

The Expressive Power of Low-Rank Adaptation

Paper Reading - 2025/02/20

Fine-Tuning

- Foundation models (e.g., LLMs, Multi-modal) : large-scale neural networks trained on diverse, extensive datasets
 - Generalizable representational frameworks, can be adapted to downstream applications through fine-tuning (e.g., GPT → ChatGPT)
 - Parameters: billions or trillions (e.g., GPT-4 1.76 trillion)
 - Fine-tuning (type of transfer learning): Computational Cost
- Parameter-efficient fine-tuning

Low-Rank Adaptation (LoRA)

- Parameter-efficient fine-tuning method
- Full fine-tuning:

$$\max_{\theta \in \Phi} \hat{L}(f_{\theta})$$

Initial Param: θ^{t-1} . Update from (θ^{t-1}) to $(\theta^{t-1} + \Delta\theta)$ over Φ .

- LoRA:

$$\max_{\Delta\theta \in \Theta_k} \hat{L}(f_{\theta^{t-1} + \Delta\theta})$$

Where Θ_k is a smaller set, $|\Theta_k| \ll |\Phi|$.

Specifically, $\Theta_k = \{\Delta\theta \in R^{m \times n} : \text{rank}(\Delta\theta) \leq k\}$

Paper

- What is the minimum rank of the LoRA adapters required to adapt a (pre-trained) model f to match the functionality of the target model \bar{f} ?

Published as a conference paper at ICLR 2024

THE EXPRESSIVE POWER OF LOW-RANK ADAPTATION

Yuchen Zeng

Department of Computer Science
University of Wisconsin-Madison
yzeng58@wisc.edu

Kangwook Lee

Department of Electrical and Computer Engineering
University of Wisconsin-Madison
kangwook.lee@wisc.edu

ABSTRACT

Case 1: Linear Model

- Problem:

Frozen Model $f_0(\mathbf{x}) = \mathbf{W}_L \cdots \mathbf{W}_1 \mathbf{x} = \left(\prod_{l=1}^L \mathbf{W}_l \right) \mathbf{x}$

Target Model $\bar{f}(\mathbf{x}) = \bar{\mathbf{W}} \mathbf{x}$

$$\bar{\mathbf{W}}, \mathbf{W}_1, \dots, \mathbf{W}_L \in \mathbb{R}^{D \times D}$$

- Find $\Delta \mathbf{W}_1, \dots, \Delta \mathbf{W}_L$ s.t. $f = \bar{f}$,

Adapted Model $f(\mathbf{x}) = (\mathbf{W}_L + \Delta \mathbf{W}_L) \cdots (\mathbf{W}_1 + \Delta \mathbf{W}_1) \mathbf{x}$

where $\text{rank}(\Delta \mathbf{W}_l) \leq R$ for all $l \in [L]$.

Case 1: Linear Model

Lemma 1. Define error matrix $\mathbf{E} := \bar{\mathbf{W}} - \prod_{l=1}^L \mathbf{W}_l$, and denote its rank by $R_{\mathbf{E}} = \text{rank}(\mathbf{E})$. For a given LoRA-rank $R \in [D]$, assume that all the weight matrices of the frozen model $(\mathbf{W}_l)_{l=1}^L$, and $\prod_{l=1}^L \mathbf{W}_l + \text{LR}_r(\mathbf{E})$ are non-singular for all $r \leq R(L-1)$. Then, we have the following:

$$\min_{\Delta \mathbf{W}_l: \text{rank}(\Delta \mathbf{W}_l) \leq R} \left\| \prod_{l=1}^L (\mathbf{W}_l + \Delta \mathbf{W}_l) - \bar{\mathbf{W}} \right\|_2 = \sigma_{RL+1}(\mathbf{E}).$$

Thus, when $R \geq \lceil \frac{R_{\mathbf{E}}}{L} \rceil$, the optimal solution satisfies $\prod_{l=1}^L (\mathbf{W}_l + \Delta \mathbf{W}_l) = \bar{\mathbf{W}}$, implying $f = \bar{f}$.

→ approximate $R \times L$ ranks of the error \mathbf{E} .

Proof:

- Solve constrained optimization:

$$\min_{\Delta \mathbf{W}_l: \text{rank}(\Delta \mathbf{W}_l) \leq R} \left\| \prod_{l=1}^L (\mathbf{W}_l + \Delta \mathbf{W}_l) - \bar{\mathbf{W}} \right\|_2.$$

By subtracting $\prod_{l=1}^L \mathbf{W}_l$ from both terms, the constrain optimization problem becomes

$$\min_{\Delta \mathbf{W}_l: \text{rank}(\Delta \mathbf{W}_l) \leq R} \left\| \underbrace{\left(\prod_{l=1}^L (\mathbf{W}_l + \Delta \mathbf{W}_l) - \prod_{l=1}^L \mathbf{W}_l \right)}_{:=\mathbf{A}} - \underbrace{\left(\bar{\mathbf{W}} - \prod_{l=1}^L \mathbf{W}_l \right)}_{:=\mathbf{E}} \right\|_2. \quad (2)$$

To perform analysis on (2), we start with the analysis of \mathbf{A} as follows:

$$\begin{aligned} \mathbf{A} &= \prod_{l=1}^L (\Delta \mathbf{W}_l + \mathbf{W}_l) - \prod_{l=1}^L \mathbf{W}_l \\ &= \Delta \mathbf{W}_L \prod_{l=1}^{L-1} (\Delta \mathbf{W}_l + \mathbf{W}_l) + \mathbf{W}_L \prod_{l=1}^{L-1} (\Delta \mathbf{W}_l + \mathbf{W}_l) - \prod_{l=1}^L \mathbf{W}_l. \end{aligned}$$

At this point, it becomes clear that this expression can be iteratively decomposed. Following this pattern, we can express \mathbf{A} as:

$$\begin{aligned}
\mathbf{A} &= \Delta \mathbf{W}_L \prod_{l=1}^{L-1} (\Delta \mathbf{W}_l + \mathbf{W}_l) + \mathbf{W}_L \Delta \mathbf{W}_{L-1} \prod_{l=1}^{L-2} (\Delta \mathbf{W}_l + \mathbf{W}_l) \\
&\quad + \dots + \left(\prod_{l=2}^L \mathbf{W}_l \right) (\Delta \mathbf{W}_1 + \mathbf{W}_1) - \prod_{l=1}^L \mathbf{W}_l \\
&= \sum_{l=1}^L \underbrace{\left[\left(\prod_{i=l+1}^L \mathbf{W}_i \right) \Delta \mathbf{W}_l \left(\prod_{i=1}^{l-1} (\mathbf{W}_i + \Delta \mathbf{W}_i) \right) \right]}_{:= \mathbf{A}_l}.
\end{aligned} \tag{3}$$

In this final form, \mathbf{A} is decomposed as $\mathbf{A} = \sum_{l=1}^L \mathbf{A}_l$. It is important to note that $\text{rank}(\mathbf{A}_l) \leq \text{rank}(\Delta \mathbf{W}_l) \leq R$. Consequently, $\text{rank}(\mathbf{A}) \leq \sum_{l=1}^L \text{rank}(\mathbf{A}_l) \leq RL$.

Then, the optimization problem (2) can be relaxed into a low-rank approximation problem

$$(2) \geq \min_{\mathbf{A}: \text{rank}(\mathbf{A}) \leq RL} \|\mathbf{A} - \mathbf{E}\|_2, \quad (4)$$

$$\min_{\Delta \mathbf{W}_l: \text{rank}(\Delta \mathbf{W}_l) \leq R} \left\| \underbrace{\left(\prod_{l=1}^L (\mathbf{W}_l + \Delta \mathbf{W}_l) - \prod_{l=1}^L \mathbf{W}_l \right)}_{:= \mathbf{A}} - \underbrace{\left(\overline{\mathbf{W}} - \prod_{l=1}^L \mathbf{W}_l \right)}_{:= \mathbf{E}} \right\|_2. \quad (2)$$

where the optimal solution is $\mathbf{A} = \text{LR}_{RL \wedge R_E}(\mathbf{E}) := \mathbf{E}'$. Therefore, if we can identify rank- R or lower matrices $(\Delta \mathbf{W}_l)_{l=1}^L$ such that

$$\underbrace{\prod_{l=1}^L (\mathbf{W}_l + \Delta \mathbf{W}_l) - \prod_{l=1}^L \mathbf{W}_l}_{:= \mathbf{A}} = \underbrace{\text{LR}_{RL \wedge R_E}(\overline{\mathbf{W}} - \prod_{l=1}^L \mathbf{W}_l)}_{:= \mathbf{E}'}, \quad (5)$$

Denote $R_{E'} = RL \wedge R_E$. To derive the explicit form of E' , we first refer to the SVD of E as

$$E = UDV^\top,$$

where U and V are orthonormal matrices and the first R_E diagonal entries of D are non-zero, with all remaining entries being zero. Based on this, E' is expressed as

$$E' = UDI_{1:RL,D}V^\top.$$

Having already derived the decomposition $A = \sum_{l=1}^L A_l$, we next aim to decompose E' as $E' = \sum_{l=1}^L E'Q_l$, where $Q_1, \dots, Q_L \in \mathbb{R}^{D \times D}$. The goal now shifts to identifying $\Delta W_l, Q_l$ such that $A_l = E'Q_l$ for each $l \in [L]$. Achieving this would complete the proof of (5).

Therefore, our goal becomes finding $\Delta W_1, \dots, \Delta W_L$ with $\text{rank}(\Delta W_l) \leq R$ for all $l \in [L]$ such that

$$A_l = \left(\prod_{i=l+1}^L W_i \right) \Delta W_l \left(\prod_{i=1}^{l-1} (W_i + \Delta W_i) \right) = E'Q_l, \quad \text{for all } l \in [L]. \quad (6)$$

One sufficient condition for achieving (6) is that the decomposed matrices $\mathbf{Q}_1, \mathbf{Q}_L$ and low-rank adapters $\Delta \mathbf{W}_1, \dots, \Delta \mathbf{W}_L$ meet the following conditions:

$$\sum_{l=1}^L \mathbf{E}' \mathbf{Q}_l = \mathbf{E}', \quad (7)$$

$$\Delta \mathbf{W}_l = \left(\prod_{i=l+1}^L \mathbf{W}_i \right)^{-1} \mathbf{E}' \mathbf{Q}_l \left(\prod_{i=1}^{l-1} (\mathbf{W}_i + \Delta \mathbf{W}_i) \right)^{-1}, \text{ for all } l \in [L] \quad (8)$$

$$\text{rank}(\Delta \mathbf{W}_l) \leq R, \text{ for all } l \in [L], \quad (9)$$

$$\text{rank}(\mathbf{W}_l + \Delta \mathbf{W}_l) = D, \text{ for all } l \in [L-1]. \quad (10)$$

$$\mathbf{A}_l = \left(\prod_{i=l+1}^L \mathbf{W}_i \right) \Delta \mathbf{W}_l \left(\prod_{i=1}^{l-1} (\mathbf{W}_i + \Delta \mathbf{W}_i) \right) = \mathbf{E}' \mathbf{Q}_l, \quad \text{for all } l \in [L].$$

- Consider

$$E = (\sigma_1 u_1 v_1^T + \cdots + \sigma_R u_R v_R^T) + (\dots) + (\dots)$$

We will show that the matrices $(Q_l)_{l=1}^L$ defined by

$$\begin{aligned} Q_l &= V I_{(R(l-1)+1) \wedge R_{E'} : Rl \wedge R_{E'}, D} V^\top, \quad \text{for all } l \in [L], \\ \Delta W_l &= \left(\prod_{i=l+1}^L W_i \right)^{-1} E' Q_l \left(\prod_{i=1}^{l-1} (W_i + \Delta W_i) \right)^{-1}, \quad \text{for all } l \in [L] \end{aligned} \quad (11)$$

When $l = 1$. We begin by examining the three conditions (8), (9) and (10) under the base case $l = 1$. We first determine Q_1 and ΔW_1 based on (11) and (8):

$$\Delta W_1 = \left(\prod_{i=2}^L W_i \right)^{-1} E' Q_1 \quad (12)$$

By the choice of ΔW_1 , we satisfy the condition (8). Moreover, it directly follows that $\text{rank}(\Delta W_1) \leq \text{rank}(Q_1) = R$, thereby fulfilling the rank constraint in (9).

Therefore, we just need to prove that $\mathbf{W}_1 + \Delta\mathbf{W}_1$ is full-rank, as required by condition (10). To compute $\text{rank}(\mathbf{W}_1 + \Delta\mathbf{W}_1)$, we proceed as follows:

$$\begin{aligned}
& \text{rank}(\mathbf{W}_1 + \Delta\mathbf{W}_1) \\
& \stackrel{(12)}{=} \text{rank}(\mathbf{W}_1 + (\prod_{i=2}^L \mathbf{W}_i)^{-1} \mathbf{E}' \mathbf{Q}_1) && \text{(Substituting for } \Delta\mathbf{W}_1) \\
& = \text{rank}((\prod_{i=1}^L \mathbf{W}_i) + \mathbf{E}' \mathbf{Q}_1) && \text{(Left multiplying with invertible } (\prod_{i=2}^L \mathbf{W}_i)^{-1}) \\
& = \text{rank}((\prod_{i=1}^L \mathbf{W}_i) + \text{LR}_{R \wedge R_{\mathbf{E}'}}(\mathbf{E})). && \text{(Simplifying)}
\end{aligned}$$

Given the assumption that $\prod_{l=1}^L \mathbf{W}_l + \text{LR}_r(\mathbf{E})$ is full rank for all $r \leq R(L-1)$, $\text{rank}(\mathbf{W}_1 + \Delta\mathbf{W}_1) = \text{rank}((\prod_{i=1}^L \mathbf{W}_i) + \text{LR}_{R \wedge R_{\mathbf{E}'}}(\mathbf{E})) = D$, satisfying the last condition (10).

When $l > 1$. Consider $l = 2, \dots, L$. We assume that for $i \in [l-1]$, we have determined matrices \mathbf{Q}_i and $\Delta \mathbf{W}_i$ based on (11) and (8), respectively, and we assume that they satisfy the conditions (8), (9), and (10).

First, under the induction assumption that $\mathbf{W}_i + \Delta \mathbf{W}_i$ is invertible for all $i \in [l-1]$, to achieve $\mathbf{A}_l = \mathbf{E}' \mathbf{Q}_l$, we set $\Delta \mathbf{W}_l$ based on (8). This definition ensures $\text{rank}(\Delta \mathbf{W}_l) \leq \text{rank}(\mathbf{Q}_l) = R$, thereby satisfying the condition (9). To prove that $\mathbf{W}_l + \Delta \mathbf{W}_l$ is full-rank (condition (10)), we focus on computing $\text{rank}(\mathbf{W}_l + \Delta \mathbf{W}_l)$. We proceed as follows:

$$\begin{aligned}
& \text{rank}(\mathbf{W}_l + \Delta \mathbf{W}_l) \\
& \stackrel{(8)}{=} \text{rank}\left(\mathbf{W}_l + \left(\prod_{i=l+1}^L \mathbf{W}_i\right)^{-1} \mathbf{E}' \mathbf{Q}_l \left(\prod_{i=1}^{l-1} (\mathbf{W}_i + \Delta \mathbf{W}_i)^{-1}\right)\right) \quad (\text{Substituting for } \Delta \mathbf{W}_l) \\
& = \text{rank}\left(\mathbf{I}_D + \left(\prod_{i=l}^L \mathbf{W}_i\right)^{-1} \mathbf{E}' \mathbf{Q}_l \left(\prod_{i=1}^{l-1} (\mathbf{W}_i + \Delta \mathbf{W}_i)\right)^{-1}\right) \quad (\text{Left multiplying invertible } \mathbf{W}_l^{-1}) \\
& = \text{rank}\left(\prod_{i=1}^{l-1} (\mathbf{W}_i + \Delta \mathbf{W}_i) + \left(\prod_{i=l}^L \mathbf{W}_i\right)^{-1} \mathbf{E}' \mathbf{Q}_l\right) \quad (\text{Right multiplying invertible } \prod_{i=1}^{l-1} (\mathbf{W}_i + \Delta \mathbf{W}_i)) \\
& = \text{rank}\left((\mathbf{W}_{l-1} + \Delta \mathbf{W}_{l-1}) \prod_{i=1}^{l-2} (\mathbf{W}_i + \Delta \mathbf{W}_i) + \left(\prod_{i=l}^L \mathbf{W}_i\right)^{-1} \mathbf{E}' \mathbf{Q}_l\right) \quad (\text{Rearranging terms})
\end{aligned}$$

$$\begin{aligned}
&\stackrel{(8)}{=} \text{rank} \left((\mathbf{W}_{l-1} + (\prod_{i=l}^L \mathbf{W}_i)^{-1} \mathbf{E}' \mathbf{Q}_{l-1} (\prod_{i=1}^{l-2} (\mathbf{W}_i + \Delta \mathbf{W}_i))^{-1}) \prod_{i=1}^{l-2} (\mathbf{W}_i + \Delta \mathbf{W}_i) \right. \\
&\quad \left. + (\prod_{i=l}^L \mathbf{W}_i)^{-1} \mathbf{E}' \mathbf{Q}_l \right) \quad \text{(Substituting for } \Delta \mathbf{W}_{l-1} \text{)} \\
&= \text{rank} \left((\prod_{i=l-1}^L \mathbf{W}_i + \mathbf{E}' \mathbf{Q}_{l-1} (\prod_{i=1}^{l-2} (\mathbf{W}_i + \Delta \mathbf{W}_i))^{-1}) \prod_{i=1}^{l-2} (\mathbf{W}_i + \Delta \mathbf{W}_i) \right. \\
&\quad \left. + \mathbf{E}' \mathbf{Q}_l \right) \quad \text{(Left multiplying } \prod_{i=l}^L \mathbf{W}_i \text{)} \\
&= \text{rank} \left((\prod_{i=l-1}^L \mathbf{W}_i \prod_{i=1}^{l-2} (\mathbf{W}_i + \Delta \mathbf{W}_i) + \mathbf{E}' \mathbf{Q}_{l-1} + \mathbf{E}' \mathbf{Q}_l) \right) \quad \text{(Rearranging terms)} \\
&= \dots \\
&= \text{rank} (\prod_{i=1}^L \mathbf{W}_i + \mathbf{E}' (\sum_{i=1}^l \mathbf{Q}_i)) \quad \text{(Taking similar steps)} \\
&= \text{rank} (\prod_{i=1}^L \mathbf{W}_i + \text{LR}_{Rl \wedge R_{\mathbf{E}'}}(\mathbf{E})). \quad \text{(Simplifying)}
\end{aligned}$$

By the assumption that $\prod_{l=1}^L \mathbf{W}_l + \text{LR}_r(\mathbf{E})$ is full-rank for $r \leq R(L-1)$ and consequently, $\text{rank}(\mathbf{W}_l + \Delta \mathbf{W}_l) = \text{rank}(\prod_{i=1}^L \mathbf{W}_i + \text{LR}_{Rl \wedge R_{\mathbf{E}'}}(\mathbf{E})) = D$, satisfying the last condition (10).

$$\min_{\Delta \mathbf{W}_l: \text{rank}(\Delta \mathbf{W}_l) \leq R} \left\| \underbrace{\left(\prod_{l=1}^L (\mathbf{W}_l + \Delta \mathbf{W}_l) - \prod_{l=1}^L \mathbf{W}_l \right)}_{:= \mathbf{A}} - \underbrace{\left(\overline{\mathbf{W}} - \prod_{l=1}^L \mathbf{W}_l \right)}_{:= \mathbf{E}} \right\|_2$$

$$(2) \geq \min_{\mathbf{A}: \text{rank}(\mathbf{A}) \leq RL} \|\mathbf{A} - \mathbf{E}\|_2$$

$$\underbrace{\prod_{l=1}^L (\mathbf{W}_l + \Delta \mathbf{W}_l) - \prod_{l=1}^L \mathbf{W}_l}_{:= \mathbf{A}} = \underbrace{\text{LR}_{RL \wedge R_E}(\overline{\mathbf{W}} - \prod_{l=1}^L \mathbf{W}_l)}_{:= \mathbf{E}'}$$

$$\mathbf{A} = \sum_{l=1}^L \mathbf{A}_l, \quad \sum_{l=1}^{\sim} \mathbf{E}' \mathbf{Q}_l = \mathbf{E}'$$

$$\mathbf{A}_l = \left(\prod_{i=l+1}^L \mathbf{W}_i \right) \Delta \mathbf{W}_l \left(\prod_{i=1}^{l-1} (\mathbf{W}_i + \Delta \mathbf{W}_i) \right) = \mathbf{E}' \mathbf{Q}_l.$$

$$\Delta \mathbf{W}_l = \left(\prod_{i=l+1}^L \mathbf{W}_i \right)^{-1} \mathbf{E}' \mathbf{Q}_l \left(\prod_{i=1}^{l-1} (\mathbf{W}_i + \Delta \mathbf{W}_i) \right)^{-1}$$

Case 2: 1-Layer ReLU FNN

- Suppose target is 1-layer ReLU FNN

Lemma 9 (Detailed version of Lemma 2). *Define error matrix $\mathbf{E} := \bar{\mathbf{W}}_1 - \prod_{l=1}^L \mathbf{W}_l$, with its rank represented by $R_{\mathbf{E}} = \text{rank}(\mathbf{E})$. Consider a LoRA-rank $R \in [D]$. Assume that the weight matrices $\mathbf{W}_1, \dots, \mathbf{W}_L \in \mathbb{R}^{D \times D}$ and $\prod_{l=1}^L \mathbf{W}_l + \text{LR}_r(\mathbf{E})$ for all $r \leq R(L-1)$ are non-singular. Let \mathbf{x} be a random input sampled from a distribution with bounded support \mathcal{X} and let $\Sigma = \mathbb{E}\mathbf{x}\mathbf{x}^\top$. Then, there exists rank- R or lower matrices $\Delta\mathbf{W}_1, \dots, \Delta\mathbf{W}_L \in \mathbb{R}^{D \times D}$ and bias vectors $\hat{\mathbf{b}}_1, \dots, \hat{\mathbf{b}}_L \in \mathbb{R}^D$ such that for any input $\mathbf{x} \in \mathcal{X}$,*

$$f(\mathbf{x}) - \bar{f}(\mathbf{x}) = \text{ReLU} \left(\left(\text{LR}_{RL \wedge R_{\mathbf{E}}}(\bar{\mathbf{W}}_1 - \prod_{l=1}^L \mathbf{W}_l) - (\bar{\mathbf{W}}_1 - \prod_{l=1}^L \mathbf{W}_l) \right) \mathbf{x} \right).$$

Therefore, when $R \geq \lceil R_{\mathbf{E}}/L \rceil$, the adapted model exactly approximates the target model, i.e., $f(\mathbf{x}) = \bar{f}(\mathbf{x})$ for all $\mathbf{x} \in \mathcal{X}$.

Furthermore, let \mathbf{x} be a random input sampled from a distribution with bounded support \mathcal{X} and let $\Sigma = \mathbb{E}\mathbf{x}\mathbf{x}^\top$. Then, the expected squared error is bounded as

$$\mathbb{E} \|f(\mathbf{x}) - \bar{f}(\mathbf{x})\|_2^2 \leq \|\Sigma\|_F \sigma_{RL \wedge R_{\mathbf{E}}+1}^2 (\bar{\mathbf{W}}_1 - \prod_{l=1}^L \mathbf{W}_l).$$

Proof:

Linearization. The main challenge here stems from the non-linearities introduced by the ReLU activation function. To remove the non-linearities in the first $L - 1$ layers of updated model f , since the input space \mathcal{X} is bounded, we can set all the entries of $\hat{\mathbf{b}}_1, \dots, \hat{\mathbf{b}}_{L-1}$ sufficiently large, thereby

activating all ReLUs in the first $L - 1$ layers of f . Consequently, we have

$$\begin{aligned} f(\mathbf{x}) &= \text{ReLU}((\mathbf{W}_L + \Delta\mathbf{W}_L)\mathbf{z}_{L-1} + \hat{\mathbf{b}}_L) \\ &= \text{ReLU}\left((\mathbf{W}_L + \Delta\mathbf{W}_L)\text{ReLU}((\mathbf{W}_{L-1} + \Delta\mathbf{W}_{L-1})\mathbf{z}_{L-2} + \hat{\mathbf{b}}_{L-1}) + \hat{\mathbf{b}}_L\right) \\ &= \text{ReLU}\left((\mathbf{W}_L + \Delta\mathbf{W}_L)((\mathbf{W}_{L-1} + \Delta\mathbf{W}_{L-1})\mathbf{z}_{L-2} + \hat{\mathbf{b}}_{L-1}) + \hat{\mathbf{b}}_L\right) \\ &= \text{ReLU}\left((\mathbf{W}_L + \Delta\mathbf{W}_L)(\mathbf{W}_{L-1} + \Delta\mathbf{W}_{L-1})\mathbf{z}_{L-2} + (\mathbf{W}_L + \Delta\mathbf{W}_L)\hat{\mathbf{b}}_{L-1} + \hat{\mathbf{b}}_L\right) \\ &= \dots \\ &= \text{ReLU}\left(\prod_{l=1}^L (\mathbf{W}_l + \Delta\mathbf{W}_l)\mathbf{x} + \left(\sum_{l=1}^{L-1} \prod_{i=l+1}^L (\mathbf{W}_i + \Delta\mathbf{W}_i)\hat{\mathbf{b}}_l\right) + \hat{\mathbf{b}}_L\right), \end{aligned}$$

which is equivalent to a single-layer ReLU neural network with weight matrix $\prod_{l=1}^L (\mathbf{W}_l + \Delta\mathbf{W}_l)$ and bias vector $(\sum_{l=1}^{L-1} \prod_{i=l+1}^L (\mathbf{W}_i + \Delta\mathbf{W}_i)\hat{\mathbf{b}}_l) + \hat{\mathbf{b}}_L$.

Parameter Alignment. To match the updated model $f(\mathbf{x})$ and target model $\bar{f}(\mathbf{x})$, we proceed as follows. For weight matrix, Lemma 7 guarantees the existence of rank- R or lower matrices $\Delta \mathbf{W}_1, \dots, \Delta \mathbf{W}_L \in \mathbb{R}^{D \times D}$ such that

$$\prod_{l=1}^L (\mathbf{W}_l + \Delta \mathbf{W}_l) = \prod_{l=1}^L \mathbf{W}_l + \text{LR}_{RL \wedge R_E}(\bar{\mathbf{W}} - \prod_{l=1}^L \mathbf{W}_l). \quad (14)$$

For the bias vector, we set $\hat{\mathbf{b}}_L = \bar{\mathbf{b}}_1 - \sum_{l=1}^{L-1} \prod_{i=l+1}^L (\mathbf{W}_i + \Delta \mathbf{W}_i) \hat{\mathbf{b}}_l$ such that $\sum_{l=1}^{L-1} \prod_{i=l+1}^L (\mathbf{W}_i + \Delta \mathbf{W}_i) \hat{\mathbf{b}}_l + \hat{\mathbf{b}}_L = \bar{\mathbf{b}}_1$. Therefore, we obtain

$$f(\mathbf{x}) - \bar{f}(\mathbf{x}) = \text{ReLU} \left(\left(\text{LR}_{RL \wedge R_E}(\bar{\mathbf{W}}_1 - \prod_{l=1}^L \mathbf{W}_l) - (\bar{\mathbf{W}}_1 - \prod_{l=1}^L \mathbf{W}_l) \right) \mathbf{x} \right).$$

Error Derivation. We compute the expected squared error as follows:

$$\begin{aligned} & \mathbb{E} \|f(\mathbf{x}) - \bar{f}(\mathbf{x})\|_2^2 \\ & \leq \mathbb{E} \left\| \left(\text{LR}_{RL \wedge R_E}(\bar{\mathbf{W}}_1 - \prod_{l=1}^L \mathbf{W}_l) - (\bar{\mathbf{W}}_1 - \prod_{l=1}^L \mathbf{W}_l) \right) \mathbf{x} \right\|_2^2 \quad (\text{ReLU is 1-Lipschitz}) \\ & \stackrel{(1)}{\leq} \left\| \text{LR}_{RL \wedge R_E}(\bar{\mathbf{W}}_1 - \prod_{l=1}^L \mathbf{W}_l) - (\bar{\mathbf{W}}_1 - \prod_{l=1}^L \mathbf{W}_l) \right\|_2^2 \mathbb{E} \|\mathbf{x}\|_2^2 \\ & = \|\Sigma\|_{\text{F}} \sigma_{RL \wedge R_E + 1}^2 (\bar{\mathbf{W}}_1 - \prod_{l=1}^L \mathbf{W}_l). \quad (\text{By the definition of } \text{LR}_{RL \wedge R_E}(\cdot)) \end{aligned}$$

This completes the proof. \square

Case 3: L-Layer ReLU FNN

- Model Partition:

Example 1. Consider the case where $\bar{L} = 2$ and $L = 4$. We view a two-layer target model \bar{f} as a composition of two one-layer ReLU FNNs. Accordingly, we partition the four-layer adapted model f into two submodels, each consisting of two layers. For each layer in the target model, we utilize two corresponding layers in the frozen/adapted model for approximation. This problem then simplifies into a one-layer FNN approximation problem, which has already been addressed in Lemma 2.

Based on this example, we introduce a ordered partition $\mathcal{P} = \{P_1, \dots, P_{\bar{L}}\}$ to partition the layers in the adapted model f , where $\bigcup_{i=1}^{\bar{L}} P_i = [L]$. Each element $P_i \in \mathcal{P}$ consists of consecutive integers. Given a partition \mathcal{P} , each element P_i specifies that the layers with index $l \in P_i$ in the adapted model will be used to approximate the i -th layer in the target model. Example 1, which uses every two layers in the adapted model to approximate each layer in the target model, can be considered as a partition represented as $\{\{1, 2\}, \{3, 4\}\}$. Similarly, we extend this simple uniform partition into general cases for \bar{L} -layer target FNN and L -layer frozen FNN:

$$\mathcal{P}^u = \{P_1^u, \dots, P_{\bar{L}}^u\} := \{\{1, \dots, M\}, \{M+1, \dots, 2M\}, \dots, \{(\bar{L}-1)M+1, \dots, L\}\},$$

where $M := \lfloor L/\bar{L} \rfloor$. The uniform partition indicates that every M layers in the adapted model are employed to approximate each layer in the target model. We use $\prod_{l \in P_i} \mathbf{W}_l$ to denote the product of the weight matrices from the layers $l \in P_i$, with the later layer positioned to the left and the earlier layer to the right in the matrix product. For example, $\prod_{l \in P_1^u} \mathbf{W}_l = \prod_{l=1}^M \mathbf{W}_l = \mathbf{W}_M \cdots \mathbf{W}_1$.

Case 3: L-Layer ReLU FNN

Theorem 5. Define the approximation error of i -th layer as $E_i = \sigma_{RM+1}(\bar{\mathbf{W}}_i - \prod_{l \in P_i^u} \mathbf{W}_l)$, and the magnitude of the parameters and the input as $\beta := \max_{i \in [\bar{L}]} \left(\sqrt{\|\Sigma\|_F} \prod_{j=1}^i \|\bar{\mathbf{W}}_j\|_F + \sum_{j=1}^i \prod_{k=j+1}^{i-1} \|\bar{\mathbf{W}}_k\|_F \|\bar{\mathbf{b}}_j\|_2 \right) \vee \sqrt{\|\Sigma\|_F}$.

Under Assumption 1, there exists rank- R or lower matrices $(\Delta \mathbf{W}_l)_{l=1}^L$ with $\Delta \mathbf{W}_l \in \mathbb{R}^{D \times D}$ and bias vectors $(\hat{\mathbf{b}}_l)_{l=1}^L$ with $\hat{\mathbf{b}}_l \in \mathbb{R}^D$ such that for input $\mathbf{x} \in \mathcal{X}$ with $\mathbb{E} \mathbf{x} \mathbf{x}^\top = \Sigma$,

$$\mathbb{E} \|f(\mathbf{x}) - \bar{f}(\mathbf{x})\|_2 \leq \beta \sum_{i=1}^{\bar{L}} \max_{k \in [\bar{L}]} (\|\bar{\mathbf{W}}_k\|_F + E_k)^{\bar{L}-i} E_i.$$

Proof:

Model Decomposition. We partition the adapted model f into \bar{L} sub-models, each defined as

$$f_i(\cdot) = \text{FNN}_{\bar{L}, D}(\cdot; (\mathbf{W}_l + \Delta \mathbf{W}_l)_{l \in P_i^u}, (\hat{\mathbf{b}}_l)_{l \in P_i^u}), \quad i \in [\bar{L}].$$

In a similar manner, we break down \bar{f} into \bar{L} sub-models, each is a one-layer FNN:

$$\bar{f}_i(\cdot) = \text{FNN}_{1, D}(\cdot; \bar{\mathbf{W}}_i, \bar{\mathbf{b}}_i), \quad i \in [\bar{L}].$$

We can then express $f(\mathbf{x})$ and $\bar{f}(\mathbf{x})$ as compositions of their respective sub-models:

$$f(\cdot) = f_{\bar{L}} \circ \cdots \circ f_1(\cdot), \quad \bar{f}(\cdot) = \bar{f}_{\bar{L}} \circ \cdots \circ \bar{f}_1(\cdot).$$

To analyze the error $\mathbb{E} \|f(\mathbf{x}) - \bar{f}(\mathbf{x})\|_2 = \mathbb{E} \|f(\mathbf{x}) - \bar{f}(\mathbf{x})\|_2$, we consider the error caused by each submodel. Let $\tilde{R}_i = \text{rank}(\bar{\mathbf{W}}_i - \prod_{l \in P_i^u} \mathbf{W}_l)$ denote the rank of the discrepancy between the target weight matrix and the frozen weight matrices, where $i \in [\bar{L}]$. By Lemma 9, we can select $\Delta \mathbf{W}_1, \dots, \Delta \mathbf{W}_L, \hat{\mathbf{b}}_1, \dots, \hat{\mathbf{b}}_L$ such that

$$f_i(\mathbf{z}) - \bar{f}_i(\mathbf{z}) = \text{ReLU} \left(\left(\text{LR}_{RL \wedge \tilde{R}_i}(\bar{\mathbf{W}}_i - \prod_{l \in P_i^u} \mathbf{W}_l) - (\bar{\mathbf{W}}_i - \prod_{l \in P_i^u} \mathbf{W}_l) \right) \mathbf{z} \right), \quad (15)$$

$$\mathbb{E} \|f_i(\mathbf{z}) - \bar{f}_i(\mathbf{z})\|_2^2 \leq \|\mathbb{E} \mathbf{z} \mathbf{z}^\top\|_{\text{F}} \sigma_{RL \wedge \tilde{R}_i + 1}^2 (\bar{\mathbf{W}}_i - \prod_{l=1}^L \mathbf{W}_l). \quad (16)$$

Error Decomposition. For submodel $i = 2, \dots, \bar{L}$, we calculate the expected error of the composition of the first i sub-models,

$$\begin{aligned}
\mathbb{E} \|\widehat{\mathbf{z}}_i - \bar{\mathbf{z}}_i\|_2 &= \mathbb{E} \|f_i(\widehat{\mathbf{z}}_{i-1}) - \bar{f}_i(\bar{\mathbf{z}}_{i-1})\|_2 & (17) \\
&= \mathbb{E} \|(f_i(\widehat{\mathbf{z}}_{i-1}) - f_i(\bar{\mathbf{z}}_{i-1})) + (f_i(\bar{\mathbf{z}}_{i-1}) - \bar{f}_i(\bar{\mathbf{z}}_{i-1}))\|_2 & \text{(Rearranging terms)} \\
&\leq \underbrace{\mathbb{E} \|f_i(\widehat{\mathbf{z}}_{i-1}) - f_i(\bar{\mathbf{z}}_{i-1})\|_2}_{A_i} + \underbrace{\mathbb{E} \|f_i(\bar{\mathbf{z}}_{i-1}) - \bar{f}_i(\bar{\mathbf{z}}_{i-1})\|_2}_{B_i}. & \text{(Applying triangle inequality)}
\end{aligned}$$

Here A_i represents the error resulting from the discrepancy between the first $i - 1$ submodels, while B_i represents the error arising from the mismatch between the i -th submodel.

Computing A_i . We start by computing the error introduced by the first $i - 1$ submodels, denoted by A_i :

$$\begin{aligned}
A_i &= \mathbb{E} \|f_i(\widehat{\mathbf{z}}_{i-1}) - f_i(\bar{\mathbf{z}}_{i-1})\|_2 = \mathbb{E} \left\| \text{ReLU}(\widehat{\mathbf{W}}_i(\widehat{\mathbf{z}}_{i-1} - \bar{\mathbf{z}}_{i-1})) \right\|_2 \\
&\leq \mathbb{E} \left\| \widehat{\mathbf{W}}_i(\widehat{\mathbf{z}}_{i-1} - \bar{\mathbf{z}}_{i-1}) \right\|_2 & \text{(ReLU is 1-Lipschitz)} \\
&\stackrel{(1)}{\leq} \left\| \widehat{\mathbf{W}}_i \right\|_{\text{F}} \mathbb{E} \|\widehat{\mathbf{z}}_{i-1} - \bar{\mathbf{z}}_{i-1}\|_2. & (18)
\end{aligned}$$

$$\widehat{\mathbf{W}}_i = \text{LR}_{RL \wedge \tilde{R}_i}(\bar{\mathbf{W}}_1 - \prod_{l \in P_i^u} \mathbf{W}_l) + \prod_{l \in P_i^u} \mathbf{W}_l, \quad \prod_{l=1}^L (\mathbf{W}_l + \Delta \mathbf{W}_l) = \prod_{l=1}^L \mathbf{W}_l + \text{LR}_{RL \wedge R_E}(\bar{\mathbf{W}} - \prod_{l=1}^L \mathbf{W}_l).$$

Here,

$$\begin{aligned}
\|\widehat{\mathbf{W}}_i\|_{\text{F}} &= \left\| \prod_{l \in P_i^u} \mathbf{W}_l + \text{LR}_{RM \wedge \tilde{R}_i}(\overline{\mathbf{W}}_i - \prod_{l \in P_i^u} \mathbf{W}_l) \right\|_{\text{F}} \\
&= \left\| \overline{\mathbf{W}}_i + \left(\prod_{l \in P_i^u} \mathbf{W}_l - \overline{\mathbf{W}}_i \right) + \text{LR}_{RM \wedge \tilde{R}_i}(\overline{\mathbf{W}}_i - \prod_{l \in P_i^u} \mathbf{W}_l) \right\|_{\text{F}} \quad (\text{Rearranging terms}) \\
&\leq \|\overline{\mathbf{W}}_i\|_{\text{F}} + \left\| \left(\prod_{l \in P_i^u} \mathbf{W}_l - \overline{\mathbf{W}}_i \right) + \text{LR}_{RM \wedge \tilde{R}_i}(\overline{\mathbf{W}}_i - \prod_{l \in P_i^u} \mathbf{W}_l) \right\|_{\text{F}} \\
&\quad (\text{Applying triangle inequality}) \\
&= \|\overline{\mathbf{W}}_i\|_{\text{F}} + \sqrt{\sum_{j=RM \wedge \tilde{R}_i+1}^D \sigma_j^2(\overline{\mathbf{W}}_i - \prod_{l \in P_i^u} \mathbf{W}_l)} \quad (19) \\
&\quad (\text{By the definition of } \overline{\mathbf{W}}_i \text{ and } \text{LR}_{RM \wedge \tilde{R}_i+1}(\cdot)) \\
&\leq \max_{k \in [\tilde{L}]} (\|\overline{\mathbf{W}}_k\|_{\text{F}} + E_i) := \alpha.
\end{aligned}$$

By combining (18) and (19), we get

$$A_i \leq \max_{k \in [\tilde{L}]} (\|\overline{\mathbf{W}}_k\|_{\text{F}} + E_i) \mathbb{E} \|\hat{\mathbf{z}}_{i-1} - \bar{\mathbf{z}}_{i-1}\|_2 \leq \alpha \mathbb{E} \|\hat{\mathbf{z}}_{i-1} - \bar{\mathbf{z}}_{i-1}\|_2. \quad (20)$$

$$E_i = \sigma_{RM+1}(\overline{\mathbf{W}}_i - \prod_{l \in P_i^u} \mathbf{W}_l)$$

Computing B_i . We proceed to compute the error associated with the i -th submodel, which we denote as B_i . It can be evaluated as follows:

$$\begin{aligned}
B_i &= \mathbb{E} \left\| f_i(\bar{\mathbf{z}}_{i-1}) - \bar{f}_i(\bar{\mathbf{z}}_{i-1}) \right\|_2 \\
&\stackrel{(15)}{=} \mathbb{E} \left\| \text{ReLU} \left(\left(\text{LR}_{RM \wedge \tilde{R}_i}(\bar{\mathbf{W}}_i - \prod_{l \in P_i^u} \mathbf{w}_l) - (\bar{\mathbf{W}}_i - \prod_{l \in P_i^u} \mathbf{w}_l) \right) \bar{\mathbf{z}}_{i-1} \right) \right\|_2 \\
&\leq \mathbb{E} \left\| \left(\text{LR}_{RM \wedge \tilde{R}_i}(\bar{\mathbf{W}}_i - \prod_{l \in P_i^u} \mathbf{w}_l) - (\bar{\mathbf{W}}_i - \prod_{l \in P_i^u} \mathbf{w}_l) \right) \bar{\mathbf{z}}_{i-1} \right\|_2 \quad (\text{ReLU is 1-Lipschitz}) \\
&\stackrel{(1)}{\leq} \left\| \text{LR}_{RM \wedge \tilde{R}_i}(\bar{\mathbf{W}}_i - \prod_{l \in P_i^u} \mathbf{w}_l) - (\bar{\mathbf{W}}_i - \prod_{l \in P_i^u} \mathbf{w}_l) \right\|_2 \mathbb{E} \|\bar{\mathbf{z}}_{i-1}\|_2 \\
&= \sigma_{RM \wedge \tilde{R}_i+1}(\bar{\mathbf{W}}_i - \prod_{l \in P_i^u} \mathbf{w}_l) \mathbb{E} \|\bar{\mathbf{z}}_{i-1}\|_2.
\end{aligned}$$

We can further simplify $\mathbb{E} \|\bar{\mathbf{z}}_{i-1}\|_2$ as :

$$\begin{aligned}
& \mathbb{E} \|\bar{\mathbf{z}}_{i-1}\|_2 \\
&= \mathbb{E} \|\text{ReLU}(\bar{\mathbf{W}}_{i-1} \bar{\mathbf{z}}_{i-2} + \bar{\mathbf{b}}_{i-1})\|_2 \\
&= \mathbb{E} \|\bar{\mathbf{W}}_{i-1} \bar{\mathbf{z}}_{i-2} + \bar{\mathbf{b}}_{i-1}\|_2 && \text{(ReLU is 1-Lipschitz)} \\
\\
&\leq \|\bar{\mathbf{W}}_{i-1}\|_F \mathbb{E} \|\bar{\mathbf{z}}_{i-2}\|_2 + \|\bar{\mathbf{b}}_{i-1}\|_2 && \text{(Applying triangle inequality and (1))} \\
&\leq \|\bar{\mathbf{W}}_{i-1}\|_F (\|\bar{\mathbf{W}}_{i-2}\|_F \mathbb{E} \|\bar{\mathbf{z}}_{i-3}\|_2 + \|\bar{\mathbf{b}}_{i-2}\|_2) + \|\bar{\mathbf{b}}_{i-1}\|_2 && \text{(Following the same steps)} \\
&\leq \prod_{j=1}^{i-1} \|\bar{\mathbf{W}}_j\|_F \mathbb{E} \|\mathbf{x}\|_2 + \sum_{j=1}^{i-1} \prod_{k=j+1}^{i-1} \|\bar{\mathbf{W}}_k\|_F \|\bar{\mathbf{b}}_j\|_2 && \text{(Repeating the same steps)} \\
&= \sqrt{\|\Sigma\|_F} \prod_{j=1}^{i-1} \|\bar{\mathbf{W}}_j\|_F + \sum_{j=1}^{i-1} \prod_{k=j+1}^{i-1} \|\bar{\mathbf{W}}_k\|_F \|\bar{\mathbf{b}}_j\|_2 \leq \beta.
\end{aligned}$$

Therefore, we obtain

$$B_i \leq \beta \sigma_{RM \wedge \tilde{R}_i+1} (\bar{\mathbf{W}}_i - \prod_{l \in P_i^u} \mathbf{W}_l).$$

Error Composition. Having established upper bounds for A_i and B_i , we next evaluate the expected error for the composition of the first i adapted submodels.

$$\begin{aligned}
\mathbb{E} \|\hat{\mathbf{z}}_i - \bar{\mathbf{z}}_i\|_2 &\stackrel{(17)}{\leq} A_i + B_i \stackrel{(20)}{\leq} \alpha \mathbb{E} \|\hat{\mathbf{z}}_{i-1} - \bar{\mathbf{z}}_{i-1}\|_2 + B_i \leq \alpha(\alpha \mathbb{E} \|\hat{\mathbf{z}}_{i-2} - \bar{\mathbf{z}}_{i-2}\|_2 + B_{i-1}) + B_i \\
&= \alpha^2 \mathbb{E} \|\hat{\mathbf{z}}_{i-2} - \bar{\mathbf{z}}_{i-2}\|_2 + \alpha B_{i-1} + B_i \leq \dots \leq \alpha^{i-1} \mathbb{E} \|\hat{\mathbf{z}}_1 - \bar{\mathbf{z}}_1\|_2 + \sum_{k=2}^i \alpha^{i-k} B_k. \tag{21}
\end{aligned}$$

To compute the overall approximation error of f , which is the composite of all submodels, we have

$$\begin{aligned}
\mathbb{E} \|f(\mathbf{x}) - \bar{f}(\mathbf{x})\|_2 &= \mathbb{E} \|f(\mathbf{x}) - \bar{f}(\mathbf{x})\|_2 = \mathbb{E} \|\hat{\mathbf{z}}_{\bar{L}} - \bar{\mathbf{z}}_{\bar{L}}\|_2 \\
&\stackrel{(21)}{\leq} \alpha^{\bar{L}-1} \mathbb{E} \|\hat{\mathbf{z}}_1 - \bar{\mathbf{z}}_1\|_2 + \sum_{i=2}^{\bar{L}} \alpha^{\bar{L}-i} B_i \\
&\stackrel{(16)}{\leq} \alpha^{\bar{L}-1} \beta \sigma_{RM \wedge \tilde{R}_{i+1}}(\bar{\mathbf{w}}_i - \prod_{l \in P_i^u} \mathbf{w}_l) + \beta \sum_{i=2}^{\bar{L}} \alpha^{\bar{L}-i} \sigma_{RM \wedge \tilde{R}_{i+1}}(\bar{\mathbf{w}}_i - \prod_{l \in P_i^u} \mathbf{w}_l) \\
&= \beta \sum_{i=1}^{\bar{L}} \alpha^{\bar{L}-i} \sigma_{RM \wedge \tilde{R}_{i+1}}(\bar{\mathbf{w}}_i - \prod_{l \in P_i^u} \mathbf{w}_l) \\
&= \beta \sum_{i=1}^{\bar{L}} \alpha^{\bar{L}-i} \sigma_{RM+1}(\bar{\mathbf{w}}_i - \prod_{l \in P_i^u} \mathbf{w}_l).
\end{aligned}$$

Substituting α with $\max_{k \in [\bar{L}]} (\|\bar{\mathbf{w}}_k\|_F + E_i)$ concludes the proof. \square

Transfer Learning:

Given:

- $\{X_i^{(0)}, Y_i^{(0)}\}_N \sim P^{(0)}$
- $\{X_i^{(1)}, Y_i^{(1)}\}_n \sim P^{(1)}, n \ll N$

Suppose:

- $y_i^{(0)} = M_L \sigma \left(M_{L-1} \sigma \left(M_{L-2} \sigma \left(\dots M_2 \sigma \left(M_1 x_i^{(0)} \right) \right) \right) \right) + \epsilon_i$
- $y_i^{(1)} = M_L \sigma \left(M_{L-1} \sigma \left(M_{L-2} \sigma \left(\dots M_2 \sigma \left((M_1 + Z) x_i^{(1)} \right) \right) \right) \right) + \epsilon_i$

Where $rank(Z) \leq k$

- Suppose have exactly $\{M_1\}_{1:L}$, without estimation error ($n \ll N$).
- $W_1 = M_1 + Z$, we can estimate without gradient method to obtain RW_1 for some orthogonal R , and hence the V in $M_1 + Z = UV^T, V^TV = I$

Under the case with Gaussian input $\mathbf{x} \in \mathcal{N}(0, \Sigma)$, the first and second order score reduce to

$$\begin{aligned} S(\mathbf{x}) &= \Sigma^{-1}\mathbf{x}, \\ T(\mathbf{x}) &= \Sigma^{-1}\mathbf{x}\mathbf{x}^T\Sigma^{-1} - \Sigma^{-1}. \end{aligned}$$

If we further let $\mathbf{z} = \mathbf{W}_1\mathbf{x}$ and assume that both $\mathbb{E}[yT(\mathbf{x})]$ and $\mathbb{E}[\nabla_{\mathbf{z}}^2 f(\mathbf{W}_1\mathbf{x})]$ are well-defined, then a second-order Stein's formula suggests that

$$\mathbb{E}[yT(\mathbf{x})] = \mathbf{W}_1^T \mathbb{E}[\nabla_{\mathbf{z}}^2 f(\mathbf{W}_1\mathbf{x})] \cdot \mathbf{W}_1. \quad (10)$$

The equation (10) serves as the basis for estimating \mathcal{S}_0 . Let $\mathbf{A} = \mathbb{E}[yT(\mathbf{x})]$ and denote its eigenvalue decomposition as $\mathbf{A} = \mathbf{W}^T \mathbf{D} \mathbf{W}$. The second-order Stein's formula suggests that there exists an orthogonal matrix $\mathbf{R} \in \mathbb{R}^{k_1 \times k_1}$ such that $\mathbf{W}_1 = \mathbf{R} \mathbf{W}$. In other words, \mathbf{W}_1 has the same row space as \mathbf{W} , regardless of the specific form of link function f . Consequently, \mathcal{S}_0 can be obtained from the eigenvalue decomposition of \mathbf{A} .

$$\begin{aligned}
M_1 + Z &= UV^T \\
(M_1 + Z)VV^T &= UV^T VV^T = UV^T = M_1 + Z \\
M_1(VV^T - I) + Z(VV^T - I) &= 0
\end{aligned}$$

To estimate Z , consider the constrained optimization:

$$\begin{aligned}
&\min_{Z: \text{rank}(Z) \leq k} \|M_1(VV^T - I) + Z(VV^T - I)\|_F^2 \\
&= \min_{Z: \text{rank}(Z) \leq k} \|Y - ZX\|_F^2
\end{aligned}$$

Where $Y = M_1(VV^T - I)$, $X = (I - VV^T)$. There is closed-form solution to this.

- (1) Estimate Z without gradient method
- (2) Parameter-efficient, $\text{rank}(Z) \leq k$, good when $n \ll N$.

- Some potential Scenarios?

$\{X_i^{(0)}, Y_i^{(0)}\}_N \sim P^{(0)}$: Population

$\{X_i^{(1)}, Y_i^{(1)}\}_n \sim P^{(1)}$: Subgroup, $n \ll N$

→ low-rank adjustment is enough

To formalize the idea:

(1): Error and noise: $\hat{M}_1 \neq M_1$, etc

(2): Deal with non-linearity to extend to $M_2 + Z_2, M_3 + Z_3, \dots$

$$\begin{aligned}
y_i^{(1)} &= (M_L + Z_L) \sigma \left((M_{L-1} + Z_{L-1}) \sigma \left((M_{L-2} + Z_{L-2}) \sigma \left(\dots (M_2 + Z_2) \sigma \left((M_1 + Z_1) x_i^{(1)} \right) \right) \right) \right) + \epsilon_i \\
&= (M_L + Z_L) \sigma \left((M_{L-1} + Z_{L-1}) \sigma \left((M_{L-2} + Z_{L-2}) \sigma \left(\dots (M_2 + Z_2) \mathbf{g}_1 \right) \right) \right) + \epsilon_i
\end{aligned}$$

To continue, we would need to deal with truncated normal:

$$g_1 = \sigma \left((M_1 + Z_1) x_i^{(1)} \right)$$

We need the second-order score:

Under the case with Gaussian input $\mathbf{x} \in \mathcal{N}(0, \Sigma)$, the first and second order score reduce to

$$\begin{aligned}
S(\mathbf{x}) &= \Sigma^{-1} \mathbf{x}, \\
T(\mathbf{x}) &= \Sigma^{-1} \mathbf{x} \mathbf{x}^\top \Sigma^{-1} - \Sigma^{-1}.
\end{aligned}$$

If we further let $\mathbf{z} = \mathbf{W}_1 \mathbf{x}$ and assume that both $\mathbb{E}[yT(\mathbf{x})]$ and $\mathbb{E}[\nabla_{\mathbf{z}}^2 f(\mathbf{W}_1 \mathbf{x})]$ are well-defined, then a second-order Stein's formula suggests that

$$\mathbb{E}[yT(\mathbf{x})] = \mathbf{W}_1^\top \mathbb{E}[\nabla_{\mathbf{z}}^2 f(\mathbf{W}_1 \mathbf{x})] \cdot \mathbf{W}_1. \quad (10)$$

- Randomly initializes $\{M_l\}_{1:L}$ (they are known), and obtain $\{Z_l\}_{1:L}$, if $\text{rank}(Z_1)$ is large enough, this can approximate any L-layer FNN

$$y_i^{(1)} = (M_L + Z_L) \sigma \left((M_{L-1} + Z_{L-1}) \sigma \left((M_{L-2} + Z_{L-2}) \sigma \left(\dots (M_2 + Z_2) \sigma \left((M_1 + Z_1) x_i^{(1)} \right) \right) \right) \right) + \epsilon_i$$