, 0], \quad \mathbf{W}_1 \mathbf{W}_2 = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}, \quad \mathbf{R}_\delta = \begin{bmatrix} \delta & 0 \\ 0 & 0 \end{bmatrix},$$

then $\mathbf{W}_1 \mathbf{W}_2 + \mathbf{R}_\delta$ is rank one and hence feasible perturbation. However, for any perturbation $\widetilde{\mathbf{W}}_1 = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \end{bmatrix}$ and $\widetilde{\mathbf{W}}_2 = [\epsilon_3, \epsilon_4]$, we have

$$(\mathbf{W}_1 + \widetilde{\mathbf{W}}_1)(\mathbf{W}_2 + \widetilde{\mathbf{W}}_2) = \begin{bmatrix} (1 + \epsilon_1)\epsilon_3 & (1 + \epsilon_1)\epsilon_4 \\ (1 + \epsilon_2)\epsilon_3 & (1 + \epsilon_2)\epsilon_4 \end{bmatrix}.$$

Hence, in order for this perturbation to be equal to $\mathbf{W}_1 \mathbf{W}_2 + \mathbf{R}_\delta$, we need ϵ_3 to be different from zero. But when ϵ_3 is different from zero, for small enough ϵ_2 , there does not exist such $\widetilde{\mathbf{W}}_1$ and $\widetilde{\mathbf{W}}_2$, or equivalently, $\mathcal{M}(\cdot, \cdot)$ is not locally open at $(\mathbf{W}_1, \mathbf{W}_2)$.

In the next sections, we use our local openness result to characterize the cases where the local optima of various training optimization problem of the form (5) are globally optimal.

4 NON-LINEAR DEEP NEURAL NETWORK WITH A PYRAMIDAL STRUCTURE:

Consider the non-linear deep neural network optimization problem with a pyramidal structure

$$\underset{\mathbf{W}}{\text{minimize}} \ell(\mathcal{F}_h(\mathbf{W})) \quad \text{with} \quad \mathcal{F}_1(\mathbf{W}) \triangleq \sigma_1(\mathbf{W}_1 \mathbf{X}); \quad \mathcal{F}_k(\mathbf{W}) \triangleq \sigma_k(\mathbf{W}_k \mathcal{F}_{k-1}(\mathbf{W})), \quad (10)$$

for $i \in [2, h]$, where $\sigma_i(\cdot)$ is the activation function applied component-wise to the entries of each layer, i.e., $\sigma_i(\mathbf{A}) = [\sigma_i(\mathbf{A}_{ij})]_{i,j}$ with $\sigma_i: \mathbb{R} \mapsto \mathbb{R}$ being continuous and strictly monotone. Here $\mathbf{W} = (\mathbf{W}_i)_{i=1}^h$ where $\mathbf{W}_i \in \mathbb{R}^{d_i \times d_{i-1}}$ is the weight matrix of layer i , and $\mathbf{X} \in \mathbb{R}^{d_0 \times n}$ is the input training data. In this section, we consider the pyramidal network structure with $d_0 > n$ and $d_i \leq d_{i-1}$ for $1 \leq i \leq h$; see Nguyen & Hein (2017) for more details on these types of networks.

First notice that when \mathbf{X} is full column rank and the functions σ_i 's are all continuous and strictly monotone, the image of the mapping \mathcal{F}_h is convex and hence every local optimum of the auxiliary optimization problem (6) is global. We

now show that when \mathbf{W}_i 's are all full row rank and the functions σ_i are all strictly monotone, the mapping \mathcal{F}_h is locally open at \mathbf{W} .

Lemma 6. *Assume the functions $\sigma_i(\cdot) : \mathbb{R} \mapsto \mathbb{R}$ are all continuous strictly monotone. Then the mapping \mathcal{F}_h defined in (10) is locally open at the point $\mathbf{W} = (\mathbf{W}_1, \dots, \mathbf{W}_h)$ if \mathbf{W}_i 's are all full row rank.*

Before proving this result, we would like to remark that many of the popular activation functions such as logsitic, tangent hyperbolic, and leaky ReLu are strictly monotone and satisfy the assumptions of this lemma.

Proof. Let us prove by induction. Since linear mappings are open, and since $\sigma_1(\cdot)$ is strictly monotone; by using the composition property of open maps, we get that \mathcal{F}_1 is open.

Assume $\mathcal{F}_{k-1} \left((\mathbf{W}_i)_{i=1}^{k-1} \right)$ is locally open at $(\mathbf{W}_i)_{i=1}^{k-1}$, then using Proposition 4, due to the full row rankness of \mathbf{W}_k , the mapping $\mathbf{W}_k \mathcal{F}_{k-1} \left((\mathbf{W}_i)_{i=1}^{k-1} \right)$ is locally open at $(\mathbf{W}_k, (\mathbf{W}_i)_{i=1}^{k-1})$. Using the composition property of open maps and strict monotonicity of $\sigma_k(\cdot)$, we get $\mathcal{F}_k \left((\mathbf{W}_i)_{i=1}^k \right)$ is locally open at $(\mathbf{W}_i)_{i=1}^k$. \square

Thus, by Observation 1, if $\bar{\mathbf{W}}$ is a local optimum of problem (10) with $\bar{\mathbf{W}}_i$'s being full row rank, then $\bar{\mathbf{Z}} = \mathcal{F}_h(\bar{\mathbf{W}})$ is a local optimum of the corresponding auxiliary problem:

$$\underset{\mathbf{Z} \in \mathcal{Z}}{\text{minimize}} \ell(\mathbf{Z})$$

where \mathcal{Z} is convex. Consequently, $\bar{\mathbf{Z}}$ is a global optimum of problem (10) when the loss function $\ell(\cdot)$ is convex. Nguyen & Hein (2017) show that every critical point \mathbf{W} of problem (10) with \mathbf{W}_i 's being full row rank is a global optimum when both $\sigma(\cdot)$ and $\ell(\cdot)$ are differentiable. Our result relaxes the differentiability assumption on both the activation and loss functions; however, we can only show all local optima are global. A popular activation function that is strictly monotonic and not differentiable is the Leaky ReLU, for which our result follows. It is also worth mentioning that Nguyen & Hein (2017) allow wide intermediate layers in parts of their result. It is not clear if this result can be extended to non-differentiable activation functions as well or not.

5 TWO-LAYER LINEAR NEURAL NETWORK

Consider the two layer linear neural network optimization problem

$$\underset{\mathbf{W}}{\text{minimize}} \frac{1}{2} \|\mathbf{W}_2 \mathbf{W}_1 \mathbf{X} - \mathbf{Y}\|^2 \quad (11)$$

where $\mathbf{W}_2 \in \mathbb{R}^{d_2 \times d_1}$ and $\mathbf{W}_1 \in \mathbb{R}^{d_1 \times d_0}$ are weight matrices, $\mathbf{X} \in \mathbb{R}^{d_0 \times n}$ is the input data, and $\mathbf{Y} \in \mathbb{R}^{d_2 \times n}$ is the target training data. Using our transformation, the corresponding auxiliary optimization problem can be written as

$$\begin{aligned} \underset{\mathbf{Z}}{\text{minimum}} & \quad \frac{1}{2} \|\mathbf{Z} \mathbf{X} - \mathbf{Y}\|^2 \\ \text{subject to} & \quad \text{rank}(\mathbf{Z}) \leq \min\{d_2, d_1, d_0\} \end{aligned} \quad (12)$$

(Lu & Kawaguchi, 2017, Theorem 2.2) shows that when \mathbf{X} is full rank, every local minimum of problem (12) is global. By using local openness, we first show that this result holds without any assumption on \mathbf{X} or \mathbf{Y} . The proof of Lemma 7 can be found in Appendix A.3

Lemma 7. *Every local minimum of problem (12) is global.*

Lemma 7 uses local openness to simplify the proof of (Lu & Kawaguchi, 2017, Theorem 2.2) and relax the full rankness assumption on \mathbf{X} . In another related work, (Kawaguchi, 2016, Theorem 2.3) shows that when $\mathbf{X} \mathbf{X}^T$ and $\mathbf{Y} \mathbf{X}^T$ are full rank, $d_2 \leq d_0$, and when $\mathbf{Y} \mathbf{X}^T (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{Y}^T$ has d_2 distinct eigenvalues, every local optimum is global and all saddle points are second order saddles. While the local/global equivalence result holds for deeper networks, the property that all saddles are second order does not hold in that case. Another result by (Yun et al., 2017, Theorem 2.2) shows that when $\mathbf{X} \mathbf{X}^T$, $\mathbf{Y} \mathbf{X}^T$, and $\mathbf{Y} \mathbf{X}^T (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{Y}^T$ are full rank, every local optimum of a linear deep network is global. Moreover, they provide necessary and sufficient conditions for a critical point to be a

global minimum. However, in their proof, the full rankness assumption of $\mathbf{Y}\mathbf{X}^T$ was not used in showing the result for non-degenerate critical points and thus can be relaxed in that case. In this section, without any assumptions on both \mathbf{X} and \mathbf{Y} , we reconstruct the proof that shows the latter result for 2-layer networks using local openness, and then show a similar result for the degenerate case. The result for the degenerate case holds when replacing the square loss error by a general convex loss function as we will see in Colorollary 9. The proofs of the theorem and corollary stated below can be found in Appendices A.1 and A.2, respectively.

Theorem 8. *Every local minimum of problem (11) is global. Moreover, every degenerate saddle point of problem (11) is a second order saddle.*

Corollary 9. *Let the square loss error in (11) be replaced by a general convex loss function $\ell(\cdot)$. Then every degenerate critical point is either a global minimum or a second order saddle.*

Baldi & Hornik (1989) and Srebro & Jaakkola (2003) show the same result when both \mathbf{X} and \mathbf{Y} are full row rank. Theorem 8 generalizes their results by relaxing the assumptions on both \mathbf{X} and \mathbf{Y} .

6 MULTI-LAYER LINEAR NEURAL NETWORK

Consider the training problem of multi-layer deep linear neural networks:

$$\underset{\mathbf{W}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{W}_h \cdots \mathbf{W}_1 \mathbf{X} - \mathbf{Y}\|^2. \quad (13)$$

Here $\mathbf{W} = (\mathbf{W}_i)_{i=1}^h$, $\mathbf{W}_i \in \mathbb{R}^{d_i \times d_{i-1}}$ are the weight matrices, $\mathbf{X} \in \mathbb{R}^{d_0 \times n}$ is the input training data, and $\mathbf{Y} \in \mathbb{R}^{d_h \times n}$ is the target training data. Based on our general framework, the corresponding auxiliary optimization problem is given by

$$\begin{aligned} \underset{\mathbf{Z} \in \mathbb{R}^{d_h \times n}}{\text{minimum}} \quad & \frac{1}{2} \|\mathbf{Z}\mathbf{X} - \mathbf{Y}\|^2 \\ \text{subject to} \quad & \text{rank}(\mathbf{Z}) \leq d_p \triangleq \min_{0 \leq i \leq h} d_i \end{aligned} \quad (14)$$

Paper Lu & Kawaguchi (2017) showed that when \mathbf{X} and \mathbf{Y} are full row rank, every local minimum of (13) is global. We now relax the full rankness assumption and reproduce similar results. However, as we will see, the local/global equivalence does not always follow if we relax the full rankness. In such cases, we will provide detailed counter examples. Before proceeding to the proof we define the following mapping:

$$\mathcal{M}_{i,j}(\mathbf{W}_i, \dots, \mathbf{W}_j) : \{\mathbf{W}_i, \dots, \mathbf{W}_j\} \rightarrow \mathcal{R}_{\mathcal{M}_{i,j}} \triangleq \{\mathbf{Z} = \mathbf{W}_i \dots \mathbf{W}_j \in \mathbb{R}^{d_i \times d_{j-1}} \mid \text{rank}(\mathbf{Z}) \leq \min_{j-1 \leq l \leq i} d_l\} \quad \text{for } i > j$$

Now we state Theorem 3.1 of Lu & Kawaguchi (2017) using our notation. The proof of the lemma stated below can be found in Appendix A.4.

Lemma 10. *If \mathbf{W} is non-degenerate, then $\mathcal{M}_{h,1}(\mathbf{W}) = \mathbf{W}_h \cdots \mathbf{W}_1$ is locally open at \mathbf{W} .*

We now demonstrate our main results for this optimization problem which shows that under a set of necessary conditions, every local minimum of problem (13) is global. Although the result for the non-degenerate case directly follows from (Yun et al., 2017, Theorem 2.2), we provide in Lemma 11 a more intuitive proof that uses local openness of \mathcal{M} . Moreover, Theorem 12 extends the result to degenerate critical points.

Lemma 11. *Every non-degenerate local minimum of (13) is global minimum.*

Proof. Suppose $\mathbf{W} = (\mathbf{W}_h, \dots, \mathbf{W}_1)$ is a non-degenerate local minimum. Then it follows by Lemma 10 that $\mathcal{M}_{h,1}$ is locally open at \mathbf{W} . Then by Lemma 1, $\mathbf{Z} = \mathcal{M}_h(\mathbf{W}_h, \dots, \mathbf{W}_1)$ is a local optimum of problem (14) which is in fact global by Lemma 7. \square

We now state the desired result for degenerate critical points.

Theorem 12. *Let $p_1^* \triangleq \underset{0 \leq i \leq h}{\text{argmin}} d_i$ and $p_2^* \triangleq \underset{j \neq p_1^*}{\text{argmin}} d_j$. If $d_{p_2^*} < \min\{d_h, d_0\}$, we can find a rank deficient \mathbf{Y} such that problem (30) has a local minimum that is not global. Otherwise, given any \mathbf{X} and \mathbf{Y} , every local minimum of problem (30) is a global minimum.*

The proof of Theorem 12 can be found in Appendix C.

REFERENCES

- M. Balcerzak, A. Wachowicz, and W. Wilczyński. Multiplying balls in the space of continuous functions on $[0, 1]$. *Studia Mathematica*, 170:203–209, 2005.
- M. Balcerzak, A. Majchrzycki, and A. Wachowicz. Openness of multiplication in some function spaces. *Taiwanese J. Math*, 17: 1115–1126, 2013.
- P. Baldi and K. Hornik. Neural networks and principal component analysis: Learning from examples without local minima. *Neural networks*, 2(1):53–58, 1989.
- E. Behrends. Products of n open subsets in the space of continuous functions on $[0, 1]$. *Studia Mathematica*, 204:73–95, 2011.
- E. Behrends. Where is matrix multiplication locally open? *Linear Algebra and its Applications*, 517:167–176, 2017.
- S. Bhojanapalli, B. Neyshabur, and N. Srebro. Global optimality of local search for low rank matrix recovery. In *Advances in Neural Information Processing Systems*, pp. 3873–3881, 2016.
- A. Blum and R. L. Rivest. Training a 3-node neural network is np-complete. In *Advances in neural information processing systems*, pp. 494–501, 1989.
- A. Choromanska, M. Henaff, M. Mathieu, G. B. Arous, and Y. LeCun. The loss surfaces of multilayer networks. In *Artificial Intelligence and Statistics*, pp. 192–204, 2015.
- R. Ge, J. D. Lee, and T. Ma. Matrix completion has no spurious local minimum. In *Advances in Neural Information Processing Systems*, pp. 2973–2981, 2016.
- I. Goodfellow and A. Courville. Deep learning. *Book in preparation for MIT Press, Cambridge*, 2016.
- M. Hardt and T. Ma. Identity matters in deep learning. *arXiv preprint arXiv:1611.04231*, 2016.
- C. Horowitz. An elementary counterexample to the open mapping principle for bilinear maps. *Proceedings of the American Mathematical Society*, 53(2):293–294, 1975.
- K. Kawaguchi. Deep learning without poor local minima. In *Advances in Neural Information Processing Systems*, pp. 586–594, 2016.
- J. D. Lee, M. Simchowitz, M. I. Jordan, and B. Recht. Gradient descent only converges to minimizers. In *Conference on Learning Theory*, pp. 1246–1257, 2016.
- H. Lu and K. Kawaguchi. Depth creates no bad local minima. *arXiv preprint arXiv:1702.08580*, 2017.
- Q. Nguyen and M. Hein. The loss surface of deep and wide neural networks. *arXiv preprint arXiv:1704.08045*, 2017.
- D. Park, A. Kyriillidis, C. Caramanis, and S. Sanghavi. Non-square matrix sensing without spurious local minima via the burer-monteiro approach. *arXiv preprint arXiv:1609.03240*, 2016.
- W. Rudin. Functional analysis, mcgraw-hill series in higher mathematics. 1973.
- M. Soltanolkotabi, A. Javanmard, and J. D. Lee. Theoretical insights into the optimization landscape of over-parameterized shallow neural networks. *arXiv preprint arXiv:1707.04926*, 2017.
- D. Soudry and Y. Carmon. No bad local minima: Data independent training error guarantees for multilayer neural networks. *arXiv preprint arXiv:1605.08361*, 2016.
- N. Srebro and T. Jaakkola. Weighted low-rank approximations. In *Proceedings of the 20th International Conference on Machine Learning (ICML-03)*, pp. 720–727, 2003.
- R. Sun. *Matrix completion via nonconvex factorization: Algorithms and theory*. PhD thesis, University of Minnesota, 2015.
- L. Wang, X. Zhang, and Q. Gu. A unified computational and statistical framework for nonconvex low-rank matrix estimation. *arXiv preprint arXiv:1610.05275*, 2016.
- B. Xie, Y. Liang, and L. Song. Diverse neural network learns true target functions. In *Artificial Intelligence and Statistics*, pp. 1216–1224, 2017.
- C. Yun, S. Sra, and A. Jadbabaie. Global optimality conditions for deep neural networks. *arXiv preprint arXiv:1707.02444*, 2017.

Appendix

A PROOF OF THEOREM 8, COROLLARY 9 AND LEMMA 10

A.1 PROOF OF THE THEOREM 8

Proof. The proof for the degenerate case is done by constructing a descent direction if the point is critical but not global. Let $(\bar{\mathbf{W}}_2, \bar{\mathbf{W}}_1)$ be a degenerate critical point, i.e. $\text{rank}(\bar{\mathbf{W}}_2 \bar{\mathbf{W}}_1) < \min\{d_2, d_1, d_0\}$. Then, based on the dimensions of $d_0, d_1,$ and $d_2,$ we have one of the following cases:

$$\begin{aligned} d_2 < d_1 \text{ then } \exists \mathbf{b} \neq \mathbf{0} \text{ such that } \mathbf{b} \in \mathcal{N}(\bar{\mathbf{W}}_2) \\ d_0 < d_1 \text{ then } \exists \mathbf{b} \neq \mathbf{0} \text{ such that } \mathbf{b} \in \mathcal{N}(\bar{\mathbf{W}}_1^T) \\ d_1 \leq d_2, \text{ and } d_1 \leq d_0 \text{ then either } \bar{\mathbf{W}}_2 \text{ is rank deficient and } \exists \mathbf{b} \neq \mathbf{0} \text{ such that } \mathbf{b} \in \mathcal{N}(\bar{\mathbf{W}}_2) \text{ or} \\ \bar{\mathbf{W}}_1 \text{ is rank deficient and } \exists \mathbf{b} \neq \mathbf{0} \text{ such that } \mathbf{b} \in \mathcal{N}(\bar{\mathbf{W}}_1^T) \end{aligned}$$

So in all cases either $\mathcal{N}(\bar{\mathbf{W}}_2) \neq \emptyset$ or $\mathcal{N}(\bar{\mathbf{W}}_1^T) \neq \emptyset$. Also, let $\Delta = \bar{\mathbf{W}}_2 \bar{\mathbf{W}}_1 \mathbf{X} - \mathbf{Y}$. If $\Delta \mathbf{X}^T = \mathbf{0}$, then by convexity of the square loss error function, the point $(\bar{\mathbf{W}}_2, \bar{\mathbf{W}}_1)$ is a global minimum of (11). Else, there exists (i, j) such that $\langle \mathbf{X}_{i,:}, \Delta_{j,:} \rangle \neq 0$. We now use first and second order optimality conditions to construct a descent direction when the current critical point is not global.

First order optimality condition: By considering perturbation in the directions $\mathbf{A} \in \mathbb{R}^{d_2 \times d_1}$ and $\mathbf{B} \in \mathbb{R}^{d_1 \times d_0}$ for the optimization problem

$$\underset{t}{\text{minimize}} \frac{1}{2} \|(\mathbf{W}_2 + t\mathbf{A})(\mathbf{W}_1 + t\mathbf{B})\mathbf{X} - \mathbf{Y}\|^2 \quad (15)$$

we obtain

$$\langle \mathbf{A} \bar{\mathbf{W}}_1 \mathbf{X} + \bar{\mathbf{W}}_2 \mathbf{B} \mathbf{X}, \Delta \rangle = 0, \quad \forall \mathbf{A} \in \mathbb{R}^{d_2 \times d_1}, \mathbf{B} \in \mathbb{R}^{d_1 \times d_0}$$

Second order optimality condition:

$$2 \langle \mathbf{A} \mathbf{B} \mathbf{X}, \Delta \rangle + \|\mathbf{A} \bar{\mathbf{W}}_1 \mathbf{X} + \bar{\mathbf{W}}_2 \mathbf{B} \mathbf{X}\|^2 \geq 0 \quad \forall \mathbf{A} \in \mathbb{R}^{d_2 \times d_1}, \mathbf{B} \in \mathbb{R}^{d_1 \times d_0}$$

Suppose $(\bar{\mathbf{W}}_2, \bar{\mathbf{W}}_1)$ is a critical point and there exists $\mathbf{b} \neq \mathbf{0}, \mathbf{b} \in \mathcal{N}(\bar{\mathbf{W}}_2)$. We define

$$\mathbf{B}_{:,l} \triangleq \begin{cases} \alpha \mathbf{b} & \text{if } l = i, \\ \mathbf{0} & \text{otherwise} \end{cases} \quad \mathbf{A}_{l,:} \triangleq \begin{cases} \mathbf{b}^T & \text{if } l = j, \\ \mathbf{0} & \text{otherwise} \end{cases}$$

where α is a scalar constant. Then, using the second order optimality condition, for $c = \|\mathbf{A} \bar{\mathbf{W}}_1 \mathbf{X}\|^2$, we get

$$\alpha \underbrace{\|\mathbf{b}\|^2}_{\neq 0} \underbrace{\langle \mathbf{X}_{i,:}, \Delta_{j,:} \rangle}_{\neq 0} + c \geq 0$$

Since this is true for every value of α , \mathbf{b} should be zero which contradicts the assumption on the choice of \mathbf{b} . Hence $\mathcal{N}(\bar{\mathbf{W}}_2) = \emptyset$.

Similarly, suppose $(\bar{\mathbf{W}}_2, \bar{\mathbf{W}}_1)$ is a critical point and there exists $\mathbf{a}^T \neq \mathbf{0}, \mathbf{a}^T \in \mathcal{N}(\bar{\mathbf{W}}_1^T)$. Let

$$\mathbf{A}_{l,:} \triangleq \begin{cases} \alpha \mathbf{a}^T & \text{if } l = j, \\ \mathbf{0} & \text{otherwise} \end{cases} \quad \mathbf{B}_{:,l} \triangleq \begin{cases} \mathbf{a} & \text{if } l = i, \\ \mathbf{0} & \text{otherwise} \end{cases}$$

where α is a constant. Then, for $c = \|\bar{\mathbf{W}}_2 \mathbf{B} \mathbf{X}\|^2$, we get

$$\alpha \underbrace{\|\mathbf{a}\|^2}_{\neq 0} \underbrace{\langle \mathbf{X}_{i,:}, \Delta_{j,:} \rangle}_{\neq 0} + c \geq 0$$

Using the same argument, we can show that $(\bar{\mathbf{W}}_2, \bar{\mathbf{W}}_1)$ is a second order saddle point of problem (11).

We now show the result for the non-degenerate case. Let $(\bar{\mathbf{W}}_2, \bar{\mathbf{W}}_1)$ be a non-degenerate local minimum, i.e. $\text{rank}(\bar{\mathbf{W}}_2 \bar{\mathbf{W}}_1) = \min\{d_2, d_1, d_0\}$. Then it follows by Lemma 10 that the matrix multiplication $\mathcal{M}(\cdot, \cdot)$ is locally open at $(\bar{\mathbf{W}}_2, \bar{\mathbf{W}}_1)$. Then by Observation 1, $\mathbf{Z} = \bar{\mathbf{W}}_2 \bar{\mathbf{W}}_1$ is a local optimum of problem (12) which is in fact global by Lemma 7. □

A.2 PROOF OF COROLLARY 9

Proof. We follow the same steps used in the proof of Theorem 8 to show the result.

First order optimality condition: By considering perturbation in the directions $\mathbf{A} \in \mathbb{R}^{d_2 \times d_1}$ and $\mathbf{B} \in \mathbb{R}^{d_1 \times d_0}$ for the optimization problem

$$\underset{t}{\text{minimize}} \ell((\mathbf{W}_2 + t\mathbf{A})(\mathbf{W}_1 + t\mathbf{B})\mathbf{X} - \mathbf{Y}) \quad (16)$$

we obtain

$$\langle \mathbf{A}\bar{\mathbf{W}}_1\mathbf{X} + \bar{\mathbf{W}}_2\mathbf{B}\mathbf{X}, \nabla \ell(\bar{\mathbf{W}}_2\bar{\mathbf{W}}_1\mathbf{X} - \mathbf{Y}) \rangle = 0 \quad \forall \mathbf{A} \in \mathbb{R}^{d_2 \times d_1}, \mathbf{B} \in \mathbb{R}^{d_1 \times d_0}$$

Second order optimality condition:

$$2\langle \mathbf{A}\mathbf{B}\mathbf{X}, \nabla \ell(\bar{\mathbf{W}}_2\bar{\mathbf{W}}_1\mathbf{X} - \mathbf{Y}) \rangle + h(\mathbf{A}\bar{\mathbf{W}}_1\mathbf{X}, \bar{\mathbf{W}}_2\mathbf{B}\mathbf{X}, \bar{\mathbf{W}}_2\bar{\mathbf{W}}_1\mathbf{X}) \geq 0 \quad \forall \mathbf{A} \in \mathbb{R}^{d_2 \times d_1}, \mathbf{B} \in \mathbb{R}^{d_1 \times d_0}$$

where $h(\cdot)$ is a function that has a tensor representation. But we only need to know that it is a function of $\mathbf{A}\bar{\mathbf{W}}_1\mathbf{X}$, $\bar{\mathbf{W}}_2\mathbf{B}\mathbf{X}$, and $\bar{\mathbf{W}}_2\bar{\mathbf{W}}_1\mathbf{X}$.

If $\nabla \ell(\bar{\mathbf{W}}_2\bar{\mathbf{W}}_1\mathbf{X} - \mathbf{Y})\mathbf{X}^T$, then by convexity of $\ell(\cdot)$, $(\bar{\mathbf{W}}_2, \bar{\mathbf{W}}_1)$ is a global minimum. Otherwise, there exists (i, j) such that $\langle \mathbf{X}_{i,:}, (\nabla \ell(\bar{\mathbf{W}}_2\bar{\mathbf{W}}_1\mathbf{X} - \mathbf{Y}))_{j,:} \rangle \neq 0$. Using the same former argument in proof of Theorem 8, we choose \mathbf{A} and \mathbf{B} such that $h(\mathbf{A}\bar{\mathbf{W}}_1\mathbf{X}, \bar{\mathbf{W}}_2\mathbf{B}\mathbf{X}, \bar{\mathbf{W}}_2\bar{\mathbf{W}}_1\mathbf{X})$ is some constant that does not depend on α , and $\langle \mathbf{A}\mathbf{B}\mathbf{X}, \nabla \ell(\bar{\mathbf{W}}_2\bar{\mathbf{W}}_1\mathbf{X} - \mathbf{Y}) \rangle = \alpha \underbrace{\langle \mathbf{X}_{i,:}, (\nabla \ell(\bar{\mathbf{W}}_2\bar{\mathbf{W}}_1\mathbf{X} - \mathbf{Y}))_{j,:} \rangle}_{\neq 0}$. Then by proper choice of α we show

that the point $(\bar{\mathbf{W}}_2, \bar{\mathbf{W}}_1)$ is a second order saddle point. □

A.3 PROOF OF LEMMA 7

Proof. Let $r_{\mathbf{X}} = \text{rank}(\mathbf{X})$ and $\mathbf{U}_{\mathbf{X}}\boldsymbol{\Sigma}_{\mathbf{X}}\mathbf{V}_{\mathbf{X}}^T$ with $\mathbf{U}_{\mathbf{X}} \in \mathbb{R}^{d_0 \times d_0}$, $\boldsymbol{\Sigma}_{\mathbf{X}} \in \mathbb{R}^{d_0 \times n}$, $\mathbf{V}_{\mathbf{X}} \in \mathbb{R}^{n \times n}$ be a singular value decomposition of \mathbf{X} . Then

$$\begin{aligned} \|\mathbf{Z}\mathbf{X} - \mathbf{Y}\|^2 &= \|\mathbf{Z}\mathbf{U}_{\mathbf{X}} \left[(\boldsymbol{\Sigma}_{\mathbf{X}})_{:,1:r_{\mathbf{X}}} \mid \mathbf{0} \right] - \mathbf{Y}\mathbf{V}_{\mathbf{X}}\|^2 \\ &= \|\mathbf{Z}\mathbf{U}_{\mathbf{X}}(\boldsymbol{\Sigma}_{\mathbf{X}})_{:,1:r_{\mathbf{X}}} - (\mathbf{Y}\mathbf{V}_{\mathbf{X}})_{:,1:r_{\mathbf{X}}}\|^2 + \underbrace{\|(\mathbf{Y}\mathbf{V}_{\mathbf{X}})_{:,r_{\mathbf{X}}+1:n}\|^2}_{\text{constant in problem (12)}}. \end{aligned}$$

Since $\mathbf{U}_{\mathbf{X}}(\boldsymbol{\Sigma}_{\mathbf{X}})_{:,1:r_{\mathbf{X}}}$ is full column rank, then the linear mapping $\mathbf{Z}\mathbf{U}_{\mathbf{X}}(\boldsymbol{\Sigma}_{\mathbf{X}})_{:,1:r_{\mathbf{X}}}$ is open, and

$$\text{rank}(\mathbf{Z}\mathbf{U}_{\mathbf{X}}(\boldsymbol{\Sigma}_{\mathbf{X}})_{:,1:r_{\mathbf{X}}}) \leq \min\{\text{rank}(\mathbf{Z}), r_{\mathbf{X}}\} = \min\{d_2, d_1, d_0, r_{\mathbf{X}}\}.$$

Consequently, every local minimum in problem (12) corresponds to a local minimum in problem

$$\begin{aligned} &\underset{\bar{\mathbf{Z}} \in \mathbb{R}^{d_2 \times r_{\mathbf{X}}}}{\text{minimum}} \quad \frac{1}{2} \|\bar{\mathbf{Z}} - \bar{\mathbf{Y}}\|^2 \\ &\text{subject to} \quad \text{rank}(\bar{\mathbf{Z}}) \leq \min\{d_2, d_1, d_0, r_{\mathbf{X}}\} \end{aligned} \quad (17)$$

where $\bar{\mathbf{Y}} = (\mathbf{Y}\mathbf{V}_{\mathbf{X}})_{:,1:r_{\mathbf{X}}}$. The result follows using (Lu & Kawaguchi, 2017, Theorem 2.2). □

A.4 PROOF OF LEMMA 10

Proof. We construct a proof by induction on h to show the desired result. When $h = 2$, we either have $d_1 < \min\{d_2, d_0\}$ or $d_1 \geq \min\{d_2, d_0\}$. In the first case,

$$d_1 = \text{rank}(\mathbf{W}_2\mathbf{W}_1) \leq \text{rank}(\mathbf{W}_1) \leq d_1 \Leftrightarrow \text{rank}(\mathbf{W}_1) = d_1, \quad \text{and} \quad d_1 = \text{rank}(\mathbf{W}_2\mathbf{W}_1) \leq \text{rank}(\mathbf{W}_2) \leq d_1 \Leftrightarrow \text{rank}(\mathbf{W}_2) = d_1.$$

Since \mathbf{W}_1 is full column rank and \mathbf{W}_2 is full column rank, then by Theorem 5, $\mathcal{M}_{2,1}(\cdot)$ is locally open at $(\mathbf{W}_2, \mathbf{W}_1)$. In the second case, either

$$d_2 = \text{rank}(\mathbf{W}_2 \mathbf{W}_1) \leq \text{rank}(\mathbf{W}_2) \leq d_2 \Leftrightarrow \text{rank}(\mathbf{W}_2) = d_2,$$

or

$$d_0 = \text{rank}(\mathbf{W}_2 \mathbf{W}_1) \leq \text{rank}(\mathbf{W}_1) \leq d_0 \Leftrightarrow \text{rank}(\mathbf{W}_1) = d_0$$

Thus, either \mathbf{W}_2 is full row rank or \mathbf{W}_1 is full row rank, then by Proposition 4, $\mathcal{M}_{2,1}(\cdot)$ is locally open at $(\mathbf{W}_2, \mathbf{W}_1)$.

Now assume the result holds for the product of h matrices $\mathcal{M}_{h,1}(\mathbf{W})$, we show it is true for $\mathcal{M}_{h+1,1}(\mathbf{W})$.

Since

$$d_p = \text{rank}(\mathbf{W}_h \dots \mathbf{W}_1) \leq \text{rank}(\mathbf{W}_{p+1} \mathbf{W}_p) \leq d_p \Leftrightarrow \text{rank}(\mathbf{W}_{p+1} \mathbf{W}_p) = d_p,$$

then using Proposition 4, we get $\mathcal{M}_{p+1,p}(\cdot)$ is locally open at $(\mathbf{W}_{p+1}, \mathbf{W}_p)$. So we can replace $\mathbf{W}_{p+1} \mathbf{W}_p$ by a new matrix \mathbf{Z}_p with rank d_p . Then by induction hypothesis, the product mapping $\mathcal{M}_{h+1,1} = \mathbf{W}_{h+1} \dots \mathbf{W}_{p+2} \mathbf{Z}_p \mathbf{W}_{p-1} \dots \mathbf{W}_1$ is locally open at \mathbf{W} . Since the composition of locally open maps is locally open, the result follows. \square

B PROOF OF THEOREM 5

In this section, we prove Theorem 5. This theorem provides a complete characterization of points $(\bar{\mathbf{W}}_1, \bar{\mathbf{W}}_2) \in \mathbb{R}^{m \times k} \times \mathbb{R}^{k \times n}$ for which the matrix multiplication mapping $\mathcal{M}(\cdot, \cdot)$ is locally open for the case of $k < \min\{m, n\}$. In particular, we show that if $\text{rank}(\bar{\mathbf{W}}_1) \neq \text{rank}(\bar{\mathbf{W}}_2)$, $\mathcal{M}(\cdot, \cdot)$ is not locally open at $(\bar{\mathbf{W}}_1, \bar{\mathbf{W}}_2)$. Else, if $\text{rank}(\bar{\mathbf{W}}_1) = \text{rank}(\bar{\mathbf{W}}_2)$, then the following statements are equivalent:

- i) $\exists \tilde{\mathbf{W}}_1 \in \mathbb{R}^{m \times k}$ such that $\tilde{\mathbf{W}}_1 \bar{\mathbf{W}}_2 = \mathbf{0}$ and $\bar{\mathbf{W}}_1 + \tilde{\mathbf{W}}_1$ is full column rank.
- ii) $\exists \tilde{\mathbf{W}}_2 \in \mathbb{R}^{k \times n}$ such that $\bar{\mathbf{W}}_1 \tilde{\mathbf{W}}_2 = \mathbf{0}$ and $\bar{\mathbf{W}}_2 + \tilde{\mathbf{W}}_2$ is full row rank.
- iii) $\dim(\mathcal{N}(\bar{\mathbf{W}}_1) \cap \mathcal{C}(\bar{\mathbf{W}}_2)) = 0$.
- iv) $\dim(\mathcal{N}(\bar{\mathbf{W}}_2^T) \cap \mathcal{C}(\bar{\mathbf{W}}_1^T)) = 0$.
- v) $\mathcal{M}(\cdot, \cdot)$ is locally open at $(\bar{\mathbf{W}}_1, \bar{\mathbf{W}}_2)$ in its range $\mathcal{R}_{\mathcal{M}} = \{\mathbf{Z} \in \mathbb{R}^{m \times n} \text{ with } \text{rank}(\mathbf{Z}) \leq \min\{m, k, n\}\}$.

To prove this result, we first show that the local openness of $\mathcal{M}(\cdot, \cdot)$ at $(\mathbf{W}_1, \mathbf{W}_2)$ is equivalent to the local openness of $\mathcal{M}(\cdot, \cdot)$ at $(\mathbf{U}^T \mathbf{W}_1, \mathbf{W}_2 \mathbf{V})$ where the columns of $\mathbf{U} \in \mathbb{R}^{m \times m}$ and the columns of $\mathbf{V} \in \mathbb{R}^{n \times n}$ are the left and right singular vectors of the product $\mathbf{W}_1 \mathbf{W}_2$, respectively. This allows us to focus our study on the local openness of the mapping \mathcal{M} to matrix pairs whose product is a diagonal matrix. We then show in Lemma 16 that when $\text{rank}(\mathbf{W}_1) = \text{rank}(\mathbf{W}_2)$, the statements *i*, *ii*, *iii*, and *iv* are equivalent. Finally, we show in Proposition 18, that these conditions hold if and only if the mapping $\mathcal{M}(\cdot, \cdot)$ is locally open at $(\mathbf{W}_1, \mathbf{W}_2)$. Before proceeding we state and prove Lemma 13 that will be used later in the proof.

Lemma 13. *Let $\mathbf{V} \in \mathbb{R}^{m \times n}$ be a matrix with $\text{rank}(\mathbf{V}) = r < m$. Then there exist an index set $\mathcal{B} = \{i_1, \dots, i_r\} \subseteq \{1, \dots, m\}$ and a matrix $\mathbf{A} \in \mathbb{R}^{(m-r) \times r}$ such that*

$$\|\mathbf{A}\|_{\infty} = \max_{i,j} |\mathbf{A}_{ij}| \leq 2^{m-r-1} \quad \text{and} \quad \mathbf{V}_{\mathcal{B}^c} = \mathbf{A} \mathbf{V}_{\mathcal{B}},$$

where $\mathbf{V}_{\mathcal{B}} \in \mathbb{R}^{r \times n}$ is a matrix with rows $\{\mathbf{V}_{i,:}\}_{i \in \mathcal{B}}$ and $\mathbf{V}_{\mathcal{B}^c} \in \mathbb{R}^{(m-r) \times n}$ is a matrix with rows $\{\mathbf{V}_{i,:}\}_{i \in \mathcal{B}^c}$.

Notice that in the above lemma, the bound on the norm of matrix \mathbf{A} is independent of the dimension n and it also does not depend on the choice of matrix \mathbf{V} .

Proof. To ease the notation, we denote the i^{th} row of V by v_i . We use induction on m to show that there exists a basis $\mathcal{B} = \{i_1, \dots, i_r\}$ and a vector $\mathbf{a}_j \in \mathbb{R}^r$ such that $\forall j \in \mathcal{B}^c$,

$$\mathbf{v}_j = \sum_{i \in \mathcal{B}} \mathbf{a}_{j,i} \mathbf{v}_i \quad \text{with } |\mathbf{a}_{j,i}| \leq 2^{m-r-1} \quad \forall i \in \mathcal{B}.$$

• *Induction Base Case* $m = r + 1$: Without loss of generality, assume $\mathcal{B} = \{1, \dots, r\}$. Since the case of $\mathbf{v}_{r+1} = 0$ trivially holds, we consider $\mathbf{v}_{r+1} \neq 0$. By the property of basis, there exists a non-zero vector $\mathbf{a}_{r+1} \in \mathbb{R}^r$ such that $\mathbf{v}_{r+1} = \sum_{i=1}^r \mathbf{a}_{r+1,i} \mathbf{v}_i$.

Let $i^* = \arg \max_{i \in \mathcal{B}} |\mathbf{a}_{r+1,i}|$. If $|\mathbf{a}_{r+1,i^*}| \leq 1$, then the induction hypothesis is true. Otherwise, when $|\mathbf{a}_{r+1,i^*}| > 1$, we have

$$\begin{aligned} \mathbf{v}_{i^*} &= \frac{1}{\underbrace{\mathbf{a}_{r+1,i^*}}_{\bar{\mathbf{a}}_{r+1,r+1}}} \mathbf{v}_{r+1} - \sum_{i=1; i \neq i^*}^r \frac{\mathbf{a}_{r+1,i}}{\underbrace{\mathbf{a}_{r+1,i^*}}_{\bar{\mathbf{a}}_{r+1,i}}} \mathbf{v}_i \\ &= \sum_{i \in \mathcal{B}^*} \bar{\mathbf{a}}_{r+1,i} \mathbf{v}_i \quad \text{where } \mathcal{B}^* = (\mathcal{B} \cup \{r+1\}) \setminus \{i^*\}, \end{aligned}$$

i.e., we remove the item i^* from \mathcal{B} and include the item $r+1$ instead. Since $|\bar{\mathbf{a}}_{r+1,i}| \leq 1$, the induction base holds.

• *Inductive Step*: Assume the induction hypothesis is true for $m > r$, we show it is also true for $m+1$. Without loss of generality we can assume that $\mathcal{B} = \{1, \dots, r\}$. By induction hypothesis,

$$\mathbf{v}_j = \sum_{i=1}^r \mathbf{a}_{j,i} \mathbf{v}_i \quad \text{with } |\mathbf{a}_{j,i}| \leq 2^{m-r-1}, \quad \forall j = \{r+1, \dots, m\}.$$

Since the case of $\mathbf{v}_{m+1} = 0$ trivially holds, we consider $\mathbf{v}_{m+1} \neq 0$. Since \mathcal{B} is a basis, there exists $\mathbf{a}_{m+1} \neq 0$ such that $\mathbf{v}_{m+1} = \sum_{i=1}^r \mathbf{a}_{m+1,i} \mathbf{v}_i$. Let $i^* = \arg \max_{i \in \mathcal{B}} |\mathbf{a}_{m+1,i}|$. If $|\mathbf{a}_{m+1,i^*}| \leq 2^{m-r}$, the induction step is done. Otherwise, for the case of $|\mathbf{a}_{m+1,i^*}| > 2^{m-r}$, we have

$$\begin{aligned} \mathbf{v}_{i^*} &= \frac{1}{\underbrace{\mathbf{a}_{m+1,i^*}}_{\bar{\mathbf{a}}_{m+1,m+1}}} \mathbf{v}_{m+1} - \sum_{i=1; i \neq i^*}^r \frac{\mathbf{a}_{m+1,i}}{\underbrace{\mathbf{a}_{m+1,i^*}}_{\bar{\mathbf{a}}_{m+1,i}}} \mathbf{v}_i \\ &= \sum_{i \in \mathcal{B}^*} \bar{\mathbf{a}}_{m+1,i} \mathbf{v}_i, \end{aligned}$$

where $\mathcal{B}^* = (\mathcal{B} \cup \{m+1\}) \setminus \{i^*\}$ and clearly $|\bar{\mathbf{a}}_{m+1,i}| \leq 1, \forall i \in \mathcal{B}^*$ according to the definition of i^* . For all $j \in \{r+1, \dots, m\}$

$$\begin{aligned} \mathbf{v}_j &= \sum_{i=1; i \neq i^*}^r \mathbf{a}_{j,i} \mathbf{v}_i + \mathbf{a}_{j,i^*} \mathbf{v}_{i^*} \\ &= \sum_{i=1; i \neq i^*}^r \mathbf{a}_{j,i} \mathbf{v}_i + \frac{\mathbf{a}_{j,i^*}}{\bar{\mathbf{a}}_{m+1,i^*}} \mathbf{v}_{m+1} - \sum_{i=1; i \neq i^*}^r \frac{\mathbf{a}_{m+1,i} \mathbf{a}_{j,i^*}}{\bar{\mathbf{a}}_{m+1,i^*}} \mathbf{v}_i \\ &= \sum_{i=1; i \neq i^*}^r \left(\underbrace{\mathbf{a}_{j,i} - \frac{\mathbf{a}_{j,i^*} \mathbf{a}_{m+1,i}}{\bar{\mathbf{a}}_{m+1,i^*}}}_{\bar{\mathbf{a}}_{j,i}} \right) \mathbf{v}_i + \underbrace{\frac{\mathbf{a}_{j,i^*}}{\bar{\mathbf{a}}_{m+1,i^*}}}_{\bar{\mathbf{a}}_{j,m+1}} \mathbf{v}_{m+1} \\ &= \sum_{i \in \mathcal{B}^*} \bar{\mathbf{a}}_{j,i} \mathbf{v}_i. \end{aligned}$$

It remains to show that $|\bar{a}_{j,i}| \leq 2^{m-r}$ for all $i \in \mathcal{B}^*$, $j \in \{r+1, \dots, m\}$. Let us first consider $i \in \mathcal{B}^* \setminus \{m+1\}$ and $j \in \{r+1, \dots, m\}$:

$$\begin{aligned} |\bar{a}_{j,i}| &\leq |\mathbf{a}_{j,i}| + \left| \mathbf{a}_{j,i^*} \frac{\mathbf{a}_{m+1,i}}{\mathbf{a}_{m+1,i^*}} \right| && \text{by triangular inequality} \\ &\leq 2^{m-r-1} + 2^{m-r-1} \left| \frac{\mathbf{a}_{m+1,i}}{\mathbf{a}_{m+1,i^*}} \right| && \text{by induction hypothesis} \\ &\leq 2^{m-r} && \text{by definition of } i^*. \end{aligned}$$

For $i = m+1$, $|\bar{a}_{j,m+1}| = \left| \frac{\mathbf{a}_{j,i^*}}{\mathbf{a}_{m+1,i^*}} \right| \leq \left| \frac{2^{m-r-1}}{\mathbf{a}_{m+1,i^*}} \right| \leq 2^{m-r}$. This concludes the inductive step and completes our proof. \square

Lemma 14. *Let $\mathbf{W}_1 \in \mathbb{R}^{m \times k}$ and $\mathbf{W}_2 \in \mathbb{R}^{k \times n}$. Assume further that $\mathbf{W}_1 \mathbf{W}_2 = \mathbf{U} \Sigma \mathbf{V}^T$ is a singular value decomposition of the matrix product $\mathbf{W}_1 \mathbf{W}_2$ with $\mathbf{U} \in \mathbb{R}^{m \times m}$, $\mathbf{V} \in \mathbb{R}^{n \times n}$, and $\Sigma \in \mathbb{R}^{m \times n}$. Then*

$$\mathcal{M}(\cdot, \cdot) \text{ is locally open at } (\mathbf{W}_1, \mathbf{W}_2) \Leftrightarrow \mathcal{M}(\cdot, \cdot) \text{ is locally open at } (\mathbf{U}^T \mathbf{W}_1, \mathbf{W}_2 \mathbf{V}).$$

Proof. We first show the direction “ \Rightarrow ”. Suppose $\mathcal{M}(\cdot, \cdot)$ is locally open at $(\mathbf{W}_1, \mathbf{W}_2)$, then by definition of local openness, for any given $\epsilon > 0$, there exists $\delta > 0$ such that

$$\mathbb{B}_\delta(\mathbf{W}_1 \mathbf{W}_2) \cap \mathcal{R}_{\mathcal{M}} \subseteq \mathcal{M}(\mathbb{B}_\epsilon(\mathbf{W}_1), \mathbb{B}_\epsilon(\mathbf{W}_2)) = \{(\mathbf{W}_1 + \mathbf{W}_1^\epsilon)(\mathbf{W}_2 + \mathbf{W}_2^\epsilon) \mid \|\mathbf{W}_1^\epsilon\| \leq \epsilon, \|\mathbf{W}_2^\epsilon\| \leq \epsilon\}. \quad (18)$$

We now show that

$$\mathbb{B}_\delta(\mathbf{U}^T \mathbf{W}_1 \mathbf{W}_2 \mathbf{V}) \cap \mathcal{R}_{\mathcal{M}} \subseteq \mathcal{M}(\mathbb{B}_\epsilon(\mathbf{U}^T \mathbf{W}_1), \mathbb{B}_\epsilon(\mathbf{W}_2 \mathbf{V})).$$

Consider $\tilde{\Sigma} \in \mathbb{B}_\delta(\mathbf{U}^T \mathbf{W}_1 \mathbf{W}_2 \mathbf{V}) \cap \mathcal{R}_{\mathcal{M}}$, i.e., $\tilde{\Sigma} = \mathbf{U}^T \mathbf{W}_1 \mathbf{W}_2 \mathbf{V} + \mathbf{R}_\delta$ with $\text{rank}(\tilde{\Sigma}) \leq \min\{m, k, n\}$ and $\|\mathbf{R}_\delta\| \leq \delta$. Since \mathbf{U} and \mathbf{V} are unitary matrices, we get $\mathbf{U} \tilde{\Sigma} \mathbf{V}^T = \mathbf{W}_1 \mathbf{W}_2 + \mathbf{U} \mathbf{R}_\delta \mathbf{V}^T$ with $\text{rank}(\mathbf{U} \tilde{\Sigma} \mathbf{V}^T) = \text{rank}(\tilde{\Sigma}) \leq \min\{m, k, n\}$ and $\|\mathbf{U} \mathbf{R}_\delta \mathbf{V}^T\| = \|\mathbf{R}_\delta\| \leq \delta$. According to (18), we have

$$\mathbf{U} \tilde{\Sigma} \mathbf{V}^T \in \mathbb{B}_\delta(\mathbf{W}_1 \mathbf{W}_2) \cap \mathcal{R}_{\mathcal{M}} \subseteq \{(\mathbf{W}_1 + \mathbf{W}_1^\epsilon)(\mathbf{W}_2 + \mathbf{W}_2^\epsilon) \mid \|\mathbf{W}_1^\epsilon\| \leq \epsilon, \|\mathbf{W}_2^\epsilon\| \leq \epsilon\}.$$

which implies,

$$\begin{aligned} \tilde{\Sigma} &\in \{(\mathbf{U}^T \mathbf{W}_1 + \mathbf{U}^T \mathbf{W}_1^\epsilon)(\mathbf{W}_2 \mathbf{V} + \mathbf{W}_2^\epsilon \mathbf{V}) \mid \|\mathbf{W}_1^\epsilon\| \leq \epsilon, \|\mathbf{W}_2^\epsilon\| \leq \epsilon\} \\ &= \{(\mathbf{U}^T \mathbf{W}_1 + \mathbf{U}^T \mathbf{W}_1^\epsilon)(\mathbf{W}_2 \mathbf{V} + \mathbf{W}_2^\epsilon \mathbf{V}) \mid \|\mathbf{U}^T \mathbf{W}_1^\epsilon\| \leq \epsilon, \|\mathbf{W}_2^\epsilon \mathbf{V}\| \leq \epsilon\}. \end{aligned}$$

Since $\tilde{\Sigma}$ was arbitrarily chosen, we get $\mathbb{B}_\delta(\mathbf{U}^T \mathbf{W}_1 \mathbf{W}_2 \mathbf{V}) \cap \mathcal{R}_{\mathcal{M}} \subseteq \mathcal{M}(\mathbb{B}_\epsilon(\mathbf{U}^T \mathbf{W}_1), \mathbb{B}_\epsilon(\mathbf{W}_2 \mathbf{V}))$.

Proving the converse direction “ \Leftarrow ” is similar. Suppose $\mathcal{M}(\cdot, \cdot)$ is locally open at $(\mathbf{U}^T \mathbf{W}_1, \mathbf{W}_2 \mathbf{V})$, then by definition of local openness, for any given $\epsilon > 0$, there exists $\delta > 0$ such that

$$\begin{aligned} \mathbb{B}_\delta(\mathbf{U}^T \mathbf{W}_1 \mathbf{W}_2 \mathbf{V}) \cap \mathcal{R}_{\mathcal{M}} &\subseteq \mathcal{M}(\mathbb{B}_\epsilon(\mathbf{U}^T \mathbf{W}_1), \mathbb{B}_\epsilon(\mathbf{W}_2 \mathbf{V})) \\ &= \{(\mathbf{U}^T \mathbf{W}_1 + \mathbf{W}_1^\epsilon)(\mathbf{W}_2 \mathbf{V} + \mathbf{W}_2^\epsilon) \mid \|\mathbf{W}_1^\epsilon\| \leq \epsilon, \|\mathbf{W}_2^\epsilon\| \leq \epsilon\}. \end{aligned} \quad (19)$$

We now show that

$$\mathbb{B}_\delta(\mathbf{W}_1 \mathbf{W}_2) \cap \mathcal{R}_{\mathcal{M}} \subseteq \mathcal{M}(\mathbb{B}_\epsilon(\mathbf{W}_1), \mathbb{B}_\epsilon(\mathbf{W}_2)).$$

Consider $\tilde{\mathbf{Z}} \in \mathbb{B}_\delta(\mathbf{W}_1 \mathbf{W}_2) \cap \mathcal{R}_{\mathcal{M}}$, i.e. $\tilde{\mathbf{Z}} = \mathbf{W}_1 \mathbf{W}_2 + \mathbf{R}_\delta$ with $\text{rank}(\tilde{\mathbf{Z}}) \leq \min\{m, k, n\}$ and $\|\mathbf{R}_\delta\| \leq \delta$. Since \mathbf{U} and \mathbf{V} are unitary matrices we get $\mathbf{U}^T \tilde{\mathbf{Z}} \mathbf{V} = \mathbf{U}^T \mathbf{W}_1 \mathbf{W}_2 \mathbf{V} + \mathbf{U}^T \mathbf{R}_\delta \mathbf{V}$ with $\text{rank}(\mathbf{U}^T \tilde{\mathbf{Z}} \mathbf{V}) = \text{rank}(\tilde{\mathbf{Z}}) \leq \min\{m, k, n\}$ and $\|\mathbf{U}^T \mathbf{R}_\delta \mathbf{V}\| = \|\mathbf{R}_\delta\| \leq \delta$. According to (19), we have

$$\mathbf{U}^T \tilde{\mathbf{Z}} \mathbf{V} \in \mathbb{B}_\delta(\mathbf{U}^T \mathbf{W}_1 \mathbf{W}_2 \mathbf{V}) \cap \mathcal{R}_{\mathcal{M}} \subseteq \mathcal{M}(\mathbb{B}_\epsilon(\mathbf{U}^T \mathbf{W}_1), \mathbb{B}_\epsilon(\mathbf{W}_2 \mathbf{V})).$$

which implies,

$$\begin{aligned} \tilde{\mathbf{Z}} &\in \{(\mathbf{W}_1 + \mathbf{U} \mathbf{W}_1^\epsilon)(\mathbf{W}_2 + \mathbf{W}_2^\epsilon \mathbf{V}^T) \mid \|\mathbf{W}_1^\epsilon\| \leq \epsilon, \|\mathbf{W}_2^\epsilon\| \leq \epsilon\} \\ &= \{(\mathbf{W}_1 + \mathbf{U} \mathbf{W}_1^\epsilon)(\mathbf{W}_2 + \mathbf{W}_2^\epsilon \mathbf{V}^T) \mid \|\mathbf{U} \mathbf{W}_1^\epsilon\| \leq \epsilon, \|\mathbf{W}_2^\epsilon \mathbf{V}^T\| \leq \epsilon\}. \end{aligned}$$

Since $\tilde{\mathbf{Z}}$ was arbitrarily chosen, we get $\mathbb{B}_\delta(\mathbf{W}_1\mathbf{W}_2) \cap \mathcal{R}_M \subseteq \mathcal{M}(\mathbb{B}_\epsilon(\mathbf{W}_1), \mathbb{B}_\epsilon(\mathbf{W}_2))$, which completes the proof. \square

Lemma 15. *Let $\mathbf{W}_1 \in \mathbb{R}^{m \times k}$ and $\mathbf{W}_2 \in \mathbb{R}^{k \times n}$. Assume further that $\mathbf{W}_1\mathbf{W}_2 = \mathbf{U}\Sigma\mathbf{V}^T$ is a singular value decomposition of the matrix product $\mathbf{W}_1\mathbf{W}_2$ with $\mathbf{U} \in \mathbb{R}^{m \times m}$, $\mathbf{V} \in \mathbb{R}^{n \times n}$, and $\Sigma \in \mathbb{R}^{m \times n}$. Define $\bar{\mathbf{W}}_1 \triangleq \mathbf{U}^T\mathbf{W}_1$ and $\bar{\mathbf{W}}_2 \triangleq \mathbf{W}_2\mathbf{V}$. Then the condition (A) below holds true if and only if the condition (B) is true. Similarly, condition (C) is true if and only if condition (D) is true.*

(A) $\exists \widehat{\mathbf{W}}_1 \in \mathbb{R}^{m \times k}$ such that $\widehat{\mathbf{W}}_1\mathbf{W}_2 = \mathbf{0}$ and $\mathbf{W}_1 + \widehat{\mathbf{W}}_1$ is full column rank.

(B) $\exists \widetilde{\mathbf{W}}_1 \in \mathbb{R}^{m \times k}$ such that $\widetilde{\mathbf{W}}_1\bar{\mathbf{W}}_2 = \mathbf{0}$ and $\bar{\mathbf{W}}_1 + \widetilde{\mathbf{W}}_1$ is full column rank.

(C) $\exists \widehat{\mathbf{W}}_2 \in \mathbb{R}^{k \times n}$ such that $\mathbf{W}_1\widehat{\mathbf{W}}_2 = \mathbf{0}$ and $\mathbf{W}_2 + \widehat{\mathbf{W}}_2$ is full row rank.

(D) $\exists \widetilde{\mathbf{W}}_2 \in \mathbb{R}^{k \times n}$ such that $\bar{\mathbf{W}}_1\widetilde{\mathbf{W}}_2 = \mathbf{0}$ and $\bar{\mathbf{W}}_2 + \widetilde{\mathbf{W}}_2$ is full row rank.

Proof. Setting $\widetilde{\mathbf{W}}_1 = \mathbf{U}^T\widehat{\mathbf{W}}_1$ and $\widetilde{\mathbf{W}}_2 = \widehat{\mathbf{W}}_2\mathbf{V}$ leads to the desired result. \square

Lemma 14 and Lemma 15 imply that for proving Theorem 5, without loss of generality, we can assume that the product $\bar{\mathbf{W}}_1\bar{\mathbf{W}}_2$ is equal to a diagonal matrix. We next show in Lemma 16 that if $k < \min\{m, n\}$ and $\text{rank}(\mathbf{W}_1) = \text{rank}(\mathbf{W}_2)$, then statements *i*, *ii*, *iii*, and *iv* in Theorem 5 are all equivalent.

Lemma 16. *Let $\mathbf{W}_1 \in \mathbb{R}^{m \times k}$, $\mathbf{W}_2 \in \mathbb{R}^{k \times n}$ with $\text{rank}(\mathbf{W}_1) = \text{rank}(\mathbf{W}_2) = r$. Assume further that $k < \min\{m, n\}$. Then, the following conditions are equivalent*

i) $\exists \widetilde{\mathbf{W}}_1 \in \mathbb{R}^{m \times k}$ such that $\widetilde{\mathbf{W}}_1\mathbf{W}_2 = \mathbf{0}$ and $\mathbf{W}_1 + \widetilde{\mathbf{W}}_1$ is full column rank.

ii) $\exists \widetilde{\mathbf{W}}_2 \in \mathbb{R}^{k \times n}$ such that $\mathbf{W}_1\widetilde{\mathbf{W}}_2 = \mathbf{0}$ and $\mathbf{W}_2 + \widetilde{\mathbf{W}}_2$ is full row rank.

iii) $\dim(\mathcal{N}(\mathbf{W}_1) \cap \mathcal{C}(\mathbf{W}_2)) = 0$.

iv) $\dim(\mathcal{N}(\mathbf{W}_2^T) \cap \mathcal{C}(\mathbf{W}_1^T)) = 0$.

Proof. To prove the desired result we show the equivalences *ii* \Leftrightarrow *iii*, and *i* \Leftrightarrow *iv*. Then we complete the proof by showing *iii* \Leftrightarrow *iv*.

We first show the direction “*ii* \Rightarrow *iii*”. Consider $\mathbf{W}_1 \in \mathbb{R}^{m \times k}$, $\mathbf{W}_2 \in \mathbb{R}^{k \times n}$ with both being rank r matrices. Suppose *ii* holds, then

$$\mathcal{C}(\widetilde{\mathbf{W}}_2) \subseteq \mathcal{N}(\mathbf{W}_1) \Rightarrow \text{rank}(\widetilde{\mathbf{W}}_2) \leq \dim(\mathcal{N}(\mathbf{W}_1)) = k - r. \quad (20)$$

Also,

$$k = \text{rank}(\mathbf{W}_2 + \widetilde{\mathbf{W}}_2) \leq \text{rank}(\mathbf{W}_2) + \text{rank}(\widetilde{\mathbf{W}}_2) = r + \text{rank}(\widetilde{\mathbf{W}}_2). \quad (21)$$

From inequalities (20) and (21), we get

$$k - r \leq \text{rank}(\widetilde{\mathbf{W}}_2) \leq k - r \Rightarrow \text{rank}(\widetilde{\mathbf{W}}_2) = k - r.$$

Note that $\dim(\mathcal{C}(\widetilde{\mathbf{W}}_2)) = \dim(\mathcal{N}(\mathbf{W}_1))$ and $\mathcal{C}(\widetilde{\mathbf{W}}_2) \subseteq \mathcal{N}(\mathbf{W}_1)$, which implies that $\mathcal{C}(\widetilde{\mathbf{W}}_2) = \mathcal{N}(\mathbf{W}_1)$.

Then since $\text{rank}(\mathbf{W}_2 + \widetilde{\mathbf{W}}_2) = \text{rank}(\mathbf{W}_2) + \text{rank}(\widetilde{\mathbf{W}}_2)$, we get

$$\emptyset = \mathcal{C}(\widetilde{\mathbf{W}}_2) \cap \mathcal{C}(\mathbf{W}_2) = \mathcal{N}(\mathbf{W}_1) \cap \mathcal{C}(\mathbf{W}_2) \Rightarrow \dim(\mathcal{N}(\mathbf{W}_1) \cap \mathcal{C}(\mathbf{W}_2)) = 0.$$

We now show the other direction “ $ii \Leftarrow iii$ ”. Without loss of generality, let $\mathbf{W}_2 = [(\mathbf{W}_2')^{k \times r} \mathbf{A}^{r \times n-r} (\mathbf{W}_2')^{k \times r}]$ where columns of \mathbf{W}_2' are linearly independent and let $\widetilde{\mathbf{W}}_2 = \epsilon [\mathbf{w}_1^1, \dots, \mathbf{w}_1^{k-r}, 0, \dots, 0] \in \mathbb{R}^{k \times n}$ be a rank $k-r$ matrix where \mathbf{w}_1^i are unit basis of $\mathcal{N}(\mathbf{W}_1)$ which yields $\mathcal{C}(\widetilde{\mathbf{W}}_2) = \mathcal{N}(\mathbf{W}_1)$. Then since $\dim(\mathcal{N}(\mathbf{W}_1) \cap \mathcal{C}(\mathbf{W}_2)) = 0$, we get $\text{rank}(\mathbf{W}_2 + \widetilde{\mathbf{W}}_2) = k$ for generic choice of ϵ . This completes the proof.

Note that by setting $\mathbf{W}_1 = \mathbf{W}_2^T$ and $\mathbf{W}_2 = \mathbf{W}_1^T$, the same proof can be used to show $i \Leftrightarrow iv$. Next, we will prove the equivalence $iii \Leftrightarrow iv$. Notice that

$$\begin{aligned} & \dim\left(\text{span}(\mathcal{N}(\mathbf{W}_1) \cup \mathcal{C}(\mathbf{W}_2))\right) \\ &= \dim(\mathcal{N}(\mathbf{W}_1)) + \dim(\mathcal{C}(\mathbf{W}_2)) - \dim(\mathcal{N}(\mathbf{W}_1) \cap \mathcal{C}(\mathbf{W}_2)) \\ &= k - r + r - \dim(\mathcal{N}(\mathbf{W}_1) \cap \mathcal{C}(\mathbf{W}_2)) \\ &= k - \dim(\mathcal{N}(\mathbf{W}_1) \cap \mathcal{C}(\mathbf{W}_2)). \end{aligned}$$

Thus,

$$\begin{aligned} \dim(\mathcal{N}(\mathbf{W}_1) \cap \mathcal{C}(\mathbf{W}_2)) \neq 0 &\Leftrightarrow \dim\left(\text{span}(\mathcal{N}(\mathbf{W}_1) \cup \mathcal{C}(\mathbf{W}_2))\right) < k \\ &\Leftrightarrow \exists \mathbf{a} \neq 0 \text{ such that } \mathbf{a} \perp \mathcal{C}(\mathbf{W}_2), \text{ and } \mathbf{a} \perp \mathcal{N}(\mathbf{W}_1) \\ &\Leftrightarrow \exists \mathbf{a} \neq 0 \text{ such that } \mathbf{a} \in \mathcal{N}(\mathbf{W}_2^T), \text{ and } \mathbf{a} \in \mathcal{C}(\mathbf{W}_1^T) \\ &\Leftrightarrow \dim(\mathcal{N}(\mathbf{W}_2^T) \cap \mathcal{C}(\mathbf{W}_1^T)) \neq 0, \end{aligned}$$

which completes the proof. \square

The next result shows that if the conditions in Lemma 16 hold, then $\text{rank}(\mathbf{W}_1) = \text{rank}(\mathbf{W}_2)$. Moreover, for $r \triangleq \text{rank}(\mathbf{W}_1 \mathbf{W}_2)$, the last $n-r$ rows of $\mathbf{U}^T \mathbf{W}_1$ and last $n-r$ columns of $\mathbf{W}_2 \mathbf{V}$ are all zeros, where the columns of $\mathbf{U} \in \mathbb{R}^{m \times m}$ and the columns of $\mathbf{V} \in \mathbb{R}^{n \times n}$ are respectively the left and right singular vectors of the product $\mathbf{W}_1 \mathbf{W}_2$.

Lemma 17. *Let $\mathbf{W}_1 \in \mathbb{R}^{m \times k}$, $\mathbf{W}_2 \in \mathbb{R}^{k \times n}$ with $k < \min\{m, n\}$ and let $r \triangleq \text{rank}(\mathbf{W}_1 \mathbf{W}_2)$. Assume further that $\mathbf{W}_1 \mathbf{W}_2 = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$ is an SVD decomposition of $\mathbf{W}_1 \mathbf{W}_2$ with $\mathbf{U} \in \mathbb{R}^{m \times m}$, and $\mathbf{V} \in \mathbb{R}^{n \times n}$, and $\mathbf{\Sigma} \in \mathbb{R}^{m \times n}$. If*

$$\begin{cases} i) \exists \widetilde{\mathbf{W}}_1 \in \mathbb{R}^{m \times k} \text{ such that } \widetilde{\mathbf{W}}_1 \mathbf{W}_2 = \mathbf{0} \text{ and } \mathbf{W}_1 + \widetilde{\mathbf{W}}_1 \text{ is full column rank.} \\ \text{and} \\ ii) \exists \widetilde{\mathbf{W}}_2 \in \mathbb{R}^{k \times n} \text{ such that } \mathbf{W}_1 \widetilde{\mathbf{W}}_2 = \mathbf{0} \text{ and } \mathbf{W}_2 + \widetilde{\mathbf{W}}_2 \text{ is full row rank.} \end{cases}$$

then

$$\text{rank}(\mathbf{W}_1) = \text{rank}(\mathbf{W}_2), \quad (\mathbf{W}_2 \mathbf{V})_{:,r+1:n} = 0, \quad \text{and} \quad (\mathbf{U}^T \mathbf{W}_1)_{r+1:n,:} = 0.$$

Proof. Suppose that ii) holds, then

$$\mathcal{C}(\widetilde{\mathbf{W}}_2) \subseteq \mathcal{N}(\mathbf{W}_1) \Rightarrow \text{rank}(\widetilde{\mathbf{W}}_2) \leq \dim(\mathcal{N}(\mathbf{W}_1)) = k - \text{rank}(\mathbf{W}_1). \quad (22)$$

Also,

$$k = \text{rank}(\mathbf{W}_2 + \widetilde{\mathbf{W}}_2) \leq \text{rank}(\mathbf{W}_2) + \text{rank}(\widetilde{\mathbf{W}}_2). \quad (23)$$

From inequalities (22) and (23), we get

$$k - \text{rank}(\mathbf{W}_2) \leq \text{rank}(\widetilde{\mathbf{W}}_2) \leq k - \text{rank}(\mathbf{W}_1) \Rightarrow \text{rank}(\mathbf{W}_2) \geq \text{rank}(\mathbf{W}_1). \quad (24)$$

Similarly, condition i) implies

$$\mathcal{C}(\widetilde{\mathbf{W}}_1^T) \subseteq \mathcal{N}(\mathbf{W}_2^T) \Rightarrow \text{rank}(\widetilde{\mathbf{W}}_1) \leq \dim(\mathcal{N}(\mathbf{W}_2^T)) = k - \text{rank}(\mathbf{W}_2). \quad (25)$$

Also,

$$k = \text{rank}(\mathbf{W}_1 + \widetilde{\mathbf{W}}_1) \leq \text{rank}(\mathbf{W}_1) + \text{rank}(\widetilde{\mathbf{W}}_1). \quad (26)$$

From inequalities (25) and (26), we get

$$k - \text{rank}(\mathbf{W}_1) \leq \text{rank}(\widetilde{\mathbf{W}}_1) \leq k - \text{rank}(\mathbf{W}_2) \Rightarrow \text{rank}(\mathbf{W}_1) \geq \text{rank}(\mathbf{W}_2). \quad (27)$$

From inequalities (24) and (27), we get $\text{rank}(\mathbf{W}_1) = \text{rank}(\mathbf{W}_2)$, which combined with Lemma 16 implies $\dim(\mathcal{N}(\mathbf{W}_1) \cap \mathcal{C}(\mathbf{W}_2)) = 0$. Let $r \triangleq \text{rank}(\mathbf{W}_1 \mathbf{W}_2)$. It directly follows from the SVD decomposition of the matrix $\mathbf{W}_1 \mathbf{W}_2$, that $U^T \mathbf{W}_1 (\mathbf{W}_2 \mathbf{V})_{:,r+1:n} = \boldsymbol{\Sigma}_{:,r+1:n} = \mathbf{0}$, or equivalently $\mathbf{W}_1 (\mathbf{W}_2 \mathbf{V})_{:,r+1:n} = \mathbf{0}$. On the other hand, since $\mathcal{C}(\mathbf{W}_2 \mathbf{V}_{:,r+1:n}) \subset \mathcal{C}(\mathbf{W}_2)$ and $\mathcal{N}(\mathbf{W}_1) \cap \mathcal{C}(\mathbf{W}_2) = \emptyset$, we conclude that $(\mathbf{W}_2 \mathbf{V})_{:,r+1:n} = \mathbf{0}$. Similarly, one can show that $(U^T \mathbf{W}_1)_{r+1:n,:} = \mathbf{0}$. \square

Proposition 18. *Let $\mathcal{M}(\mathbf{W}_1, \mathbf{W}_2) = \mathbf{W}_1 \mathbf{W}_2$ be the matrix product mapping with $\mathbf{W}_1 \in \mathbb{R}^{m \times k}$, $\mathbf{W}_2 \in \mathbb{R}^{k \times n}$, and $k \leq \min\{m, n\}$. Then, $\mathcal{M}(\cdot, \cdot)$ is locally open in its range $\mathcal{R}_{\mathcal{M}} \triangleq \{\mathbf{Z} \in \mathbb{R}^{m \times n} : \text{rank}(\mathbf{Z}) \leq k\}$ at the point $(\bar{\mathbf{W}}_1, \bar{\mathbf{W}}_2)$ if and only if the following two conditions are satisfied:*

$$\begin{cases} i) \exists \widetilde{\mathbf{W}}_1 \in \mathbb{R}^{m \times k} \text{ such that } \widetilde{\mathbf{W}}_1 \bar{\mathbf{W}}_2 = \mathbf{0} \text{ and } \bar{\mathbf{W}}_1 + \widetilde{\mathbf{W}}_1 \text{ is full column rank.} \\ \quad \text{and} \\ ii) \exists \widetilde{\mathbf{W}}_2 \text{ such that } \bar{\mathbf{W}}_1 \widetilde{\mathbf{W}}_2 = \mathbf{0} \text{ and } \bar{\mathbf{W}}_2 + \widetilde{\mathbf{W}}_2 \text{ is full row rank.} \end{cases}$$

Proof. First of all, according to Lemma 14 and Lemma 15, without loss of generality we can assume that the matrix product $\bar{\mathbf{W}}_1 \bar{\mathbf{W}}_2$ is of diagonal form.

Let us start by first proving the ‘‘only if’’ direction. Notice that the result clearly holds when $\text{rank}(\bar{\mathbf{W}}_1) = \text{rank}(\bar{\mathbf{W}}_2) = k$ by choosing $\widetilde{\mathbf{W}}_1 = \widetilde{\mathbf{W}}_2 = \mathbf{0}$. Moreover, the mapping $\mathcal{M}(\cdot, \cdot)$ cannot be locally open if only one of the matrices $\bar{\mathbf{W}}_1$ or $\bar{\mathbf{W}}_2$ is rank deficient. To see this, let us assume that $\bar{\mathbf{W}}_1$ is full column rank, while $\bar{\mathbf{W}}_2$ is rank deficient. Assume further that the mapping $\mathcal{M}(\cdot, \cdot)$ is locally open at $(\bar{\mathbf{W}}_1, \bar{\mathbf{W}}_2)$, it follows from the definition of openness that the mapping $\mathcal{M}^1(\mathbf{W}_1, \mathbf{W}_2^1) \triangleq \mathbf{W}_1 \mathbf{W}_2^1$ is locally open at $(\bar{\mathbf{W}}_1, \bar{\mathbf{W}}_2^1)$ where $\bar{\mathbf{W}}_2^1 \triangleq (\bar{\mathbf{W}}_2)_{:,1:k}$ only contains the first k columns of $\bar{\mathbf{W}}_2$. Since the range of the mapping \mathcal{M}^1 at $(\bar{\mathbf{W}}_1, \bar{\mathbf{W}}_2^1)$ is the entire space $\mathbb{R}^{m \times k}$, Proposition 4 implies that

$$\begin{cases} \exists \widetilde{\mathbf{W}}_1 \text{ such that } \widetilde{\mathbf{W}}_1 \bar{\mathbf{W}}_2^1 = \mathbf{0} \text{ and } \bar{\mathbf{W}}_1 + \widetilde{\mathbf{W}}_1 \text{ is full row rank.} \\ \quad \text{or} \\ \exists \widetilde{\mathbf{W}}_2^1 \text{ such that } \bar{\mathbf{W}}_1 \widetilde{\mathbf{W}}_2^1 = \mathbf{0} \text{ and } \bar{\mathbf{W}}_2^1 + \widetilde{\mathbf{W}}_2^1 \text{ is full rank.} \end{cases}$$

Moreover, since $\bar{\mathbf{W}}_1 \in \mathbb{R}^{m \times k}$ and $m > k$, it is impossible for $\widetilde{\mathbf{W}}_1 + \bar{\mathbf{W}}_1$ to be full row rank. On the other hand, since $\bar{\mathbf{W}}_1$ is full column rank, $\widetilde{\mathbf{W}}_1 \bar{\mathbf{W}}_2^1 = \mathbf{0}$ implies that $\widetilde{\mathbf{W}}_2^1 = \mathbf{0}$; and hence $\bar{\mathbf{W}}_2^1 + \widetilde{\mathbf{W}}_2^1$ is not full column rank. Hence none of the above two conditions can hold and consequently, $\mathcal{M}(\cdot, \cdot)$ cannot be open at the point $(\bar{\mathbf{W}}_1, \bar{\mathbf{W}}_2^1)$ in this case. Similarly, we can show that when $\bar{\mathbf{W}}_1$ is rank deficient and $\bar{\mathbf{W}}_2$ is full row rank, the mapping $\mathcal{M}(\cdot, \cdot)$ cannot be locally open. Hence, if $\bar{\mathbf{W}}_1$ and $\bar{\mathbf{W}}_2$ are not both full rank, then they both should be rank deficient.

Assume that the matrices $\bar{\mathbf{W}}_1$ and $\bar{\mathbf{W}}_2$ are both rank deficient and $\mathcal{M}(\cdot, \cdot)$ is locally open at $(\bar{\mathbf{W}}_1, \bar{\mathbf{W}}_2)$. It follows that $\mathcal{M}^1(\mathbf{W}_1, \mathbf{W}_2^1) \triangleq \mathbf{W}_1 \mathbf{W}_2^1$ is locally open at $(\bar{\mathbf{W}}_1, \bar{\mathbf{W}}_2^1)$. By Proposition 4, and since there does not exist $\widetilde{\mathbf{W}}_1$ such that $\bar{\mathbf{W}}_1 + \widetilde{\mathbf{W}}_1$ is full row rank, there should exist $\widetilde{\mathbf{W}}_2^1$ such that $\bar{\mathbf{W}}_1 \widetilde{\mathbf{W}}_2^1 = \mathbf{0}$ and $\bar{\mathbf{W}}_2^1 + \widetilde{\mathbf{W}}_2^1$ is full rank. Defining $\widetilde{\mathbf{W}}_2 \triangleq \begin{bmatrix} \widetilde{\mathbf{W}}_2^1 & | & \mathbf{0} \end{bmatrix}$, we satisfy the desired condition *ii*).

Similarly, by looking at the transpose of the mapping \mathcal{M} , we can show that condition *i*) is true when \mathcal{M} is locally open.

We now prove the ‘‘if’’ direction. Suppose *i*) and *ii*) hold, i.e.,

$$\begin{cases} \exists \widetilde{\mathbf{W}}_1 \text{ such that } \widetilde{\mathbf{W}}_1 \bar{\mathbf{W}}_2 = \mathbf{0} \text{ and } \bar{\mathbf{W}}_1 + \widetilde{\mathbf{W}}_1 \text{ is full column rank.} \\ \quad \text{and} \\ \exists \widetilde{\mathbf{W}}_2 \text{ such that } \bar{\mathbf{W}}_1 \widetilde{\mathbf{W}}_2 = \mathbf{0} \text{ and } \bar{\mathbf{W}}_2 + \widetilde{\mathbf{W}}_2 \text{ is full row rank.} \end{cases}$$

Let $\Sigma = \bar{W}_1 \bar{W}_2 = [\Sigma_{:,1:r} \quad 0]$ be a rank r matrix. Lemma 17 implies that $\text{rank}(\bar{W}_1) = \text{rank}(\bar{W}_2)$, and the last $n - r$ columns of \bar{W}_2 are all zeros. We need to show that for any given $\epsilon > 0$, there exists $\delta > 0$, such that

$$\mathbb{B}_\delta(\bar{W}_1 \bar{W}_2) \cap \mathcal{R}_{\mathcal{M}} \subseteq \mathcal{M}(\mathbb{B}_\epsilon(\bar{W}_1), \mathbb{B}_\epsilon(\bar{W}_2)).$$

Consider a perturbed matrix $\tilde{\Sigma} \in \mathbb{B}_\delta(\Sigma) \cap \mathcal{R}_{\mathcal{M}}$, we show that $\tilde{\Sigma} \in \mathcal{M}(\mathbb{B}_\epsilon(\bar{W}_1), \mathbb{B}_\epsilon(\bar{W}_2))$. Without loss of generality, and by permuting the columns of $\tilde{\Sigma}$ if necessary, $\tilde{\Sigma}$ can be expressed as

$$\tilde{\Sigma} = \left[\begin{array}{c|c|c} \underbrace{\Sigma_{:,1:r} + R_\delta^1}_{m \times r} & \underbrace{R_\delta^2}_{m \times (k-r)} & \underbrace{(\Sigma_{:,1:r} + R_\delta^1)A_1 + R_\delta^2 A_2}_{m \times (n-k)} \end{array} \right].$$

Here $A_1 \in \mathbb{R}^{r \times (n-k)}$ and $A_2 \in \mathbb{R}^{(k-r) \times (n-k)}$ exist since $\text{rank}(\tilde{\Sigma}) \leq k$. Moreover, $\|\tilde{\Sigma} - \Sigma\| \leq \delta$ implies that the perturbed matrix

$$R_\delta \triangleq [R_\delta^1 \mid R_\delta^2 \mid (\Sigma_{:,1:r} + R_\delta^1)A_1 + R_\delta^2 A_2]$$

has norm less than or equal δ , i.e. $\|R_\delta\| \leq \delta$.

Since $\text{rank}(\bar{W}_2 + \tilde{W}_2) = k$, there exist a unitary basis set $\{\tilde{w}_2^1, \dots, \tilde{w}_2^{k-r}\}$ for \tilde{W}_2 such that $\text{span}\{\tilde{w}_2^1, \dots, \tilde{w}_2^{k-r}\} \cap \mathcal{C}(\bar{W}_2) = 0$. Define

$$\tilde{W}_2^1 \triangleq \frac{\epsilon}{n2^{n+1}} \begin{bmatrix} k \times r & k \times (k-r) \\ 0 & \tilde{w}_2^1 \dots \tilde{w}_2^{k-r} \end{bmatrix}, \quad (28)$$

and let us form the matrix $\bar{W}_2^1 \in \mathbb{R}^{k \times k}$ using the first k columns of \tilde{W}_2 . Since the last $n - r$ columns of the matrix \bar{W}_2 are zero, $\bar{W}_2^1 + \tilde{W}_2^1$ is a full rank $k \times k$ matrix and $\bar{W}_1 \bar{W}_2^1 = 0$. Let us define

$$\bar{W}_1^0 \triangleq [R_\delta^1 \mid R_\delta^2] (\bar{W}_2^1 + \tilde{W}_2^1)^{-1},$$

and

$$\bar{W}_2^0 \triangleq [\tilde{W}_2^1 \mid (\bar{W}_2^1 + \tilde{W}_2^1)_{:,1:r} A_1 + (\bar{W}_2^1 + \tilde{W}_2^1)_{:,r+1:k} A_2].$$

Using this definition, we have

$$\begin{aligned} & (\bar{W}_1 + \bar{W}_1^0)(\bar{W}_2 + \bar{W}_2^0) \\ &= \left[(\bar{W}_1 + \bar{W}_1^0)(\bar{W}_2 + \bar{W}_2^0)_{:,1:k} \mid (\bar{W}_1 + \bar{W}_1^0)(\bar{W}_2 + \bar{W}_2^0)_{:,k+1:n} \right] \\ &= \left[(\bar{W}_1 + \bar{W}_1^0)(\bar{W}_2^1 + \tilde{W}_2^1) \mid (\bar{W}_1 + \bar{W}_1^0)(\bar{W}_2 + \bar{W}_2^0)_{:,k+1:n} \right] \\ &= \left[\underbrace{\tilde{\Sigma}_{:,1:k} + \underbrace{\bar{W}_1 \tilde{W}_2^1}_{=0}}_{m \times k} + [R_\delta^1 \mid R_\delta^2] (\bar{W}_2^1 + \tilde{W}_2^1)^{-1} (\bar{W}_2^1 + \tilde{W}_2^1) \mid \underbrace{\mathbf{0}}_{m \times (n-k)} \right] \\ &+ \left[\underbrace{\mathbf{0}}_{m \times k} \mid (\bar{W}_1 + \bar{W}_1^0) \left[\begin{array}{cc} (\bar{W}_2^1 + \tilde{W}_2^1)_{:,1:r} & (\bar{W}_2^1 + \tilde{W}_2^1)_{:,r+1:k} \end{array} \right] \left[\begin{array}{c} A_1 \\ A_2 \end{array} \right] \right] \\ &= \bar{W}_1 \bar{W}_2 + [R_\delta^1 \mid R_\delta^2 \mid (\Sigma_{:,1:r} + R_\delta^1)A_1 + R_\delta^2 A_2] \\ &= \bar{W}_1 \bar{W}_2 + R_\delta \\ &= \tilde{\Sigma}. \end{aligned} \quad (29)$$

□

To complete the proof, it remains to show that for any $\epsilon > 0$, we can choose δ small enough such that $\|\bar{W}_1^0\| \leq \epsilon$ and $\|\bar{W}_2^0\| \leq \epsilon$. In other words, we will show $\tilde{\Sigma} \in \mathcal{M}(\mathbb{B}_\epsilon(\bar{W}_1), \mathbb{B}_\epsilon(\bar{W}_2))$.

Let \tilde{r} , with $k \geq \tilde{r} \geq r$, be the rank of $\tilde{\Sigma}$. According to Lemma 13 and by possibly permuting the columns, $\tilde{\Sigma}$ can be expressed as

$$\tilde{\Sigma} = [\tilde{\Sigma}_1 \mid \tilde{\Sigma}_1 \bar{A}],$$

where $\tilde{\Sigma}_1 \in \mathbb{R}^{m \times \tilde{r}}$ is full column rank, and \bar{A} has a bounded norm $\|\bar{A}\| \leq n2^{n-\tilde{r}-1}$. Notice that for given \bar{W}_1^0 and \bar{W}_2^0 satisfying (29), permuting the columns of $\tilde{\Sigma}$ corresponds to permuting the columns of $(\bar{W}_2 + \bar{W}_2^0)$. If we can show that the first r columns are not among the permuted ones, then using the fact that \bar{W}_2 has only its first r columns non-zero, it follows that the permutation of the columns of $\tilde{\Sigma}$ corresponds to the same permutation of the columns of \bar{W}_2^0 . Moreover, if the first r columns are not among the permuted ones, then without loss of generality we can express the perturbed matrix

$$\tilde{\Sigma} = \left[\underbrace{\Sigma_{:,1:r} + R_\delta^1}_{m \times r} \mid \underbrace{R_\delta^2}_{m \times (k-r)} \mid \underbrace{(\Sigma_{:,1:r} + R_\delta^1)\bar{A}_1 + R_\delta^2\bar{A}_2}_{m \times (n-k)} \right],$$

and the perturbation matrix

$$R_\delta = \left[\underbrace{R_\delta^1}_{m \times r} \mid \underbrace{R_\delta^2}_{m \times (k-r)} \mid \underbrace{(\Sigma_{:,1:r} + R_\delta^1)\bar{A}_1 + R_\delta^2\bar{A}_2}_{m \times (n-k)} \right],$$

where $\begin{bmatrix} \bar{A}_1 \\ \bar{A}_2 \end{bmatrix} = \bar{A}$ has a bounded norm.

We now show that the first r columns of $\tilde{\Sigma}$ before permutation $\Sigma_{:,1:r} + R_\delta^1 \subseteq \tilde{\Sigma}_1$. Assume the contrary, then there exists at least a column $\Sigma_{:,j} + (R_\delta^1)_{:,j}$ that is not a column of $\tilde{\Sigma}_1$, which implies $\Sigma_{:,j} + (R_\delta^1)_{:,j}$ is a column of $\tilde{\Sigma}_1 \bar{A}$. Without loss of generality let $\Sigma_{:,j} + (R_\delta^1)_{:,j} = \tilde{\Sigma}_1 \bar{A}_{:,1}$. It follows that

$$\Sigma_{j,j} + (R_\delta^1)_{j,j} = (\tilde{\Sigma}_1)_{j,:} \bar{A}_{:,1}.$$

But since $\Sigma_{j,j} + (R_\delta^1)_{j,j}$ is a non-zero perturbed singular value, and since elements of $(\tilde{\Sigma}_1)_{j,:}$ are all of order δ , then by choosing δ sufficiently small, we get $\|\bar{A}\| > 2^{n-\tilde{r}-1}$, which contradicts the bound we have on \bar{A} .

We now obtain an upper-bound on $\|\bar{W}_2^0\|$. Since the norm of \bar{A} is bounded, the norm of \bar{A}_2 is also bounded by some constant $K \triangleq n2^n > n2^{n-\tilde{r}-1}$. Hence,

$$\begin{aligned} \delta &\geq \|R_\delta\| \geq \|(\Sigma_{:,1:r} + R_\delta^1)\bar{A}_1 + R_\delta^2\bar{A}_2\| \\ &\geq \|(\Sigma_{:,1:r} + R_\delta^1)\bar{A}_1\| - \|R_\delta^2\bar{A}_2\| \\ &\geq \|(\Sigma_{:,1:r} + R_\delta^1)\bar{A}_1\| - K\delta \\ &\geq \frac{\sigma_{\min}}{2} \|\bar{A}_1\| - K\delta, \end{aligned}$$

where σ_{\min} is the minimum singular value of the full column rank matrix $\Sigma_{:,1:r}$ which is bounded away from zero. Here, we have chosen $\delta < \sigma_{\min}/2$ so that $\|(\Sigma_{:,1:r} + R_\delta^1)\bar{A}_1\| \leq \frac{\sigma_{\min}}{2} \|\bar{A}_1\|$. Rearranging the terms, we obtain

$$\|\bar{A}_1\| \leq \frac{2(1+K)\delta}{\sigma_{\min}}.$$

Thus, for some constant $C \triangleq \|\bar{W}_2^1\|$, we obtain

$$\begin{aligned}
\|\bar{\mathbf{W}}_2^0\|^2 &\leq \|\widetilde{\mathbf{W}}_2^1\|^2 + \|\bar{\mathbf{W}}_2^1\|^2 \|\bar{\mathbf{A}}_1\|^2 + \|\widetilde{\mathbf{W}}_2^1\|^2 \|\bar{\mathbf{A}}_2\|^2 && \text{(by triangular inequality and Cauchy Shwarz)} \\
&\leq \frac{\epsilon^2}{4n^2 2^{2n}} + \frac{\epsilon^2 K^2}{4n^2 2^{2n}} + \delta^2 C^2 \left(\frac{2+2K}{\sigma_{\min}} \right)^2 \\
&\leq \frac{\epsilon^2}{4K^2} + \frac{\epsilon^2 K^2}{4K^2} + \delta^2 C^2 \left(\frac{2+2K}{\sigma_{\min}} \right)^2 \\
&\leq \epsilon^2/2 + \delta^2 C^2 \left(\frac{2+2K}{\sigma_{\min}} \right)^2.
\end{aligned}$$

For a given $\epsilon > 0$, choose

$$\delta \leq \min \left\{ \frac{\epsilon}{1 + \max \left\{ \|(\bar{\mathbf{W}}_2^1 + \widetilde{\mathbf{W}}_2^1)^{-1}\|, \sqrt{2} C \left(\frac{2+2K}{\sigma_{\min}} \right) \right\}}, \sigma_{\min}/2 \right\}.$$

This choice of δ leads to $\|\bar{\mathbf{W}}_2^0\| \leq \epsilon$. Moreover,

$$\begin{aligned}
\|\bar{\mathbf{W}}_1^0\| &\leq \|\mathbf{R}_\delta\| \|(\mathbf{W}_2^1 + \widetilde{\mathbf{W}}_2^1)^{-1}\| \\
&\leq \delta \|(\mathbf{W}_2^1 + \widetilde{\mathbf{W}}_2^1)^{-1}\| \\
&\leq \frac{\epsilon \|(\mathbf{W}_2^1 + \widetilde{\mathbf{W}}_2^1)^{-1}\|}{1 + \|(\mathbf{W}_2^1 + \widetilde{\mathbf{W}}_2^1)^{-1}\|} \\
&\leq \epsilon,
\end{aligned}$$

which completes the proof.

We now use Proposition 18, Lemma 16, and Lemma 17 to complete the proof of Theorem 5 restated below.

Theorem 5 [Restated]: Let $\mathcal{M}(\mathbf{W}_1, \mathbf{W}_2) = \mathbf{W}_1 \mathbf{W}_2$ denote the matrix multiplication mapping with $\mathbf{W}_1 \in \mathbb{R}^{m \times k}$ and $\mathbf{W}_2 \in \mathbb{R}^{k \times n}$. Assume $k < \min\{m, n\}$. Then if $\text{rank}(\bar{\mathbf{W}}_1) \neq \text{rank}(\bar{\mathbf{W}}_2)$, $\mathcal{M}(\cdot, \cdot)$ is not locally open at $(\bar{\mathbf{W}}_1, \bar{\mathbf{W}}_2)$. Else, if $\text{rank}(\bar{\mathbf{W}}_1) = \text{rank}(\bar{\mathbf{W}}_2)$, then the following statements are equivalent:

- i) $\exists \widetilde{\mathbf{W}}_1 \in \mathbb{R}^{m \times k}$ such that $\widetilde{\mathbf{W}}_1 \bar{\mathbf{W}}_2 = \mathbf{0}$ and $\bar{\mathbf{W}}_1 + \widetilde{\mathbf{W}}_1$ is full column rank.
- ii) $\exists \widetilde{\mathbf{W}}_2 \in \mathbb{R}^{k \times n}$ such that $\bar{\mathbf{W}}_1 \widetilde{\mathbf{W}}_2 = \mathbf{0}$ and $\bar{\mathbf{W}}_2 + \widetilde{\mathbf{W}}_2$ is full row rank.
- iii) $\dim(\mathcal{N}(\bar{\mathbf{W}}_1) \cap \mathcal{C}(\bar{\mathbf{W}}_2)) = 0$.
- iv) $\dim(\mathcal{N}(\bar{\mathbf{W}}_2^T) \cap \mathcal{C}(\bar{\mathbf{W}}_1^T)) = 0$.
- v) $\mathcal{M}(\cdot, \cdot)$ is locally open at $(\bar{\mathbf{W}}_1, \bar{\mathbf{W}}_2)$ in its range $\mathcal{R}_{\mathcal{M}} = \{\mathbf{Z} \in \mathbb{R}^{m \times n} \text{ with } \text{rank}(\mathbf{Z}) \leq \min\{m, k, n\}\}$.

Proof. First of all, if $\mathcal{M}(\cdot, \cdot)$ is locally open at $(\bar{\mathbf{W}}_1, \bar{\mathbf{W}}_2)$, according to Proposition 18, the conditions *i*) and *ii*) must hold; and hence $\text{rank}(\bar{\mathbf{W}}_1) = \text{rank}(\bar{\mathbf{W}}_2)$ due to Lemma 17. Thus, $\mathcal{M}(\cdot, \cdot)$ cannot be locally open if $\text{rank}(\bar{\mathbf{W}}_1) \neq \text{rank}(\bar{\mathbf{W}}_2)$. On the other hand, when $\text{rank}(\bar{\mathbf{W}}_1) = \text{rank}(\bar{\mathbf{W}}_2)$, the conditions *i*), *ii*), *iii*), and *iv*) are equivalent due to Lemma 16. Moreover, these conditions imply local openness according to Proposition 18. □

C PROOF OF THEOREM 12

Consider the training problem of a multi-layer deep linear neural network:

$$\underset{\mathbf{W}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{W}_h \cdots \mathbf{W}_1 \mathbf{X} - \mathbf{Y}\|^2. \quad (30)$$

Here $\mathbf{W} = (\mathbf{W}_i)_{i=1}^h$, $\mathbf{W}_i \in \mathbb{R}^{d_i \times d_{i-1}}$ are the weight matrices, $\mathbf{X} \in \mathbb{R}^{d_0 \times n}$ is the input training data, and $\mathbf{Y} \in \mathbb{R}^{d_h \times n}$ is the target training data. Based on our general framework, the corresponding auxiliary optimization problem is given by

$$\begin{aligned} & \underset{\mathbf{Z} \in \mathbb{R}^{d_h \times d_0}}{\text{minimum}} \quad \frac{1}{2} \|\mathbf{Z}\mathbf{X} - \mathbf{Y}\|^2 \\ & \text{subject to} \quad \text{rank}(\mathbf{Z}) \leq d_p \triangleq \min_{0 \leq i \leq h} d_i \end{aligned} \quad (31)$$

Let $p_1^* \triangleq \underset{0 \leq i \leq h}{\text{argmin}} d_i$ and $p_2^* \triangleq \underset{j \neq p_1^*}{\text{argmin}} d_j$. In this section we show that if $d_{p_2^*} < \min\{d_h, d_0\}$, we can find a rank deficient \mathbf{Y} such that problem (30) has a local minimum that is not global. Otherwise, given any \mathbf{X} and \mathbf{Y} , every local minimum of problem (30) is a global minimum. We start with a Lemma that will be essential in our main proof.

Lemma 19. *Consider a degenerate point $\bar{\mathbf{W}} = (\bar{\mathbf{W}}_h, \dots, \bar{\mathbf{W}}_1)$ with $\mathcal{N}(\bar{\mathbf{W}}_i)$ and $\mathcal{N}(\bar{\mathbf{W}}_i^T)$ for $h-1 \leq i \leq 2$ all non-empty. If*

$$\mathcal{N}(\bar{\mathbf{W}}_h) \text{ is non-empty} \quad \text{or} \quad \mathcal{N}(\bar{\mathbf{W}}_1^T) \text{ is non-empty,}$$

then $\bar{\mathbf{W}}$ is either a global minimum or a saddle point of problem (30).

Proof. Suppose that $\mathcal{N}(\bar{\mathbf{W}}_h)$ is non-empty. Let $\Delta = \bar{\mathbf{W}}_h \cdots \bar{\mathbf{W}}_1 \mathbf{X} - \mathbf{Y}$. If $\Delta \mathbf{X}^T = \mathbf{0}$, then by convexity of the square loss error function, the point $\bar{\mathbf{W}} = (\bar{\mathbf{W}}_h, \dots, \bar{\mathbf{W}}_1)$ is a global minimum of (30). Else, there exist (i, j) such that $\langle \mathbf{X}_{i,:}, \Delta_{j,:} \rangle \neq 0$. We define the set $\mathcal{K} \triangleq \{k \mid 3 \leq k \leq h, \mathcal{N}(\bar{\mathbf{W}}_k) \perp \mathcal{N}((\bar{\mathbf{W}}_{k-1} \bar{\mathbf{W}}_{k-2} \cdots \bar{\mathbf{W}}_2)^T)\}$. We split the rest of the proof into two cases that correspond to \mathcal{K} being empty and non-empty.

Case a: Assume \mathcal{K} is non-empty. We define $k^* \triangleq \underset{k \in \mathcal{K}}{\text{maximum}} k$.

By definition of the set \mathcal{K} and choice of k^* , the null space $\mathcal{N}(\bar{\mathbf{W}}_{k^*})$ is orthogonal to the null-space $\mathcal{N}((\bar{\mathbf{W}}_{k^*-1} \cdots \bar{\mathbf{W}}_2)^T)$. This implies there exists a non-zero $\mathbf{b} \in \mathbb{R}^{d_{k^*-1}}$ such that $\mathbf{b} \in \mathcal{N}(\bar{\mathbf{W}}_{k^*}) \cap \mathcal{C}(\bar{\mathbf{W}}_{k^*-1} \cdots \bar{\mathbf{W}}_2)$. By considering perturbation in directions $\mathbf{A} = (\mathbf{A}_h, \dots, \mathbf{A}_1)$, $\mathbf{A}_i \in \mathbb{R}^{d_i \times d_{i-1}}$ for the optimization problem

$$\underset{t}{\text{minimize}} \quad g(t) \triangleq \frac{1}{2} \|(\bar{\mathbf{W}}_h + t\mathbf{A}_h) \cdots (\bar{\mathbf{W}}_1 + t\mathbf{A}_1) \mathbf{X} - \mathbf{Y}\|^2, \quad (32)$$

we examine the optimality conditions for a specific direction $\bar{\mathbf{A}}$.

Let

$$(\bar{\mathbf{A}}_h)_{l,:} \triangleq \begin{cases} \alpha_h \mathbf{p}_h^T & \text{if } l = j, \\ \mathbf{0} & \text{otherwise} \end{cases} \quad (\bar{\mathbf{A}}_1)_{l,:} \triangleq \begin{cases} \alpha_1 \mathbf{b}_1 & \text{if } l = i, \\ \mathbf{0} & \text{otherwise} \end{cases} \quad \bar{\mathbf{A}}_k \triangleq \begin{cases} \mathbf{b}_k \mathbf{p}_k^T & \text{if } k^* + 1 \leq k \leq h-1 \\ \mathbf{b}_k \mathbf{b}^T & \text{if } k = k^* \\ \mathbf{0} & \text{if } 2 \leq k \leq k^* - 1, \end{cases}$$

where α_h and α_1 are scalar constants, $\mathbf{b}_1 \in \mathbb{R}^{d_1}$ such that $\bar{\mathbf{W}}_{k^*-1} \cdots \bar{\mathbf{W}}_2 \mathbf{b}_1 = \mathbf{b}$, and

$$\mathbf{p}_k \in \mathcal{N}((\bar{\mathbf{W}}_{k-1} \cdots \bar{\mathbf{W}}_2)^T), \quad \mathbf{b}_{k-1} \in \mathcal{N}(\bar{\mathbf{W}}_k), \quad \text{and } \langle \mathbf{p}_k, \mathbf{b}_{k-1} \rangle \neq 0 \quad \forall k^* + 1 \leq k \leq h. \quad (33)$$

Notice that such \mathbf{p}_k and \mathbf{b}_{k-1} exist from the definition of \mathcal{K} and choice of k^* . For this particular choice of $\bar{\mathbf{A}} = (\bar{\mathbf{A}}_h, \dots, \bar{\mathbf{A}}_1)$, we obtain

$$\bar{\mathbf{W}}_{k+1} \bar{\mathbf{A}}_k = \mathbf{0} \quad \text{for } k^* \leq k \leq h-1 \quad \text{and} \quad \bar{\mathbf{A}}_k \bar{\mathbf{W}}_{k-1} \cdots \bar{\mathbf{W}}_2 = \mathbf{0} \quad \text{for } k^* + 1 \leq k \leq h. \quad (34)$$

We now show that $(\bar{\mathbf{A}}_h, \dots, \bar{\mathbf{A}}_1)$ is in fact a descent direction. Before proceeding we define some notation that will ease the expressions of the optimality conditions. Let \mathcal{V} be an index set that is a subset of $\{1, \dots, h\}$. We define the function $f(\bar{\mathbf{A}}^{\mathcal{V}}, \bar{\mathbf{W}}^{-\mathcal{V}})$ which is the matrix product attained from $\bar{\mathbf{W}}_h \cdots \bar{\mathbf{W}}_1 \mathbf{X}$ by replacing matrices $\bar{\mathbf{W}}_v$ by

matrices $\bar{\mathbf{A}}_v$ for every $v \in \mathcal{V}$. For instance, if $h = 5$ and $\mathcal{V} = \{2, 3, 5\}$, then $f(\bar{\mathbf{A}}^\mathcal{V}, \bar{\mathbf{W}}^{-\mathcal{V}}) = \bar{\mathbf{A}}_5 \bar{\mathbf{W}}_4 \bar{\mathbf{A}}_3 \bar{\mathbf{A}}_2 \bar{\mathbf{W}}_1 \mathbf{X}$. We now determine index sets \mathcal{V} , with $|\mathcal{V}| \geq 1$, that correspond to non-zero $f(\bar{\mathbf{A}}^\mathcal{V}, \bar{\mathbf{W}}^{-\mathcal{V}})$. First note by definition of $\bar{\mathbf{A}}$, if $\mathcal{V} \cap \{k^* - 1, \dots, 2\} \neq \emptyset$, then $f(\bar{\mathbf{A}}^\mathcal{V}, \bar{\mathbf{W}}^{-\mathcal{V}}) = \mathbf{0}$. Also by (34), for any $k^* \leq v \leq h - 1$, if $v \in \mathcal{V}$ then either $\{k^*, \dots, h\} \in \mathcal{V}$ or $f(\bar{\mathbf{A}}^\mathcal{V}, \bar{\mathbf{W}}^{-\mathcal{V}}) = \mathbf{0}$. This directly imply that $\bar{\mathbf{A}}_h \cdots \bar{\mathbf{A}}_{k^*} \bar{\mathbf{W}}_{k^*-1} \cdots \bar{\mathbf{W}}_1 \mathbf{X}$ and $\bar{\mathbf{A}}_h \cdots \bar{\mathbf{A}}_{k^*} \bar{\mathbf{W}}_{k^*-1} \cdots \bar{\mathbf{W}}_2 \bar{\mathbf{A}}_1 \mathbf{X}$ are the only terms that can take non-zero values. Using the definition equation (32) we obtain

$$g(t) = \frac{1}{2} \|t^{h-k^*+1} \bar{\mathbf{A}}_h \cdots \bar{\mathbf{A}}_{k^*} \bar{\mathbf{W}}_{k^*-1} \cdots \bar{\mathbf{W}}_1 \mathbf{X} + t^{h-k^*+2} \bar{\mathbf{A}}_h \cdots \bar{\mathbf{A}}_{k^*} \bar{\mathbf{W}}_{k^*-1} \cdots \bar{\mathbf{W}}_2 \bar{\mathbf{A}}_1 \mathbf{X} + \Delta\|^2.$$

It follows that

$$\left. \frac{\partial^r g(t)}{\partial t^r} \right|_{t=0} = 0 \quad \text{for all } r \leq h - k^*.$$

and

$$\left. \frac{\partial^{h-k^*+1} g(t)}{\partial t^{h-k^*+1}} \right|_{t=0} = (h - k^* + 1)! \langle \bar{\mathbf{A}}_h \cdots \bar{\mathbf{A}}_{k^*} \bar{\mathbf{W}}_{k^*-1} \cdots \bar{\mathbf{W}}_1 \mathbf{X}, \Delta \rangle.$$

If $\langle \bar{\mathbf{A}}_h \cdots \bar{\mathbf{A}}_{k^*} \bar{\mathbf{W}}_{k^*-1} \cdots \bar{\mathbf{W}}_1 \mathbf{X}, \Delta \rangle \neq 0$, then by properly choosing the sign of α_h such that $\langle \bar{\mathbf{A}}_h \cdots \bar{\mathbf{A}}_{k^*} \bar{\mathbf{W}}_{k^*-1} \cdots \bar{\mathbf{W}}_1 \mathbf{X}, \Delta \rangle < 0$, we get a descent direction. Otherwise, we examine

$$\left. \frac{\partial^{h-k^*+2} g(t)}{\partial t^{h-k^*+2}} \right|_{t=0} = (h - k^* + 2)! \langle \bar{\mathbf{A}}_h \cdots \bar{\mathbf{A}}_{k^*} \bar{\mathbf{W}}_{k^*-1} \cdots \bar{\mathbf{W}}_2 \bar{\mathbf{A}}_1 \mathbf{X}, \Delta \rangle + h \langle \bar{\mathbf{A}}_h \cdots \bar{\mathbf{A}}_{k^*} \bar{\mathbf{W}}_{k^*-1} \cdots \bar{\mathbf{W}}_1 \mathbf{X} \rangle.$$

where $h(\cdot)$ is a function of $\bar{\mathbf{A}}_h \cdots \bar{\mathbf{A}}_{k^*} \bar{\mathbf{W}}_{k^*-1} \cdots \bar{\mathbf{W}}_1 \mathbf{X}$.

We now evaluate the term $\langle \bar{\mathbf{A}}_h \cdots \bar{\mathbf{A}}_{k^*} \bar{\mathbf{W}}_{k^*-1} \cdots \bar{\mathbf{W}}_2 \bar{\mathbf{A}}_1 \mathbf{X}, \Delta \rangle$. Since $(\bar{\mathbf{A}}_h)_{l,:} = \mathbf{0}$ for all $l \neq j$ and $(\bar{\mathbf{A}}_1)_{:,l} = \mathbf{0}$ for all $l \neq i$, we only need to compute the (j, i) index $(\bar{\mathbf{A}}_h \cdots \bar{\mathbf{A}}_{k^*} \bar{\mathbf{W}}_{k^*-1} \cdots \bar{\mathbf{W}}_2 \bar{\mathbf{A}}_1)_{(j,i)}$ as all other indices are zero. For some constant $c = \mathbf{p}_h^T \mathbf{b}_{h-1} \mathbf{p}_{h-1}^T \mathbf{b}_{h-2} \cdots \mathbf{p}_{k^*+1}^T \mathbf{b}_{k^*} \mathbf{b}^T \mathbf{b}$, we obtain

$$\begin{aligned} (\bar{\mathbf{A}}_h \cdots \bar{\mathbf{A}}_{k^*} \bar{\mathbf{W}}_{k^*-1} \cdots \bar{\mathbf{W}}_2 \bar{\mathbf{A}}_1)_{(j,i)} &= \alpha_h \alpha_1 \mathbf{p}_h^T \mathbf{b}_{h-1} \mathbf{p}_{h-1}^T \mathbf{b}_{h-2} \cdots \mathbf{p}_{k^*+1}^T \mathbf{b}_{k^*} \mathbf{b}^T \bar{\mathbf{W}}_{k^*-1} \cdots \bar{\mathbf{W}}_2 \mathbf{b}_1 \\ &= \alpha_h \alpha_1 \mathbf{p}_h^T \mathbf{b}_{h-1} \mathbf{p}_{h-1}^T \mathbf{b}_{h-2} \cdots \mathbf{p}_{k^*+1}^T \mathbf{b}_{k^*} \mathbf{b}^T \mathbf{b} \\ &= \alpha_h \alpha_1 c, \end{aligned}$$

where c is non-zero by our choice of \mathbf{b} , \mathbf{p}_k and \mathbf{b}_{k-1} for $h \leq k \leq k^* - 1$ as defined in (33). For a fixed $\alpha_h \neq 0$, $h(\bar{\mathbf{A}}_h \cdots \bar{\mathbf{A}}_{k^*} \bar{\mathbf{W}}_{k^*-1} \cdots \bar{\mathbf{W}}_1 \mathbf{X})$ is a constant scalar we denote by c_α . Then by properly choosing α_1 such that

$$\underbrace{\alpha_h}_{\neq 0} \alpha_1 \underbrace{c}_{\neq 0} \underbrace{\langle \mathbf{X}_{i,:}, \Delta_{j,:} \rangle}_{\neq 0} + c_\alpha < 0,$$

we get a descent direction. This completes the first case.

Case b: Assume \mathcal{K} is empty. We consider

$$(\bar{\mathbf{A}}_h)_{l,:} \triangleq \begin{cases} \alpha_h \mathbf{p}_h^T & \text{if } l = j, \\ \mathbf{0} & \text{otherwise} \end{cases} \quad (\bar{\mathbf{A}}_1)_{:,l} \triangleq \begin{cases} \alpha_1 \mathbf{b}_1 & \text{if } l = i, \\ \mathbf{0} & \text{otherwise} \end{cases} \quad \bar{\mathbf{A}}_k \triangleq \begin{cases} \mathbf{b}_k \mathbf{p}_k^T & \text{if } 3 \leq k \leq h - 1 \\ \mathbf{b}_k \mathbf{b}_1^T & \text{if } k = 2, \end{cases}$$

where α_h and α_1 are scalar constants, $\mathbf{b}_1 \in \mathcal{N}(\bar{\mathbf{W}}_2)$, and

$$\mathbf{p}_k \in \mathcal{N}((\bar{\mathbf{W}}_{k-1} \cdots \bar{\mathbf{W}}_2)^T), \quad \mathbf{b}_{k-1} \in \mathcal{N}(\bar{\mathbf{W}}_k), \quad \text{and } \langle \mathbf{p}_k, \mathbf{b}_{k-1} \rangle \neq 0 \quad \forall 3 \leq k \leq h. \quad (35)$$

For this particular choice of $\bar{\mathbf{A}} = (\bar{\mathbf{A}}_h, \dots, \bar{\mathbf{A}}_1)$, we obtain

$$\bar{\mathbf{W}}_{k+1} \bar{\mathbf{A}}_k = \mathbf{0} \quad \text{for } 2 \leq k \leq h - 1 \quad \text{and} \quad \bar{\mathbf{A}}_k \bar{\mathbf{W}}_{k-1} \cdots \bar{\mathbf{W}}_2 = \mathbf{0} \quad \text{for } 3 \leq k \leq h.$$

We now determine index sets \mathcal{V} , with $|\mathcal{V}| \geq 1$, that correspond to non-zero $f(\bar{\mathbf{A}}^\mathcal{V}, \bar{\mathbf{W}}^{-\mathcal{V}})$. By (35), for any $2 \leq v \leq h-1$, if $v \in \mathcal{V}$ then either $\{2, \dots, h\} \in \mathcal{V}$ or $f(\bar{\mathbf{A}}^\mathcal{V}, \bar{\mathbf{W}}^{-\mathcal{V}}) = 0$. This directly imply that $\bar{\mathbf{A}}_h \cdots \bar{\mathbf{A}}_2 \bar{\mathbf{W}}_1 \mathbf{X}$ and $\bar{\mathbf{A}}_h \cdots \bar{\mathbf{A}}_1 \mathbf{X}$ are the only terms that can take non-zero values. Using the definition of equation (32) we obtain

$$g(t) = \frac{1}{2} \|t^{h-1} \bar{\mathbf{A}}_h \cdots \bar{\mathbf{A}}_2 \bar{\mathbf{W}}_1 \mathbf{X} + t^h \bar{\mathbf{A}}_h \cdots \bar{\mathbf{A}}_1 \mathbf{X} + \Delta\|^2.$$

It follows that

$$\left. \frac{\partial^r g(t)}{\partial t^r} \right|_{t=0} = 0 \quad \text{for all } r \leq h-2.$$

and

$$\left. \frac{\partial^{h-1} g(t)}{\partial t^{h-1}} \right|_{t=0} = (h-1)! \langle \bar{\mathbf{A}}_h \cdots \bar{\mathbf{A}}_2 \bar{\mathbf{W}}_1 \mathbf{X}, \Delta \rangle.$$

If $\langle \bar{\mathbf{A}}_h \cdots \bar{\mathbf{A}}_2 \bar{\mathbf{W}}_1 \mathbf{X}, \Delta \rangle \neq 0$, then by properly choosing the sign of α_h such that $\langle \bar{\mathbf{A}}_h \cdots \bar{\mathbf{A}}_2 \bar{\mathbf{W}}_1 \mathbf{X}, \Delta \rangle < 0$, we get a descent direction. Otherwise, we examine

$$\left. \frac{\partial^h g(t)}{\partial t^h} \right|_{t=0} = \langle \bar{\mathbf{A}}_h \cdots \bar{\mathbf{A}}_1 \mathbf{X}, \Delta \rangle + h \langle \bar{\mathbf{A}}_h \cdots \bar{\mathbf{A}}_2 \bar{\mathbf{W}}_1 \mathbf{X}, \Delta \rangle,$$

where $h(\cdot)$ is a function of $\bar{\mathbf{A}}_h \cdots \bar{\mathbf{A}}_2 \bar{\mathbf{W}}_1 \mathbf{X}$. We now evaluate the term $\langle \bar{\mathbf{A}}_h \cdots \bar{\mathbf{A}}_1 \mathbf{X}, \Delta \rangle$. Since $(\bar{\mathbf{A}}_h)_{l,:} = \mathbf{0}$ for all $l \neq j$ and $(\bar{\mathbf{A}}_1)_{:,l} = \mathbf{0}$ for all $l \neq i$, we only need to compute the (j, i) index $(\bar{\mathbf{A}}_h \cdots \bar{\mathbf{A}}_1)_{(j,i)}$ as all other indices are zero. For some constant $c = \mathbf{p}_h^T \mathbf{b}_{h-1} \mathbf{p}_{h-1}^T \mathbf{b}_{h-2} \cdots \mathbf{p}_3^T \mathbf{b}_2 \mathbf{b}_1^T \mathbf{b}_1$, we obtain

$$\begin{aligned} (\bar{\mathbf{A}}_h \cdots \bar{\mathbf{A}}_1)_{(j,i)} &= \alpha_h \alpha_1 \mathbf{p}_h^T \mathbf{b}_{h-1} \mathbf{p}_{h-1}^T \mathbf{b}_{h-2} \cdots \mathbf{p}_3^T \mathbf{b}_2 \mathbf{b}_1^T \mathbf{b}_1 \\ &= \alpha_h \alpha_1 c, \end{aligned}$$

where c is non-zero by our choice of \mathbf{b} , \mathbf{p}_k and \mathbf{b}_{k-1} for $3 \leq k \leq h$ as defined in (35). For a fixed $\alpha_h \neq 0$, $h(\bar{\mathbf{A}}_h \cdots \bar{\mathbf{A}}_2 \bar{\mathbf{W}}_1 \mathbf{X})$ is a constant scalar we denote by c_α . Then by properly choosing α_1 such that

$$\underbrace{\alpha_h}_{\neq 0} \alpha_1 \underbrace{c}_{\neq 0} \underbrace{\langle \mathbf{X}_{i,:}, \Delta_{j,:} \rangle}_{\neq 0} + c_\alpha < 0,$$

we get a descent direction. This completes the second case.

Now if $\mathcal{N}(\bar{\mathbf{W}}_1^T)$ is non-empty, we define the set

$$\mathcal{K} \triangleq \{k \mid 1 \leq k \leq h-2, \mathcal{N}(\bar{\mathbf{W}}_{h-1} \cdots \bar{\mathbf{W}}_{k+1}) \perp \mathcal{N}(\bar{\mathbf{W}}_k^T)\},$$

and use a similar proof scheme to show the result. More specifically, we split the proof into two cases that correspond to \mathcal{K} being empty and non-empty.

Case a: Assume \mathcal{K} is non-empty. We define $k^* \triangleq \underset{k \in \mathcal{K}}{\text{minimum}} k$.

By definition of the set \mathcal{K} and choice of k^* , the null space $\mathcal{N}(\bar{\mathbf{W}}_{k^*}^T)$ is orthogonal to the null-space $\mathcal{N}(\bar{\mathbf{W}}_{h-1} \cdots \bar{\mathbf{W}}_{k^*+1})$. This implies there exists a non-zero $\mathbf{p} \in \mathbb{R}^{d_{k^*}}$ such that $\mathbf{p} \in \mathcal{N}(\bar{\mathbf{W}}_{k^*}^T) \cap \mathcal{C}((\bar{\mathbf{W}}_{h-1} \cdots \bar{\mathbf{W}}_{k^*+1})^T)$. By considering perturbation in directions $\mathbf{A} = (\mathbf{A}_h, \dots, \mathbf{A}_1)$, $\mathbf{A}_i \in \mathbb{R}^{d_i \times d_{i-1}}$ for the optimization problem

$$\underset{t}{\text{minimize}} \quad g(t) \triangleq \frac{1}{2} \|(\bar{\mathbf{W}}_h + t\mathbf{A}_h) \cdots (\bar{\mathbf{W}}_1 + t\mathbf{A}_1) \mathbf{X} - \mathbf{Y}\|^2, \quad (36)$$

we examine the optimality conditions for a specific direction $\bar{\mathbf{A}}$.

Let

$$(\bar{\mathbf{A}}_h)_{l,:} \triangleq \begin{cases} \alpha_h \mathbf{p}_h^T & \text{if } l = j, \\ \mathbf{0} & \text{otherwise} \end{cases} \quad (\bar{\mathbf{A}}_1)_{:,l} \triangleq \begin{cases} \alpha_1 \mathbf{b}_1 & \text{if } l = i, \\ \mathbf{0} & \text{otherwise} \end{cases} \quad \bar{\mathbf{A}}_k \triangleq \begin{cases} \mathbf{b}_k \mathbf{p}_k^T & \text{if } 2 \leq k \leq k^* - 1 \\ \mathbf{p}_k \mathbf{p}_k^T & \text{if } k = k^* \\ \mathbf{0} & \text{if } k^* + 1 \leq k \leq h - 1, \end{cases}$$

where α_h and α_1 are scalar constants, $\mathbf{p}_h \in \mathbb{R}^{d_{h-1}}$ such that $\mathbf{p}_h^T \bar{\mathbf{W}}_{h-1} \cdots \bar{\mathbf{W}}_{k^*+1} = \mathbf{p}^T$, and

$$\mathbf{p}_k \in \mathcal{N}(\bar{\mathbf{W}}_{k-1}^T), \quad \mathbf{b}_{k-1} \in \mathcal{N}(\bar{\mathbf{W}}_{h-1} \cdots \bar{\mathbf{W}}_k), \quad \text{and } \langle \mathbf{p}_k, \mathbf{b}_{k-1} \rangle \neq 0 \quad \forall 2 \leq k \leq k^*. \quad (37)$$

Notice that such \mathbf{p}_k and \mathbf{b}_{k-1} exist from the definition of \mathcal{K} and choice of k^* . For this particular choice of $\bar{\mathbf{A}} = (\bar{\mathbf{A}}_h, \dots, \bar{\mathbf{A}}_1)$, we obtain

$$\bar{\mathbf{A}}_k \bar{\mathbf{W}}_{k-1} = \mathbf{0} \quad \text{for } 2 \leq k \leq k^* \quad \text{and} \quad \bar{\mathbf{W}}_{h-1} \cdots \bar{\mathbf{W}}_{k+1} \bar{\mathbf{A}}_k = \mathbf{0} \quad \text{for } 1 \leq k \leq k^* - 1. \quad (38)$$

The same argument used above can be used to show that $(\bar{\mathbf{A}}_h, \dots, \bar{\mathbf{A}}_1)$ is actually a descent direction. This completes the first case.

Case b: Assume \mathcal{K} is empty. We consider

$$(\bar{\mathbf{A}}_h)_{l,:} \triangleq \begin{cases} \alpha_h \mathbf{p}_h^T & \text{if } l = j, \\ \mathbf{0} & \text{otherwise} \end{cases} \quad (\bar{\mathbf{A}}_1)_{:,l} \triangleq \begin{cases} \alpha_1 \mathbf{b}_1 & \text{if } l = i, \\ \mathbf{0} & \text{otherwise} \end{cases} \quad \bar{\mathbf{A}}_k \triangleq \begin{cases} \mathbf{b}_k \mathbf{p}_k^T & \text{if } 2 \leq k \leq h - 2 \\ \mathbf{p}_h \mathbf{p}_k^T & \text{if } k = h - 1, \end{cases}$$

where α_h and α_1 are scalar constants, $\mathbf{p}_h \in \mathcal{N}(\bar{\mathbf{W}}_{h-1}^T)$, and

$$\mathbf{p}_k \in \mathcal{N}(\bar{\mathbf{W}}_{k-1}^T), \quad \mathbf{b}_{k-1} \in \mathcal{N}(\bar{\mathbf{W}}_{h-1} \cdots \bar{\mathbf{W}}_k), \quad \text{and } \langle \mathbf{p}_k, \mathbf{b}_{k-1} \rangle \neq 0 \quad \forall 2 \leq k \leq h - 1. \quad (39)$$

For this particular choice of $\bar{\mathbf{A}} = (\bar{\mathbf{A}}_h, \dots, \bar{\mathbf{A}}_1)$, we obtain

$$\bar{\mathbf{A}}_k \bar{\mathbf{W}}_{k-1} = \mathbf{0} \quad \text{for } 2 \leq k \leq h - 1 \quad \text{and} \quad \bar{\mathbf{W}}_{h-1} \cdots \bar{\mathbf{W}}_{k+1} \bar{\mathbf{A}}_k = \mathbf{0} \quad \text{for } 1 \leq k \leq h - 2. \quad (40)$$

The same argument used above can be used to show that $(\bar{\mathbf{A}}_h, \dots, \bar{\mathbf{A}}_1)$ is actually a descent direction. This completes the first case and thus completes the proof. \square

Note that following the same steps of the proof in Lemma 19, we get the same result when replacing the square loss error by a general convex and differentiable function $\ell(\cdot)$. We are now ready to prove the main result restated below.

Theorem 12 Let $p_1^* \triangleq \underset{0 \leq i \leq h}{\operatorname{argmin}} d_i$ and $p_2^* \triangleq \underset{j \neq p_1^*}{\operatorname{argmin}} d_j$. If $d_{p_2^*} < \min(d_h, d_0)$, we can find a rank deficient \mathbf{Y} such that problem (30) has a local minimum that is not global. Otherwise, given any \mathbf{X} and \mathbf{Y} , every local minimum of problem (30) is a global minimum.

Proof. Suppose $d_{p_2^*} < \min\{d_h, d_0\}$, we define

$$p_2 \triangleq \max(p_1^*, p_2^*) \quad \text{and} \quad p_1 \triangleq \min(p_1^*, p_2^*).$$

Let $\mathbf{X} \triangleq \mathbf{I}$,

$$(\bar{\mathbf{Y}})_{(i,j)} \triangleq \begin{cases} 1 & \text{if } (i,j) = (d_h, d_0) \\ 0 & \text{otherwise} \end{cases}, \quad \bar{\mathbf{W}}_k \triangleq \begin{cases} \begin{bmatrix} \mathbf{I}_{d_k} & \mathbf{0} \end{bmatrix} & \text{if } d_k \leq d_{k-1}, \\ \begin{bmatrix} \mathbf{I}_{d_{k-1}} \\ \mathbf{0} \end{bmatrix} & \text{if } d_k > d_{k-1}, \end{cases}$$

for $k \in \{h, \dots, p_2 + 1\} \cup \{p_1, \dots, 1\}$, and $\bar{\mathbf{W}}_k = \mathbf{0}$ for $k \in \{p_2, \dots, p_1 + 1\}$. Since $\bar{\mathbf{W}}_h \cdots \bar{\mathbf{W}}_{p_2+1}$ and $\bar{\mathbf{W}}_{p_1} \cdots \bar{\mathbf{W}}_1$ are both full rank, then using Lemma 10, the matrix products \mathcal{M}_{h,p_2+1} and $\mathcal{M}_{p_1,1}$ are locally open at $(\bar{\mathbf{W}}_h, \dots, \bar{\mathbf{W}}_{p_2+1})$ and $(\bar{\mathbf{W}}_{p_1}, \dots, \bar{\mathbf{W}}_1)$, respectively. Moreover, using Proposition 4 and the composition property

of open maps, the matrix product mapping \mathcal{M}_{p_2, p_1+1} is locally open at $(\bar{\mathbf{W}}_{p_2}, \dots, \bar{\mathbf{W}}_{p_1+1})$. It follows by Observation 1 that if $\bar{\mathbf{W}}$ is a local minimum of

$$\underset{\mathbf{W}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{W}_h \cdots \mathbf{W}_1 - \bar{\mathbf{Y}}\|^2. \quad (41)$$

then $(\bar{\mathbf{Z}}_3, \bar{\mathbf{Z}}_2, \bar{\mathbf{Z}}_1)$ is a local minimum of

$$\underset{\mathbf{Z}_3 \in \mathbb{R}^{d_h \times d_{p_2}}, \mathbf{Z}_2 \in \mathbb{R}^{d_{p_2} \times d_{p_1}}, \mathbf{Z}_1 \in \mathbb{R}^{d_{p_1} \times d_0}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{Z}_3 \mathbf{Z}_2 \mathbf{Z}_1 - \bar{\mathbf{Y}}\|^2. \quad (42)$$

where

$$\bar{\mathbf{Z}}_3 = \bar{\mathbf{W}}_h \cdots \bar{\mathbf{W}}_{p_2+1} = \begin{bmatrix} \mathbf{I}_{d_{p_2}} \\ \mathbf{0} \end{bmatrix}, \quad \bar{\mathbf{Z}}_2 = \mathbf{0}, \quad \text{and} \quad \bar{\mathbf{Z}}_1 = \bar{\mathbf{W}}_{p_1} \cdots \bar{\mathbf{W}}_1 = \begin{bmatrix} \mathbf{I}_{d_{p_1}} & \mathbf{0} \end{bmatrix}.$$

The point $(\bar{\mathbf{Z}}_3, \bar{\mathbf{Z}}_2, \bar{\mathbf{Z}}_1)$ is obviously not global, we show using optimality conditions that the point is a local minimum. By considering perturbations in the directions $\bar{\mathbf{A}} = (\bar{\mathbf{A}}_3, \bar{\mathbf{A}}_2, \bar{\mathbf{A}}_1)$ for the optimization problem

$$\begin{aligned} \underset{t}{\text{minimize}} \quad g(t) &\triangleq \frac{1}{2} \|(\bar{\mathbf{Z}}_3 + t\bar{\mathbf{A}}_3)(\bar{\mathbf{Z}}_2 + t\bar{\mathbf{A}}_2)(\bar{\mathbf{Z}}_1 + t\bar{\mathbf{A}}_1) - \bar{\mathbf{Y}}\|^2 \\ &= \frac{1}{2} \|t(\bar{\mathbf{Z}}_3 + t\bar{\mathbf{A}}_3)\bar{\mathbf{A}}_2(\bar{\mathbf{Z}}_1 + t\bar{\mathbf{A}}_1) - \bar{\mathbf{Y}}\|^2. \end{aligned} \quad (43)$$

It follows that

$$\begin{aligned} \left. \frac{\partial g(t)}{\partial t} \right|_{t=0} &= \langle \bar{\mathbf{Z}}_3 \bar{\mathbf{A}}_2 \bar{\mathbf{Z}}_1, -\bar{\mathbf{Y}} \rangle \\ &= -(\bar{\mathbf{Z}}_3 \bar{\mathbf{A}}_2 \bar{\mathbf{Z}}_1)_{d_h, d_0} \bar{\mathbf{Y}}_{d_h, d_0} \\ &= 0, \end{aligned} \quad (44)$$

where the last equality holds since the last row (d_h^{th} row) of $\bar{\mathbf{Z}}_3$ is zero. Also,

$$\begin{aligned} \left. \frac{\partial^2 g(t)}{\partial t^2} \right|_{t=0} &= 2 \langle \bar{\mathbf{Z}}_3 \bar{\mathbf{A}}_2 \bar{\mathbf{A}}_1 + \bar{\mathbf{A}}_3 \bar{\mathbf{A}}_2 \bar{\mathbf{Z}}_1, -\bar{\mathbf{Y}} \rangle + \|\bar{\mathbf{Z}}_3 \bar{\mathbf{A}}_2 \bar{\mathbf{Z}}_1\|^2 \\ &= -2(\bar{\mathbf{Z}}_3 \bar{\mathbf{A}}_2 \bar{\mathbf{A}}_1)_{d_h, d_0} \bar{\mathbf{Y}}_{d_h, d_0} - 2(\bar{\mathbf{A}}_3 \bar{\mathbf{A}}_2 \bar{\mathbf{Z}}_1)_{d_h, d_0} \bar{\mathbf{Y}}_{d_h, d_0} + \|\bar{\mathbf{Z}}_3 \bar{\mathbf{A}}_2 \bar{\mathbf{Z}}_1\|^2 \\ &= \|\bar{\mathbf{A}}_2\|^2, \end{aligned} \quad (45)$$

where the last equality holds since the last row (d_h^{th} row) of $\bar{\mathbf{Z}}_3$ and the last column (d_0^{th} column) of $\bar{\mathbf{Z}}_1$ are both zeros. Then for $\|\bar{\mathbf{A}}_2\| \neq 0$, it follows from the second-order optimality condition that the point is a local minimum, and if $\|\bar{\mathbf{A}}_2\| = 0$ we get

$$g(t) = \frac{1}{2} \|\bar{\mathbf{Y}}\|^2 = \frac{1}{2}$$

which implies $(\bar{\mathbf{Z}}_3, \bar{\mathbf{Z}}_2, \bar{\mathbf{Z}}_1)$ is a local optimum that is not global.

Note that the same method used to construct the example above can be used to find a local minimum that is not global whenever the $\text{rank}(\mathbf{Y}) \leq \min\{d_h - d_{p_2}, d_0 - d_{p_1}\}$. When \mathbf{Y} is full rank, we know from the results of Lu & Kawaguchi (2017); Yun et al. (2017) that every local minimum is global. To have a complete characterization of problems for which every local minimum is global, it remains to either prove or disprove the statement when \mathbf{Y} is a rank deficient matrix with $\text{rank}(\mathbf{Y}) > \min\{d_h - d_{p_2}, d_0 - d_{p_1}\}$. We now provide a counterexample that disproves the statement. In particular, we construct a three layer network with input \mathbf{X} and output \mathbf{Y} with $\text{rank}(\mathbf{Y}) > \min\{d_h - d_{p_2}, d_0 - d_{p_1}\}$, and then find a local minimum $(\bar{\mathbf{W}}_3, \bar{\mathbf{W}}_2, \bar{\mathbf{W}}_1)$ that is not global. Let $\mathbf{X} = \mathbf{I}$,

$$\mathbf{Y} \triangleq \begin{bmatrix} 1 & 0 & -1 \\ 0 & 4 & 0 \\ -1 & 0 & 1 \end{bmatrix}, \quad \bar{\mathbf{W}}_3 \triangleq \begin{bmatrix} 1 & -1 \\ -1 & -1 \\ 1 & -1 \end{bmatrix}, \quad \bar{\mathbf{W}}_2 \triangleq \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \quad \text{and} \quad \bar{\mathbf{W}}_1 \triangleq \bar{\mathbf{W}}_3^T.$$

Obviously $(\bar{\mathbf{W}}_3, \bar{\mathbf{W}}_2, \bar{\mathbf{W}}_1)$ is not a global minimum. We define $\Delta \triangleq \bar{\mathbf{W}}_3 \bar{\mathbf{W}}_2 \bar{\mathbf{W}}_1 - \mathbf{Y}$. Then we get

$$\bar{\mathbf{W}}_3^T \Delta = \Delta \bar{\mathbf{W}}_1^T = \mathbf{0}. \quad (46)$$

By considering perturbations in the directions $\mathbf{A} = (\mathbf{A}_3, \mathbf{A}_2, \mathbf{A}_1)$ for the optimization problem

$$\text{minimize}_t g(t, \mathbf{A}) \triangleq \frac{1}{2} \|(\bar{\mathbf{W}}_3 + t\mathbf{A}_3)(\bar{\mathbf{W}}_3 + t\mathbf{A}_2)(\bar{\mathbf{W}}_1 + t\mathbf{A}_1) - \mathbf{Y}\|^2,$$

it follows that

$$\left. \frac{\partial g(t, \mathbf{A})}{\partial t} \right|_{t=0} = \langle \bar{\mathbf{W}}_3 \bar{\mathbf{W}}_2 \mathbf{A}_1 + \bar{\mathbf{W}}_3 \mathbf{A}_2 \bar{\mathbf{W}}_1 + \mathbf{A}_3 \bar{\mathbf{W}}_2 \bar{\mathbf{W}}_1, \Delta \rangle = 0,$$

where the last equality is directly implied from (46). Also

$$\begin{aligned} g^{(2)}(0, \mathbf{A}) &\triangleq \left. \frac{\partial^2 g(t, \mathbf{A})}{\partial t^2} \right|_{t=0} = 2 \langle \mathbf{A}_3 \mathbf{A}_2 \bar{\mathbf{W}}_1 + \mathbf{A}_3 \bar{\mathbf{W}}_2 \mathbf{A}_1 + \bar{\mathbf{W}}_3 \mathbf{A}_2 \mathbf{A}_1, \Delta \rangle + \|\bar{\mathbf{W}}_3 \bar{\mathbf{W}}_2 \mathbf{A}_1 + \bar{\mathbf{W}}_3 \mathbf{A}_2 \bar{\mathbf{W}}_1 + \mathbf{A}_3 \bar{\mathbf{W}}_2 \bar{\mathbf{W}}_1\|^2 \\ &= 2 \langle \mathbf{A}_3 \bar{\mathbf{W}}_2 \mathbf{A}_1, \Delta \rangle + \|\bar{\mathbf{W}}_3 \bar{\mathbf{W}}_2 \mathbf{A}_1 + \bar{\mathbf{W}}_3 \mathbf{A}_2 \bar{\mathbf{W}}_1 + \mathbf{A}_3 \bar{\mathbf{W}}_2 \bar{\mathbf{W}}_1\|^2. \end{aligned}$$

which is a quadratic function of \mathbf{A} we denote

$$f_{\mathbf{A}} \triangleq \frac{1}{2} \mathbf{a}^T \mathbf{H}_{\mathbf{A}} \mathbf{a}.$$

Here $\mathbf{a} \in \mathbb{R}^{16 \times 1}$ is a vectorization of matrices \mathbf{A}_3 , \mathbf{A}_2 , and \mathbf{A}_1 , and $\mathbf{H}_{\mathbf{A}}$ is the hessian of $f_{\mathbf{A}}$. By computing the eigenvalues of $\mathbf{H}_{\mathbf{A}}$ we get that $\mathbf{H}_{\mathbf{A}} \succeq 0$ which directly implies

$$g^{(2)}(0, \mathbf{A}) \geq 0 \quad \forall \mathbf{A}.$$

Moreover, let \mathbf{a}_{opt} be the optimal solution set of the problem

$$\text{minimize}_{\mathbf{a}} f_{\mathbf{A}}.$$

Then $\mathbf{a}_{\text{opt}} = \{\mathbf{a} \mid \mathbf{a} \in \mathcal{N}(\mathbf{H}_{\mathbf{a}})\}$. We notice that for any $\bar{\mathbf{a}} \in \mathbf{a}_{\text{opt}}$, the corresponding direction $\bar{\mathbf{A}}$ has

$$\bar{\mathbf{W}}_3 \bar{\mathbf{W}}_2 \bar{\mathbf{A}}_1 + \bar{\mathbf{W}}_3 \bar{\mathbf{A}}_2 \bar{\mathbf{W}}_1 + \bar{\mathbf{A}}_3 \bar{\mathbf{W}}_2 \bar{\mathbf{W}}_1 = \mathbf{0} \quad \text{and} \quad \langle \bar{\mathbf{A}}_3 \bar{\mathbf{A}}_2 \bar{\mathbf{A}}_1, \Delta \rangle = 0.$$

Then, it follows that

$$g^{(3)}(0, \bar{\mathbf{A}}) \triangleq \left. \frac{\partial^3 g(t, \bar{\mathbf{A}})}{\partial t^3} \right|_{t=0} = 6 \langle \bar{\mathbf{A}}_3 \bar{\mathbf{A}}_2 \bar{\mathbf{A}}_1, \Delta \rangle = 0,$$

and

$$g^{(4)}(0, \bar{\mathbf{A}}) \triangleq \left. \frac{\partial^4 g(t, \bar{\mathbf{A}})}{\partial t^4} \right|_{t=0} = 12 \|\bar{\mathbf{A}}_3 \bar{\mathbf{A}}_2 \bar{\mathbf{W}}_1 + \bar{\mathbf{A}}_3 \bar{\mathbf{W}}_2 \bar{\mathbf{A}}_1 + \bar{\mathbf{W}}_3 \bar{\mathbf{A}}_2 \bar{\mathbf{A}}_1\|^2 \geq 0.$$

If $\bar{\mathbf{A}}_3 \bar{\mathbf{A}}_2 \bar{\mathbf{W}}_1 + \bar{\mathbf{A}}_3 \bar{\mathbf{W}}_2 \bar{\mathbf{A}}_1 + \bar{\mathbf{W}}_3 \bar{\mathbf{A}}_2 \bar{\mathbf{A}}_1 \neq \mathbf{0}$, then using the fourth order optimality conditions $(\bar{\mathbf{W}}_3, \bar{\mathbf{W}}_2, \bar{\mathbf{W}}_1)$ is a local minimum. Otherwise, we get

$$g^{(5)}(0, \bar{\mathbf{A}}) \triangleq \left. \frac{\partial^5 g(t, \bar{\mathbf{A}})}{\partial t^5} \right|_{t=0} = 0,$$

and

$$g^{(6)}(0, \bar{\mathbf{A}}) \triangleq \left. \frac{\partial^6 g(t, \bar{\mathbf{A}})}{\partial t^6} \right|_{t=0} = \|\bar{\mathbf{A}}_3 \bar{\mathbf{A}}_2 \bar{\mathbf{A}}_1\|^2 \geq 0,$$

which also implies that $(\bar{\mathbf{W}}_3, \bar{\mathbf{W}}_2, \bar{\mathbf{W}}_1)$ is a local minimum.

We now show that if $d_{p_2^*} \geq \min\{d_h, d_0\}$, every local minimum of (30) is global. In particular, we show that for any \mathbf{X} and \mathbf{Y} , if $\bar{\mathbf{W}}$ is not a global minimum, we can construct a descent direction.

First notice that if for some $1 \leq i \leq h-1$, $\bar{\mathbf{W}}_i$ is full column rank, then using Proposition 4, $\mathcal{M}_{i+1,i}(\cdot)$ is locally open at $(\bar{\mathbf{W}}_{i+1}, \bar{\mathbf{W}}_i)$ and $\bar{\mathbf{W}}_{i+1} \bar{\mathbf{W}}_i \in \mathbb{R}^{d_{i+1} \times d_{i-1}}$. Using Observation 1, we conclude that any local minimum of problem (30) is a local minimum of the problem obtained by replacing $\bar{\mathbf{W}}_{i+1} \bar{\mathbf{W}}_i$ by $\bar{\mathbf{Z}}_{i+1,i} \in \mathbb{R}^{d_{i+1} \times d_{i-1}}$. By a similar

argument, we conclude that if $\bar{\mathbf{W}}_i$ is a full row rank for some $2 \leq i \leq h$, any local minimum of problem (30) is a local minimum of the problem obtained by replacing $\bar{\mathbf{W}}_i \bar{\mathbf{W}}_{i-1}$ by $\bar{\mathbf{Z}}_{i,i-1} \in \mathbb{R}^{d_i \times d_{i-2}}$. Thus, if $\bar{\mathbf{W}} = (\bar{\mathbf{W}}_h, \dots, \bar{\mathbf{W}}_1)$ is a local minimum of problem (30), the new point $\bar{\mathbf{Z}} = (\bar{\mathbf{Z}}_{h'}, \dots, \bar{\mathbf{Z}}_1)$, where $\bar{\mathbf{Z}}_i \in \mathbb{R}^{d'_i \times d'_{i-1}}$ and $h' \leq h$, is a local minimum of the problem attained by applying the replacements discussed above. If $h' = 1$, we get the desired result from Lemma 7. Else, if $h' = 2$, the auxiliary problem becomes a two layer linear network for which Theorem 8 provides the desired result. When $h' > 2$, examine $d'_{h'}, d'_{h'-1}, d'_1$ and d'_0 . If $d'_{h'} > d'_{h'-1}$ and $d'_0 > d'_1$, then $d_{p_2^*} < \min\{d_h, d_0\}$ which contradicts our assumption. It follows by construction of $\bar{\mathbf{Z}}_i$, that either $d'_{h'} \leq d'_{h'-1}$ and $\bar{\mathbf{Z}}'_{h'}$ is not full row rank or $d'_0 \leq d'_1$ and $\bar{\mathbf{Z}}'_1$ is not full column rank; thus at least one of the null spaces $\mathcal{N}((\bar{\mathbf{Z}}'_{h'})^T)$, $\mathcal{N}(\bar{\mathbf{Z}}'_1)$ is non empty. Moreover, $\bar{\mathbf{Z}}_i$ has non-empty right and left null spaces for $2 \leq i \leq h - 1$. The result follows using Lemma 19.

□