
Supplementary Materials of *Deep Neural Networks with Multi-Branch Architectures Are Intrinsically Less Non-Convex*

Hongyang Zhang
Carnegie Mellon University

Junru Shao
Carnegie Mellon University

Ruslan Salakhutdinov
Carnegie Mellon University

A Supplementary Experiments

A.1 Performance of Multi-Branch Architecture

In this section, we test the classification accuracy of the multi-branch architecture on the CIFAR-10 dataset. We use a 9-layer VGG network [4] as our sub-network in each branch, which is memory-efficient for practitioners to fit many branches into GPU memory simultaneously. The detailed network setup of VGG-9 is in Table 1, where the width of VGG-9 is either 16 or 32. We test the performance of varying numbers of branches in the overall architecture from 4 to 32, with cross-entropy loss. Figure 1 presents the test accuracy on CIFAR-10 as the number of branches increases. It shows that the test accuracy improves monotonously with the increasing number of parallel branches/paths.

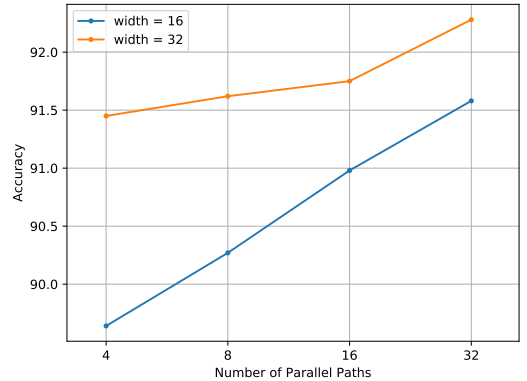


Figure 1: Test accuracy of using VGG-9 as the sub-networks in the multi-branch architecture.

Table 1: Network architecture of VGG-9. Here w is the width of the network, which controls the number of filters in each convolution layer. All convolution layers have a kernel of size 3, and zero padding of size 1. All layers followed by the batch normalization have no bias term. All max pooling layers have a stride of 2.

Layer	Weight	Activation	Input size	Output size
Input	N / A	N / A	N / A	$3 \times 32 \times 32$
Conv1	$3 \times 3 \times 3 \times w$	BN + ReLU	$3 \times 32 \times 32$	$w \times 32 \times 32$
Conv2	$3 \times 3 \times w \times w$	BN + ReLU	$w \times 32 \times 32$	$w \times 32 \times 32$
MaxPool	N / A	N / A	$w \times 32 \times 32$	$w \times 16 \times 16$
Conv3	$3 \times 3 \times w \times 2w$	BN + ReLU	$w \times 16 \times 16$	$2w \times 16 \times 16$
Conv4	$3 \times 3 \times 2w \times 2w$	BN + ReLU	$2w \times 16 \times 16$	$2w \times 16 \times 16$
MaxPool	N / A	N / A	$2w \times 16 \times 16$	$2w \times 8 \times 8$
Conv5	$3 \times 3 \times 2w \times 4w$	BN + ReLU	$2w \times 8 \times 8$	$4w \times 8 \times 8$
Conv6	$3 \times 3 \times 4w \times 4w$	BN + ReLU	$4w \times 8 \times 8$	$4w \times 8 \times 8$
Conv7	$3 \times 3 \times 4w \times 4w$	BN + ReLU	$4w \times 8 \times 8$	$4w \times 8 \times 8$
MaxPool	N / A	N / A	$4w \times 8 \times 8$	$4w \times 4 \times 4$
Flatten	N / A	N / A	$4w \times 4 \times 4$	$64w$
FC1	$64w \times 4w$	BN + ReLU	$64w$	$4w$
FC2	$4w \times 10$	Softmax	$4w$	10

A.2 Strong Duality of Deep Linear Neural Networks

We compare the optima of primal problem (4) and dual problem (5) by numerical experiments for three-layer linear neural networks ($H = 3$). The data are generated as follows. We construct the output matrix $\mathbf{Y} \in \mathbb{R}^{100 \times 100}$ by drawing the entries of \mathbf{Y} from i.i.d. standard Gaussian distribution and the input matrix $\mathbf{X} \in \mathbb{R}^{100 \times 100}$ by the identity matrix. The d_{\min} varies from 5 to 50. Both primal and dual problems are solved by numerical algorithms. Given the non-convex nature of primal problem, we rerun the algorithm by multiple initializations and choose the best solution that we obtain. The results are shown in Figure 2. We can easily see that the optima of primal and dual problems almost match. The small gap is due to the numerical inaccuracy.

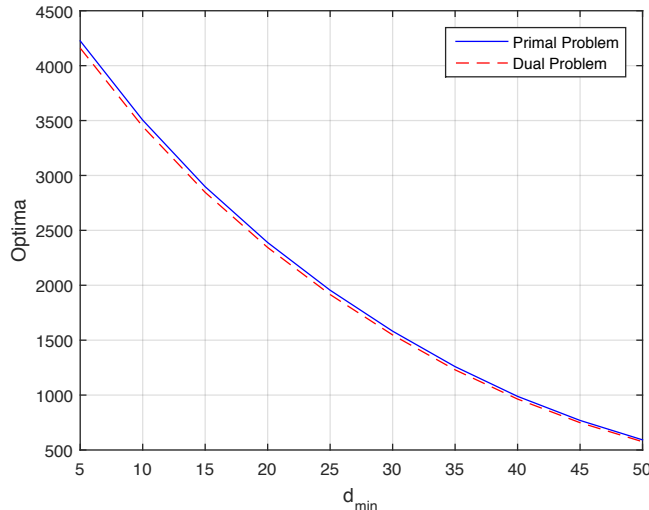


Figure 2: Comparison of optima between primal and dual problems.

We also compare the ℓ_2 distance between the solution $\mathbf{W}_H^* \mathbf{W}_{H-1}^* \dots \mathbf{W}_1^*$ of primal problem and the solution $\text{svd}_{d_{\min}}(\tilde{\mathbf{Y}} - \mathbf{\Lambda}^*)$ of dual problem in Table 2. We see that the solutions are close to each other.

Table 2: Comparison of the ℓ_2 distance between the solutions of primal and dual problems.

d_{\min}	5	10	15	20	25	30	35	40	45	50
ℓ_2 distance ($\times 10^{-10}$)	1.95	1.26	7.89	3.80	3.14	1.92	1.04	3.92	6.53	8.00

B Proofs of Theorem 1: Duality Gap of Multi-Branch Neural Networks

The lower bound $0 \leq \frac{\inf(\mathbf{P}) - \sup(\mathbf{D})}{\Delta_{\text{worst}}}$ is obvious by the weak duality. So we only need to prove the upper bound $\frac{\inf(\mathbf{P}) - \sup(\mathbf{D})}{\Delta_{\text{worst}}} \leq \frac{2}{I}$.

Consider the subset of \mathbb{R}^2 :

$$\mathcal{Y}_i := \left\{ \mathbf{y}_i \in \mathbb{R}^2 : \mathbf{y}_i = \frac{1}{I} \left[h_i(\mathbf{w}_{(i)}), \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{P}} \left(1 - \frac{y \cdot f_i(\mathbf{w}_{(i)}; \mathbf{x})}{\tau} \right) \right], \mathbf{w}_{(i)} \in \mathcal{W}_i \right\}, \quad i \in [I].$$

Define the vector summation

$$\mathcal{Y} := \mathcal{Y}_1 + \mathcal{Y}_2 + \dots + \mathcal{Y}_I.$$

Since $f_i(\mathbf{w}_{(i)}; \mathbf{x})$ and $h_i(\mathbf{w}_{(i)})$ are continuous w.r.t. $\mathbf{w}_{(i)}$ and \mathcal{W}_i 's are compact, the set

$$\{(\mathbf{w}_{(i)}, h_i(\mathbf{w}_{(i)}), f_i(\mathbf{w}_{(i)}; \mathbf{x})) : \mathbf{w}_{(i)} \in \mathcal{W}_i\}$$

is compact as well. So \mathcal{Y} , $\text{conv}(\mathcal{Y})$, \mathcal{Y}_i , and $\text{conv}(\mathcal{Y}_i)$, $i \in [I]$ are all compact sets. According to the definition of \mathcal{Y} and the standard duality argument [3], we have

$$\inf(\mathbf{P}) = \min \{w : \text{there exists } (r, w) \in \mathcal{Y} \text{ such that } r \leq K\},$$

and

$$\sup(\mathbf{D}) = \min \{w : \text{there exists } (r, w) \in \text{conv}(\mathcal{Y}) \text{ such that } r \leq K\}.$$

Technique (a): Shapley-Folkman Lemma. We are going to apply the following Shapley-Folkman lemma.

Lemma 1 (Shapley-Folkman, [5]). *Let $\mathcal{Y}_i, i \in [I]$ be a collection of subsets of \mathbb{R}^m . Then for every $\mathbf{y} \in \text{conv}(\sum_{i=1}^I \mathcal{Y}_i)$, there is a subset $\mathcal{I}(\mathbf{y}) \subseteq [I]$ of size at most m such that*

$$\mathbf{y} \in \left[\sum_{i \notin \mathcal{I}(\mathbf{y})} \mathcal{Y}_i + \sum_{i \in \mathcal{I}(\mathbf{y})} \text{conv}(\mathcal{Y}_i) \right].$$

We apply Lemma 1 to prove Theorem 1 with $m = 2$. Let $(\bar{r}, \bar{w}) \in \text{conv}(\mathcal{Y})$ be such that

$$\bar{r} \leq K, \quad \text{and} \quad \bar{w} = \sup(\mathbf{D}).$$

Applying the above Shapley-Folkman lemma to the set $\mathcal{Y} = \sum_{i=1}^I \mathcal{Y}_i$, we have that there are a subset $\bar{\mathcal{I}} \subseteq [I]$ of size 2 and vectors

$$(\bar{r}_i, \bar{w}_i) \in \text{conv}(\mathcal{Y}_i), \quad i \in \bar{\mathcal{I}} \quad \text{and} \quad \bar{\mathbf{w}}_{(i)} \in \mathcal{W}_i, \quad i \notin \bar{\mathcal{I}},$$

such that

$$\frac{1}{I} \sum_{i \notin \bar{\mathcal{I}}} h_i(\bar{\mathbf{w}}_{(i)}) + \sum_{i \in \bar{\mathcal{I}}} \bar{r}_i = \bar{r} \leq K, \quad (1)$$

$$\frac{1}{I} \sum_{i \notin \bar{\mathcal{I}}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{P}} \left(1 - \frac{y \cdot f_i(\bar{\mathbf{w}}_{(i)}; \mathbf{x})}{\tau} \right) + \sum_{i \in \bar{\mathcal{I}}} \bar{w}_i = \sup(\mathbf{D}). \quad (2)$$

Representing elements of the convex hull of $\mathcal{Y}_i \subseteq \mathbb{R}^2$ by Carathéodory theorem, we have that for each $i \in \bar{\mathcal{I}}$, there are vectors $\mathbf{w}_{(i)}^1, \mathbf{w}_{(i)}^2, \mathbf{w}_{(i)}^3 \in \mathcal{W}_i$ and scalars $a_i^1, a_i^2, a_i^3 \in \mathbb{R}$ such that

$$\sum_{j=1}^3 a_i^j = 1, \quad a_i^j \geq 0, \quad j = 1, 2, 3,$$

$$\bar{r}_i = \frac{1}{I} \sum_{j=1}^3 a_i^j h_i(\mathbf{w}_{(i)}^j), \quad \bar{w}_i = \frac{1}{I} \sum_{j=1}^3 a_i^j \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{P}} \left(1 - \frac{y \cdot f_i(\mathbf{w}_{(i)}^j; \mathbf{x})}{\tau} \right).$$

Recall that we define

$$\hat{f}_i(\tilde{\mathbf{w}}) := \inf_{\mathbf{w}_{(i)} \in \mathcal{W}_i} \left\{ \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{P}} \left(1 - \frac{y \cdot f_i(\mathbf{w}_{(i)}; \mathbf{x})}{\tau} \right) : h_i(\mathbf{w}_{(i)}) \leq h_i(\tilde{\mathbf{w}}) \right\}, \quad (3)$$

$$\tilde{f}_i(\tilde{\mathbf{w}}) := \inf_{a^j, \mathbf{w}_{(i)}^j \in \mathcal{W}_i} \left\{ \sum_{j=1}^{p_i+2} a^j \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{P}} \left(1 - \frac{y \cdot f_i(\mathbf{w}_{(i)}^j; \mathbf{x})}{\tau} \right) : \tilde{\mathbf{w}} = \sum_{j=1}^{p_i+2} a^j \mathbf{w}_{(i)}^j, \sum_{j=1}^{p_i+2} a^j = 1, a^j \geq 0 \right\},$$

and $\Delta_i := \sup_{\mathbf{w} \in \mathcal{W}_i} \{ \hat{f}_i(\mathbf{w}) - \tilde{f}_i(\mathbf{w}) \} \geq 0$. We have for $i \in \bar{\mathcal{I}}$,

$$\bar{r}_i \geq \frac{1}{I} h_i \left(\sum_{j=1}^3 a_i^j \mathbf{w}_{(i)}^j \right), \quad (\text{because } h_i(\cdot) \text{ is convex}) \quad (4)$$

and

$$\begin{aligned} \bar{w}_i &\geq \frac{1}{I} \tilde{f}_i \left(\sum_{j=1}^3 a_i^j \mathbf{w}_{(i)}^j \right) \quad (\text{by the definition of } \tilde{f}_i(\cdot)) \\ &\geq \frac{1}{I} \hat{f}_i \left(\sum_{j=1}^3 a_i^j \mathbf{w}_{(i)}^j \right) - \frac{1}{I} \Delta_i. \quad (\text{by the definition of } \Delta_i) \end{aligned} \quad (5)$$

Thus, by Eqns. (1) and (4), we have

$$\frac{1}{I} \sum_{i \notin \bar{\mathcal{I}}} h_i(\bar{\mathbf{w}}_{(i)}) + \frac{1}{I} \sum_{i \in \bar{\mathcal{I}}} h_i \left(\sum_{j=1}^3 a_i^j \mathbf{w}_{(i)}^j \right) \leq K, \quad (6)$$

and by Eqns. (2) and (5), we have

$$\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{P}} \left[\frac{1}{I} \sum_{i \notin \bar{\mathcal{I}}} \left(1 - \frac{y \cdot f_i(\bar{\mathbf{w}}_{(i)}; \mathbf{x})}{\tau} \right) \right] + \frac{1}{I} \sum_{i \in \bar{\mathcal{I}}} \hat{f}_i \left(\sum_{j=1}^3 a_i^j \mathbf{w}_{(i)}^j \right) \leq \sup(\mathbf{D}) + \frac{1}{I} \sum_{i \in \bar{\mathcal{I}}} \Delta_i. \quad (7)$$

Given any $\epsilon > 0$ and $i \in \bar{\mathcal{I}}$, we can find a vector $\bar{\mathbf{w}}_{(i)} \in \mathcal{W}_i$ such that

$$h_i(\bar{\mathbf{w}}_{(i)}) \leq h_i \left(\sum_{j=1}^3 a_i^j \mathbf{w}_{(i)}^j \right) \text{ and } \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{P}} \left(1 - \frac{y \cdot f_i(\bar{\mathbf{w}}_{(i)}; \mathbf{x})}{\tau} \right) \leq \hat{f}_i \left(\sum_{j=1}^3 a_i^j \mathbf{w}_{(i)}^j \right) + \epsilon, \quad (8)$$

where the first inequality holds because \mathcal{W}_i is convex and the second inequality holds by the definition (3) of $\hat{f}_i(\cdot)$. Therefore, Eqns. (6) and (8) imply that

$$\frac{1}{I} \sum_{i=1}^I h_i(\bar{\mathbf{w}}_{(i)}) \leq K.$$

Namely, $(\bar{\mathbf{w}}_{(1)}, \dots, \bar{\mathbf{w}}_{(I)})$ is a feasible solution of problem (2). Also, Eqns. (7) and (8) yield

$$\begin{aligned} \inf(\mathbf{P}) &\leq \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{P}} \left[\frac{1}{I} \sum_{i=1}^I \left(1 - \frac{y \cdot f_i(\bar{\mathbf{w}}_{(i)}; \mathbf{x})}{\tau} \right) \right] \\ &\leq \sup(\mathbf{D}) + \frac{1}{I} \sum_{i \in \bar{\mathcal{I}}} (\Delta_i + \epsilon) \\ &\leq \sup(\mathbf{D}) + \frac{2}{I} \Delta_{\text{worst}} + 2\epsilon, \end{aligned}$$

where the last inequality holds because $|\bar{\mathcal{I}}| = 2$. Finally, letting $\epsilon \rightarrow 0$ leads to the desired result.

C Proofs of Theorem 2: Strong Duality of Deep Linear Neural Networks

Let $\tilde{\mathbf{Y}} = \mathbf{Y}\mathbf{X}^\dagger\mathbf{X}$. We note that by Pythagorean theorem, for every \mathbf{Y} ,

$$\frac{1}{2} \|\mathbf{Y} - \mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X}\|_F^2 = \frac{1}{2} \|\tilde{\mathbf{Y}} - \mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X}\|_F^2 + \underbrace{\frac{1}{2} \|\mathbf{Y} - \tilde{\mathbf{Y}}\|_F^2}_{\text{independent of } \mathbf{W}_1, \dots, \mathbf{W}_H}.$$

So we can focus on the following optimization problem instead of problem (4):

$$\min_{\mathbf{W}_1, \dots, \mathbf{W}_H} \frac{1}{2} \|\tilde{\mathbf{Y}} - \mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X}\|_F^2 + \frac{\gamma}{H} \left[\|\mathbf{W}_1 \mathbf{X}\|_{S_H}^H + \sum_{i=2}^H \|\mathbf{W}_i\|_{S_H}^H \right]. \quad (9)$$

Technique (b): Variational Form. Our work is inspired by a variational form of problem (9) given by the following lemma.

Lemma 2. *If $(\mathbf{W}_1^*, \dots, \mathbf{W}_H^*)$ is optimal to problem*

$$\min_{\mathbf{W}_1, \dots, \mathbf{W}_H} F(\mathbf{W}_1, \dots, \mathbf{W}_H) := \frac{1}{2} \|\tilde{\mathbf{Y}} - \mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X}\|_F^2 + \gamma \|\mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X}\|_*, \quad (10)$$

*then $(\mathbf{W}_1^{**}, \dots, \mathbf{W}_H^{**})$ is optimal to problem (9), where $\mathbf{U}\Sigma\mathbf{V}^T$ is the skinny SVD of $\mathbf{W}_H^* \mathbf{W}_{H-1}^* \cdots \mathbf{W}_1^* \mathbf{X}$, $\mathbf{W}_i^{**} = [\Sigma^{1/H}, \mathbf{0}; \mathbf{0}, \mathbf{0}] \in \mathbb{R}^{d_i \times d_{i-1}}$ for $i = 2, 3, \dots, H-1$, $\mathbf{W}_H^{**} = [\mathbf{U}\Sigma^{1/H}, \mathbf{0}] \in \mathbb{R}^{d_H \times d_{H-2}}$ and $\mathbf{W}_1^{**} = [\Sigma^{1/H} \mathbf{V}^T; \mathbf{0}] \mathbf{X}^\dagger \in \mathbb{R}^{d_1 \times d_0}$. Furthermore, problems (9) and (10) have the same optimal objective function value.*

Proof of Lemma 2. Let $\mathbf{U}\Sigma\mathbf{V}^T$ be the skinny SVD of matrix $\mathbf{W}_H\mathbf{W}_{H-1}\cdots\mathbf{W}_1\mathbf{X} =: \mathbf{Z}$. We notice that

$$\begin{aligned} \|\mathbf{Z}\|_* &= \|\mathbf{W}_H\mathbf{W}_{H-1}\cdots\mathbf{W}_1\mathbf{X}\|_* \\ &\leq \|\mathbf{W}_1\mathbf{X}\|_{\mathcal{S}_H} \prod_{i=2}^H \|\mathbf{W}_i\|_{\mathcal{S}_H} \quad (\text{by the generalized Hölder's inequality}) \\ &\leq \frac{1}{H} \left[\|\mathbf{W}_1\mathbf{X}\|_{\mathcal{S}_H}^H + \sum_{i=2}^H \|\mathbf{W}_i\|_{\mathcal{S}_H}^H \right]. \quad (\text{by the inequality of mean}) \end{aligned}$$

Hence, on one hand, for every $(\mathbf{W}_1, \dots, \mathbf{W}_H)$,

$$\begin{aligned} \min_{\mathbf{W}_1, \dots, \mathbf{W}_H} F(\mathbf{W}_1, \dots, \mathbf{W}_H) &\leq \frac{1}{2} \|\tilde{\mathbf{Y}} - \mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X}\|_F^2 + \gamma \|\mathbf{W}_H \mathbf{W}_{H-1} \cdots \mathbf{W}_1 \mathbf{X}\|_* \\ &\leq \frac{1}{2} \|\tilde{\mathbf{Y}} - \mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X}\|_F^2 + \frac{\gamma}{H} \left[\|\mathbf{W}_1 \mathbf{X}\|_{\mathcal{S}_H}^H + \sum_{i=2}^H \|\mathbf{W}_i\|_{\mathcal{S}_H}^H \right], \end{aligned}$$

which yields

$$\min_{\mathbf{W}_1, \dots, \mathbf{W}_H} F(\mathbf{W}_1, \dots, \mathbf{W}_H) \leq \min_{\mathbf{W}_1, \dots, \mathbf{W}_H} \frac{1}{2} \|\tilde{\mathbf{Y}} - \mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X}\|_F^2 + \frac{\gamma}{H} \left[\|\mathbf{W}_1 \mathbf{X}\|_{\mathcal{S}_H}^H + \sum_{i=2}^H \|\mathbf{W}_i\|_{\mathcal{S}_H}^H \right].$$

On the other hand, suppose $(\mathbf{W}_1^*, \dots, \mathbf{W}_H^*)$ is optimal to problem (10), and let $\mathbf{U}\Sigma\mathbf{V}^T$ be the skinny SVD of matrix $\mathbf{W}_H^* \mathbf{W}_{H-1}^* \cdots \mathbf{W}_1^* \mathbf{X}$. We choose $(\mathbf{W}_1^{**}, \dots, \mathbf{W}_H^{**})$ such that

$$\mathbf{W}_H^{**} = [\mathbf{U}\Sigma^{\frac{1}{H}}, \mathbf{0}], \quad \mathbf{W}_1^{**} \mathbf{X} = [\Sigma^{\frac{1}{H}} \mathbf{V}^T; \mathbf{0}], \quad \mathbf{W}_i^{**} = [\Sigma^{\frac{1}{H}}, \mathbf{0}; \mathbf{0}, \mathbf{0}], \quad i = 2, \dots, H-1.$$

We pad $\mathbf{0}$ around \mathbf{W}_i^{**} so as to adapt to the dimensionality of each \mathbf{W}_i^{**} . Notice that

$$\begin{aligned} \|\mathbf{W}_H^* \mathbf{W}_{H-1}^* \cdots \mathbf{W}_1^* \mathbf{X}\|_* &= \|\mathbf{W}_H^{**} \mathbf{W}_{H-1}^{**} \cdots \mathbf{W}_1^{**} \mathbf{X}\|_* \\ &= \frac{1}{H} \left[\|\mathbf{W}_1^{**} \mathbf{X}\|_{\mathcal{S}_H}^H + \sum_{i=2}^H \|\mathbf{W}_i^{**}\|_{\mathcal{S}_H}^H \right]. \end{aligned}$$

Since $\mathbf{W}_H^* \mathbf{W}_{H-1}^* \cdots \mathbf{W}_1^* \mathbf{X} = \mathbf{W}_H^{**} \mathbf{W}_{H-1}^{**} \cdots \mathbf{W}_1^{**} \mathbf{X}$, for every $\tilde{\mathbf{Y}}$,

$$\|\tilde{\mathbf{Y}} - \mathbf{W}_H^* \mathbf{W}_{H-1}^* \cdots \mathbf{W}_1^* \mathbf{X}\|_F = \|\tilde{\mathbf{Y}} - \mathbf{W}_H^{**} \mathbf{W}_{H-1}^{**} \cdots \mathbf{W}_1^{**} \mathbf{X}\|_F.$$

Hence

$$\begin{aligned} \min_{\mathbf{W}_1, \dots, \mathbf{W}_H} F(\mathbf{W}_1, \dots, \mathbf{W}_H) &= F(\mathbf{W}_1^*, \dots, \mathbf{W}_H^*) = F(\mathbf{W}_1^{**}, \dots, \mathbf{W}_H^{**}) \\ &= \frac{1}{2} \|\tilde{\mathbf{Y}} - \mathbf{W}_H^{**} \cdots \mathbf{W}_1^{**} \mathbf{X}\|_F^2 + \frac{\gamma}{H} \left[\|\mathbf{W}_1^{**} \mathbf{X}\|_{\mathcal{S}_H}^H + \sum_{i=2}^H \|\mathbf{W}_i^{**}\|_{\mathcal{S}_H}^H \right] \\ &\geq \min_{\mathbf{W}_1, \dots, \mathbf{W}_H} \frac{1}{2} \|\tilde{\mathbf{Y}} - \mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X}\|_F^2 + \frac{\gamma}{H} \left[\|\mathbf{W}_1 \mathbf{X}\|_{\mathcal{S}_H}^H + \sum_{i=2}^H \|\mathbf{W}_i\|_{\mathcal{S}_H}^H \right], \end{aligned}$$

which yields the other direction of the inequality and hence completes the proof. \square

Technique (c): Reduction to Low-Rank Approximation. We now reduce problem (10) to the classic problem of low-rank approximation of the form $\min_{\mathbf{W}_1, \dots, \mathbf{W}_H} \frac{1}{2} \|\hat{\mathbf{Y}} - \mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X}\|_F^2$, which has the following nice properties.

Lemma 3. For any $\hat{\mathbf{Y}} \in \text{Row}(\mathbf{X})$, every global minimum $(\mathbf{W}_1^*, \dots, \mathbf{W}_H^*)$ of function

$$f(\mathbf{W}_1, \dots, \mathbf{W}_H) = \frac{1}{2} \|\hat{\mathbf{Y}} - \mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X}\|_F^2$$

obeys $\mathbf{W}_H^* \cdots \mathbf{W}_1^* \mathbf{X} = \text{svd}_{d_{\min}}(\hat{\mathbf{Y}})$. Here $\hat{\mathbf{Y}} \in \text{Row}(\mathbf{X})$ means the row vectors of $\hat{\mathbf{Y}}$ belongs to the row space of \mathbf{X} .

Proof of Lemma 3. Note that the optimal solution to $\min_{\mathbf{W}_H, \dots, \mathbf{W}_1} \frac{1}{2} \|\widehat{\mathbf{Y}} - \mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X}\|_F^2$ is equal to the optimal solution to the low-rank approximation problem $\min_{\text{rank}(\mathbf{Z}) \leq d_{\min}} \frac{1}{2} \|\widehat{\mathbf{Y}} - \mathbf{Z}\|_F^2$ when $\widehat{\mathbf{Y}} \in \text{Row}(\mathbf{X})$, which has a closed-form solution $\text{svd}_{d_{\min}}(\widehat{\mathbf{Y}})$.¹ \square

We now reduce $F(\mathbf{W}_1, \dots, \mathbf{W}_H)$ to the form of $\frac{1}{2} \|\widehat{\mathbf{Y}} - \mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X}\|_F^2$ for some $\widehat{\mathbf{Y}}$ plus an extra additive term that is independent of $(\mathbf{W}_1, \dots, \mathbf{W}_H)$. To see this, denote by $K(\cdot) = \gamma \|\cdot\|_*$. We have

$$\begin{aligned} F(\mathbf{W}_1, \dots, \mathbf{W}_H) &= \frac{1}{2} \|\widetilde{\mathbf{Y}} - \mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X}\|_F^2 + K^{**}(\mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X}) \\ &= \max_{\Lambda} \frac{1}{2} \|\widetilde{\mathbf{Y}} - \mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X}\|_F^2 + \langle \Lambda, \mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X} \rangle - K^*(\Lambda) \\ &= \max_{\Lambda} \frac{1}{2} \|\widetilde{\mathbf{Y}} - \Lambda - \mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X}\|_F^2 - \frac{1}{2} \|\Lambda\|_F^2 - K^*(\Lambda) + \langle \widetilde{\mathbf{Y}}, \Lambda \rangle \\ &=: \max_{\Lambda} L(\mathbf{W}_1, \dots, \mathbf{W}_H, \Lambda), \end{aligned}$$

where we define $L(\mathbf{W}_1, \dots, \mathbf{W}_H, \Lambda) := \frac{1}{2} \|\widetilde{\mathbf{Y}} - \Lambda - \mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X}\|_F^2 - \frac{1}{2} \|\Lambda\|_F^2 - K^*(\Lambda) + \langle \widetilde{\mathbf{Y}}, \Lambda \rangle$ as the Lagrangian of problem (10). The first equality holds because $K(\cdot)$ is closed and convex w.r.t. the argument $\mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X}$ so $K(\cdot) = K^{**}(\cdot)$, and the second equality is by the definition of conjugate function. One can check that $L(\mathbf{W}_1, \dots, \mathbf{W}_H, \Lambda) = \min_{\mathbf{M}} L'(\mathbf{W}_1, \dots, \mathbf{W}_H, \mathbf{M}, \Lambda)$, where $L'(\mathbf{W}_1, \dots, \mathbf{W}_H, \mathbf{M}, \Lambda)$ is the Lagrangian of the constraint optimization problem $\min_{\mathbf{W}_1, \dots, \mathbf{W}_H, \mathbf{M}} \frac{1}{2} \|\widetilde{\mathbf{Y}} - \mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X}\|_F^2 + K(\mathbf{M})$, s.t. $\mathbf{M} = \mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X}$. With a little abuse of notation, we call $L(\mathbf{A}, \mathbf{B}, \Lambda)$ the Lagrangian of the unconstrained problem (10) as well.

The remaining analysis is to choose a proper $\Lambda^* \in \text{Row}(\mathbf{X})$ such that $(\mathbf{W}_1^*, \dots, \mathbf{W}_H^*, \Lambda^*)$ is a primal-dual saddle point of $L(\mathbf{W}_1, \dots, \mathbf{W}_H, \Lambda)$, so that the problem $\min_{\mathbf{W}_1, \dots, \mathbf{W}_H} L(\mathbf{W}_1, \dots, \mathbf{W}_H, \Lambda^*)$ and problem (10) have the same optimal solution $(\mathbf{W}_1^*, \dots, \mathbf{W}_H^*)$. For this, we introduce the following condition, and later we will show that the condition holds.

Condition 1. For a solution $(\mathbf{W}_1^*, \dots, \mathbf{W}_H^*)$ to optimization problem (10), there exists an

$$\Lambda^* \in \partial_{\mathbf{Z}} K(\mathbf{Z})|_{\mathbf{Z}=\mathbf{W}_H^* \cdots \mathbf{W}_1^* \mathbf{X}} \cap \text{Row}(\mathbf{X})$$

such that

$$\begin{aligned} \mathbf{W}_{i+1}^{*T} \cdots \mathbf{W}_H^{*T} (\mathbf{W}_H^* \cdots \mathbf{W}_1^* \mathbf{X} + \Lambda^* - \widetilde{\mathbf{Y}}) \mathbf{X}^T \mathbf{W}_1^{*T} \cdots \mathbf{W}_{i-1}^{*T} &= \mathbf{0}, \quad i = 2, \dots, H-1, \\ \mathbf{W}_2^{*T} \cdots \mathbf{W}_H^{*T} (\mathbf{W}_H^* \mathbf{W}_{H-1}^* \cdots \mathbf{W}_1^* \mathbf{X} + \Lambda^* - \widetilde{\mathbf{Y}}) \mathbf{X}^T &= \mathbf{0}, \\ (\mathbf{W}_H^* \mathbf{W}_{H-1}^* \cdots \mathbf{W}_1^* \mathbf{X} + \Lambda^* - \widetilde{\mathbf{Y}}) \mathbf{X}^T \mathbf{W}_1^{*T} \cdots \mathbf{W}_{H-1}^{*T} &= \mathbf{0}. \end{aligned} \tag{11}$$

We note that if we set Λ to be the Λ^* in (11), then $\nabla_{\mathbf{W}_i} L(\mathbf{W}_1^*, \dots, \mathbf{W}_H^*, \Lambda^*) = \mathbf{0}$ for every i . So $(\mathbf{W}_1^*, \dots, \mathbf{W}_H^*)$ is either a saddle point, a local minimizer, or a global minimizer of $L(\mathbf{W}_1, \dots, \mathbf{W}_H, \Lambda^*)$ as a function of $(\mathbf{W}_1, \dots, \mathbf{W}_H)$ for the fixed Λ^* . The following lemma states that if it is a global minimizer, then strong duality holds.

Lemma 4. Let $(\mathbf{W}_1^*, \dots, \mathbf{W}_H^*)$ be a global minimizer of $F(\mathbf{W}_1, \dots, \mathbf{W}_H)$. If there exists a dual certificate Λ^* satisfying Condition 1 and the pair $(\mathbf{W}_1^*, \dots, \mathbf{W}_H^*)$ is a global minimizer of $L(\mathbf{W}_1, \dots, \mathbf{W}_H, \Lambda^*)$ for the fixed Λ^* , then strong duality holds. Moreover, we have the relation $\mathbf{W}_H^* \cdots \mathbf{W}_1^* \mathbf{X} = \text{svd}_{d_{\min}}(\widetilde{\mathbf{Y}} - \Lambda^*)$.

Proof of Lemma 4. By the assumption of the lemma, $(\mathbf{W}_1^*, \dots, \mathbf{W}_H^*)$ is a global minimizer of

$$L(\mathbf{W}_1, \dots, \mathbf{W}_H, \Lambda^*) = \frac{1}{2} \|\widetilde{\mathbf{Y}} - \Lambda^* - \mathbf{W}_H \mathbf{W}_{H-1} \cdots \mathbf{W}_1 \mathbf{X}\|_F^2 + c(\Lambda^*),$$

where $c(\Lambda^*)$ is a function of Λ^* that is independent of \mathbf{W}_i for all i 's. Namely, $(\mathbf{W}_1^*, \dots, \mathbf{W}_H^*)$ globally minimizes $L(\mathbf{W}_1, \dots, \mathbf{W}_H, \Lambda)$ when Λ is fixed to Λ^* . Furthermore, $\Lambda^* \in \partial_{\mathbf{Z}} K(\mathbf{Z})|_{\mathbf{Z}=\mathbf{W}_H^* \cdots \mathbf{W}_1^* \mathbf{X}}$

¹Note that the low-rank approximation problem might have non-unique solution. However, we will use in this paper the abuse of language $\text{svd}_{d_{\min}}(\widehat{\mathbf{Y}})$ as the non-uniqueness issue does not lead to any issue in our developments.

implies that $\mathbf{W}_H^* \mathbf{W}_{H-1}^* \cdots \mathbf{W}_1^* \mathbf{X} \in \partial_{\Lambda} K^*(\Lambda)|_{\Lambda=\Lambda^*}$ by the convexity of function $K(\cdot)$, meaning that $\mathbf{0} \in \partial_{\Lambda} L(\mathbf{W}_1^*, \dots, \mathbf{W}_H^*, \Lambda)$. So $\Lambda^* = \operatorname{argmax}_{\Lambda} L(\mathbf{W}_1^*, \dots, \mathbf{W}_H^*, \Lambda)$ due to the concavity of function $L(\mathbf{W}_1^*, \dots, \mathbf{W}_H^*, \Lambda)$ w.r.t. variable Λ . Thus $(\mathbf{W}_1^*, \dots, \mathbf{W}_H^*, \Lambda^*)$ is a primal-dual saddle point of $L(\mathbf{W}_1, \dots, \mathbf{W}_H, \Lambda)$.

We now prove the strong duality. By the fact that $F(\mathbf{W}_1^*, \dots, \mathbf{W}_H^*) = \max_{\Lambda} L(\mathbf{W}_1^*, \dots, \mathbf{W}_H^*, \Lambda)$ and that $\Lambda^* = \operatorname{argmax}_{\Lambda} L(\mathbf{W}_1^*, \dots, \mathbf{W}_H^*, \Lambda)$, for every $\mathbf{W}_1, \dots, \mathbf{W}_H$, we have

$$F(\mathbf{W}_1^*, \dots, \mathbf{W}_H^*) = L(\mathbf{W}_1^*, \dots, \mathbf{W}_H^*, \Lambda^*) \leq L(\mathbf{W}_1, \dots, \mathbf{W}_H, \Lambda^*),$$

where the inequality holds because $(\mathbf{W}_1^*, \dots, \mathbf{W}_H^*, \Lambda^*)$ is a primal-dual saddle point of L . Notice that we also have

$$\begin{aligned} \min_{\mathbf{W}_1, \dots, \mathbf{W}_H} \max_{\Lambda} L(\mathbf{W}_1, \dots, \mathbf{W}_H, \Lambda) &= F(\mathbf{W}_1^*, \dots, \mathbf{W}_H^*) \\ &\leq \min_{\mathbf{W}_1, \dots, \mathbf{W}_H} L(\mathbf{W}_1, \dots, \mathbf{W}_H, \Lambda^*) \\ &\leq \max_{\Lambda} \min_{\mathbf{W}_1, \dots, \mathbf{W}_H} L(\mathbf{W}_1, \dots, \mathbf{W}_H, \Lambda). \end{aligned}$$

On the other hand, by weak duality,

$$\min_{\mathbf{W}_1, \dots, \mathbf{W}_H} \max_{\Lambda} L(\mathbf{W}_1, \dots, \mathbf{W}_H, \Lambda) \geq \max_{\Lambda} \min_{\mathbf{W}_1, \dots, \mathbf{W}_H} L(\mathbf{W}_1, \dots, \mathbf{W}_H, \Lambda).$$

Therefore,

$$\min_{\mathbf{W}_1, \dots, \mathbf{W}_H} \max_{\Lambda} L(\mathbf{W}_1, \dots, \mathbf{W}_H, \Lambda) = \max_{\Lambda} \min_{\mathbf{W}_1, \dots, \mathbf{W}_H} L(\mathbf{W}_1, \dots, \mathbf{W}_H, \Lambda),$$

i.e., strong duality holds. Hence,

$$\begin{aligned} \mathbf{W}_H^* \mathbf{W}_{H-1}^* \cdots \mathbf{W}_1^* &= \operatorname{argmin}_{\mathbf{W}_H \mathbf{W}_{H-1} \cdots \mathbf{W}_1} L(\mathbf{W}_1, \dots, \mathbf{W}_H, \Lambda^*) \\ &= \operatorname{argmin}_{\mathbf{W}_H \mathbf{W}_{H-1} \cdots \mathbf{W}_1} \frac{1}{2} \|\tilde{\mathbf{Y}} - \Lambda^* - \mathbf{W}_H \mathbf{W}_{H-1} \cdots \mathbf{W}_1 \mathbf{X}\|_F^2 - \frac{1}{2} \|\Lambda^*\|_F^2 - K^*(\Lambda^*) + \langle \tilde{\mathbf{Y}}, \Lambda^* \rangle \\ &= \operatorname{argmin}_{\mathbf{W}_H \mathbf{W}_{H-1} \cdots \mathbf{W}_1} \frac{1}{2} \|\tilde{\mathbf{Y}} - \Lambda^* - \mathbf{W}_H \mathbf{W}_{H-1} \cdots \mathbf{W}_1 \mathbf{X}\|_F^2 \\ &= \operatorname{svd}_{d_{\min}}(\tilde{\mathbf{Y}} - \Lambda^*). \end{aligned}$$

The proof of Lemma 4 is completed. \square

Technique (d): Dual Certificate. We now construct dual certificate Λ^* such that all of conditions in Lemma 4 hold. We note that Λ^* should satisfy the followings by Lemma 4:

- (a) $\Lambda^* \in \partial K(\mathbf{W}_H^* \mathbf{W}_{H-1}^* \cdots \mathbf{W}_1^* \mathbf{X}) \cap \operatorname{Row}(\mathbf{X})$; (by Condition 1)
- (b) Equations (11); (by Condition 1) (12)
- (c) $\mathbf{W}_H^* \mathbf{W}_{H-1}^* \cdots \mathbf{W}_1^* \mathbf{X} = \operatorname{svd}_{d_{\min}}(\tilde{\mathbf{Y}} - \Lambda^*)$. (by the global optimality and Lemma 3)

Before proceeding, we denote by $\tilde{\mathbf{A}} := \mathbf{W}_H^* \cdots \mathbf{W}_{\min+1}^*$, $\tilde{\mathbf{B}} := \mathbf{W}_{\min}^* \cdots \mathbf{W}_1^* \mathbf{X}$, where \mathbf{W}_{\min}^* is a matrix among $\{\mathbf{W}_i^*\}_{i=1}^{H-1}$ which has d_{\min} rows, and let

$$\mathcal{T} := \{\tilde{\mathbf{A}} \mathbf{C}_1^T + \mathbf{C}_2 \tilde{\mathbf{B}} : \mathbf{C}_1 \in \mathbb{R}^{n \times d_{\min}}, \mathbf{C}_2 \in \mathbb{R}^{d_H \times d_{\min}}\}$$

be a matrix space. Denote by \mathcal{U} the left singular space of $\tilde{\mathbf{A}} \tilde{\mathbf{B}}$ and \mathcal{V} the right singular space. Then the linear space \mathcal{T} can be equivalently represented as $\mathcal{T} = \mathcal{U} + \mathcal{V}$. Therefore, $\mathcal{T}^{\perp} = (\mathcal{U} + \mathcal{V})^{\perp} = \mathcal{U}^{\perp} \cap \mathcal{V}^{\perp}$. With this, we note that: (b) $\mathbf{W}_H^* \mathbf{W}_{H-1}^* \cdots \mathbf{W}_1^* \mathbf{X} + \Lambda^* - \tilde{\mathbf{Y}} \in \operatorname{Null}(\tilde{\mathbf{A}}^T) = \operatorname{Col}(\tilde{\mathbf{A}})^{\perp}$ and $\mathbf{W}_H^* \mathbf{W}_{H-1}^* \cdots \mathbf{W}_1^* \mathbf{X} + \Lambda^* - \tilde{\mathbf{Y}} \in \operatorname{Row}(\tilde{\mathbf{B}})^{\perp}$ (so $\mathbf{W}_H^* \mathbf{W}_{H-1}^* \cdots \mathbf{W}_1^* \mathbf{X} + \Lambda^* - \tilde{\mathbf{Y}} \in \mathcal{T}^{\perp}$) imply Equations (11) since either $\mathbf{W}_{i+1}^{*T} \cdots \mathbf{W}_H^{*T} (\mathbf{W}_H^* \mathbf{W}_{H-1}^* \cdots \mathbf{W}_1^* \mathbf{X} + \Lambda^* - \tilde{\mathbf{Y}}) = \mathbf{0}$ or $(\mathbf{W}_H^* \mathbf{W}_{H-1}^* \cdots \mathbf{W}_1^* \mathbf{X} + \Lambda^* - \tilde{\mathbf{Y}}) \mathbf{X}^T \mathbf{W}_1^{*T} \cdots \mathbf{W}_{i-1}^{*T} = \mathbf{0}$ for all i 's. And (c) for an orthogonal decomposition $\tilde{\mathbf{Y}} - \Lambda^* = \mathbf{W}_H^* \mathbf{W}_{H-1}^* \cdots \mathbf{W}_1^* \mathbf{X} + \mathbf{E}$ where $\mathbf{W}_H^* \mathbf{W}_{H-1}^* \cdots \mathbf{W}_1^* \mathbf{X} \in \mathcal{T}$ and $\mathbf{E} \in \mathcal{T}^{\perp}$, we have that

$$\|\mathbf{E}\| \leq \sigma_{d_{\min}}(\mathbf{W}_H^* \mathbf{W}_{H-1}^* \cdots \mathbf{W}_1^* \mathbf{X})$$

and condition (b) together imply $\mathbf{W}_H^* \mathbf{W}_{H-1}^* \cdots \mathbf{W}_1^* \mathbf{X} = \text{svd}_{d_{\min}}(\tilde{\mathbf{Y}} - \mathbf{\Lambda}^*)$ by Lemma 3. Therefore, the dual conditions in (12) are implied by

- (1) $\mathbf{\Lambda}^* \in \partial K(\mathbf{W}_H^* \mathbf{W}_{H-1}^* \cdots \mathbf{W}_1^* \mathbf{X}) \cap \text{Row}(\mathbf{X})$;
- (2) $\mathcal{P}_{\mathcal{T}}(\tilde{\mathbf{Y}} - \mathbf{\Lambda}^*) = \mathbf{W}_H^* \mathbf{W}_{H-1}^* \cdots \mathbf{W}_1^* \mathbf{X}$;
- (3) $\|\mathcal{P}_{\mathcal{T}^\perp}(\tilde{\mathbf{Y}} - \mathbf{\Lambda}^*)\| \leq \sigma_{d_{\min}}(\mathbf{W}_H^* \mathbf{W}_{H-1}^* \cdots \mathbf{W}_1^* \mathbf{X})$.

It thus suffices to construct a dual certificate $\mathbf{\Lambda}^*$ such that conditions (1), (2) and (3) hold, because conditions (1), (2) and (3) are stronger than conditions (a), (b) and (c). Let $r = \text{rank}(\tilde{\mathbf{Y}})$ and $\bar{r} = \min\{r, d_{\min}\}$. To proceed, we need the following lemma.

Lemma 5 ([6]). *Suppose $\tilde{\mathbf{Y}} \in \text{Row}(\mathbf{X})$. Let $(\mathbf{W}_1^*, \dots, \mathbf{W}_H^*)$ be the solution to problem (10) and let $\text{Udiag}(\sigma_1(\tilde{\mathbf{Y}}), \dots, \sigma_r(\tilde{\mathbf{Y}})) \mathbf{V}^T$ denote the skinny SVD of $\tilde{\mathbf{Y}} \in \text{Row}(\mathbf{X})$. We have $\mathbf{W}_H^* \mathbf{W}_{H-1}^* \cdots \mathbf{W}_1^* \mathbf{X} = \text{Udiag}((\sigma_1(\tilde{\mathbf{Y}}) - \gamma)_+, \dots, (\sigma_{\bar{r}}(\tilde{\mathbf{Y}}) - \gamma)_+, 0, \dots, 0) \mathbf{V}^T$.*

Recall that the sub-differential of the nuclear norm of a matrix \mathbf{Z} is

$$\partial_{\mathbf{Z}} \|\mathbf{Z}\|_* = \{\mathbf{U}_{\mathbf{Z}} \mathbf{V}_{\mathbf{Z}}^T + \mathbf{T}_{\mathbf{Z}} : \mathbf{T}_{\mathbf{Z}} \in \mathcal{T}^\perp, \|\mathbf{T}_{\mathbf{Z}}\| \leq 1\},$$

where $\mathbf{U}_{\mathbf{Z}} \mathbf{\Sigma}_{\mathbf{Z}} \mathbf{V}_{\mathbf{Z}}^T$ is the skinny SVD of the matrix \mathbf{Z} . So with Lemma 5, the sub-differential of (scaled) nuclear norm at optimizer $\mathbf{W}_H^* \mathbf{W}_{H-1}^* \cdots \mathbf{W}_1^* \mathbf{X}$ is given by

$$\partial(\gamma \|\mathbf{W}_H^* \mathbf{W}_{H-1}^* \cdots \mathbf{W}_1^* \mathbf{X}\|_*) = \{\gamma \mathbf{U}_{:,1:\bar{r}} \mathbf{V}_{:,1:\bar{r}}^T + \mathbf{T} : \mathbf{T} \in \mathcal{T}^\perp, \|\mathbf{T}\| \leq \gamma\}. \quad (13)$$

To construct the dual certificate, we set

$$\mathbf{\Lambda}^* = \underbrace{\gamma \mathbf{U}_{:,1:\bar{r}} \mathbf{V}_{:,1:\bar{r}}^T}_{\text{Component in space } \mathcal{T}} + \underbrace{\mathbf{U}_{:,\bar{r}+1:r} \text{diag}(\gamma, \dots, \gamma) \mathbf{V}_{:,\bar{r}+1:r}^T}_{\text{Component } \mathbf{T} \text{ in space } \mathcal{T}^\perp \text{ with } \|\mathbf{T}\| \leq \gamma} \in \text{Row}(\mathbf{X}),$$

where $\mathbf{\Lambda}^* \in \text{Row}(\mathbf{X})$ because $\mathbf{V}^T \in \text{Row}(\mathbf{X})$ (This is because \mathbf{V}^T is the right singular matrix of $\tilde{\mathbf{Y}}$ and $\tilde{\mathbf{Y}} \in \text{Row}(\mathbf{X})$). So condition (1) is satisfied according to (13). To see condition (2), $\mathcal{P}_{\mathcal{T}}(\tilde{\mathbf{Y}} - \mathbf{\Lambda}^*) = \mathcal{P}_{\mathcal{T}} \tilde{\mathbf{Y}} - \gamma \mathbf{U}_{:,1:\bar{r}} \mathbf{V}_{:,1:\bar{r}}^T = \text{Udiag}((\sigma_1(\tilde{\mathbf{Y}}) - \gamma)_+, \dots, (\sigma_{\bar{r}}(\tilde{\mathbf{Y}}) - \gamma)_+, 0, \dots, 0) \mathbf{V}^T = \mathbf{W}_H^* \mathbf{W}_{H-1}^* \cdots \mathbf{W}_1^* \mathbf{X}$, where the last equality is by Lemma 5 and the assumption $\sigma_{\min}(\tilde{\mathbf{Y}}) > \gamma$. As for condition (3), note that

$$\begin{aligned} \left\| \mathcal{P}_{\mathcal{T}^\perp}(\tilde{\mathbf{Y}} - \mathbf{\Lambda}^*) \right\| &= \left\| \mathbf{U}_{:,\bar{r}+1:r} \text{diag}(\sigma_{\bar{r}+1}(\tilde{\mathbf{Y}}) - \gamma, \dots, \sigma_r(\tilde{\mathbf{Y}}) - \gamma) \mathbf{V}_{:,\bar{r}+1:r}^T \right\| \\ &= \begin{cases} 0, & \text{if } \bar{r} = r, \\ \sigma_{d_{\min}+1}(\tilde{\mathbf{Y}}) - \gamma, & \text{otherwise.} \end{cases} \end{aligned}$$

By Lemma 5, $\sigma_{d_{\min}}(\mathbf{W}_H^* \mathbf{W}_{H-1}^* \cdots \mathbf{W}_1^* \mathbf{X}) \geq \|\mathcal{P}_{\mathcal{T}^\perp}(\tilde{\mathbf{Y}} - \mathbf{\Lambda}^*)\|$. So the proof of strong duality is completed, where the dual problem is given in Section D.

To see the relation between the solutions of primal and dual problems, it is a direct result of Lemmas 2 and 4.

D Dual Problem of Deep Linear Neural Network

In this section, we derive the dual problem of non-convex program (4). Denote by $G(\mathbf{W}_1, \dots, \mathbf{W}_H)$ the objective function of problem (4). Let $K(\cdot) = \gamma \|\cdot\|_*$, and let $\tilde{\mathbf{Y}} = \mathbf{Y} \mathbf{X}^\dagger \mathbf{X}$ be the projection of \mathbf{Y} on the row span of \mathbf{X} .

We note that

$$\begin{aligned}
 & \min_{\mathbf{W}_1, \dots, \mathbf{W}_H} G(\mathbf{W}_1, \dots, \mathbf{W}_H) - \frac{1}{2} \|\mathbf{Y} - \tilde{\mathbf{Y}}\|_F^2 \\
 &= \min_{\mathbf{W}_1, \dots, \mathbf{W}_H} \frac{1}{2} \|\mathbf{Y} - \mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X}\|_F^2 - \frac{1}{2} \|\mathbf{Y} - \tilde{\mathbf{Y}}\|_F^2 + K(\mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X}) \\
 &= \min_{\mathbf{W}_1, \dots, \mathbf{W}_H} \frac{1}{2} \|\tilde{\mathbf{Y}} - \mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X}\|_F^2 + K^{**}(\mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X}) \\
 &= \min_{\mathbf{W}_1, \dots, \mathbf{W}_H} \max_{\text{Row}(\mathbf{\Lambda}) \subseteq \text{Row}(\mathbf{X})} \frac{1}{2} \|\tilde{\mathbf{Y}} - \mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X}\|_F^2 + \langle \mathbf{\Lambda}, \mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X} \rangle - K^*(\mathbf{\Lambda}) \\
 &= \min_{\mathbf{W}_1, \dots, \mathbf{W}_H} \max_{\text{Row}(\mathbf{\Lambda}) \subseteq \text{Row}(\mathbf{X})} \frac{1}{2} \|\tilde{\mathbf{Y}} - \mathbf{\Lambda} - \mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X}\|_F^2 - \frac{1}{2} \|\mathbf{\Lambda}\|_F^2 - K^*(\mathbf{\Lambda}) + \langle \tilde{\mathbf{Y}}, \mathbf{\Lambda} \rangle,
 \end{aligned}$$

where the second equality holds since $K(\cdot)$ is closed and convex w.r.t. the argument $\mathbf{W}_H \mathbf{W}_{H-1} \cdots \mathbf{W}_1 \mathbf{X}$ and the third equality is by the definition of conjugate function of nuclear norm. Therefore, the dual problem is given by

$$\begin{aligned}
 & \max_{\text{Row}(\mathbf{\Lambda}) \subseteq \text{Row}(\mathbf{X})} \min_{\mathbf{W}_1, \dots, \mathbf{W}_H} \frac{1}{2} \|\tilde{\mathbf{Y}} - \mathbf{\Lambda} - \mathbf{W}_H \cdots \mathbf{W}_1 \mathbf{X}\|_F^2 - \frac{1}{2} \|\mathbf{\Lambda}\|_F^2 - K^*(\mathbf{\Lambda}) + \langle \tilde{\mathbf{Y}}, \mathbf{\Lambda} \rangle + \frac{1}{2} \|\mathbf{Y} - \tilde{\mathbf{Y}}\|_F^2 \\
 &= \max_{\text{Row}(\mathbf{\Lambda}) \subseteq \text{Row}(\mathbf{X})} \frac{1}{2} \sum_{i=d_{\min}+1}^{\min\{d_H, n\}} \sigma_i^2(\tilde{\mathbf{Y}} - \mathbf{\Lambda}) - \frac{1}{2} \|\tilde{\mathbf{Y}} - \mathbf{\Lambda}\|_F^2 - K^*(\mathbf{\Lambda}) + \frac{1}{2} \|\mathbf{Y}\|_F^2 \\
 &= \max_{\text{Row}(\mathbf{\Lambda}) \subseteq \text{Row}(\mathbf{X})} -\frac{1}{2} \|\tilde{\mathbf{Y}} - \mathbf{\Lambda}\|_{d_{\min}}^2 - K^*(\mathbf{\Lambda}) + \frac{1}{2} \|\mathbf{Y}\|_F^2,
 \end{aligned}$$

where $\|\cdot\|_{d_{\min}}^2 = \sum_{i=1}^{d_{\min}} \sigma_i^2(\cdot)$. We note that

$$K^*(\mathbf{\Lambda}) = \begin{cases} 0, & \|\mathbf{\Lambda}\| \leq \gamma; \\ +\infty, & \|\mathbf{\Lambda}\| > \gamma. \end{cases}$$

So the dual problem is given by

$$\max_{\text{Row}(\mathbf{\Lambda}) \subseteq \text{Row}(\mathbf{X})} -\frac{1}{2} \|\tilde{\mathbf{Y}} - \mathbf{\Lambda}\|_{d_{\min}}^2 + \frac{1}{2} \|\mathbf{Y}\|_F^2, \quad \text{s.t.} \quad \|\mathbf{\Lambda}\| \leq \gamma. \quad (14)$$

Problem (14) can be solved efficiently due to their convexity. In particular, Grussler et al. [1] provided a computationally efficient algorithm to compute the proximal operators of functions $\frac{1}{2} \|\cdot\|_r^2$. Hence, the Douglas-Rachford algorithm can find the global minimum up to an ϵ error in function value in time $\text{poly}(1/\epsilon)$ [2].

References

- [1] C. Grussler, A. Rantzer, and P. Giselsson. Low-rank optimization with convex constraints. *arXiv:1606.01793*, 2016.
- [2] B. He and X. Yuan. On the $O(1/n)$ convergence rate of the Douglas–Rachford alternating direction method. *SIAM Journal on Numerical Analysis*, 50(2):700–709, 2012.
- [3] T. L. Magnanti, J. F. Shapiro, and M. H. Wagner. Generalized linear programming solves the dual. *Management Science*, 22(11):1195–1203, 1976.
- [4] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [5] R. M. Starr. Quasi-equilibria in markets with non-convex preferences. *Econometrica: Journal of the Econometric Society*, pages 25–38, 1969.
- [6] M. Udell, C. Horn, R. Zadeh, and S. Boyd. Generalized low rank models. *Foundations and Trends in Machine Learning*, 9(1):1–118, 2016.