

Linear Convergence Analysis of Neural Collapse under the MSE Loss

July 5, 2024

1 Preliminaries and Main Results

Without loss of generality, let $\mathbf{H}_i := [\mathbf{h}_{1,i}, \dots, \mathbf{h}_{K,i}] \in \mathbb{R}^{d \times K}$ for $i \in [n]$ and $\mathbf{H} := [\mathbf{H}_1, \dots, \mathbf{H}_n] \in \mathbb{R}^{d \times N}$ be the matrix with the columns being organized unconstrained features, associated with the label matrix $\mathbf{Y} = \mathbf{1}_n^T \otimes \mathbf{I}_K \in \mathbb{R}^{K \times N}$. Let

$$F(\mathbf{W}, \mathbf{H}) := \frac{1}{nK} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}(\mathbf{W}^T \mathbf{h}_{k,i} + \mathbf{b}, \mathbf{y}_k) + \frac{\lambda_W}{2} \|\mathbf{W}\|_F^2 + \frac{\lambda_H}{2} \|\mathbf{H}\|_F^2 + \frac{\lambda_b}{2} \|\mathbf{b}\|^2. \quad (1)$$

where $\mathbf{y}_k \in \mathbb{R}^K$ is a membership matrix with all the entries being 0 but the k -th one being 1, $\lambda > 0$ is the regularized parameter, and $\mathcal{L} : \mathbb{R}^d \times \mathbb{R}^K \rightarrow \mathbb{R}_+$ is a loss function. Let

$$\min_{\mathbf{W}, \mathbf{H} \in \mathbb{R}^{d \times K}} f(\mathbf{W}, \mathbf{H}) := \frac{1}{K} \sum_{k=1}^K \mathcal{L}(\mathbf{W}^T \mathbf{h}_k + \mathbf{b}, \mathbf{y}_k) + \frac{\lambda_W}{2} \|\mathbf{W}\|_F^2 + \frac{\lambda_H}{2} \|\mathbf{H}\|_F^2 + \frac{\lambda_b}{2} \|\mathbf{b}\|^2. \quad (2)$$

We define $\mathbf{Z} := \mathbf{W}^T \mathbf{H} + \mathbf{b} \mathbf{1}_K^T$. This implies $\mathbf{z}_k = \mathbf{W}^T \mathbf{h}_k$ for all $k \in [K]$. To simplify our development, let

$$g_k(\mathbf{z}) = \frac{1}{K} \mathcal{L}(\mathbf{z}, \mathbf{y}_k), \quad g(\mathbf{Z}) = \sum_{k=1}^K g_k(\mathbf{z}_k). \quad (3)$$

Consider a optimization problem

$$v^* = \min_{\mathbf{x} \in \mathcal{E}} F(\mathbf{x}), \quad (4)$$

where \mathcal{E} is a finite-dimensional Euclidean space and $f : \mathcal{E} \rightarrow (\infty, \infty)$ is a continuously differentiable function. Let $\mathcal{X} \subseteq \mathcal{E}$ denote the set of optimal solutions of Problem (4).

Definition 1 (Error bound condition). *We say that an error bound condition holds for Problem (4) if there exists a constant $\kappa > 0$ such that for all $\mathbf{x} \in \mathbb{R}^n$ with $\text{dist}(\mathbf{x}, \mathcal{X}) \leq \delta$,*

$$\text{dist}(\mathbf{x}, \mathcal{X}) \leq \kappa \|\nabla F(\mathbf{x})\|. \quad (5)$$

It follows from [1, 2, 3] that the error bound condition can be used to analyze the convergence rate of first-order methods.

Fact 1 (cf. [1, 2, 3]). Suppose that the optimal solution of Problem (4) is non-empty, i.e., $\mathcal{X} \neq \emptyset$, and the error bound holds for Problem (4). Suppose in addition that the sequence $\{\mathbf{x}^k\}_{k \geq k_1}$ for an index $k_1 \geq 0$ satisfies the following properties:

(A1). (Sufficient Decrease) There exists a constant $\kappa_1 > 0$ such that

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

(A2). (Cost-to-Go Estimate) There exists a constant $\kappa_2 > 0$ such that

$$F(\mathbf{x}^{k+1}) - v^* \leq \kappa_2 \left(d^2(\mathbf{x}^k, \mathcal{X}) + \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2 \right).$$

(A3). (Safeguard) There exists a constant $\kappa_3 > 0$ such that

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

Then, the sequence $\{F(\mathbf{x}^k)\}_{k \geq 0}$ converges Q -linearly to v^* and $\{\mathbf{x}^k\}_{k \geq 0}$ converges R -linearly to some $\mathbf{x}^* \in \mathcal{X}$.

Despite the fact that F in Problem (4) is non-convex, we can verify that the sequence generated by the gradient descent method for solving Problem (4) satisfies (A1)-(A3) in Fact (1).

Notation. For a matrix $\mathbf{A} \in \mathbb{R}^{d \times K}$, we denote by \mathbf{a}_k by its k -th column and by \mathbf{a}^i by its i -th row. Let $\mathbf{P} = \mathbf{I}_K - \frac{1}{K} \mathbf{1}_K \mathbf{1}_K^T$ be a projection matrix and $\mathbf{P}^\perp = \frac{1}{K} \mathbf{1}_K \mathbf{1}_K^T$ be its complement. Let $\lambda_{\min} = \min\{\lambda_W, \lambda_H\}$ and $\lambda_{\max} = \max\{\lambda_W, \lambda_H\}$.

2 Neural Collapse with Regularized MSE Loss

Suppose that $\mathcal{L}(\cdot, \cdot)$ is the mean squared error (MSE) loss:

$$\mathcal{L}(\mathbf{z}, \mathbf{y}_k) = \frac{1}{2} \|\mathbf{z} - \mathbf{y}_k\|^2. \quad (6)$$

Substituting this into Problem (1) and Problem (2) respectively yields

$$\min_{\mathbf{W} \in \mathbb{R}^{d \times K}, \mathbf{H} \in \mathbb{R}^{d \times N}} F(\mathbf{W}, \mathbf{H}) = \frac{1}{2N} \|\mathbf{W}^T \mathbf{H} - \mathbf{Y}\|_F^2 + \frac{\lambda_W}{2} \|\mathbf{W}\|_F^2 + \frac{\lambda_H}{2} \|\mathbf{H}\|_F^2, \quad (7)$$

$$\min_{\mathbf{W}, \boldsymbol{\Theta} \in \mathbb{R}^{d \times K}} f(\mathbf{W}, \boldsymbol{\Theta}) = \frac{1}{2K} \|\mathbf{W}^T \boldsymbol{\Theta} - \mathbf{I}_K\|_F^2 + \frac{\lambda_W}{2} \|\mathbf{W}\|_F^2 + \frac{n\lambda_H}{2} \|\boldsymbol{\Theta}\|_F^2. \quad (8)$$

Note that

$$F(\mathbf{W}, \mathbf{H}) = \frac{1}{n} \sum_{i=1}^n f(\mathbf{W}, \mathbf{H}_i). \quad (9)$$

We first study the following optimization problems:

$$\min_{\mathbf{W}, \boldsymbol{\Theta} \in \mathbb{R}^{d \times K}} f(\mathbf{W}, \boldsymbol{\Theta}) = \frac{1}{2K} \|\mathbf{W}^T \boldsymbol{\Theta} - \boldsymbol{\Sigma}\|_F^2 + \frac{\lambda_W}{2} \|\mathbf{W}\|_F^2 + \frac{n\lambda_H}{2} \|\boldsymbol{\Theta}\|_F^2, \quad (10)$$

where $\boldsymbol{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_K)$ is a diagonal matrix with $\sigma_1 \geq \dots \geq \sigma_K \geq 0$. It is worth noting that Problem (8) is a special case of this problem by taking $\sigma_k = 1$ for all $k \in [K]$.

Lemma 1. *The optimal solution set of Problem (10) takes the form of*

$$\left\{ (\mathbf{W}, \mathbf{\Theta}) : \mathbf{W} = \frac{\sqrt[4]{n\lambda_H}}{\sqrt[4]{\lambda_W}} \mathbf{U} \left(\max \left\{ \mathbf{\Sigma} - K\sqrt{\lambda_W\lambda_H} \mathbf{I}_K, \mathbf{0} \right\} \right)^{\frac{1}{2}}, \mathbf{\Theta} = \frac{\sqrt{\lambda_W}}{\sqrt{n\lambda_H}} \mathbf{W}, \mathbf{U} \in \mathcal{O}^{d \times K} \right\}. \quad (11)$$

Proof. According to the first-order optimality condition, we have

$$\begin{cases} \nabla_{\mathbf{W}} f(\mathbf{W}, \mathbf{\Theta}) = \mathbf{\Theta} (\mathbf{W}^T \mathbf{\Theta} - \mathbf{\Sigma})^T + \lambda_W K \mathbf{W} = \mathbf{0}, \\ \nabla_{\mathbf{\Theta}} f(\mathbf{W}, \mathbf{\Theta}) = \mathbf{W} (\mathbf{W}^T \mathbf{\Theta} - \mathbf{\Sigma}) + n\lambda_H K \mathbf{\Theta} = \mathbf{0}. \end{cases} \quad (12)$$

Using $\mathbf{W} \nabla_{\mathbf{W}}^T f(\mathbf{W}, \mathbf{\Theta}) - \nabla_{\mathbf{\Theta}} f(\mathbf{W}, \mathbf{\Theta}) \mathbf{\Theta}^T = \mathbf{0}$ yields

$$\mathbf{\Theta} \mathbf{\Theta}^T = \frac{\lambda_W}{n\lambda_H} \mathbf{W} \mathbf{W}^T. \quad (13)$$

This, together with (12), implies

$$\begin{cases} \frac{\lambda_W}{n\lambda_H} \mathbf{W} \mathbf{W}^T \mathbf{W} - \mathbf{\Theta} \mathbf{\Sigma} + \lambda_W K \mathbf{W} = \mathbf{0}, \\ \mathbf{W} \mathbf{W}^T \mathbf{\Theta} - \mathbf{W} \mathbf{\Sigma} + n\lambda_H K \mathbf{\Theta} = \mathbf{0}. \end{cases} \quad (14a)$$

$$(14b)$$

According to $\sqrt{n\lambda_H} \times (14a) - \sqrt{\lambda_W} \times (14b) = \mathbf{0}$, we have

$$\left(\frac{\sqrt{\lambda_W}}{\sqrt{n\lambda_H}} \mathbf{W} \mathbf{W}^T + K\sqrt{n\lambda_W\lambda_H} \mathbf{I}_d \right) \left(\sqrt{\lambda_W} \mathbf{W} - \sqrt{n\lambda_H} \mathbf{\Theta} \right) + \left(\sqrt{\lambda_W} \mathbf{W} - \sqrt{n\lambda_H} \mathbf{\Theta} \right) \mathbf{\Sigma} = \mathbf{0}.$$

Since $\sqrt{\lambda_W}/\sqrt{n\lambda_H} \mathbf{W} \mathbf{W}^T + K\sqrt{n\lambda_W\lambda_H} \mathbf{I}_d$ is positive definite and $\mathbf{\Theta}$ is diagonal with non-negative diagonal entries, we have $\sqrt{\lambda_W} \mathbf{W} = \sqrt{n\lambda_H} \mathbf{\Theta}$. For all $(\mathbf{W}, \mathbf{\Theta})$ satisfying the first-order optimality condition, we compute

$$\begin{aligned} f(\mathbf{W}, \mathbf{\Theta}) &= \frac{1}{2K} \|\mathbf{W}^T \mathbf{\Theta} - \mathbf{\Theta}\|_F^2 + \frac{\lambda_W}{2} \|\mathbf{W}\|_F^2 + \frac{n\lambda_H}{2} \|\mathbf{\Theta}\|_F^2 \\ &= \frac{1}{2K} \left\| \frac{\sqrt{\lambda_W}}{\sqrt{n\lambda_H}} \mathbf{W}^T \mathbf{W} - \mathbf{\Theta} \right\|_F^2 + \lambda_W \|\mathbf{W}\|_F^2 \\ &= \frac{1}{2K} \left\| \frac{\sqrt{\lambda_W}}{\sqrt{n\lambda_H}} \mathbf{W}^T \mathbf{W} - \left(\mathbf{\Theta} - K\sqrt{n\lambda_W\lambda_H} \mathbf{I}_K \right) \right\|_F^2 + \sqrt{n\lambda_W\lambda_H} \text{tr}(\mathbf{\Sigma}) - \frac{nK^2}{2} \lambda_W \lambda_H. \end{aligned}$$

To obtain the global minimum of $f(\mathbf{W}, \mathbf{\Theta})$, we have

$$\mathbf{W}^T \mathbf{W} = \frac{\sqrt{n\lambda_H}}{\sqrt{\lambda_W}} \max \left\{ \mathbf{\Sigma} - K\sqrt{n\lambda_W\lambda_H} \mathbf{I}_K, \mathbf{0} \right\},$$

which implies

$$\mathbf{W} = \frac{\sqrt[4]{n\lambda_H}}{\sqrt[4]{\lambda_W}} \mathbf{U} \left(\max \left\{ \mathbf{\Sigma} - K\sqrt{n\lambda_W\lambda_H} \mathbf{I}_K, \mathbf{0} \right\} \right)^{1/2}, \text{ for all } \mathbf{U} \in \mathcal{O}^{d \times K}.$$

This, together with $\sqrt{\lambda_W} \mathbf{W} = \sqrt{n\lambda_H} \mathbf{\Theta}$, gives (11). \square

Proposition 1. *The optimal solution set of Problem (8) takes the form of*

$$\mathcal{X}_f = \left\{ (\mathbf{W}, \mathbf{\Theta}) : \mathbf{W} = \frac{\sqrt[4]{n\lambda_H}}{\sqrt[4]{\lambda_W}} (\max \{c, 0\})^{\frac{1}{2}} \mathbf{U}, \mathbf{\Theta} = \frac{\sqrt{\lambda_W}}{\sqrt{n\lambda_H}} \mathbf{W}, \mathbf{U} \in \mathcal{O}^{d \times K} \right\}, \quad (15)$$

where $c := 1 - K\sqrt{n\lambda_W\lambda_H}$.

Corollary 1. *The optimal solution set of Problem (7) takes the form of*

$$\mathcal{X}_F = \left\{ (\mathbf{W}, \mathbf{H}) : \mathbf{W} = \frac{\sqrt[4]{n\lambda_H}}{\sqrt[4]{\lambda_W}} (\max\{c, 0\})^{\frac{1}{2}} \mathbf{U}, \mathbf{H}_i = \frac{\sqrt{\lambda_W}}{\sqrt{n\lambda_H}} \mathbf{W}, \forall i \in [n], \mathbf{U} \in \mathcal{O}^{d \times K} \right\}, \quad (16)$$

where $c := 1 - K\sqrt{n\lambda_W\lambda_H}$.

Proof. Suppose that $(\mathbf{W}^*, \boldsymbol{\Theta}^*) \in \mathcal{X}_f$ is an optimal solution of Problem (8). According to (9), we have

$$\min F(\mathbf{W}, \mathbf{H}) \geq \frac{1}{n} \sum_{i=1}^n \min f(\mathbf{W}, \boldsymbol{\Theta}) = f(\mathbf{W}^*, \boldsymbol{\Theta}^*),$$

where the equality holds if $\mathbf{W} = \mathbf{W}^*$ and $\mathbf{H}_i = \boldsymbol{\Theta}^*$ for all $i \in [n]$. This, together with (15), implies (16). \square

Note that when $\lambda_W\lambda_H \geq 1/(nK^2)$, we have $c \leq 0$. Then the optimal solution set of Problem (7) is $\{(\mathbf{0}, \mathbf{0})\}$. To avoid this trivial case, it suffices to consider $\lambda_W\lambda_H < 1/(nK^2)$.

Lemma 2. *Suppose that $\lambda_W\lambda_H < 1/(nK^2)$ and the error bound holds for Problem (2), i.e., there exist constants $\delta, \kappa > 0$ such that*

$$\text{dist}((\mathbf{W}, \boldsymbol{\Theta}), \mathcal{X}_f) \leq \kappa \|\nabla f(\mathbf{W}, \boldsymbol{\Theta})\|_F \quad (17)$$

for all $(\mathbf{W}, \boldsymbol{\Theta}) \in \mathbb{R}^{d \times K}$ satisfying $\text{dist}((\mathbf{W}, \boldsymbol{\Theta}), \mathcal{X}_f) \leq \delta$. Then for all (\mathbf{W}, \mathbf{H}) satisfying

$$\text{dist}((\mathbf{W}, \mathbf{H}), \mathcal{X}_F) \leq \delta_F := \min \left\{ \frac{\sqrt[4]{n\lambda_H}}{\sqrt[4]{\lambda_W}}, \frac{\sqrt[4]{\lambda_W}}{\sqrt[4]{n\lambda_H}}, \frac{\sqrt{2}}{8} nK^2 \lambda_H \sqrt{n\lambda_W\lambda_H} \delta, \frac{\sqrt{2}}{2} \delta \right\}, \quad (18)$$

it holds that

$$\text{dist}((\mathbf{W}, \mathbf{H}), \mathcal{X}_F) \leq \kappa_F \|\nabla F(\mathbf{W}, \mathbf{H})\|_F, \quad (19)$$

where

$$\kappa_F := \frac{\sqrt{2n}\kappa \max \left\{ 1, \frac{n\lambda_H}{\kappa} + \frac{9}{nK\lambda_H\sqrt{\lambda_W\lambda_H}} \right\}}{\min \left\{ \sqrt{\frac{n}{2}} \left(1 + \frac{2\lambda_W}{n\lambda_H} \right)^{-1/2}, \frac{\sqrt{2}}{4} \min \left\{ 1, \frac{\sqrt{n\lambda_H}}{\sqrt{\lambda_W}} \right\} \right\}}. \quad (20)$$

Proof. For ease of exposition, let $c := 1 - K\sqrt{n\lambda_W\lambda_H}$ and $\rho := \sqrt[4]{n\lambda_H}/\sqrt[4]{\lambda_W}$. According to (16) in Corollary 1 and $\lambda_W\lambda_H < 1/(nK^2)$, we compute

$$\begin{aligned} \text{dist}^2((\mathbf{W}, \mathbf{H}), \mathcal{X}_F) &= \min_{\mathbf{U} \in \mathcal{O}^{d \times K}} \left\{ \|\mathbf{W} - \rho\sqrt{c}\mathbf{U}\|_F^2 + \sum_{i=1}^n \left\| \mathbf{H}_i - \frac{\sqrt{c}}{\rho} \mathbf{U} \right\|_F^2 \right\} \\ &\leq \left(1 + \frac{2\lambda_W}{n\lambda_H} \right) \min_{\mathbf{U} \in \mathcal{O}^{d \times K}} \|\mathbf{W} - \rho\sqrt{c}\mathbf{U}\|_F^2 + 2 \sum_{i=1}^n \left\| \mathbf{H}_i - \frac{1}{\rho^2} \mathbf{W} \right\|_F^2, \end{aligned} \quad (21)$$

where the inequality follows from $\|\mathbf{A} + \mathbf{B}\|_F^2 \leq 2\|\mathbf{A}\|_F^2 + 2\|\mathbf{B}\|_F^2$. Moreover, given some $(\mathbf{W}, \boldsymbol{\Theta})$, it follows from (15) in Proposition 1 that

$$\begin{aligned} \text{dist}^2((\mathbf{W}, \boldsymbol{\Theta}), \mathcal{X}_f) &= \min_{\mathbf{U} \in \mathcal{O}^{d \times K}} \left\{ \|\mathbf{W} - \rho\sqrt{c}\mathbf{U}\|_F^2 + \left\| \boldsymbol{\Theta} - \frac{\sqrt{c}}{\rho}\mathbf{U} \right\|_F^2 \right\} \\ &\geq \min_{\mathbf{U} \in \mathcal{O}^{d \times K}} \left\{ \|\mathbf{W} - \rho\sqrt{c}\mathbf{U}\|_F^2 + \alpha \left\| \boldsymbol{\Theta} - \frac{\sqrt{c}}{\rho}\mathbf{U} \right\|_F^2 \right\} \\ &\geq \frac{1}{2} \min_{\mathbf{U} \in \mathcal{O}^{d \times K}} \|\mathbf{W} - \rho\sqrt{c}\mathbf{U}\|_F^2 + \frac{\alpha}{2} \left\| \boldsymbol{\Theta} - \frac{1}{\rho^2}\mathbf{W} \right\|_F^2, \end{aligned} \quad (22)$$

where the first inequality is due to $\alpha := \min\{1, n\lambda_H/\lambda_W\}/2 \leq 1$, and the second inequality follows from $\|\mathbf{A} + \mathbf{B}\|_F^2 \geq \|\mathbf{A}\|_F^2/2 - \|\mathbf{B}\|_F^2$ and $\alpha\lambda_W/(n\lambda_H) \leq 1/2$. This, together with (21), yields

$$\begin{aligned} \min \left\{ \frac{n}{2} \left(1 + \frac{2\lambda_W}{n\lambda_H} \right)^{-1}, \frac{\alpha}{4} \right\} \text{dist}^2((\mathbf{W}, \mathbf{H}), \mathcal{X}_F) &\leq \frac{n}{2} \min_{\mathbf{U} \in \mathcal{O}^{d \times K}} \|\mathbf{W} - \rho\sqrt{c}\mathbf{U}\|_F^2 \\ + \frac{\alpha}{2} \sum_{i=1}^n \left\| \mathbf{H}_i - \frac{1}{\rho^2}\mathbf{W} \right\|_F^2 &\leq \sum_{i=1}^n \text{dist}^2((\mathbf{W}, \mathbf{H}_i), \mathcal{X}_f). \end{aligned} \quad (23)$$

Noting that $f(\mathbf{W}, \boldsymbol{\Theta})$ is strongly convex w.r.t. $\boldsymbol{\Theta}$ with constant $n\lambda_H$, let $\theta(\mathbf{W}) := \arg \min_{\boldsymbol{\Theta}} f(\mathbf{W}, \boldsymbol{\Theta})$ denote the unique minimizer. Using the first-order optimality condition, we obtain

$$\theta(\mathbf{W}) = (\mathbf{W}\mathbf{W}^T + nK\lambda_H\mathbf{I})^{-1}\mathbf{W}. \quad (24)$$

Using the strongly convexity again, we have for arbitrary $\boldsymbol{\Theta} \in \mathbb{R}^{d \times K}$,

$$n\lambda_H\|\boldsymbol{\Theta} - \theta(\mathbf{W})\|_F^2 \leq \langle \nabla_{\boldsymbol{\Theta}} f(\mathbf{W}, \boldsymbol{\Theta}) - \nabla_{\boldsymbol{\Theta}} f(\mathbf{W}, \theta(\mathbf{W})), \boldsymbol{\Theta} - \theta(\mathbf{W}) \rangle.$$

This, together with the Cauchy–Schwarz inequality and $\nabla_{\boldsymbol{\Theta}} f(\mathbf{W}, \theta(\mathbf{W})) = \mathbf{0}$, implies $\|\nabla_{\boldsymbol{\Theta}} f(\mathbf{W}, \boldsymbol{\Theta})\|_F \geq n\lambda_H\|\boldsymbol{\Theta} - \theta(\mathbf{W})\|_F$. Then, it holds for arbitrary \mathbf{H} that for all $i \in [n]$,

$$\|\nabla_{\mathbf{H}_i} f(\mathbf{W}, \mathbf{H}_i)\|_F \geq n\lambda_H\|\mathbf{H}_i - \theta(\mathbf{W})\|_F. \quad (25)$$

Let $(\mathbf{W}^*, \mathbf{H}^*) \in \mathcal{X}_F$ be such that $\text{dist}((\mathbf{W}, \mathbf{H}), \mathcal{X}_F) = \|(\mathbf{W}, \mathbf{H}) - (\mathbf{W}^*, \mathbf{H}^*)\|_F \leq \delta_F$. It follows from $\text{dist}((\mathbf{W}, \mathbf{H}), \mathcal{X}_F) \leq \delta_F$ that

$$\|\mathbf{W} - \mathbf{W}^*\|_F^2 + \sum_{i=1}^n \|\mathbf{H}_i - \mathbf{H}_i^*\|_F^2 \leq \delta_F^2. \quad (26)$$

This, together with (16) and (24), gives

$$\|\mathbf{W}\| \leq \|\mathbf{W} - \mathbf{W}^*\|_F + \|\mathbf{W}^*\| \leq \delta_F + \rho\sqrt{c}, \quad \|\mathbf{H}_i\| \leq \delta_F + \frac{\sqrt{c}}{\rho}, \quad \forall i \in [n], \quad (27)$$

$$\|\theta(\mathbf{W})\| \leq \left\| (\mathbf{W}\mathbf{W}^T + nK\lambda_H\mathbf{I})^{-1} \right\| \|\mathbf{W}\| \leq \frac{1}{nK\lambda_H} (\delta_F + \rho\sqrt{c}). \quad (28)$$

Using (12), we compute

$$\begin{aligned} \|\nabla_{\mathbf{W}} f(\mathbf{W}, \mathbf{H}_i) - \nabla_{\mathbf{W}} f(\mathbf{W}, \theta(\mathbf{W}))\|_F &= \|\mathbf{H}_i\mathbf{H}_i^T\mathbf{W} - \theta(\mathbf{W})\theta(\mathbf{W})^T\mathbf{W} + \theta(\mathbf{W}) - \mathbf{H}_i\|_F \\ &\leq ((\|\mathbf{H}_i\| + \|\theta(\mathbf{W})\|) \|\mathbf{W}\| + 1) \|\mathbf{H}_i - \theta(\mathbf{W})\|_F \\ &\leq \frac{9}{K\sqrt{n\lambda_W\lambda_H}} \|\mathbf{H}_i - \theta(\mathbf{W})\|_F, \end{aligned} \quad (29)$$

where the last inequality follows from (18), (27), and (28). Next, we bound

$$\begin{aligned}
\left\| \sum_{i=1}^n \nabla_{\mathbf{W}} f(\mathbf{W}, \mathbf{H}_i) \right\|_F &= \left\| \sum_{i=1}^n \nabla_{\mathbf{W}} f(\mathbf{W}, \theta(\mathbf{W})) + \nabla_{\mathbf{W}} f(\mathbf{W}, \mathbf{H}_i) - \nabla_{\mathbf{W}} f(\mathbf{W}, \theta(\mathbf{W})) \right\|_F \\
&\geq n \|\nabla_{\mathbf{W}} f(\mathbf{W}, \theta(\mathbf{W}))\|_F - \sum_{i=1}^n \|\nabla_{\mathbf{W}} f(\mathbf{W}, \mathbf{H}_i) - \nabla_{\mathbf{W}} f(\mathbf{W}, \theta(\mathbf{W}))\|_F \\
&\geq n \|\nabla_{\mathbf{W}} f(\mathbf{W}, \theta(\mathbf{W}))\|_F - \frac{9}{K\sqrt{n\lambda_W\lambda_H}} \sum_{i=1}^n \|\mathbf{H}_i - \theta(\mathbf{W})\|_F,
\end{aligned}$$

where the last inequality follows from (29). Substituting (25) into this inequality yields

$$\left\| \sum_{i=1}^n \nabla_{\mathbf{W}} f(\mathbf{W}, \mathbf{H}_i) \right\|_F + \frac{9}{nK\lambda_H\sqrt{n\lambda_W\lambda_H}} \sum_{i=1}^n \|\nabla_{\mathbf{H}_i} f(\mathbf{W}, \mathbf{H}_i)\|_F \geq n \|\nabla_{\mathbf{W}} f(\mathbf{W}, \theta(\mathbf{W}))\|_F. \quad (30)$$

According to (24), we compute

$$\begin{aligned}
\|\theta(\mathbf{W}) - \theta(\mathbf{W}^*)\|_F &= \left\| \left((\mathbf{W}\mathbf{W}^T + nK\lambda_H\mathbf{I})^{-1} - (\mathbf{W}^*\mathbf{W}^{*T} + nK\lambda_H\mathbf{I})^{-1} \right) \mathbf{W}^* \right. \\
&\quad \left. + (\mathbf{W}\mathbf{W}^T + nK\lambda_H\mathbf{I})^{-1} (\mathbf{W} - \mathbf{W}^*) \right\|_F \\
&\leq \left(\frac{1}{(nK\lambda_H)^2} (\|\mathbf{W}^*\| + \|\mathbf{W}\|) \|\mathbf{W}^*\| + \frac{1}{nK\lambda_H} \right) \|\mathbf{W}^* - \mathbf{W}\|_F \\
&\leq \left(\frac{1}{(nK\lambda_H)^2} (\delta_F + 2\rho\sqrt{c}) \rho\sqrt{c} + \frac{1}{nK\lambda_H} \right) \delta_F \leq \frac{\sqrt{2}}{2} \delta, \quad (31)
\end{aligned}$$

where the first inequality follows from $(\mathbf{W}\mathbf{W}^T + nK\lambda_H\mathbf{I})^{-1} - (\mathbf{W}^*\mathbf{W}^{*T} + nK\lambda_H\mathbf{I})^{-1} = (\mathbf{W}\mathbf{W}^T + nK\lambda_H\mathbf{I})^{-1}(\mathbf{W}^*\mathbf{W}^{*T} + nK\lambda_H\mathbf{I} - \mathbf{W}\mathbf{W}^T - nK\lambda_H\mathbf{I})(\mathbf{W}^*\mathbf{W}^{*T} + nK\lambda_H\mathbf{I})^{-1} = (\mathbf{W}\mathbf{W}^T + nK\lambda_H\mathbf{I})^{-1}(\mathbf{W}^*(\mathbf{W}^* - \mathbf{W})^T - (\mathbf{W}^* - \mathbf{W})\mathbf{W}^T)(\mathbf{W}^*\mathbf{W}^{*T} + nK\lambda_H\mathbf{I})^{-1}$, $\|(\mathbf{W}\mathbf{W}^T + nK\lambda_H\mathbf{I})^{-1}\| \leq 1/(nK\lambda_H)$, and $\|(\mathbf{W}^*\mathbf{W}^{*T} + nK\lambda_H\mathbf{I})^{-1}\| \leq 1/(nK\lambda_H)$, the second inequality uses (26) and (27), and the last inequality is due to (18) and $c \leq 1$. This, together with $\theta(\mathbf{W}^*) = \mathbf{H}^*$ implies $\|\theta(\mathbf{W}) - \mathbf{H}^*\|_F = \|\theta(\mathbf{W}) - \theta(\mathbf{W}^*)\|_F \leq \sqrt{2}\delta/2$. Using this and (18), we have

$$\text{dist}^2((\mathbf{W}, \theta(\mathbf{W})), \mathcal{X}_f) \leq \|\mathbf{W} - \mathbf{W}^*\|_F^2 + \|\theta(\mathbf{W}) - \mathbf{H}^*\|_F^2 \leq \delta^2. \quad (32)$$

Using $\nabla_{\Theta} f(\mathbf{W}, \theta(\mathbf{W})) = \mathbf{0}$, we have

$$\begin{aligned}
\|\nabla_{\mathbf{W}} f(\mathbf{W}, \theta(\mathbf{W}))\|_F &= \|\nabla f(\mathbf{W}, \theta(\mathbf{W}))\|_F \geq \kappa^{-1} \text{dist}((\mathbf{W}, \theta(\mathbf{W})), \mathcal{X}_f) \\
&\geq \kappa^{-1} \text{dist}((\mathbf{W}, \mathbf{H}_i), \mathcal{X}_f) - \kappa^{-1} \|\mathbf{H}_i - \theta(\mathbf{W})\|_F, \quad \forall i \in [n],
\end{aligned}$$

where the first inequality follows from (17) and (32). Substituting (25) into this inequality yields

$$\|\nabla_{\mathbf{W}} f(\mathbf{W}, \theta(\mathbf{W}))\|_F + \frac{n\lambda_H}{\kappa} \|\nabla_{\mathbf{H}_i} f(\mathbf{W}, \mathbf{H}_i)\|_F \geq \frac{1}{\kappa} \text{dist}((\mathbf{W}, \mathbf{H}_i), \mathcal{X}_f).$$

Summing up this inequality from $i = 1$ to $i = n$ and using (30) gives

$$\left\| \sum_{i=1}^n \nabla_{\mathbf{W}} f(\mathbf{W}, \mathbf{H}_i) \right\|_F + \left(\frac{n\lambda_H}{\kappa} + \frac{9}{nK\lambda_H\sqrt{n\lambda_W\lambda_H}} \right) \sum_{i=1}^n \|\nabla_{\mathbf{H}_i} f(\mathbf{W}, \mathbf{H}_i)\|_F \geq \frac{1}{\kappa} \sum_{i=1}^n \text{dist}((\mathbf{W}, \mathbf{H}_i), \mathcal{X}_f).$$

This, together with (9), implies

$$\left(\sum_{i=1}^n \text{dist}((\mathbf{W}, \mathbf{H}_i), \mathcal{X}_f) \right)^2 \leq 2n\kappa^2 \max \left\{ 1, \left(\frac{n\lambda_H}{\kappa} + \frac{9}{nK\lambda_H\sqrt{\lambda_W\lambda_H}} \right)^2 \right\} \|\nabla F(\mathbf{W}, \mathbf{H})\|_F^2$$

Combining this with (23) yields

$$\text{dist}((\mathbf{W}, \mathbf{H}), \mathcal{X}_F) \leq \frac{\sqrt{2n\kappa} \max \left\{ 1, \frac{n\lambda_H}{\kappa} + \frac{9}{nK\lambda_H\sqrt{\lambda_W\lambda_H}} \right\}}{\min \left\{ \sqrt{\frac{n}{2}} \left(1 + \frac{2\lambda_W}{n\lambda_H} \right)^{-1/2}, \frac{\sqrt{2}}{4} \min \left\{ 1, \frac{\sqrt{n\lambda_H}}{\sqrt{\lambda_W}} \right\} \right\}} \|\nabla F(\mathbf{W}, \mathbf{H})\|_F.$$

□

According to this lemma, it suffices to consider the error bound of Problem (8).

Theorem 1. Suppose that $\lambda_W\lambda_H < 1/(nK^2)$. For all $(\mathbf{W}, \mathbf{\Theta}) \in \mathbb{R}^{d \times K} \times \mathbb{R}^{d \times K}$ satisfying

$$\text{dist}((\mathbf{W}, \mathbf{\Theta}), \mathcal{X}_f) \leq \delta := \frac{1}{2} \min \left\{ \frac{\sqrt[4]{\lambda_W}}{\sqrt[4]{n\lambda_H}}, \frac{\sqrt[4]{n\lambda_H}}{\sqrt[4]{\lambda_W}} \right\} \left(1 - K\sqrt{n\lambda_W\lambda_H} \right)^{1/2}, \quad (33)$$

it holds that

$$\text{dist}((\mathbf{W}, \mathbf{\Theta}), \mathcal{X}_f) \leq \kappa \|\nabla f(\mathbf{W}, \mathbf{H})\|_F, \quad (34)$$

where κ is a constant that depends on λ_W and λ_H .

Proof. For ease of exposition, let

$$c := 1 - K\sqrt{n\lambda_W\lambda_H}, \quad \rho := \sqrt[4]{n\lambda_H}/\sqrt[4]{\lambda_W}, \quad \lambda_{\min} := \min\{n\lambda_H, \lambda_W\}, \quad \lambda_{\max} := \max\{n\lambda_H, \lambda_W\}.$$

By Proposition 1 and the condition $\lambda_W\lambda_H < 1/(nK^2)$, we obtain

$$\mathcal{X}_f = \left\{ (\mathbf{W}, \mathbf{\Theta}) : \mathbf{W} = \rho\sqrt{c}\mathbf{U}, \mathbf{\Theta} = \frac{\sqrt{c}}{\rho}\mathbf{U}, \mathbf{U} \in \mathcal{O}^{d \times K} \right\}, \quad (35)$$

Next, we calculate

$$\begin{aligned} \text{dist}^2((\mathbf{W}, \mathbf{\Theta}), \mathcal{X}_f) &= \min_{\mathbf{U} \in \mathcal{O}^{d \times K}} \left\{ \|\mathbf{W} - \rho\sqrt{c}\mathbf{U}\|_F^2 + \left\| \mathbf{\Theta} - \frac{\sqrt{c}}{\rho}\mathbf{U} \right\|_F^2 \right\} \\ &\leq 2 \left\| \mathbf{\Theta} - \frac{1}{\rho^2}\mathbf{W} \right\|_F^2 + \left(1 + \frac{2\lambda_W}{n\lambda_H} \right) \min_{\mathbf{U} \in \mathcal{O}^{d \times K}} \|\mathbf{W} - \rho\sqrt{c}\mathbf{U}\|_F^2, \end{aligned} \quad (36)$$

where the inequality follows from $\|\mathbf{A} + \mathbf{B}\|_F^2 \leq 2\|\mathbf{A}\|_F^2 + 2\|\mathbf{B}\|_F^2$ for any \mathbf{A}, \mathbf{B} of the same size. Then, we bound each term above in turn. Let $(\mathbf{W}^*, \mathbf{\Theta}^*) \in \mathcal{X}_f$ be such that $\text{dist}((\mathbf{W}, \mathbf{\Theta}), \mathcal{X}_f) = \|(\mathbf{W}, \mathbf{\Theta}) - (\mathbf{W}^*, \mathbf{\Theta}^*)\|_F \leq \delta$. According to (33) and (36), we have

$$\|\mathbf{W}\| \leq \|\mathbf{W} - \mathbf{W}^*\| + \|\mathbf{W}^*\| \leq \delta + \rho\sqrt{c} \leq 2\rho\sqrt{c}, \quad \|\mathbf{\Theta}\| \leq \|\mathbf{\Theta} - \mathbf{\Theta}^*\| + \|\mathbf{\Theta}^*\| \leq \frac{2\sqrt{c}}{\rho}. \quad (37)$$

It directly follows from (12) that

$$\|\nabla_{\mathbf{W}} f(\mathbf{W}, \mathbf{\Theta})\|_F = \|(\mathbf{\Theta}\mathbf{\Theta}^T + \lambda_W K \mathbf{I}) \mathbf{W} - \mathbf{\Theta}\|_F, \quad (38)$$

$$\|\nabla_{\mathbf{\Theta}} f(\mathbf{W}, \mathbf{\Theta})\|_F = \|(\mathbf{W}\mathbf{W}^T + n\lambda_H K \mathbf{I}) \mathbf{\Theta} - \mathbf{W}\|_F. \quad (39)$$

Summing up $\sqrt{n\lambda_H}\|\nabla_{\mathbf{W}}f(\mathbf{W}, \boldsymbol{\Theta})\|_F + \sqrt{\lambda_W}\|\nabla_{\boldsymbol{\Theta}}f(\mathbf{W}, \boldsymbol{\Theta})\|_F$ yields

$$\begin{aligned} & \sqrt{n\lambda_H}\|\nabla_{\mathbf{W}}f(\mathbf{W}, \boldsymbol{\Theta})\|_F + \sqrt{\lambda_W}\|\nabla_{\boldsymbol{\Theta}}f(\mathbf{W}, \boldsymbol{\Theta})\|_F \\ &= \sqrt{n\lambda_H}\|(\boldsymbol{\Theta}\boldsymbol{\Theta}^T + \lambda_W K\mathbf{I})\mathbf{W} - \boldsymbol{\Theta}\|_F + \sqrt{\lambda_W}\|(\mathbf{W}\mathbf{W}^T + n\lambda_H K\mathbf{I})\boldsymbol{\Theta} - \mathbf{W}\|_F. \end{aligned} \quad (40)$$

This, together with (37), yields

$$\begin{aligned} & 2 \max\left\{\rho, \frac{1}{\rho}\right\} \sqrt{c\lambda_{\max}} (\|\nabla_{\mathbf{W}}f(\mathbf{W}, \boldsymbol{\Theta})\|_F + \|\nabla_{\boldsymbol{\Theta}}f(\mathbf{W}, \boldsymbol{\Theta})\|_F) \\ & \geq 2 \max\left\{\rho, \frac{1}{\rho}\right\} \sqrt{c} \left(\sqrt{n\lambda_H}\|\nabla_{\mathbf{W}}f(\mathbf{W}, \boldsymbol{\Theta})\|_F + \sqrt{\lambda_W}\|\nabla_{\boldsymbol{\Theta}}f(\mathbf{W}, \boldsymbol{\Theta})\|_F \right) \\ & \geq \sqrt{n\lambda_H}\|(\boldsymbol{\Theta}\boldsymbol{\Theta}^T + \lambda_W K\mathbf{I})\mathbf{W} - \boldsymbol{\Theta}\|_F \|\mathbf{W}\| + \sqrt{\lambda_W}\|(\mathbf{W}\mathbf{W}^T + n\lambda_H K\mathbf{I})\boldsymbol{\Theta} - \mathbf{W}\|_F \|\boldsymbol{\Theta}\| \\ & \geq \sqrt{\lambda_{\min}} (\|(\boldsymbol{\Theta}\boldsymbol{\Theta}^T + \lambda_W K\mathbf{I})\mathbf{W}\mathbf{W}^T - \boldsymbol{\Theta}\mathbf{W}^T\|_F + \|(\mathbf{W}\mathbf{W}^T + n\lambda_H K\mathbf{I})\boldsymbol{\Theta}\boldsymbol{\Theta}^T - \mathbf{W}\boldsymbol{\Theta}^T\|_F) \\ & \geq K\sqrt{\lambda_{\min}} \|\lambda_W \mathbf{W}\mathbf{W}^T - n\lambda_H \boldsymbol{\Theta}\boldsymbol{\Theta}^T\|_F, \end{aligned}$$

where the last inequality uses the triangle inequality. This implies

$$\|\lambda_W \mathbf{W}\mathbf{W}^T - n\lambda_H \boldsymbol{\Theta}\boldsymbol{\Theta}^T\|_F \leq \kappa_1 (\|\nabla_{\mathbf{W}}f(\mathbf{W}, \boldsymbol{\Theta})\|_F + \|\nabla_{\boldsymbol{\Theta}}f(\mathbf{W}, \boldsymbol{\Theta})\|_F), \quad (41)$$

where

$$\kappa_1 := \frac{2 \max\{\rho, 1/\rho\} \sqrt{c\lambda_{\max}}}{K\sqrt{\lambda_{\min}}}. \quad (42)$$

By letting $\boldsymbol{\Phi} := \sqrt{\lambda_W}\mathbf{W} - \sqrt{n\lambda_H}\boldsymbol{\Theta}$, we obtain

$$\begin{aligned} & \sqrt{\lambda_{\max}} (\|\nabla_{\mathbf{W}}f(\mathbf{W}, \boldsymbol{\Theta})\|_F + \|\nabla_{\boldsymbol{\Theta}}f(\mathbf{W}, \boldsymbol{\Theta})\|_F) \\ & \geq \sqrt{n\lambda_H}\|(\boldsymbol{\Theta}\boldsymbol{\Theta}^T + \lambda_W K\mathbf{I})\mathbf{W} - \boldsymbol{\Theta}\|_F + \sqrt{\lambda_W}\|(\mathbf{W}\mathbf{W}^T + n\lambda_H K\mathbf{I})\boldsymbol{\Theta} - \mathbf{W}\|_F \\ & \geq \left\| \sqrt{n\lambda_H}(\boldsymbol{\Theta}\boldsymbol{\Theta}^T + \lambda_W K\mathbf{I})\mathbf{W} - \sqrt{\lambda_W}(\mathbf{W}\mathbf{W}^T + n\lambda_H K\mathbf{I})\boldsymbol{\Theta} + \boldsymbol{\Phi} \right\|_F \\ & = \left\| \sqrt{n\lambda_H} \left(\boldsymbol{\Theta}\boldsymbol{\Theta}^T - \frac{\lambda_W}{n\lambda_H} \mathbf{W}\mathbf{W}^T \right) \mathbf{W} + \left(\frac{\sqrt{\lambda_W}}{\sqrt{n\lambda_H}} \mathbf{W}\mathbf{W}^T + (K\sqrt{n\lambda_H\lambda_W} + 1)\mathbf{I} \right) \boldsymbol{\Phi} \right\|_F \\ & \geq \left\| \left(\frac{\sqrt{\lambda_W}}{\sqrt{n\lambda_H}} \mathbf{W}\mathbf{W}^T + (K\sqrt{n\lambda_H\lambda_W} + 1)\mathbf{I} \right) \boldsymbol{\Phi} \right\|_F - \sqrt{n\lambda_H}\|\mathbf{W}\| \left\| \boldsymbol{\Theta}\boldsymbol{\Theta}^T - \frac{\lambda_W}{n\lambda_H} \mathbf{W}\mathbf{W}^T \right\|_F \\ & \geq \left(1 + K\sqrt{n\lambda_H\lambda_W} \right) \|\boldsymbol{\Phi}\|_F - \frac{2\rho\sqrt{c}}{\sqrt{n\lambda_H}} \|\lambda_W \mathbf{W}\mathbf{W}^T - n\lambda_H \boldsymbol{\Theta}\boldsymbol{\Theta}^T\|_F, \end{aligned}$$

where the first inequality uses (40), and the last inequality is due to (37). This, together with (41), implies

$$\left\| \boldsymbol{\Theta} - \frac{1}{\rho^2} \mathbf{W} \right\|_F = \frac{1}{\sqrt{n\lambda_H}} \left\| \sqrt{\lambda_W}\mathbf{W} - \sqrt{n\lambda_H}\boldsymbol{\Theta} \right\|_F \leq \kappa_2 (\|\nabla_{\mathbf{W}}f(\mathbf{W}, \boldsymbol{\Theta})\|_F + \|\nabla_{\boldsymbol{\Theta}}f(\mathbf{W}, \boldsymbol{\Theta})\|_F), \quad (43)$$

where

$$\kappa_2 := \frac{\sqrt{\lambda_{\max}} + 2\rho\kappa_1\sqrt{c}/\sqrt{n\lambda_H}}{\sqrt{n\lambda_H} (1 + K\sqrt{n\lambda_H\lambda_W})}. \quad (44)$$

Using (39), we compute

$$\begin{aligned}
\|\nabla_{\Theta} f(\mathbf{W}, \Theta)\|_F &= \left\| (\mathbf{W}\mathbf{W}^T + n\lambda_H K \mathbf{I}) \left(\Theta - \frac{1}{\rho^2} \mathbf{W} + \frac{1}{\rho^2} \mathbf{W} \right) - \mathbf{W} \right\|_F \\
&\geq \frac{1}{\rho^2} \left\| (\mathbf{W}\mathbf{W}^T + n\lambda_H K \mathbf{I}) \mathbf{W} - \rho^2 \mathbf{W} \right\|_F - \left\| (\mathbf{W}\mathbf{W}^T + n\lambda_H K \mathbf{I}) \left(\Theta - \frac{1}{\rho^2} \mathbf{W} \right) \right\|_F \\
&= \frac{1}{\rho^2} \left\| \mathbf{W} (\mathbf{W}^T \mathbf{W} - \rho^2 c \mathbf{I}) \right\|_F - \left\| (\mathbf{W}\mathbf{W}^T + n\lambda_H K \mathbf{I}) \left(\Theta - \frac{1}{\rho^2} \mathbf{W} \right) \right\|_F \\
&\geq \frac{\sigma_{\min}(\mathbf{W})}{\rho^2} \left\| \mathbf{W}^T \mathbf{W} - \rho^2 c \mathbf{I} \right\|_F - (\|\mathbf{W}\|^2 + n\lambda_H K) \left\| \Theta - \frac{1}{\rho^2} \mathbf{W} \right\|_F,
\end{aligned}$$

where the second equality follows from $\rho^2 - n\lambda_H K = \sqrt{n\lambda_H}/\sqrt{\lambda_W} - n\lambda_H K = \rho^2 c$. This, together with (37) and (43), yields

$$\frac{\sigma_{\min}(\mathbf{W})}{\rho^2} \left\| \mathbf{W}^T \mathbf{W} - \rho^2 c \mathbf{I} \right\|_F \leq (\kappa_2(4\rho^2 c + n\lambda_H K) + 1) (\|\nabla_{\mathbf{W}} f(\mathbf{W}, \Theta)\|_F + \|\nabla_{\Theta} f(\mathbf{W}, \Theta)\|_F). \quad (45)$$

Let $\mathbf{W} = \mathbf{U}_1 \mathbf{\Lambda}_1 \mathbf{V}_1^T$ be the thin singular value decomposition of \mathbf{W} , where $\mathbf{U}_1 \in \mathcal{O}^{d \times K}$, $\mathbf{V}_1 \in \mathcal{O}^K$, and $\mathbf{\Lambda}_1 \in \mathbb{R}^{K \times K}$ is a diagonal matrix. Now, we compute

$$\begin{aligned}
\left\| \mathbf{W}^T \mathbf{W} - \rho^2 c \mathbf{I} \right\|_F &= \left\| \mathbf{\Lambda}_1^2 - \rho^2 c \mathbf{I} \right\|_F = \left\| (\mathbf{\Lambda}_1 - \rho\sqrt{c} \mathbf{I})(\mathbf{\Lambda}_1 + \rho\sqrt{c} \mathbf{I}) \right\|_F \\
&\geq \rho\sqrt{c} \left\| \mathbf{\Lambda}_1 - \rho\sqrt{c} \mathbf{I} \right\|_F = \rho\sqrt{c} \left\| \mathbf{U}_1 (\mathbf{\Lambda}_1 - \rho\sqrt{c} \mathbf{I}) \mathbf{V}_1^T \right\|_F \\
&= \rho\sqrt{c} \left\| \mathbf{W} - \rho\sqrt{c} \mathbf{U}_1 \mathbf{V}_1^T \right\|_F \geq \rho\sqrt{c} \min_{\mathbf{U} \in \mathcal{O}^{d \times K}} \left\| \mathbf{W} - \rho\sqrt{c} \mathbf{U} \right\|_F. \quad (46)
\end{aligned}$$

Let $\mathbf{U}^* \in \mathcal{O}^{d \times K}$ be such that

$$\left\| \mathbf{W} - \rho\sqrt{c} \mathbf{U}^* \right\|_F = \min_{\mathbf{U} \in \mathcal{O}^{d \times K}} \left\| \mathbf{W} - \rho\sqrt{c} \mathbf{U} \right\|_F.$$

Using Weyl's inequality, we obtain

$$\sigma_{\min}(\mathbf{W}) \geq \sigma_{\min}(\rho\sqrt{c} \mathbf{U}^*) - \left\| \mathbf{W} - \rho\sqrt{c} \mathbf{U}^* \right\|_F \geq \rho\sqrt{c} - \delta \geq \frac{1}{2} \rho\sqrt{c},$$

where the second inequality follows from the equality in (36) and (33), and the last inequality is due to $\delta \leq \rho\sqrt{c}/2$ by (33). Substituting this and (46) into (45) yields

$$\min_{\mathbf{U} \in \mathcal{O}^{d \times K}} \left\| \mathbf{W} - \rho\sqrt{c} \mathbf{U} \right\|_F \leq \frac{2}{c} (\kappa_2(4\rho^2 c + n\lambda_H K) + 1) (\|\nabla_{\mathbf{W}} f(\mathbf{W}, \Theta)\|_F + \|\nabla_{\Theta} f(\mathbf{W}, \Theta)\|_F)$$

This, together with (43) and (36), yields

$$\begin{aligned}
\text{dist}^2((\mathbf{W}, \Theta), \mathcal{X}_f) &\leq 2\kappa_2^2 (\|\nabla_{\mathbf{W}} f(\mathbf{W}, \Theta)\|_F + \|\nabla_{\Theta} f(\mathbf{W}, \Theta)\|_F)^2 + \left(1 + \frac{2\lambda_W}{n\lambda_H} \right) \\
&\quad \frac{4}{c^2} (\kappa_2(4\rho^2 c + n\lambda_H K) + 1)^2 (\|\nabla_{\mathbf{W}} f(\mathbf{W}, \Theta)\|_F + \|\nabla_{\Theta} f(\mathbf{W}, \Theta)\|_F)^2 \\
&\leq \left(4\kappa_2^2 + \frac{8}{c^2} (\kappa_2(4\rho^2 c + n\lambda_H K) + 1)^2 \right) \|\nabla f(\mathbf{W}, \Theta)\|_F^2.
\end{aligned}$$

Then, we complete the proof. \square

3 MSE Loss with a Bias Term

We consider the following problems:

$$\min_{\mathbf{W}, \mathbf{H}, \mathbf{b}} F(\mathbf{W}, \mathbf{H}, \mathbf{b}) = \frac{1}{2N} \|\mathbf{W}^T \mathbf{H} + \mathbf{b} \mathbf{1}_N^T - \mathbf{Y}\|_F^2 + \frac{\lambda_W}{2} \|\mathbf{W}\|_F^2 + \frac{\lambda_H}{2} \|\mathbf{H}\|_F^2 + \frac{\lambda_b}{2} \|\mathbf{b}\|^2, \quad (47)$$

$$\min_{\mathbf{W}, \mathbf{\Theta}, \mathbf{b}} f(\mathbf{W}, \mathbf{\Theta}, \mathbf{b}) = \frac{1}{2K} \|\mathbf{W}^T \mathbf{\Theta} + \mathbf{b} \mathbf{1}_K^T - \mathbf{I}_K\|_F^2 + \frac{\lambda_W}{2} \|\mathbf{W}\|_F^2 + \frac{n\lambda_H}{2} \|\mathbf{\Theta}\|_F^2 + \frac{\lambda_b}{2} \|\mathbf{b}\|^2, \quad (48)$$

where $\lambda_W, \lambda_H, \lambda_b > 0$ are the penalties for \mathbf{W} , \mathbf{H} , and \mathbf{b} , respectively. One can easily verify

$$F(\mathbf{W}, \mathbf{H}, \mathbf{b}) = \frac{1}{n} \sum_{i=1}^n f(\mathbf{W}, \mathbf{H}_i, \mathbf{b}). \quad (49)$$

To proceed, let

$$\bar{\mathbf{V}} := \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{6}} & \cdots & \frac{1}{\sqrt{(K-1)K}} \\ -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{6}} & \cdots & \frac{1}{\sqrt{(K-1)K}} \\ 0 & -\frac{\sqrt{2}}{\sqrt{3}} & \cdots & \frac{1}{\sqrt{(K-1)K}} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & -\frac{\sqrt{K-1}}{\sqrt{K}} \end{bmatrix} \in \mathcal{O}^{K \times (K-1)}. \quad (50)$$

Proposition 2. *The optimal solution set of Problem (48) can be characterized as follows:*

(i) *If $\lambda_W \lambda_H \geq 1/(nK^2)$, we have*

$$\mathcal{X}_f = \left\{ \left(\mathbf{0}, \mathbf{0}, \frac{1}{K(1+\lambda_b)} \mathbf{1}_K \right) \right\}.$$

(ii) *If $\lambda_b^2/(nK^2(1+\lambda_b)^2) \leq \lambda_W \lambda_H < 1/(nK^2)$, we have*

$$\mathcal{X}_f = \left\{ \left(\mathbf{W}, \mathbf{\Theta}, \frac{1}{K(1+\lambda_b)} \mathbf{1}_K \right) : \mathbf{W} = \frac{\sqrt[4]{n\lambda_H}}{\sqrt[4]{\lambda_W}} \sqrt{c} \mathbf{U} \bar{\mathbf{V}}^T, \mathbf{\Theta} = \frac{\sqrt{\lambda_W}}{\sqrt{n\lambda_H}} \mathbf{W}, \mathbf{U} \in \mathcal{O}^{d \times (K-1)} \right\}, \quad (51)$$

where $c := 1 - K\sqrt{n\lambda_W\lambda_H}$.

(iii) *If $\lambda_W \lambda_H < \lambda_b^2/(nK^2(1+\lambda_b)^2)$, we have*

$$\mathcal{X} = \left\{ \left(\mathbf{W}, \mathbf{\Theta}, \frac{\sqrt{n\lambda_W\lambda_H}}{\lambda_b} \mathbf{1}_K \right) : \mathbf{W} = \frac{\sqrt[4]{n\lambda_H}}{\sqrt[4]{\lambda_W}} \mathbf{U} \begin{bmatrix} \sqrt{c'} \bar{\mathbf{V}}^T \\ \sqrt{c'/K} \mathbf{1}^T \end{bmatrix}, \mathbf{\Theta} = \frac{\sqrt{\lambda_W}}{\sqrt{n\lambda_H}} \mathbf{W}, \mathbf{U} \in \mathcal{O}^{d \times K} \right\},$$

where $c' := 1 - K(1+\lambda_b)\sqrt{n\lambda_W\lambda_H}/\lambda_b$.

Proof. Using the first-order optimality of \mathbf{b} , we obtain

$$\mathbf{b}^* = \frac{1}{K(1+\lambda_b)} (\mathbf{I} - \mathbf{W}^T \mathbf{\Theta}) \mathbf{1}. \quad (52)$$

Substituting this back into Problem (48) yields

$$\begin{aligned} & \frac{1}{2K} \|\mathbf{W}^T \mathbf{\Theta} - \mathbf{I}\|_F^2 - \frac{1}{2K^2(1+\lambda_b)} \|(\mathbf{W}^T \mathbf{\Theta} - \mathbf{I}) \mathbf{1}\|^2 + \frac{\lambda_W}{2} \|\mathbf{W}\|_F^2 + \frac{n\lambda_H}{2} \|\mathbf{H}\|_F^2 \\ &= \frac{1}{2K} \|(\mathbf{W}^T \mathbf{\Theta} - \mathbf{I})(\mathbf{I} - \alpha \mathbf{1} \mathbf{1}^T)\|_F^2 + \frac{\lambda_W}{2} \|\mathbf{W}\|_F^2 + \frac{n\lambda_H}{2} \|\mathbf{\Theta}\|_F^2, \end{aligned}$$

where $\alpha := \frac{1+\sqrt{1-1/(1+\lambda_b)}}{K}$. Then, it suffices to consider

$$\min_{\mathbf{W}, \boldsymbol{\Theta} \in \mathbb{R}^{d \times K}} h(\mathbf{W}, \boldsymbol{\Theta}) := \frac{1}{2K} \|(\mathbf{W}^T \boldsymbol{\Theta} - \mathbf{I})(\mathbf{I} - \alpha \mathbf{1}\mathbf{1}^T)\|_F^2 + \frac{\lambda_W}{2} \|\mathbf{W}\|_F^2 + \frac{n\lambda_H}{2} \|\boldsymbol{\Theta}\|_F^2. \quad (53)$$

Noting that $\mathbf{I} - \alpha \mathbf{1}\mathbf{1}^T = \mathbf{P} + (1 - K\alpha)\mathbf{P}^\perp$, we have

$$\begin{aligned} h(\mathbf{W}, \boldsymbol{\Theta}) &= \frac{1}{2K} \|(\mathbf{W}^T \boldsymbol{\Theta} - \mathbf{I})(\mathbf{P} + (1 - K\alpha)\mathbf{P}^\perp)\|_F^2 + \frac{\lambda_W}{2} \|\mathbf{W}\|_F^2 + \frac{n\lambda_H}{2} \|\boldsymbol{\Theta}\|_F^2 \\ &= \frac{1}{2K} \|(\mathbf{W}^T \boldsymbol{\Theta} - \mathbf{I})\mathbf{P}\|_F^2 + \frac{\lambda_b}{2K(1+\lambda_b)} \|(\mathbf{W}^T \boldsymbol{\Theta} - \mathbf{I})\mathbf{P}^\perp\|_F^2 + \frac{\lambda_W}{2} \|\mathbf{W}\|_F^2 + \frac{n\lambda_H}{2} \|\boldsymbol{\Theta}\|_F^2, \end{aligned}$$

where the second equality follows from $(1 - K\alpha)^2 = 1/(1 + \lambda_b)$. For ease of exposition, let

$$\mathbf{W}_1 := \mathbf{W}\mathbf{P}, \mathbf{W}_2 := \mathbf{W}\mathbf{P}^\perp, \boldsymbol{\Theta}_1 := \boldsymbol{\Theta}\mathbf{P}, \boldsymbol{\Theta}_2 := \boldsymbol{\Theta}\mathbf{P}^\perp. \quad (54)$$

Since \mathbf{P} and \mathbf{P}^\perp are both projection matrices, we have

$$\begin{aligned} h(\mathbf{W}, \boldsymbol{\Theta}) &\geq \frac{1}{2K} \|\mathbf{P}(\mathbf{W}^T \boldsymbol{\Theta} - \mathbf{I})\mathbf{P}\|_F^2 + \frac{\lambda_b}{2K(1+\lambda_b)} \|\mathbf{P}^\perp(\mathbf{W}^T \boldsymbol{\Theta} - \mathbf{I})\mathbf{P}^\perp\|_F^2 + \frac{\lambda_W}{2} \|\mathbf{W}\|_F^2 \\ &\quad + \frac{n\lambda_H}{2} \|\boldsymbol{\Theta}\|_F^2 = \frac{1}{2K} \|\mathbf{W}_1^T \boldsymbol{\Theta}_1 - \mathbf{P}\|_F^2 + \frac{\lambda_W}{2} \|\mathbf{W}_1\|_F^2 + \frac{n\lambda_H}{2} \|\boldsymbol{\Theta}_1\|_F^2 + \\ &\quad \frac{\lambda_b}{2K(1+\lambda_b)} \|\mathbf{W}_2^T \boldsymbol{\Theta}_2 - \mathbf{P}^\perp\|_F^2 + \frac{\lambda_W}{2} \|\mathbf{W}_2\|_F^2 + \frac{n\lambda_H}{2} \|\boldsymbol{\Theta}_2\|_F^2 \end{aligned} \quad (55)$$

where the first inequality becomes equality if and only if $\mathbf{W}_2^T \boldsymbol{\Theta}_1 = \mathbf{0}$ and $\mathbf{W}_1^T \boldsymbol{\Theta}_2 = \mathbf{0}$. This yields

$$\begin{aligned} h(\mathbf{W}, \boldsymbol{\Theta}) &\geq \min_{\mathbf{W}_1, \boldsymbol{\Theta}_1} \left\{ \frac{1}{2K} \|\mathbf{W}_1^T \boldsymbol{\Theta}_1 - \mathbf{P}\|_F^2 + \frac{\lambda_W}{2} \|\mathbf{W}_1\|_F^2 + \frac{n\lambda_H}{2} \|\boldsymbol{\Theta}_1\|_F^2 \right\} + \\ &\quad \frac{\lambda_b}{1+\lambda_b} \min_{\mathbf{W}_2, \boldsymbol{\Theta}_2} \left\{ \frac{1}{2K} \|\mathbf{W}_2^T \boldsymbol{\Theta}_2 - \mathbf{P}^\perp\|_F^2 + \frac{\lambda_W(1+\lambda_b)}{2\lambda_b} \|\mathbf{W}_2\|_F^2 + \frac{n\lambda_H(1+\lambda_b)}{2\lambda_b} \|\boldsymbol{\Theta}_2\|_F^2 \right\}, \end{aligned} \quad (56)$$

where the inequality becomes equality if and only if there exists $\mathbf{W}^*, \boldsymbol{\Theta}^*$ such that the optimal solutions $(\mathbf{W}_1^*, \boldsymbol{\Theta}_1^*)$ and $(\mathbf{W}_2^*, \boldsymbol{\Theta}_2^*)$ of the above optimization problems satisfy $\mathbf{W}_1^* = \mathbf{W}^*\mathbf{P}$, $\mathbf{W}_2^* = \mathbf{W}^*\mathbf{P}^\perp$, $\boldsymbol{\Theta}_1^* = \boldsymbol{\Theta}^*\mathbf{P}$, $\boldsymbol{\Theta}_2^* = \boldsymbol{\Theta}^*\mathbf{P}^\perp$. Now, we can optimize the above two optimization problems, respectively. One can verify that the matrix \mathbf{P} admits the eigenvalue decomposition $\mathbf{P} = \mathbf{V}\boldsymbol{\Sigma}\mathbf{V}^T$, where $\boldsymbol{\Sigma} = \begin{bmatrix} \mathbf{I}_{K-1} & \mathbf{0} \\ \mathbf{0} & 0 \end{bmatrix}$, and $\mathbf{V} = [\bar{\mathbf{V}} \ \mathbf{1}_K/\sqrt{K}] \in \mathcal{O}^K$ with $\bar{\mathbf{V}}\bar{\mathbf{V}}^T = \mathbf{P}$ and $\mathbf{v}\mathbf{v}^T = \mathbf{P}^\perp$. Then, we can verify that $(\mathbf{W}_1^*, \boldsymbol{\Theta}_1^*)$ is an optimal solution to the first optimization problem in (56) if and only if $(\mathbf{W}_1^*\mathbf{V}, \boldsymbol{\Theta}_1^*\mathbf{V})$ is an optimal solution to

$$\min_{\mathbf{W}_1, \boldsymbol{\Theta}_1} \left\{ \frac{1}{2K} \|\mathbf{W}_1^T \boldsymbol{\Theta}_1 - \boldsymbol{\Sigma}\|_F^2 + \frac{\lambda_W}{2} \|\mathbf{W}_1\|_F^2 + \frac{n\lambda_H}{2} \|\boldsymbol{\Theta}_1\|_F^2 \right\}$$

This, together with Lemma 1, implies that the solution set of the first optimization problem in (56) with variables $\mathbf{W}_1, \boldsymbol{\Theta}_1$ is as follows:

$$\mathcal{X}_1 := \left\{ (\mathbf{W}\mathbf{V}^T, \boldsymbol{\Theta}\mathbf{V}^T) : \mathbf{W} = \frac{\sqrt[4]{n\lambda_H}}{\sqrt[4]{\lambda_W}} \mathbf{U} \begin{bmatrix} (\max\{c, 0\})^{1/2} \mathbf{I}_{K-1} & \mathbf{0} \\ \mathbf{0} & 0 \end{bmatrix}, \boldsymbol{\Theta} = \frac{\sqrt{\lambda_W}}{\sqrt{n\lambda_H}} \mathbf{W}, \mathbf{U} \in \mathcal{O}^{d \times K} \right\}.$$

where $c := 1 - K\sqrt{n\lambda_W\lambda_H}$. Note that $\mathbf{P}^\perp = \mathbf{I} - \mathbf{P} = \mathbf{V}(\mathbf{I} - \boldsymbol{\Sigma})\mathbf{V}^T$. By the same argument, we show that the solution set of the second optimization problem in (56) with variables $\mathbf{W}_2, \boldsymbol{\Theta}_2$ is as follows:

$$\mathcal{X}_2 := \left\{ (\mathbf{W}\mathbf{V}^T, \boldsymbol{\Theta}\mathbf{V}^T) : \mathbf{W} = \frac{\sqrt[4]{n\lambda_H}}{\sqrt[4]{\lambda_W}} \mathbf{U} \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & (\max\{c', 0\}, 0)^{1/2} \end{bmatrix}, \boldsymbol{\Theta} = \frac{\sqrt{\lambda_W}}{\sqrt{n\lambda_H}} \mathbf{W}, \mathbf{U} \in \mathcal{O}^{d \times K} \right\},$$

where $c' := 1 - K(1 + \lambda_b)\sqrt{n\lambda_W\lambda_H}/\lambda_b$. Now, we discuss the optimal solution set of Problem (53) denoted by \mathcal{X} case by case.

Case 1. If $\lambda_W\lambda_H \geq \frac{1}{nK^2}$, then $c = c' = 0$. Therefore, we obtain $\mathcal{X}_1 = \{(\mathbf{0}, \mathbf{0})\}$ and $\mathcal{X}_2 = \{(\mathbf{0}, \mathbf{0})\}$. Then, one can verify that the inequalities become equalities in (55) and (56) for any $(\mathbf{W}_1^*, \boldsymbol{\Theta}_1^*) \in \mathcal{X}_1$ and $(\mathbf{W}_2^*, \boldsymbol{\Theta}_2^*) \in \mathcal{X}_2$. This directly implies $\mathcal{X} = \{(\mathbf{0}, \mathbf{0})\}$.

Case 2. If $\frac{\lambda_b^2}{nK^2(1+\lambda_b)^2} \leq \lambda_W\lambda_H < \frac{1}{nK^2}$, then $c > 0$ and $c' = 0$. Therefore, we obtain $\mathcal{X}_2 = \{(\mathbf{0}, \mathbf{0})\}$ and

$$\mathcal{X}_1 = \left\{ (\mathbf{W}, \boldsymbol{\Theta}) : \mathbf{W} = \frac{\sqrt[4]{n\lambda_H}}{\sqrt[4]{\lambda_W}} \sqrt{c} \mathbf{U} \bar{\mathbf{V}}^T, \boldsymbol{\Theta} = \frac{\sqrt{\lambda_W}}{\sqrt{n\lambda_H}} \mathbf{W}, \mathbf{U} \in \mathcal{O}^{d \times (K-1)} \right\}.$$

Then, we can verify that the inequalities become equalities in (55) and (56) for any $(\mathbf{W}_1^*, \boldsymbol{\Theta}_1^*) \in \mathcal{X}_1$ and $(\mathbf{W}_2^*, \boldsymbol{\Theta}_2^*) \in \mathcal{X}_2$. This directly implies $\mathcal{X} = \mathcal{X}_1$.

Case 3. If $\lambda_W\lambda_H < \frac{\lambda_b^2}{K^2(1+\lambda_b)^2}$, then $c > 0$ and $c' > 0$. Therefore, we obtain that \mathcal{X}_1 takes the same form as above and

$$\mathcal{X}_2 = \left\{ (\mathbf{W}, \boldsymbol{\Theta}) : \mathbf{W} = \frac{\sqrt[4]{n\lambda_H}\sqrt{c'}}{\sqrt[4]{\lambda_W}\sqrt{K}} \mathbf{u} \mathbf{1}^T, \boldsymbol{\Theta} = \frac{\sqrt{\lambda_W}}{\sqrt{n\lambda_H}} \mathbf{W}, \|\mathbf{u}\| = 1 \right\}.$$

For any $(\mathbf{W}_1^*, \boldsymbol{\Theta}_1^*) \in \mathcal{X}_1$ and $(\mathbf{W}_2^*, \boldsymbol{\Theta}_2^*) \in \mathcal{X}_2$, we can verify that $\mathbf{W}_1^{*T} \mathbf{W}_2^* = \mathbf{0}$ holds if and only if $\mathbf{U}^T \mathbf{u} = 0$ due to $\bar{\mathbf{V}}^T \bar{\mathbf{V}} = \mathbf{I}$ and $\|\mathbf{1}_K\| = \sqrt{K}$. As a result, the inequality in (55) becomes equality. It follows from $\mathbf{P} = \bar{\mathbf{V}} \bar{\mathbf{V}}^T$ and $\mathbf{P}^\perp = \mathbf{1}\mathbf{1}^T/K$ that the inequality in (56) becomes equality. Using (54), $\mathbf{U}^T \mathbf{u} = 0$, $\mathbf{U} \in \mathcal{O}^{d \times (K-1)}$, and $\|\mathbf{u}\| = 1$, we have

$$\mathcal{X} = \left\{ (\mathbf{W}, \boldsymbol{\Theta}) : \mathbf{W} = \frac{\sqrt[4]{n\lambda_H}}{\sqrt[4]{\lambda_W}} \mathbf{U} \begin{bmatrix} \sqrt{c} \bar{\mathbf{V}}^T \\ \sqrt{c'/K} \mathbf{1}_K^T \end{bmatrix}, \boldsymbol{\Theta} = \frac{\sqrt{\lambda_W}}{\sqrt{n\lambda_H}} \mathbf{W}, \mathbf{U} \in \mathcal{O}^{d \times K} \right\}.$$

Combining the above cases with (52) yields that desired result. \square

Theorem 2. Suppose that $\lambda_b^2/(nK^2(1+\lambda_b)^2) < \lambda_W\lambda_H < 1/nK^2$. For all $(\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b})$ satisfying

$$\text{dist}((\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b}), \mathcal{X}_f) \leq \delta := \frac{1}{2} \min \left\{ \frac{\sqrt[4]{\lambda_W}}{\sqrt[4]{n\lambda_H}}, \frac{\sqrt[4]{n\lambda_H}}{\sqrt[4]{\lambda_W}} \right\} \left(1 - K\sqrt{n\lambda_W\lambda_H}\right)^{1/2}, \quad (57)$$

it holds that

$$\text{dist}((\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b}), \mathcal{X}_f) \leq \kappa \|\nabla f(\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b})\|_F, \quad (58)$$

where κ is a constant that depends on λ_W , λ_H , and λ_b .

Proof. To simplify our notation, let

$$\alpha := \frac{1 + \sqrt{1 - 1/(1 + \lambda_b)}}{K}, \quad \beta := \frac{\lambda_b}{1 + \lambda_b}, \quad c := 1 - K\sqrt{n\lambda_W\lambda_H}, \quad \rho := \frac{\sqrt[4]{n\lambda_H}}{\sqrt[4]{\lambda_W}}. \quad (59)$$

According to the proof of Proposition 2 and $\lambda_b^2/(nK^2(1+\lambda_b)^2) \leq \lambda_W \lambda_H < 1/nK^2$, the optimal solution set of Problem (53) takes the form of

$$\mathcal{X}_h := \left\{ (\mathbf{W}, \boldsymbol{\Theta}) : \mathbf{W} = \rho\sqrt{c}\mathbf{U}\bar{\mathbf{V}}^T, \boldsymbol{\Theta} = \frac{1}{\rho^2}\mathbf{W}, \mathbf{U} \in \mathcal{O}^{d \times (K-1)} \right\}. \quad (60)$$

We claim that for all $(\mathbf{W}, \boldsymbol{\Theta})$ satisfying

$$\text{dist}((\mathbf{W}, \boldsymbol{\Theta}), \mathcal{X}_h) \leq \delta_1 := \frac{1}{2} \min \left\{ \rho, \frac{1}{\rho} \right\} \sqrt{c}, \quad (61)$$

it holds that

$$\text{dist}((\mathbf{W}, \boldsymbol{\Theta}), \mathcal{X}_h) \leq \kappa_1 \|\nabla h(\mathbf{W}, \boldsymbol{\Theta})\|_F, \quad (62)$$

where δ_1, κ_1 are positive constants that depend on λ_W and λ_H . Let $(\mathbf{W}^*, \boldsymbol{\Theta}^*) \in \mathcal{X}_h$ be such that $\text{dist}((\mathbf{W}, \boldsymbol{\Theta}), \mathcal{X}_h) = \|(\mathbf{W}, \boldsymbol{\Theta}) - (\mathbf{W}^*, \boldsymbol{\Theta}^*)\|_F \leq \delta_1$. According to (60) and $\delta_1 \leq \min\{\rho, 1/\rho\}\sqrt{c}$, we have

$$\|\mathbf{W}\| \leq \|\mathbf{W} - \mathbf{W}^*\| + \|\mathbf{W}^*\| \leq \delta_1 + \rho\sqrt{c} \leq 2\rho\sqrt{c}, \quad \|\boldsymbol{\Theta}\| \leq \|\boldsymbol{\Theta} - \boldsymbol{\Theta}^*\| + \|\boldsymbol{\Theta}^*\| \leq \frac{2\sqrt{c}}{\rho}. \quad (63)$$

Since $f(\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b})$ is strongly convex w.r.t. \mathbf{b} with constant λ_b , we obtain

$$\lambda_b \|\mathbf{b} - \mathbf{b}^*\|^2 \leq \langle \nabla_{\mathbf{b}} f(\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b}) - \nabla_{\mathbf{b}} f(\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b}^*), \mathbf{b} - \mathbf{b}^* \rangle,$$

where \mathbf{b}^* is defined in (52). This, together with $\nabla_{\mathbf{b}} f(\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b}^*) = \mathbf{0}$ and the Cauchy–Schwarz inequality, implies

$$\|\mathbf{b} - \mathbf{b}^*\| \leq \frac{1}{\lambda_b} \|\nabla_{\mathbf{b}} f(\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b})\| \quad (64)$$

We compute

$$\begin{cases} \nabla_{\mathbf{W}} f(\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b}) := \frac{1}{K} \boldsymbol{\Theta} (\mathbf{W}^T \boldsymbol{\Theta} + \mathbf{b} \mathbf{1}_K^T - \mathbf{I}_K)^T + \lambda_W \mathbf{W}, \\ \nabla_{\boldsymbol{\Theta}} f(\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b}) := \frac{1}{K} \mathbf{W} (\mathbf{W}^T \boldsymbol{\Theta} + \mathbf{b} \mathbf{1}_K^T - \mathbf{I}_K) + n \lambda_H \mathbf{H}. \end{cases} \quad (65)$$

It follows from (52) and (53) that $h(\mathbf{W}, \boldsymbol{\Theta}) = f(\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b}^*)$, which implies $\nabla_{\mathbf{W}} h(\mathbf{W}, \boldsymbol{\Theta}) = \nabla_{\mathbf{W}} f(\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b}^*)$ and $\nabla_{\boldsymbol{\Theta}} h(\mathbf{W}, \boldsymbol{\Theta}) = \nabla_{\boldsymbol{\Theta}} f(\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b}^*)$. Therefore, we have

$$\begin{aligned} \|\nabla_{\mathbf{W}} f(\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b}) - \nabla_{\mathbf{W}} h(\mathbf{W}, \boldsymbol{\Theta})\|_F &= \|\nabla_{\mathbf{W}} f(\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b}) - \nabla_{\mathbf{W}} f(\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b}^*)\|_F \\ &= \frac{1}{K} \|\boldsymbol{\Theta} \mathbf{1}_K (\mathbf{b} - \mathbf{b}^*)\|_F \leq \frac{2\sqrt{c}}{\rho\sqrt{K}} \|\mathbf{b} - \mathbf{b}^*\|, \end{aligned}$$

where the inequality follows from (63). This implies

$$\begin{aligned} \|\nabla_{\mathbf{W}} h(\mathbf{W}, \boldsymbol{\Theta})\|_F &\leq \|\nabla_{\mathbf{W}} f(\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b})\|_F + \frac{2\sqrt{c}}{\rho\sqrt{K}} \|\mathbf{b} - \mathbf{b}^*\| \\ &\leq \|\nabla_{\mathbf{W}} f(\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b})\|_F + \frac{2\sqrt{c}}{\lambda_b \rho \sqrt{K}} \|\nabla_{\mathbf{b}} f(\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b})\|, \end{aligned} \quad (66)$$

where the last inequality uses (64). By the same argument, we obtain

$$\|\nabla_{\boldsymbol{\Theta}} h(\mathbf{W}, \boldsymbol{\Theta})\|_F \leq \|\nabla_{\boldsymbol{\Theta}} f(\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b})\|_F + \frac{2\rho\sqrt{c}}{\lambda_b \sqrt{K}} \|\nabla_{\mathbf{b}} f(\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b})\| \quad (67)$$

According to (51) and (60), we have for all $(\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b})$,

$$\begin{aligned} \text{dist}^2((\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b}), \mathcal{X}_f) &= \text{dist}^2((\mathbf{W}, \boldsymbol{\Theta}), \mathcal{X}_h) + \|\mathbf{b} - \mathbf{b}^*\|^2 \leq \kappa_1^2 \|\nabla h(\mathbf{W}, \boldsymbol{\Theta})\|_F^2 + \|\mathbf{b} - \mathbf{b}^*\|^2 \\ &\leq 2\kappa_1^2 \left(\|\nabla_{\mathbf{W}} f(\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b})\|_F^2 + \frac{4c}{\lambda_b^2 \rho^2 K} \|\nabla_{\mathbf{b}} f(\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b})\|^2 \right) + \\ &\quad 2\kappa_1^2 \left(\|\nabla_{\boldsymbol{\Theta}} f(\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b})\|_F^2 + \frac{4\rho^2 c}{\lambda_b^2 K} \|\nabla_{\mathbf{b}} f(\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b})\|^2 \right) + \frac{1}{\lambda_b^2} \|\nabla_{\mathbf{b}} f(\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b})\|^2 \\ &\leq 2\kappa_1^2 \max \left\{ 1, \frac{4c}{\lambda_b^2 K} \left(\frac{1}{\rho^2} + \rho^2 \right) + \frac{1}{2\kappa_1^2 \lambda_b^2} \right\} \|\nabla f(\mathbf{W}, \boldsymbol{\Theta}, \mathbf{b})\|_F^2, \end{aligned}$$

where the first inequality follows from (57) and (62), where the second inequality uses (66), (67), and (64).

The rest of the proof is devoted to proving our claim. For ease of exposition, let

$$\mathbf{W}_1 := \mathbf{W}\mathbf{P}, \quad \mathbf{W}_2 := \mathbf{W}\mathbf{P}^\perp, \quad \boldsymbol{\Theta}_1 := \boldsymbol{\Theta}\mathbf{P}, \quad \boldsymbol{\Theta}_2 := \boldsymbol{\Theta}\mathbf{P}^\perp, \quad (68)$$

and

$$\gamma := K\sqrt{n\lambda_W\lambda_H} - \beta, \quad \lambda_{\max} := \max\{\lambda_W, n\lambda_H\}, \quad \lambda_{\min} := \min\{\lambda_W, n\lambda_H\}. \quad (69)$$

Since $(\mathbf{I} - \alpha\mathbf{1}\mathbf{1}^T)^2 = \mathbf{P} + \beta\mathbf{P}^\perp$, we compute

$$\begin{cases} K\nabla_{\mathbf{W}} h(\mathbf{W}, \boldsymbol{\Theta}) = \boldsymbol{\Theta}(\mathbf{P} + \beta\mathbf{P}^\perp)(\mathbf{W}^T\boldsymbol{\Theta} - \mathbf{I})^T + K\lambda_W\mathbf{W}, \\ K\nabla_{\boldsymbol{\Theta}} h(\mathbf{W}, \boldsymbol{\Theta}) = \mathbf{W}(\mathbf{W}^T\boldsymbol{\Theta} - \mathbf{I})(\mathbf{P} + \beta\mathbf{P}^\perp) + nK\lambda_H\boldsymbol{\Theta}. \end{cases} \quad (70)$$

According to (60), we compute

$$\begin{aligned} \text{dist}^2((\mathbf{W}, \boldsymbol{\Theta}), \mathcal{X}_h) &= \min_{\mathbf{U} \in \mathcal{O}^{d \times (K-1)}} \left\{ \|\mathbf{W} - \rho\sqrt{c}\mathbf{U}\bar{\mathbf{V}}^T\|_F^2 + \left\| \boldsymbol{\Theta} - \frac{\sqrt{c}}{\rho}\mathbf{U}\bar{\mathbf{V}}^T \right\|_F^2 \right\} \\ &= \min_{\mathbf{U} \in \mathcal{O}^{d \times (K-1)}} \left\{ \|\mathbf{W}_1 - \rho\sqrt{c}\mathbf{U}\bar{\mathbf{V}}^T\|_F^2 + \left\| \boldsymbol{\Theta}_1 - \frac{\sqrt{c}}{\rho}\mathbf{U}\bar{\mathbf{V}}^T \right\|_F^2 \right\} + \|\mathbf{W}_2\|_F^2 + \|\boldsymbol{\Theta}_2\|_F^2 \end{aligned} \quad (71)$$

$$\leq \|\mathbf{W}_2\|_F^2 + \|\boldsymbol{\Theta}_2\|_F^2 + 2 \left\| \boldsymbol{\Theta}_1 - \frac{1}{\rho^2}\mathbf{W}_1 \right\|_F^2 + \left(1 + \frac{2\lambda_W}{n\lambda_H} \right) \min_{\mathbf{U} \in \mathcal{O}^{d \times (K-1)}} \|\mathbf{W}_1 - \rho\sqrt{c}\mathbf{U}\bar{\mathbf{V}}^T\|_F^2, \quad (72)$$

where the second equality is due to $\bar{\mathbf{V}}^T\mathbf{P}^\perp = \mathbf{0}$. Then, we bound each term above in turn. According to (70), we compute

$$\begin{aligned} K\nabla_{\mathbf{W}} h(\mathbf{W}, \boldsymbol{\Theta})\mathbf{P} &= \boldsymbol{\Theta}(\mathbf{P} + \beta\mathbf{P}^\perp)(\mathbf{W}^T\boldsymbol{\Theta} - \mathbf{I})^T\mathbf{P} + K\lambda_W\mathbf{W}\mathbf{P} \\ &= (\boldsymbol{\Theta}_1\boldsymbol{\Theta}_1^T + \beta\boldsymbol{\Theta}_2\boldsymbol{\Theta}_2^T + K\lambda_W\mathbf{I})\mathbf{W}_1 - \boldsymbol{\Theta}_1, \end{aligned} \quad (73)$$

and

$$\begin{aligned} K\nabla_{\mathbf{W}} h(\mathbf{W}, \boldsymbol{\Theta})\mathbf{P}^\perp &= \boldsymbol{\Theta}(\mathbf{P} + \beta\mathbf{P}^\perp)(\mathbf{W}^T\boldsymbol{\Theta} - \mathbf{I})^T\mathbf{P}^\perp + K\lambda_W\mathbf{W}\mathbf{P}^\perp \\ &= (\boldsymbol{\Theta}_1\boldsymbol{\Theta}_1^T + \beta\boldsymbol{\Theta}_2\boldsymbol{\Theta}_2^T + K\lambda_W\mathbf{I})\mathbf{W}_2 - \beta\boldsymbol{\Theta}_2, \end{aligned}$$

which implies

$$K\|\nabla_{\mathbf{W}} h(\mathbf{W}, \boldsymbol{\Theta})\mathbf{P}^\perp\|_F \geq K\lambda_W\|\mathbf{W}_2\|_F - \beta\|\boldsymbol{\Theta}_2\|_F. \quad (74)$$

Using the same computation, we obtain

$$\begin{cases} K\nabla_{\Theta}h(\mathbf{W}, \Theta)\mathbf{P} = (\mathbf{W}\mathbf{W}^T + nK\lambda_H\mathbf{I})\Theta_1 - \mathbf{W}_1, \\ K\nabla_{\Theta}h(\mathbf{W}, \Theta)\mathbf{P}^\perp = (\beta\mathbf{W}\mathbf{W}^T + nK\lambda_H\mathbf{I})\Theta_2 - \beta\mathbf{W}_2. \end{cases} \quad (75)$$

Using the same argument in (73), we show

$$K\|\nabla_{\Theta}h(\mathbf{W}, \Theta)\mathbf{P}^\perp\|_F \geq nK\lambda_H\|\Theta_2\|_F - \beta\|\mathbf{W}_2\|_F. \quad (76)$$

It follows from (74), (76), and $\gamma := K\sqrt{n\lambda_W\lambda_H} - \beta$ that

$$\begin{aligned} & K\sqrt{n\lambda_H}\|\nabla_{\mathbf{W}}h(\mathbf{W}, \Theta)\mathbf{P}^\perp\|_F + K\sqrt{\lambda_W}\|\nabla_{\Theta}h(\mathbf{W}, \Theta)\mathbf{P}^\perp\|_F \\ & \geq \gamma \left(\sqrt{\lambda_W}\|\mathbf{W}_2\|_F + \sqrt{n\lambda_H}\|\mathbf{H}_2\|_F \right). \end{aligned} \quad (77)$$

This implies

$$\begin{aligned} \|\mathbf{W}_2\|_F^2 + \|\mathbf{H}_2\|_F^2 & \leq \frac{2K^2\lambda_{\max}}{\gamma^2\lambda_{\min}} \left(\|\nabla_{\mathbf{W}}h(\mathbf{W}, \Theta)\mathbf{P}^\perp\|_F^2 + \|\nabla_{\Theta}h(\mathbf{W}, \Theta)\mathbf{P}^\perp\|_F^2 \right) \\ & \leq \frac{2K^2\lambda_{\max}}{\gamma^2\lambda_{\min}} \|\nabla h(\mathbf{W}, \Theta)\|_F^2. \end{aligned} \quad (78)$$

Let $(\mathbf{W}^*, \Theta^*) \in \mathcal{X}_h$ be such that $\text{dist}((\mathbf{W}, \mathbf{H}), \mathcal{X}_h) = \|(\mathbf{W}, \mathbf{H}) - (\mathbf{W}^*, \mathbf{H}^*)\|_F \leq \delta_1$. According to (71) and $\delta_1 \leq \min\{\rho, 1/\rho\}\sqrt{c}$, we have

$$\|\mathbf{W}_1\| \leq \|\mathbf{W}_1 - \mathbf{W}_1^*\| + \|\mathbf{W}_1^*\| \leq 2\rho\sqrt{c}, \quad \|\Theta_1\| \leq \|\Theta_1 - \Theta_1^*\| + \|\Theta_1^*\| \leq \frac{2\sqrt{c}}{\rho}, \quad (79)$$

$$\|\mathbf{W}_2\| \leq \delta_1, \quad \|\Theta_2\| \leq \delta_1. \quad (80)$$

According to (73), we compute

$$\begin{aligned} K\|\nabla_{\mathbf{W}}h(\mathbf{W}, \Theta)\mathbf{P}\|_F & \geq \|(\Theta_1\Theta_1^T + K\lambda_W\mathbf{I})\mathbf{W}_1 - \Theta_1\|_F - \|\Theta_2\Theta_2^T\mathbf{W}_1\|_F \\ & \geq \|(\Theta_1\Theta_1^T + K\lambda_W\mathbf{I})\mathbf{W}_1 - \Theta_1\|_F - c\|\Theta_2\|_F, \end{aligned}$$

where the first inequality follows from the triangular inequality and $\beta < 1$, where the second inequality uses $\|\Theta_2\Theta_2^T\mathbf{W}_1\|_F \leq \|\mathbf{W}_1\|\|\Theta_2\|\|\Theta_2\|_F \leq c$ due to (80) and (61). By the same argument, we compute

$$K\|\nabla_{\Theta}h(\mathbf{W}, \Theta)\mathbf{P}\|_F \geq \|(\mathbf{W}_1\mathbf{W}_1^T + nK\lambda_H\mathbf{I})\Theta_1 - \mathbf{W}_1\|_F - c\|\mathbf{W}_2\|_F.$$

Summing up $K\sqrt{n\lambda_H}\|\nabla_{\mathbf{W}}h(\mathbf{W}, \Theta)\mathbf{P}\|_F + K\sqrt{\lambda_W}\|\nabla_{\Theta}h(\mathbf{W}, \Theta)\mathbf{P}\|_F$ yields

$$\begin{aligned} & K \left(\sqrt{n\lambda_H}\|\nabla_{\mathbf{W}}h(\mathbf{W}, \Theta)\mathbf{P}\|_F + \sqrt{\lambda_W}\|\nabla_{\Theta}h(\mathbf{W}, \Theta)\mathbf{P}\|_F \right) + c \left(\sqrt{n\lambda_H}\|\Theta_2\|_F + \sqrt{\lambda_W}\|\mathbf{W}_2\|_F \right) \\ & \geq \sqrt{n\lambda_H} \|(\Theta_1\Theta_1^T + K\lambda_W\mathbf{I})\mathbf{W}_1 - \Theta_1\|_F + \sqrt{\lambda_W} \|(\mathbf{W}_1\mathbf{W}_1^T + nK\lambda_H\mathbf{I})\Theta_1 - \mathbf{W}_1\|_F. \end{aligned}$$

Substituting (77) into the above inequality yields

$$\begin{aligned} & 2K \max \left\{ 1, \frac{c}{\gamma} \right\} \sqrt{\lambda_{\max}} (\|\nabla_{\mathbf{W}}h(\mathbf{W}, \Theta)\|_F + \|\nabla_{\Theta}h(\mathbf{W}, \Theta)\|_F) \\ & \geq \sqrt{n\lambda_H} \|(\Theta_1\Theta_1^T + K\lambda_W\mathbf{I})\mathbf{W}_1 - \Theta_1\|_F + \sqrt{\lambda_W} \|(\mathbf{W}_1\mathbf{W}_1^T + nK\lambda_H\mathbf{I})\Theta_1 - \mathbf{W}_1\|_F. \end{aligned} \quad (81)$$

Using this and (79), we obtain

$$\begin{aligned}
& 4K \max \left\{ 1, \frac{c}{\gamma} \right\} \max \left\{ \rho, \frac{1}{\rho} \right\} \frac{\sqrt{\lambda_{\max} c}}{\sqrt{\lambda_{\min}}} (\|\nabla_{\mathbf{W}} h(\mathbf{W}, \boldsymbol{\Theta})\|_F + \|\nabla_{\boldsymbol{\Theta}} h(\mathbf{W}, \boldsymbol{\Theta})\|_F) \\
& \geq \|(\boldsymbol{\Theta}_1 \boldsymbol{\Theta}_1^T + K\lambda_W \mathbf{I}) \mathbf{W}_1 - \boldsymbol{\Theta}_1\|_F \|\mathbf{W}_1\| + \|(\mathbf{W}_1 \mathbf{W}_1^T + nK\lambda_H \mathbf{I}) \boldsymbol{\Theta}_1 - \mathbf{W}_1\|_F \|\boldsymbol{\Theta}_1\| \\
& \geq \|(\boldsymbol{\Theta}_1 \boldsymbol{\Theta}_1^T + K\lambda_W \mathbf{I}) \mathbf{W}_1 \mathbf{W}_1^T - \boldsymbol{\Theta}_1 \mathbf{W}_1^T\|_F + \|(\mathbf{W}_1 \mathbf{W}_1^T + nK\lambda_H \mathbf{I}) \boldsymbol{\Theta}_1 \boldsymbol{\Theta}_1^T - \mathbf{W}_1 \boldsymbol{\Theta}_1^T\|_F \\
& \geq K\|\lambda_W \mathbf{W}_1 \mathbf{W}_1^T - n\lambda_H \boldsymbol{\Theta}_1 \boldsymbol{\Theta}_1^T\|_F.
\end{aligned}$$

Therefore, we have

$$\|\lambda_W \mathbf{W}_1 \mathbf{W}_1^T - n\lambda_H \boldsymbol{\Theta}_1 \boldsymbol{\Theta}_1^T\|_F \leq \kappa_1 (\|\nabla_{\mathbf{W}} h(\mathbf{W}, \boldsymbol{\Theta})\|_F + \|\nabla_{\boldsymbol{\Theta}} h(\mathbf{W}, \boldsymbol{\Theta})\|_F), \quad (82)$$

where

$$\kappa_1 := 4 \max \left\{ 1, \frac{c}{\gamma} \right\} \max \left\{ \rho, \frac{1}{\rho} \right\} \frac{\sqrt{\lambda_{\max} c}}{\sqrt{\lambda_{\min}}}. \quad (83)$$

According to (81), we compute

$$\begin{aligned}
& 2K \max \left\{ 1, \frac{c}{\gamma} \right\} \sqrt{\lambda_{\max}} (\|\nabla_{\mathbf{W}} h(\mathbf{W}, \boldsymbol{\Theta})\|_F + \|\nabla_{\boldsymbol{\Theta}} h(\mathbf{W}, \boldsymbol{\Theta})\|_F) \\
& \geq \|\sqrt{n\lambda_H} (\boldsymbol{\Theta}_1 \boldsymbol{\Theta}_1^T + K\lambda_W \mathbf{I}) \mathbf{W}_1 - \sqrt{\lambda_W} (\mathbf{W}_1 \mathbf{W}_1^T + nK\lambda_H \mathbf{I}) \boldsymbol{\Theta}_1 + \sqrt{\lambda_W} \mathbf{W}_1 - \sqrt{n\lambda_H} \boldsymbol{\Theta}_1\|_F \\
& = \left\| \frac{\sqrt{\lambda_W}}{\sqrt{n\lambda_H}} (\mathbf{W}_1 \mathbf{W}_1^T + nK\lambda_H \mathbf{I}) (\sqrt{\lambda_W} \mathbf{W}_1 - \sqrt{n\lambda_H} \boldsymbol{\Theta}_1) + \sqrt{n\lambda_H} \left(\boldsymbol{\Theta}_1 \boldsymbol{\Theta}_1^T - \frac{\lambda_W}{n\lambda_H} \mathbf{W}_1 \mathbf{W}_1^T \right) \mathbf{W}_1 \right. \\
& \quad \left. + (\sqrt{\lambda_W} \mathbf{W}_1 - \sqrt{n\lambda_H} \boldsymbol{\Theta}_1) \right\|_F \\
& \geq \left\| \frac{\sqrt{\lambda_W}}{\sqrt{n\lambda_H}} \left(\mathbf{W}_1 \mathbf{W}_1^T + nK\lambda_H \mathbf{I} + \frac{\sqrt{n\lambda_H}}{\sqrt{\lambda_W}} \mathbf{I} \right) (\sqrt{\lambda_W} \mathbf{W}_1 - \sqrt{n\lambda_H} \boldsymbol{\Theta}_1) \right\|_F - \\
& \quad \sqrt{n\lambda_H} \|\mathbf{W}_1\| \left\| \boldsymbol{\Theta}_1 \boldsymbol{\Theta}_1^T - \frac{\lambda_W}{n\lambda_H} \mathbf{W}_1 \mathbf{W}_1^T \right\|_F \\
& \geq \left(1 + K\sqrt{n\lambda_W \lambda_H} \right) \left\| \sqrt{\lambda_W} \mathbf{W}_1 - \sqrt{n\lambda_H} \boldsymbol{\Theta}_1 \right\|_F - \frac{2\rho\sqrt{c}}{\sqrt{n\lambda_H}} \|n\lambda_H \boldsymbol{\Theta}_1 \boldsymbol{\Theta}_1^T - \lambda_W \mathbf{W}_1 \mathbf{W}_1^T\|_F,
\end{aligned}$$

where the last inequality follows from (79). This, together with (82), yields

$$\left\| \boldsymbol{\Theta}_1 - \frac{\sqrt{\lambda_W}}{\sqrt{n\lambda_H}} \mathbf{W}_1 \right\|_F = \frac{\left\| \sqrt{\lambda_W} \mathbf{W}_1 - \sqrt{n\lambda_H} \boldsymbol{\Theta}_1 \right\|_F}{\sqrt{n\lambda_H}} \leq \kappa_2 (\|\nabla_{\mathbf{W}} h(\mathbf{W}, \boldsymbol{\Theta})\|_F + \|\nabla_{\boldsymbol{\Theta}} h(\mathbf{W}, \boldsymbol{\Theta})\|_F), \quad (84)$$

where

$$\kappa_2 := \frac{1}{\sqrt{n\lambda_H} (1 + K\sqrt{n\lambda_W \lambda_H})} \left(\frac{2\kappa_1 \rho \sqrt{c}}{\sqrt{n\lambda_H}} + 2K \max \left\{ 1, \frac{c}{\gamma} \right\} \sqrt{\lambda_{\max}} \right).$$

It follows from (61) and (71) that $\|\mathbf{W}_1 - \mathbf{W}_1^*\|_F \leq \delta_1$. This, together with Weyl's inequality, implies

$$\sigma_{K-1}(\mathbf{W}_1) \geq \sigma_{K-1}(\mathbf{W}_1^*) - \|\mathbf{W}_1 - \mathbf{W}_1^*\|_F \geq \rho\sqrt{c} - \frac{1}{2}\rho\sqrt{c} \geq \frac{1}{2}\rho\sqrt{c}, \quad (85)$$

where the second inequality follows from (61). Using (81) again, we obtain

$$\begin{aligned}
& 2K \max \left\{ 1, \frac{c}{\gamma} \right\} \frac{\sqrt{\lambda_{\max}}}{\sqrt{\lambda_W}} (\|\nabla_{\mathbf{W}} h(\mathbf{W}, \boldsymbol{\Theta})\|_F + \|\nabla_{\boldsymbol{\Theta}} h(\mathbf{W}, \boldsymbol{\Theta})\|_F) \\
& \geq \left\| (\mathbf{W}_1 \mathbf{W}_1^T + nK\lambda_H \mathbf{I}) \left(\boldsymbol{\Theta}_1 - \frac{\sqrt{\lambda_W}}{\sqrt{n\lambda_H}} \mathbf{W}_1 + \frac{\sqrt{\lambda_W}}{\sqrt{n\lambda_H}} \mathbf{W}_1 \right) - \mathbf{W}_1 \right\|_F \\
& \geq \frac{\sqrt{\lambda_W}}{\sqrt{n\lambda_H}} \left\| \mathbf{W}_1 \left(\mathbf{W}_1^T \mathbf{W}_1 + nK\lambda_H \mathbf{I} - \frac{\sqrt{n\lambda_H}}{\sqrt{\lambda_W}} \mathbf{I} \right) \right\|_F - \left\| (\mathbf{W}_1 \mathbf{W}_1^T + nK\lambda_H \mathbf{I}) \left(\boldsymbol{\Theta}_1 - \frac{\sqrt{\lambda_W}}{\sqrt{n\lambda_H}} \mathbf{W}_1 \right) \right\|_F \\
& \geq \frac{\sqrt{c}}{2\rho} \left\| \mathbf{W}_1^T \mathbf{W}_1 - \frac{c\sqrt{n\lambda_H}}{\sqrt{\lambda_W}} \mathbf{P} \right\|_F - (\|\mathbf{W}_1\|^2 + nK\lambda_H) \left\| \boldsymbol{\Theta}_1 - \frac{\sqrt{\lambda_W}}{\sqrt{n\lambda_H}} \mathbf{W}_1 \right\|_F
\end{aligned}$$

where the last inequality follows from $c := 1 - K\sqrt{n\lambda_W\lambda_H}$, and $\left\| \mathbf{W}_1 (\mathbf{W}_1^T \mathbf{W}_1 - c\sqrt{n\lambda_H} \mathbf{I} / \sqrt{\lambda_W}) \right\|_F = \left\| \mathbf{W}_1 (\mathbf{W}_1^T \mathbf{W}_1 - c\sqrt{n\lambda_H} \mathbf{P} / \sqrt{\lambda_W}) \right\|_F \geq \sigma_{K-1}(\mathbf{W}_1) \left\| \mathbf{W}_1^T \mathbf{W}_1 - c\sqrt{n\lambda_H} \mathbf{P} / \sqrt{\lambda_W} \right\|_F$ due to $\mathbf{W} = \mathbf{W}_1 \mathbf{P}$ and (85). This, together with (79) and (84), yields

$$\left\| \mathbf{W}_1^T \mathbf{W}_1 - \frac{c\sqrt{n\lambda_H}}{\sqrt{\lambda_W}} \mathbf{P} \right\|_F \leq \kappa_3 (\|\nabla_{\mathbf{W}} h(\mathbf{W}, \boldsymbol{\Theta})\|_F + \|\nabla_{\boldsymbol{\Theta}} h(\mathbf{W}, \boldsymbol{\Theta})\|_F), \quad (86)$$

where

$$\kappa_3 := \frac{2\rho}{\sqrt{c}} \left(2K \max \left\{ 1, \frac{c}{\gamma} \right\} \frac{\sqrt{\lambda_{\max}}}{\sqrt{\lambda_W}} + (4\rho^2 c + nK\lambda_H) \kappa_2 \right).$$

Let $\mathbf{W}_1 = \mathbf{U}_1 \boldsymbol{\Sigma}_1 \mathbf{V}_1^T$ be the thin singular value decomposition of \mathbf{W}_1 , where $\mathbf{U}_1 = [\bar{\mathbf{U}}_1 \mathbf{u}_1] \in \mathcal{O}^{d \times K}$ with $\bar{\mathbf{U}}_1 \in \mathbb{R}^{d \times (K-1)}$ and $\mathbf{u}_1 \in \mathbb{R}^d$, $\mathbf{V}_1 = [\bar{\mathbf{V}}_1 \mathbf{v}_1] \in \mathcal{O}^K$ with $\bar{\mathbf{V}}_1 \in \mathbb{R}^{K \times (K-1)}$ and $\mathbf{v}_1 \in \mathbb{R}^K$, and $\boldsymbol{\Sigma}_1 = \text{diag}(\sigma_1, \dots, \sigma_K)$ is a diagonal matrix with $\sigma_1 \geq \dots \geq \sigma_K \geq 0$. It follows from (85) that $\mathbf{W}_1 = \mathbf{W} \mathbf{P}$ is of rank $K-1$, which implies $\sigma_K = 0$. Therefore, we have $\mathbf{W}_1 = \bar{\mathbf{U}}_1 \text{diag}(\sigma_1, \dots, \sigma_{K-1}) \bar{\mathbf{V}}_1^T$. This, together with $\mathbf{W}_1 \mathbf{1}_K = \mathbf{W} \mathbf{P} \mathbf{1}_K = \mathbf{0}$, yields $\bar{\mathbf{V}}_1^T \mathbf{1}_K = \mathbf{0}$. This directly implies $\mathbf{v}_1 = \mathbf{1}_K / \sqrt{K}$. Using this and $\mathbf{V}_1 \mathbf{V}_1^T = \mathbf{I}$, we obtain

$$\bar{\mathbf{V}}_1 \bar{\mathbf{V}}_1^T = \mathbf{I} - \mathbf{v}_1 \mathbf{v}_1^T = \mathbf{P}. \quad (87)$$

Let $\boldsymbol{\Sigma}_2 = \begin{bmatrix} \mathbf{I}_{K-1} & \mathbf{0} \\ \mathbf{0} & 0 \end{bmatrix}$. Noting that $\rho^2 = \sqrt{n\lambda_H} / \sqrt{\lambda_W}$, we compute

$$\begin{aligned}
\left\| \mathbf{W}_1^T \mathbf{W}_1 - \rho^2 c \mathbf{P} \right\|_F &= \left\| \mathbf{V}_1 (\boldsymbol{\Sigma}_1^2 - \rho^2 c \boldsymbol{\Sigma}_2) \mathbf{V}_1^T \right\|_F = \left\| \boldsymbol{\Sigma}_1^2 - \rho^2 c \boldsymbol{\Sigma}_2 \right\|_F \\
&= \left\| (\boldsymbol{\Sigma}_1 - \rho\sqrt{c} \boldsymbol{\Sigma}_2) (\boldsymbol{\Sigma}_1 + \rho\sqrt{c} \boldsymbol{\Sigma}_2) \right\|_F \geq \rho\sqrt{c} \left\| \boldsymbol{\Sigma}_1 - \rho\sqrt{c} \boldsymbol{\Sigma}_2 \right\|_F \\
&= \rho\sqrt{c} \left\| \mathbf{U}_1 (\boldsymbol{\Sigma}_1 - \rho\sqrt{c} \boldsymbol{\Sigma}_2) \mathbf{V}_1^T \right\|_F \\
&= \rho\sqrt{c} \left\| \mathbf{W}_1 - \rho\sqrt{c} \bar{\mathbf{U}}_1 \bar{\mathbf{V}}_1^T \right\|_F \geq \rho\sqrt{c} \min_{\mathbf{U} \in \mathcal{O}^{d \times (K-1)}} \left\| \mathbf{W}_1 - \rho\sqrt{c} \mathbf{U} \bar{\mathbf{V}}^T \right\|_F,
\end{aligned}$$

where the last inequality uses the fact that there exists a $\mathbf{Q} \in \mathcal{O}^{K-1}$ such that $\bar{\mathbf{V}}_1 = \bar{\mathbf{V}} \mathbf{Q}$. Substituting this back into (86) yields

$$\min_{\mathbf{U} \in \mathcal{O}^{d \times (K-1)}} \left\| \mathbf{W}_1 - \rho\sqrt{c} \mathbf{U} \bar{\mathbf{V}}^T \right\|_F \leq \frac{\kappa_3}{\rho\sqrt{c}} (\|\nabla_{\mathbf{W}} h(\mathbf{W}, \boldsymbol{\Theta})\|_F + \|\nabla_{\boldsymbol{\Theta}} h(\mathbf{W}, \boldsymbol{\Theta})\|_F). \quad (88)$$

This, together with (72), (78), and (84), yields

$$\text{dist}^2((\mathbf{W}, \boldsymbol{\Theta}), \mathcal{X}_h) \leq \left(\frac{2K^2 \lambda_{\max}}{\gamma^2 \lambda_{\min}} + 2\kappa_2^2 + \left(1 + \frac{2\lambda_W}{n\lambda_H} \right) \frac{\kappa_3^2}{\rho^2 c} \right) \|\nabla h(\mathbf{W}, \boldsymbol{\Theta})\|_F^2.$$

Then, we prove the claim. \square

4 Cross-Entropy Loss

Suppose that $\mathcal{L}(\cdot, \cdot)$ is the cross-entropy (CE) loss:

$$\mathcal{L}(\mathbf{z}, \mathbf{y}_k) = -\log \left(\frac{\exp(z_k)}{\sum_{\ell=1}^K \exp(z_\ell)} \right). \quad (89)$$

We consider the following problems:

$$\min_{\mathbf{W} \in \mathbb{R}^{d \times K}, \mathbf{H} \in \mathbb{R}^{d \times N}} F(\mathbf{W}, \mathbf{H}) = \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}(\mathbf{W}^T \mathbf{h}_{k,i}, \mathbf{y}_k) + \frac{\lambda_W}{2} \|\mathbf{W}\|_F^2 + \frac{\lambda_H}{2} \|\mathbf{H}\|_F^2, \quad (90)$$

$$\min_{\mathbf{W} \in \mathbb{R}^{d \times K}, \mathbf{\Theta} \in \mathbb{R}^{d \times K}} f(\mathbf{W}, \mathbf{\Theta}) = \frac{1}{K} \sum_{k=1}^K \mathcal{L}(\mathbf{W}^T \mathbf{\theta}_k, \mathbf{y}_k) + \frac{\lambda_W}{2} \|\mathbf{W}\|_F^2 + \frac{n\lambda_H}{2} \|\mathbf{\Theta}\|_F^2, \quad (91)$$

where λ_W, λ_H are the penalties for \mathbf{W} and \mathbf{H} , respectively. One can easily verify

$$F(\mathbf{W}, \mathbf{H}) = \frac{1}{n} \sum_{i=1}^n f(\mathbf{W}, \mathbf{H}_i). \quad (92)$$

Lemma 3. *The optimal solution set of Problem (91) takes the form of*

$$\mathcal{X}_f = \left\{ (\mathbf{W}, \mathbf{\Theta}) : \mathbf{W} = \frac{\sqrt[4]{n\lambda_H}}{\sqrt[4]{\lambda_W}} (\max\{c, 0\})^{\frac{1}{2}} \mathbf{U}, \mathbf{\Theta} = \frac{\sqrt{\lambda_W}}{\sqrt{n\lambda_H}} \mathbf{W}, \mathbf{U} \in \mathcal{O}^{d \times K} \right\}, \quad (93)$$

where $c := 1 - K\sqrt{n\lambda_W\lambda_H}$.

Proof.

□

References

- [1] Z.-Q. Luo and P. Tseng. Error bounds and convergence analysis of feasible descent methods: a general approach. *Annals of Operations Research*, 46(1):157–178, 1993.
- [2] P. Tseng and S. Yun. A coordinate gradient descent method for nonsmooth separable minimization. *Mathematical Programming*, 117(1):387–423, 2009.
- [3] Z. Zhou and A. M.-C. So. A unified approach to error bounds for structured convex optimization problems. *Mathematical Programming*, 165(2):689–728, 2017.