

# PFL

January 11, 2025

## 1 Introduction

xxx (Collins et al., 2021)

---

### Algorithm 1

**Input:** Participation rate  $r$ , step size  $\eta$ , number of local updates for the head  $\tau_w$ , for the shortcut  $\tau_s$  and for the representation  $\tau_b$ , number of communication rounds  $T$ .

```
1: Initialize  $\mathbf{B}^0, \mathbf{w}_1^0, \dots, \mathbf{w}_n^0, \mathbf{s}_1^0, \dots, \mathbf{s}_n^0$ 
2: for  $t = 0, 1, 2, \dots, T - 1$  do
3:   Server receives a batch of clients  $\mathcal{I}^t$  of size  $rn$ 
4:   Server sends current representation  $\phi^t$  to clients in  $\mathcal{I}^t$ 
5:   for each client  $i$  in  $\mathcal{I}^t$  do
6:     Client  $i$  initializes  $\mathbf{w}_i^{t,0} \leftarrow \mathbf{w}_i^{t-1, \tau_h}$ 
7:     Client updates its head for  $\tau_h$  steps:
8:     for  $\tau = 1$  to  $\tau_w$  do
9:        $\mathbf{w}_i^{t,\tau} \leftarrow \text{GRD} \left( f_i \left( \mathbf{w}_i^{t,\tau-1}, \mathbf{B}^{t-1}, \mathbf{s}_i^{t-1, \tau_s} \right), \mathbf{w}_i^{t,\tau-1}, \eta \right)$ 
10:    end for
11:    Client  $i$  initializes  $\mathbf{s}_i^{t,0} \leftarrow \mathbf{s}_i^{t-1, \tau_s}$ 
12:    Client  $i$  updates its shortcut for  $\tau_s$  steps:
13:    for  $\tau = 1$  to  $\tau_s$  do
14:       $\mathbf{s}_i^{t,\tau} \leftarrow \text{GRD} \left( f_i \left( \mathbf{w}_i^{t-1}, \mathbf{B}^{t-1}, \mathbf{s}_i^{t,\tau-1} \right), \mathbf{s}_i^{t,\tau-1}, \eta \right)$ 
15:    end for
16:    Client  $i$  initializes  $\mathbf{B}_i^{t,0} \leftarrow \mathbf{B}^{t-1}$ 
17:    Client  $i$  updates its representation for  $\tau_b$  steps:
18:    for  $\tau = 1$  to  $\tau_b$  do
19:       $\mathbf{B}_i^{t,\tau} \leftarrow \text{GRD} \left( f_i \left( \mathbf{w}_i^{t, \tau_w}, \mathbf{B}_i^{t, \tau-1}, \mathbf{s}_i^{t, \tau_s} \right), \mathbf{B}_i^{t, \tau-1}, \eta \right)$ 
20:    end for
21:    Client  $i$  sends updated representation  $\mathbf{B}_i^{t, \tau_b}$  to server
22:  end for
23:  for each client  $j$  not in  $\mathcal{I}^t$  do
24:    Set  $\mathbf{w}_i^{t, \tau_w} \leftarrow \mathbf{w}_i^{t-1, \tau_w}$  and  $\mathbf{s}_i^{t, \tau_s} \leftarrow \mathbf{s}_i^{t-1, \tau_s}$ 
25:  end for
26:  Server computes new representation:  $\mathbf{B}^t = \frac{1}{rn} \sum_{i \in \mathcal{I}^t} \mathbf{B}_i^{t, \tau_b}$ 
27: end for
```

---

## 1.1 Preliminaries

First, we establish the notations that will be used throughout our proof. Let  $\mathbf{S} := [\mathbf{s}_1, \dots, \mathbf{s}_{rn}] \in \mathbb{R}^{d \times rn}$  represent the personalized layers, and let  $\mathbf{W} := [\mathbf{w}_1, \dots, \mathbf{w}_{rn}] \in \mathbb{R}^{k \times rn}$  denote the personalized heads, which follow the global representation  $\mathbf{B}$ .

...  
The global objective can be rewritten as

$$\min_{\mathbf{B} \in \mathbb{R}^{d \times k}, \mathbf{W} \in \mathbb{R}^{k \times rn}, \hat{\mathbf{S}} \in \mathbb{R}^{d \times rn}} \left\{ F(\hat{\mathbf{B}}, \mathbf{W}, \hat{\mathbf{S}}) := \frac{1}{2rnm} \mathbb{E}_{\mathcal{A}, \mathcal{I}} \left\| \mathbf{Y} - \mathcal{A}((1 - \alpha) \mathbf{W}_{\mathcal{I}}^{\top} \hat{\mathbf{B}}^{\top} + \alpha \mathbf{S}_{\mathcal{I}}^{\top}) \right\|_2^2 \right\}, \quad (1)$$

where  $\mathbf{Y} = \mathcal{A}((1 - \alpha) \mathbf{W}_{\mathcal{I}}^* \hat{\mathbf{B}}^{*\top} + \alpha \hat{\mathbf{S}}_{\mathcal{I}}^{*\top}) \in \mathbb{R}^{rnm}$ . Then we give the update rules of our algorithm:

$$\bar{\mathbf{W}}^{t+1} = \arg \min_{\mathbf{W} \in \mathbb{R}^{k \times rn}} \frac{1}{2rnm} \left\| \mathcal{A}^t \left( (1 - \alpha) \left( \mathbf{W}^{*\top} \hat{\mathbf{B}}^{*\top} - \mathbf{W}^{\top} \hat{\mathbf{B}}^{t\top} \right) + \alpha \left( \mathbf{S}^{*\top} - \bar{\mathbf{S}}^{t\top} \right) \right) \right\|_2^2, \quad (2)$$

$$\mathbf{W}^{t+1} = (1 - \lambda) \mathbf{W}^t + \lambda \bar{\mathbf{W}}^{t+1}, \quad (3)$$

$$\bar{\mathbf{B}}^{t+1} = \hat{\mathbf{B}}^t - \frac{\eta}{rnm} \left( (\mathcal{A}^t)^{\dagger} \mathcal{A}^t \left( (1 - \alpha) \left( \mathbf{W}^{t+1\top} \hat{\mathbf{B}}^{t\top} - \mathbf{W}^{*\top} \hat{\mathbf{B}}^{*\top} \right) + \alpha \left( \mathbf{S}^{t\top} - \mathbf{S}^{*\top} \right) \right) \right)^{\top} \mathbf{W}_{\mathcal{I}^t}^{t+1\top}, \quad (4)$$

$$\hat{\mathbf{B}}^{t+1}, \mathbf{R}^{t+1} = \text{QR}(\bar{\mathbf{B}}^{t+1}), \quad (5)$$

$$\tilde{\mathbf{S}}^{t+1} = \arg \min_{\mathbf{S} \in \mathbb{R}^{d \times rn}} \frac{1}{2rnm} \left\| \mathcal{A}^t \left( (1 - \alpha) \left( \mathbf{W}^{*\top} \hat{\mathbf{B}}^{*\top} - \mathbf{W}^{t+1\top} \hat{\mathbf{B}}^{t+1\top} \right) + \alpha \left( \mathbf{S}^{*\top} - \mathbf{S}^{\top} \right) \right) \right\|_2^2 + \frac{\beta}{2} \|\mathbf{S}\|_{\text{F}}^2, \quad (6)$$

$$\bar{\mathbf{S}}^{t+1} = \hat{\mathbf{B}}_{\perp}^{t+1} \hat{\mathbf{B}}_{\perp}^{t+1\top} \left( \tilde{\mathbf{S}}^{t+1} \right), \quad (7)$$

$$\mathbf{S}^{t+1} = (1 - \lambda) \mathbf{S}^t + \lambda \bar{\mathbf{S}}^{t+1}. \quad (8)$$

## 1.2 Auxiliary Lemmas

We first consider the update for  $\mathbf{W}$ . According to the update rule of (2),  $\mathbf{W}^{t+1}$  minimizes the function of  $\tilde{F}(\hat{\mathbf{B}}^t, \mathbf{W}, \bar{\mathbf{S}}^t) := \frac{1}{2rnm} \left\| \mathcal{A} \left( (1 - \alpha) \left( \mathbf{W}^{*\top} \hat{\mathbf{B}}^{*\top} - \mathbf{W}^{\top} \hat{\mathbf{B}}^{t\top} \right) + \alpha \left( \mathbf{S}^{*\top} - \bar{\mathbf{S}}^{t\top} \right) \right) \right\|_2^2$ .

Let  $\mathcal{W}_p^{t+1}$  be the  $p$ -th column of  $\mathbf{W}^{t+1\top}$ ,  $\mathcal{W}_p^*$  denote the  $p$ -th column of  $\mathbf{W}^{*\top}$ ,  $\mathcal{S}_l^t$  denote the  $l$ -th column of  $\bar{\mathbf{S}}^{t\top}$ ,  $\mathcal{S}_l^*$  denote the  $l$ -th column of  $\mathbf{S}^{*\top}$  and  $\hat{\mathbf{b}}_p^t$  be the  $p$ -th column of  $\hat{\mathbf{B}}^t$ , then for any  $p \in [k]$ ,  $l \in [d]$ , we have

$$\begin{aligned} \mathbf{0} &= \nabla_{\mathcal{W}_p} \tilde{F}(\hat{\mathbf{B}}^t, \mathbf{W}^{t+1}, \bar{\mathbf{S}}^t) \\ &= \frac{1 - \alpha}{rnm} \sum_{i=1}^{rn} \sum_{j=1}^m \left( \langle \mathbf{A}_{i,j}, (1 - \alpha) \left( \mathbf{W}^{t+1\top} \hat{\mathbf{B}}^{t\top} - \mathbf{W}^{*\top} \hat{\mathbf{B}}^{*\top} \right) + \alpha \left( \bar{\mathbf{S}}^{t\top} - \mathbf{S}^{*\top} \right) \rangle \right) \mathbf{A}_{i,j} \hat{\mathbf{b}}_p^t \\ &= \frac{1 - \alpha}{rnm} \sum_{i=1}^{rn} \sum_{j=1}^m \left( (1 - \alpha) \langle \mathbf{A}_{i,j}, \mathbf{W}^{t+1\top} \hat{\mathbf{B}}^{t\top} - \mathbf{W}^{*\top} \hat{\mathbf{B}}^{*\top} \rangle + \alpha \langle \mathbf{A}_{i,j}, \bar{\mathbf{S}}^{t\top} - \mathbf{S}^{*\top} \rangle \right) \mathbf{A}_{i,j} \hat{\mathbf{b}}_p^t \end{aligned}$$

$$= \frac{1-\alpha}{rnm} \sum_{i=1}^{rn} \sum_{j=1}^m \left( (1-\alpha) \left( \sum_{q=1}^k \hat{\mathbf{b}}_q^{t\top} \mathbf{A}_{i,j}^\top \mathcal{W}_q^{t+1} - \sum_{q=1}^k \hat{\mathbf{b}}_q^{*\top} \mathbf{A}_{i,j}^\top \mathcal{W}_q^* \right) + \alpha \left( \sum_{l=1}^d \mathbf{e}_l^\top \mathbf{A}_{i,j}^\top \mathcal{S}_l^t - \sum_{l=1}^d \mathbf{e}_l^\top \mathbf{A}_{i,j}^\top \mathcal{S}_l^* \right) \right) \mathbf{A}_{i,j} \hat{\mathbf{b}}_p^t, \quad (9)$$

which means

$$\begin{aligned} & \frac{1}{m} \sum_{q=1}^k \left( \sum_{i=1}^{rn} \sum_{j=1}^m \mathbf{A}_{i,j} \hat{\mathbf{b}}_p^t \hat{\mathbf{b}}_q^{t\top} \mathbf{A}_{i,j}^\top \right) (1-\alpha) \mathcal{W}_q^{t+1} \\ &= \frac{1}{m} \sum_{q=1}^k \left( \sum_{i=1}^{rn} \sum_{j=1}^m \mathbf{A}_{i,j} \hat{\mathbf{b}}_p^t \hat{\mathbf{b}}_q^{*\top} \mathbf{A}_{i,j}^\top \right) (1-\alpha) \mathcal{W}_q^* + \frac{1}{m} \sum_{l=1}^d \left( \sum_{i=1}^{rn} \sum_{j=1}^m \mathbf{A}_{i,j} \hat{\mathbf{b}}_p^t \mathbf{e}_l^\top \mathbf{A}_{i,j}^\top \right) \alpha (\mathcal{S}_l^* - \mathcal{S}_l^t). \end{aligned} \quad (10)$$

Then, define  $\mathbf{G}_{pq} := \frac{1}{m} \sum_{i=1}^{rn} \sum_{j=1}^m \mathbf{A}_{i,j} \hat{\mathbf{b}}_p^t \hat{\mathbf{b}}_q^{t\top} \mathbf{A}_{i,j}^\top$ ,  $\mathbf{C}_{pq} := \frac{1}{m} \sum_{i=1}^{rn} \sum_{j=1}^m \mathbf{A}_{i,j} \hat{\mathbf{b}}_p^t \hat{\mathbf{b}}_q^{*\top} \mathbf{A}_{i,j}^\top$  and  $\mathbf{D}_{pq} := \frac{1}{m} \sum_{i=1}^{rn} \sum_{j=1}^m \langle \hat{\mathbf{b}}_p^t, \hat{\mathbf{b}}_q^* \rangle \mathbf{I}_{rn}$ , for all  $p, q \in [k]$ , and define  $\mathbf{E}_{pl} := \frac{1}{m} \sum_{i=1}^{rn} \sum_{j=1}^m \mathbf{A}_{i,j} \hat{\mathbf{b}}_p^t \mathbf{e}_l^\top \mathbf{A}_{i,j}^\top$ , for all  $p \in [k], l \in [d]$ . Further, we define block matrices  $\mathbf{G}, \mathbf{C}, \mathbf{D} \in \mathbb{R}^{rnk \times rnk}$  and  $\mathbf{E} \in \mathbb{R}^{rnk \times rnd}$ , which are formed by  $\mathbf{G}_{pq}, \mathbf{C}_{pq}, \mathbf{D}_{pq}$  and  $\mathbf{E}_{pl}$ , respectively. In detail, take  $\mathbf{G}$  and  $\mathbf{E}$  for example,

$$\mathbf{G} := \begin{bmatrix} \mathbf{G}_{11} & \cdots & \mathbf{G}_{1k} \\ \vdots & \ddots & \vdots \\ \mathbf{G}_{k1} & \cdots & \mathbf{G}_{kk} \end{bmatrix}, \mathbf{E} := \begin{bmatrix} \mathbf{E}_{11} & \cdots & \mathbf{E}_{1d} \\ \vdots & \ddots & \vdots \\ \mathbf{E}_{k1} & \cdots & \mathbf{E}_{kd} \end{bmatrix}. \quad (11)$$

Then we define  $\widetilde{\mathcal{W}}^{t+1} := \text{vec}(\mathbf{W}^{t+1\top}) \in \mathbb{R}^{rnk}$ ,  $\widetilde{\mathcal{W}}^* := \text{vec}(\mathbf{W}^{*\top}) \in \mathbb{R}^{rnk}$ ,  $\widetilde{\mathcal{S}}^t := \text{vec}(\bar{\mathbf{S}}^{t\top}) \in \mathbb{R}^{rnd}$  and  $\widetilde{\mathcal{S}}^* := \text{vec}(\mathbf{S}^{*\top}) \in \mathbb{R}^{rnd}$ . From (10) we reach,

$$\begin{aligned} (1-\alpha)\widetilde{\mathcal{W}}^{t+1} &= (1-\alpha)\mathbf{G}^{-1}\mathbf{C}\widetilde{\mathcal{W}}^* + \alpha\mathbf{G}^{-1}\mathbf{E}(\widetilde{\mathcal{S}}^* - \widetilde{\mathcal{S}}^t) \\ &= (1-\alpha)\mathbf{D}\widetilde{\mathcal{W}}^* - (1-\alpha)\mathbf{G}^{-1}(\mathbf{GD} - \mathbf{C})\widetilde{\mathcal{W}}^* + \alpha\mathbf{G}^{-1}\mathbf{E}(\widetilde{\mathcal{S}}^* - \widetilde{\mathcal{S}}^t), \end{aligned} \quad (12)$$

where  $\mathbf{G}$  is invertible will be proved in the following lemma. Here, we consider  $\mathbf{G}_{pq}$ ,

$$\begin{aligned} \mathbf{G}_{pq} &= \frac{1}{m} \sum_{i=1}^{rn} \sum_{j=1}^m \mathbf{A}_{i,j} \hat{\mathbf{b}}_p \hat{\mathbf{b}}_q^\top \mathbf{A}_{i,j}^\top \\ &= \frac{1}{m} \sum_{i=1}^{rn} \sum_{j=1}^m \mathbf{e}_i (\mathbf{x}_i^j)^\top \hat{\mathbf{b}}_p \hat{\mathbf{b}}_q^\top \mathbf{x}_i^j \mathbf{e}_i^\top, \end{aligned} \quad (13)$$

meaning that  $\mathbf{G}_{pq}$  is diagonal with diagonal entries

$$(\mathbf{G}_{pq})_{ii} = \frac{1}{m} \sum_{j=1}^m (\mathbf{x}_i^j)^\top \hat{\mathbf{b}}_p \hat{\mathbf{b}}_q^\top \mathbf{x}_i^j = \hat{\mathbf{b}}_p^\top \left( \frac{1}{m} \sum_{j=1}^m \mathbf{x}_i^j (\mathbf{x}_i^j)^\top \right) \hat{\mathbf{b}}_q. \quad (14)$$

Define  $\mathbf{\Pi}^i := \frac{1}{m} \sum_{j=1}^m \mathbf{x}_i^j (\mathbf{x}_i^j)^\top$  for all  $i \in [rn]$ , then  $\mathbf{C}_{pq}$  is diagonal with entries  $(\mathbf{C}_{pq})_{ii} = \hat{\mathbf{b}}_p^\top \mathbf{\Pi}^i \hat{\mathbf{b}}_q^*$ , and  $\mathbf{E}_{pl}$  is diagonal with entries  $(\mathbf{E}_{pl})_{ii} = \hat{\mathbf{b}}_p^\top \mathbf{\Pi}^i \mathbf{e}_l$ . Note that  $\mathbf{D}_{pq} = \langle \hat{\mathbf{b}}_p, \hat{\mathbf{b}}_q^* \rangle \mathbf{I}_{rn}$  is also diagonal, then we define

$$\mathbf{G}^i := \left[ \hat{\mathbf{b}}_p^\top \mathbf{\Pi}^i \hat{\mathbf{b}}_q \right]_{1 \leq p, q \leq k+d} = \hat{\mathbf{B}}^\top \mathbf{\Pi}^i \hat{\mathbf{B}}, \quad \mathbf{C}^i := \left[ \hat{\mathbf{b}}_p^\top \mathbf{\Pi}^i \hat{\mathbf{b}}_q^* \right]_{1 \leq p, q \leq k+d} = \hat{\mathbf{B}}^\top \mathbf{\Pi}^i \hat{\mathbf{B}}^*, \quad (15)$$

$$\mathbf{D}^i := \left[ \langle \hat{\mathbf{b}}_p, \hat{\mathbf{b}}_q^* \rangle \right]_{1 \leq p, q \leq k+d} = \hat{\mathbf{B}}^\top \hat{\mathbf{B}}^*, \quad \mathbf{E}^i := \left[ \hat{\mathbf{b}}_p^\top \mathbf{\Pi}^i \mathbf{e}_l \right]_{1 \leq p \leq k, 1 \leq l \leq d} = \hat{\mathbf{B}}^\top \mathbf{\Pi}^i, \quad (16)$$

where  $\mathbf{G}^i$ ,  $\mathbf{C}^i$  and  $\mathbf{D}^i$  are the  $k \times k$  matrices that formed by taking the  $i$ -th diagonal entry of each block  $\mathbf{G}_{pq}$ ,  $\mathbf{C}_{pq}$  and  $\mathbf{D}_{pq}$ , respectively. Similarly,  $\mathbf{E}^i$  is the  $k \times d$  matrix that formed by taking the  $i$ -th diagonal entry of each block  $\mathbf{E}_{pl}$ . Then we can decouple the term of  $\mathbf{G}^{-1} (\mathbf{G}\mathbf{D} - \mathbf{C}) \widetilde{\mathcal{W}}^*$  in (12) into  $i$  vectors, defined as

$$\mathbf{f}_i := (\mathbf{G}^i)^{-1} (\mathbf{G}^i \mathbf{D}^i - \mathbf{C}^i) \mathbf{w}_i^*, \quad (17)$$

where  $\mathbf{w}_i^* \in \mathbb{R}^k$  is the vector formed by taking the  $((p-1)rn + i)$ -th elements of  $\widetilde{\mathcal{W}}^*$  for  $p = 1, \dots, k$ , which indeed is the  $i$ -th column of  $\mathbf{W}^*$ . Similarly, we can decouple  $\mathbf{G}^{-1} \mathbf{E} (\hat{\mathcal{S}}^* - \hat{\mathcal{S}}^t)$  into  $i$  vectors, defined as

$$\mathbf{h}_i = (\mathbf{G}^i)^{-1} \mathbf{E}^i (\mathbf{s}_i^* - \bar{\mathbf{s}}_i^t), \quad (18)$$

where  $\bar{\mathbf{s}}_i^t \in \mathbb{R}^d$  and  $\mathbf{s}_i^* \in \mathbb{R}^d$  are vectors formed by taking the  $((l-1)rn + i)$ -th elements of  $\widetilde{\mathcal{S}}^t$  and  $\widetilde{\mathcal{S}}^*$ , respectively.

Next, we consider the vector  $\mathbf{w}_i^{t+1}$  formed by taking the  $((p-1)rn + i)$ -th elements of  $\widetilde{\mathcal{W}}^{t+1}$  for  $p = 1, \dots, k$ , which is also the  $i$ -th column of  $\mathbf{W}^{t+1}$  from (12) we have

$$(1 - \alpha) \bar{\mathbf{w}}_i^{t+1} = (1 - \alpha) \mathbf{D}^i \mathbf{w}_i^* - (1 - \alpha) (\mathbf{G}^i)^{-1} (\mathbf{G}^i \mathbf{D}^i - \mathbf{C}^i) \mathbf{w}_i^* + \alpha (\mathbf{G}^i)^{-1} \mathbf{E}^i (\mathbf{s}_i^* - \bar{\mathbf{s}}_i^t), \quad (19)$$

where we preliminarily obtain the update rule of each column of  $\mathbf{W}^{t+1}$ . Next, we focus on the update for  $\mathbf{S}$ , after which we can further rewrite the update rule of  $\mathbf{W}$  in a simpler form.

According to (6),  $\widetilde{\mathcal{S}}^{t+1}$  minimizes

$$\Phi(\hat{\mathbf{B}}^{t+1}, \bar{\mathbf{W}}^{t+1}, \widetilde{\mathcal{S}}) := \frac{1}{2rnm} \left\| \mathcal{A} \left( (1 - \alpha) (\mathbf{W}^{* \top} \hat{\mathbf{B}}^{* \top} - \bar{\mathbf{W}}^{t+1 \top} \hat{\mathbf{B}}^{t+1 \top}) + \alpha (\mathbf{S}^{* \top} - \mathbf{S}^\top) \right) \right\|_2^2 + \frac{\beta}{2} \left\| \widetilde{\mathcal{S}} \right\|_{\text{F}}^2. \quad (20)$$

Then via a similar process from (9) to (19), we can obtain

$$\alpha \widetilde{\mathbf{s}}_i^{t+1} = (\mathbf{\Pi}^i + \beta \mathbf{I}_d)^{-1} \mathbf{\Pi}^i \left( \alpha \mathbf{s}_i^* + (1 - \alpha) \hat{\mathbf{B}}^* \mathbf{w}_i^* - (1 - \alpha) \hat{\mathbf{B}}^{t+1} \bar{\mathbf{w}}_i^{t+1} \right), \quad (21)$$

further, we have

$$\alpha \widetilde{\mathbf{s}}_i^{t+1} = \Delta_i^{t+1} + \hat{\mathbf{B}}_\perp^{t+1} \hat{\mathbf{B}}_\perp^{t+1 \top} \left( \alpha \mathbf{s}_i^* + (1 - \alpha) \hat{\mathbf{B}}^* \mathbf{w}_i^* - (1 - \alpha) \hat{\mathbf{B}}^{t+1} \bar{\mathbf{w}}_i^{t+1} \right) \quad (22)$$

$$= \Delta_i^{t+1} + \hat{\mathbf{B}}_\perp^{t+1} \hat{\mathbf{B}}_\perp^{t+1 \top} \left( \alpha \mathbf{s}_i^* + (1 - \alpha) \hat{\mathbf{B}}^* \mathbf{w}_i^* \right), \quad (23)$$

where

$$\Delta_i^{t+1} := \hat{\mathbf{B}}_{\perp}^{t+1} \hat{\mathbf{B}}_{\perp}^{t+1\top} \left( (\boldsymbol{\Pi}^i + \beta \mathbf{I}_d)^{-1} \boldsymbol{\Pi}^i - \mathbf{I}_d \right) \left( \alpha \mathbf{s}_i^* + (1 - \alpha) \hat{\mathbf{B}}^* \mathbf{w}_i^* - (1 - \alpha) \hat{\mathbf{B}}^{t+1} \bar{\mathbf{w}}_i^{t+1} \right). \quad (24)$$

From (23), we have

$$\begin{aligned} \alpha \mathbf{s}_i^* - \alpha \bar{\mathbf{s}}_i^t &= \alpha \mathbf{s}_i^* - \hat{\mathbf{B}}_{\perp}^t \hat{\mathbf{B}}_{\perp}^{t\top} \left( \alpha \mathbf{s}_i^* + (1 - \alpha) \hat{\mathbf{B}}^* \mathbf{w}_i^* - (1 - \alpha) \hat{\mathbf{B}}^t \mathbf{w}_i^* \right) - \alpha \Delta_i^t \\ &= \alpha \left( \mathbf{I}_d - \hat{\mathbf{B}}_{\perp}^t \hat{\mathbf{B}}_{\perp}^{t\top} \right) \mathbf{s}_i^* - (1 - \alpha) \hat{\mathbf{B}}_{\perp}^t \hat{\mathbf{B}}_{\perp}^{t\top} \hat{\mathbf{B}}^* \mathbf{w}_i^* - \alpha \Delta_i^t \\ &= \alpha \hat{\mathbf{B}}^t \hat{\mathbf{B}}^{t\top} \mathbf{s}_i^* - (1 - \alpha) \hat{\mathbf{B}}_{\perp}^t \hat{\mathbf{B}}_{\perp}^{t\top} \hat{\mathbf{B}}^* \mathbf{w}_i^* - \alpha \Delta_i^t, \end{aligned} \quad (25)$$

Then consider

$$\begin{aligned} & - (1 - \alpha) (\mathbf{G}^i)^{-1} (\mathbf{G}^i \mathbf{D}^i - \mathbf{C}^i) \mathbf{w}_i^* + \alpha (\mathbf{G}^i)^{-1} \mathbf{E}^i (\mathbf{s}_i^* - \bar{\mathbf{s}}_i^t) \\ &= - (1 - \alpha) (\mathbf{G}^i)^{-1} \left( \hat{\mathbf{B}}^{t\top} \boldsymbol{\Pi}^i \hat{\mathbf{B}}^t \hat{\mathbf{B}}^{t\top} \hat{\mathbf{B}}^* - \hat{\mathbf{B}}^{t\top} \boldsymbol{\Pi}^i \hat{\mathbf{B}}^* \right) \mathbf{w}_i^* + (\mathbf{G}^i)^{-1} \mathbf{E}^i \left( \alpha \hat{\mathbf{B}}^t \hat{\mathbf{B}}^{t\top} \mathbf{s}_i^* - (1 - \alpha) \hat{\mathbf{B}}_{\perp}^t \hat{\mathbf{B}}_{\perp}^{t\top} \hat{\mathbf{B}}^* \mathbf{w}_i^* - \alpha \Delta_i^t \right) \\ &= (1 - \alpha) (\mathbf{G}^i)^{-1} \hat{\mathbf{B}}^{t\top} \boldsymbol{\Pi}^i \hat{\mathbf{B}}_{\perp}^t \hat{\mathbf{B}}_{\perp}^{t\top} \hat{\mathbf{B}}^* \mathbf{w}_i^* + \alpha (\mathbf{G}^i)^{-1} \hat{\mathbf{B}}^{t\top} \boldsymbol{\Pi}^i \hat{\mathbf{B}}^t \hat{\mathbf{B}}^{t\top} \mathbf{s}_i^* \\ &\quad - (1 - \alpha) (\mathbf{G}^i)^{-1} \hat{\mathbf{B}}^{t\top} \boldsymbol{\Pi}^i \hat{\mathbf{B}}_{\perp}^t \hat{\mathbf{B}}_{\perp}^{t\top} \hat{\mathbf{B}}^* \mathbf{w}_i^* - \alpha (\mathbf{G}^i)^{-1} \mathbf{E}^i \Delta_i^t \end{aligned} \quad (26)$$

$$= \alpha (\mathbf{G}^i)^{-1} \hat{\mathbf{B}}^{t\top} \boldsymbol{\Pi}^i \hat{\mathbf{B}}^t \hat{\mathbf{B}}^{t\top} \mathbf{s}_i^* - \alpha (\mathbf{G}^i)^{-1} \mathbf{E}^i \Delta_i^t. \quad (27)$$

Let  $\bar{\Delta}_i^t := -\alpha (\mathbf{G}^i)^{-1} \mathbf{E}^i \Delta_i^t$ , and we can rewrite (19) as

$$(1 - \alpha) \bar{\mathbf{w}}_i^{t+1} = \hat{\mathbf{B}}^{t\top} \left( (1 - \alpha) \hat{\mathbf{B}}^* \mathbf{w}_i^* + \alpha \mathbf{s}_i^* \right) + \bar{\Delta}_i^t \quad (28)$$

Next, we focus on bounding  $\|\mathbf{H}\|_2$ .

**Lemma 1** Let  $\delta_k = c \frac{k^{3/2} \sqrt{\log(rn)}}{\sqrt{m}}$ ,  $\delta_d = c_3 \frac{\sqrt{d \log(rn)}}{\sqrt{m}}$ ,  $\delta = \frac{\delta_d}{1 - \delta_k}$  for some absolute constant  $c, c_2$ , then

$$\frac{1}{\sqrt{rn}} \|\mathbf{H}\|_2 \leq (1 + \delta) C_s \quad (29)$$

with probability at least  $1 - e^{-120k^3 \log(rn)}$ .

*Proof:* Recall that  $\mathbf{H} := [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_{rn}]$  and

$$\mathbf{h}_i = (\mathbf{G}^i)^{-1} \mathbf{E}^i (\hat{\mathbf{s}}_i^* - \hat{\mathbf{s}}_i^t) = \hat{\mathbf{B}}^{t\top} (\hat{\mathbf{s}}_i^* - \hat{\mathbf{s}}_i^t) - (\mathbf{G}^i)^{-1} (\mathbf{G}^i \hat{\mathbf{B}}^{t\top} - \mathbf{E}^i) (\hat{\mathbf{s}}_i^* - \hat{\mathbf{s}}_i^t), \quad (30)$$

then we focus on the term of  $\mathbf{G}^i \hat{\mathbf{B}}^{t\top} - \mathbf{E}^i$ , for which we have

$$\mathbf{G}^i \hat{\mathbf{B}}^{t\top} - \mathbf{E}^i = \hat{\mathbf{B}}^{t\top} \left( \frac{1}{m} \mathbf{X}_i^{\top} \mathbf{X}_i \right) (\hat{\mathbf{B}}^t \hat{\mathbf{B}}^{t\top} - \mathbf{I}_d). \quad (31)$$

Let  $\mathbf{U} := \frac{1}{\sqrt{m}} \mathbf{X}_i (\hat{\mathbf{B}}^t \hat{\mathbf{B}}^{t\top} - \mathbf{I}_d)$  and  $\mathbf{V} := \frac{1}{\sqrt{m}} \mathbf{X}_i \hat{\mathbf{B}}^t$ , then we have the  $j$ -th row of  $\mathbf{U}$  and  $\mathbf{V}$  as the following, respectively:

$$\mathbf{u}_j = \frac{1}{\sqrt{m}} (\hat{\mathbf{B}}^t \hat{\mathbf{B}}^{t\top} - \mathbf{I}_d) \mathbf{x}_i^j, \quad \mathbf{v}_j = \frac{1}{\sqrt{m}} \hat{\mathbf{B}}^{t\top} \mathbf{x}_i^j. \quad (32)$$

Note that  $\mathbf{u}_j$  is  $\frac{1}{\sqrt{m}}(\hat{\mathbf{B}}^t \hat{\mathbf{B}}^{t\top} - \mathbf{I}_d)$ -sub-gaussian and  $\mathbf{v}_j$  is  $\frac{1}{\sqrt{m}}\hat{\mathbf{B}}^t$ -sub-gaussian, therefore we can argue similarly as the derivatives for Theorem 4.4.5 in (Vershynin, 2018). First, let  $\mathcal{S}^{d-1}$  be the  $d$ -dimension unit sphere and  $\mathcal{S}^{k-1}$  be the  $k$ -dimension unit sphere, then let  $\mathcal{N}_d$  be the  $\frac{1}{4}$ -th net on  $\mathcal{S}^{d-1}$  and  $\mathcal{N}_k$  be the  $\frac{1}{4}$ -th net on  $\mathcal{S}^{k-1}$ , such that  $|\mathcal{N}_d| \leq 9^d$  and  $|\mathcal{N}_k| \leq 9^k$ , which exists according to Corollary 4.2.13 in (Vershynin, 2018). Next, by leveraging inequality 4.13 in (Vershynin, 2018), we have

$$\begin{aligned} \left\| \left( \hat{\mathbf{B}}^t \hat{\mathbf{B}}^{t\top} - \mathbf{I}_d \right) \left( \frac{d}{m} \mathbf{X}_i^\top \mathbf{X}_i \right) \hat{\mathbf{B}}^t \right\|_2 &= \left\| \mathbf{U}^\top \mathbf{V} \right\|_2 \leq 2 \max_{\mathbf{z} \in \mathcal{N}_d, \mathbf{y} \in \mathcal{N}_k} \mathbf{z}^\top \mathbf{U}^\top \mathbf{V} \mathbf{y} \\ &= 2 \max_{\mathbf{z} \in \mathcal{N}_d, \mathbf{y} \in \mathcal{N}_k} \mathbf{z}^\top \left( \sum_{j=1}^m \mathbf{u}_j \mathbf{v}_j^\top \right) \mathbf{y} \\ &= 2 \max_{\mathbf{z} \in \mathcal{N}_d, \mathbf{y} \in \mathcal{N}_k} \sum_{j=1}^m \langle \mathbf{z}, \mathbf{u}_j \rangle \langle \mathbf{v}_j, \mathbf{y} \rangle. \end{aligned} \quad (33)$$

By the definition of sub-gaussianity,  $\langle \mathbf{z}, \mathbf{u}_j \rangle$  is sub-gaussian with norm  $\frac{1}{\sqrt{m}} \left\| \hat{\mathbf{B}}^t \hat{\mathbf{B}}^{t\top} - \mathbf{I}_d \right\|_2 \leq \frac{2}{\sqrt{m}}$  and  $\langle \mathbf{v}_j, \mathbf{y} \rangle$  is sub-gaussian with norm  $\frac{1}{\sqrt{m}} \left\| \hat{\mathbf{B}}^t \right\|_2 = \frac{1}{\sqrt{m}}$ . Therefore,  $\langle \mathbf{z}, \mathbf{u}_j \rangle \langle \mathbf{v}_j, \mathbf{y} \rangle$  is sub-exponential with norm at most  $\frac{c}{m}$  for some absolute constant  $c$ , for all  $j \in [m]$ . Also, for any  $j \in [m]$  and any  $\mathbf{z} \in \mathcal{N}_d, \mathbf{y} \in \mathcal{N}_k$ , we have

$$\mathbb{E}[\langle \mathbf{z}, \mathbf{u}_j \rangle \langle \mathbf{v}_j, \mathbf{y} \rangle] = \mathbb{E}[\mathbf{z}^\top (\hat{\mathbf{B}}^t \hat{\mathbf{B}}^{t\top} - \mathbf{I}_d) \frac{d}{m} \mathbf{X}_i^\top \mathbf{X}_i \hat{\mathbf{B}}^t] = 0. \quad (34)$$

Thus, we obtain a sum of  $m$  mean-zero, independent sub-exponential random variables, for which we apply Bernstein's inequality, for any  $\mathbf{z} \in \mathcal{N}_d, \mathbf{y} \in \mathcal{N}_k$ ,

$$\mathbb{P} \left( \sum_{j=1}^m \langle \mathbf{z}, \mathbf{u}_j \rangle \langle \mathbf{v}_j, \mathbf{y} \rangle \geq s \right) \leq e^{-c' m \min(s^2, s)}. \quad (35)$$

Union bounding over all  $\mathbf{z} \in \mathcal{N}_d, \mathbf{y} \in \mathcal{N}_k$ , we obtain

$$\mathbb{P} \left( \left\| \left( \hat{\mathbf{B}}^t \hat{\mathbf{B}}^{t\top} - \mathbf{I}_d \right) \left( \frac{1}{m} \mathbf{X}_i^\top \mathbf{X}_i \right) \hat{\mathbf{B}}^t \right\|_2 \geq 2s \right) \leq 9^{d+k} e^{-c' m \min(s^2, s)}. \quad (36)$$

Here, let  $s = \max(\varepsilon, \varepsilon^2)$  for some  $\varepsilon > 0$ , then we have  $\min(s^2, s) = \varepsilon^2$ . Then we reach

$$\mathbb{P} \left( \left\| \left( \hat{\mathbf{B}}^t \hat{\mathbf{B}}^{t\top} - \mathbf{I}_d \right) \left( \frac{1}{m} \mathbf{X}_i^\top \mathbf{X}_i \right) \hat{\mathbf{B}}^t \right\|_2 \geq 2 \max(\varepsilon, \varepsilon^2) \right) \leq 9^{d+k} e^{-c' m \varepsilon^2}. \quad (37)$$

Further, let  $\varepsilon = \sqrt{\frac{c_2 d \log(rn)}{m}}$  for some constant  $c_2$ . Then conditioned on  $\varepsilon \leq 1$ , we have

$$\mathbb{P} \left( \left\| \left( \hat{\mathbf{B}}^t \hat{\mathbf{B}}^{t\top} - \mathbf{I}_d \right) \left( \frac{1}{m} \mathbf{X}_i^\top \mathbf{X}_i \right) \hat{\mathbf{B}}^t \right\|_2 \geq c_3 \sqrt{\frac{d \log(rn)}{m}} \right) \leq 9^{d+k} e^{-c_4 d \log(rn)} \leq e^{-110 d \log(rn)}, \quad (38)$$

for a large enough constant  $c_1$ . According to (30),

$$\begin{aligned}\|\mathbf{h}_i\|_2 &\leq \left\| \hat{\mathbf{B}}^{t\top} \right\|_2 \left\| \hat{\mathbf{s}}_i^* - \hat{\mathbf{s}}_i^t \right\|_2 + \left\| (\mathbf{G}^i)^{-1} \right\|_2 \left\| \mathbf{G}^i \hat{\mathbf{B}}^{t\top} - \mathbf{E}^i \right\|_2 \left\| \hat{\mathbf{s}}_i^* - \hat{\mathbf{s}}_i^t \right\|_2 \\ &= \left( 1 + \left\| (\mathbf{G}^i)^{-1} \right\|_2 \left\| \mathbf{G}^i \hat{\mathbf{B}}^{t\top} - \mathbf{E}^i \right\|_2 \right) \left\| \hat{\mathbf{s}}_i^* - \hat{\mathbf{s}}_i^t \right\|_2.\end{aligned}\quad (39)$$

From (38), we know that

$$\mathbb{P} \left( \left\| \mathbf{G}^i \hat{\mathbf{B}}^{t\top} - \mathbf{E}^i \right\|_2 \geq \delta_d \right) \leq e^{-110d \log(rn)}, \quad (40)$$

and from equation (43) in (Collins et al., 2021) we have

$$\mathbb{P} \left( \left\| (\mathbf{G}^i)^{-1} \right\|_2 \geq \frac{1}{1 - \delta_k} \right) \leq e^{-121k^3 \log(rn)} \quad (41)$$

Therefore, we obtain

$$\|\mathbf{h}_i\|_2 \leq (1 + \delta) \left\| \hat{\mathbf{s}}_i^* - \hat{\mathbf{s}}_i^t \right\|_2 \quad (42)$$

with probability at least  $1 - e^{-110d \log(rn)} - e^{-121k^3 \log(rn)}$ . Finally we take a union bound over  $i \in [rn]$ , leading to

$$\begin{aligned}\mathbb{P} \left( \frac{1}{rn} \|\mathbf{H}\|_2^2 \geq (1 + \delta)^2 \left\| \mathbf{s}_1^* - \mathbf{s}_1^t \right\|_2^2 \right) &\leq \mathbb{P} \left( \frac{1}{rn} \sum_{i=1}^{rn} \|\mathbf{h}_i\|_2^2 \geq (1 + \delta)^2 \left\| \mathbf{s}_1^* - \mathbf{s}_1^t \right\|_2^2 \right) \\ &\leq rn \mathbb{P} \left( \|\mathbf{h}_1\|_2^2 \geq (1 + \delta)^2 \left\| \mathbf{s}_1^* - \mathbf{s}_1^t \right\|_2^2 \right) \\ &\leq rn \mathbb{P} \left( \|\mathbf{h}_1\|_2^2 \geq (1 + \delta)^2 \left\| \mathbf{s}_1^* - \mathbf{s}_1^t \right\|_2^2 \right) \\ &\leq e^{-120k^3 \log(rn)}\end{aligned}\quad (43)$$

$$\leq e^{-120k^3 \log(rn)} \quad (44)$$

and thus completing the proof.  $\square$

**Lemma 2** Let  $\delta'_k = c_4 k \frac{\sqrt{d}}{\sqrt{rnm}}$  for some absolute constant  $c_4$ . Then for any  $t$ ,

$$\frac{1}{rn} \left\| \left( \frac{1}{m} \mathcal{A}^\dagger \mathcal{A} (\mathbf{Q}^{t\top}) - \mathbf{Q}^{t\top} \right)^\top (1 - \alpha) \mathbf{W}^{t+1\top} \right\|_2 \leq \delta'_k \text{dist}(\hat{\mathbf{B}}^t, \hat{\mathbf{B}}^*) \quad (45)$$

with probability at least  $1 - e^{-110d} - e^{-110k^2 \log(rn)}$ .

*Proof:* Let  $\mathbf{Q}^t = (1 - \alpha)(\hat{\mathbf{B}}^t \mathbf{W}^{t+1} - \hat{\mathbf{B}}^* \mathbf{W}^*) + \alpha(\hat{\mathbf{S}}^t - \hat{\mathbf{S}}^*)$ . We first consider the bound of the columns of  $\mathbf{Q}$ . Let  $\mathbf{q}_i \in \mathbb{R}^d$  be the  $i$ -th column of  $\mathbf{Q}$ , for all  $i \in [rn]$  we have

$$\begin{aligned}\mathbf{q}_i &= (1 - \alpha) \left( \hat{\mathbf{B}}^t \mathbf{w}_i^{t+1} - \hat{\mathbf{B}}^* \mathbf{w}_i^* \right) + \alpha \left( \hat{\mathbf{s}}_i^t - \hat{\mathbf{s}}_i^* \right) \\ &= (1 - \alpha) \hat{\mathbf{B}}^t \hat{\mathbf{B}}^{t\top} \hat{\mathbf{B}}^* \mathbf{w}_i^* - (1 - \alpha) \hat{\mathbf{B}}^t \mathbf{f}_i - \alpha \hat{\mathbf{B}}^t \mathbf{h}_i - (1 - \alpha) \hat{\mathbf{B}}^* \mathbf{w}_i^* + \alpha \hat{\mathbf{s}}_i^t - \alpha \hat{\mathbf{s}}_i^* \\ &= (1 - \alpha) \left( \hat{\mathbf{B}}^t \hat{\mathbf{B}}^{t\top} - \mathbf{I}_d \right) \hat{\mathbf{B}}^* \mathbf{w}_i^* - \hat{\mathbf{B}}^t \mathbf{k}_i + \alpha \hat{\mathbf{s}}_i^t - \alpha \hat{\mathbf{s}}_i^*\end{aligned}\quad (46)$$

Thus,

$$\begin{aligned}\|\mathbf{q}_i\|_2 &= \left\| (1-\alpha) \left( \hat{\mathbf{B}}^t \hat{\mathbf{B}}^{t\top} - \mathbf{I}_d \right) \hat{\mathbf{B}}^* \mathbf{w}_i^* - \hat{\mathbf{B}}^t \mathbf{k}_i + \alpha \hat{\mathbf{s}}_i^t - \alpha \hat{\mathbf{s}}_i^* \right\|_2 \\ &\leq \left\| (1-\alpha) \left( \hat{\mathbf{B}}^t \hat{\mathbf{B}}^{t\top} - \mathbf{I}_d \right) \hat{\mathbf{B}}^* \right\|_2 \|\mathbf{w}_i^*\|_2 + \|\mathbf{k}_i\|_2 + \alpha \|\hat{\mathbf{s}}_i^t - \hat{\mathbf{s}}_i^*\|_2 \\ &\leq (1-\alpha)\sqrt{k} \text{dist}(\hat{\mathbf{B}}^t, \hat{\mathbf{B}}^*) + \alpha C_s \text{dist}(\hat{\mathbf{B}}^t, \hat{\mathbf{B}}^*) + (\alpha C_s + (1-\alpha)\sqrt{k}) \text{dist}(\hat{\mathbf{B}}^t, \hat{\mathbf{B}}^*) \quad (47)\end{aligned}$$

$$\leq 2 \left( (1-\alpha)\sqrt{k} + \alpha C_s \right) \text{dist}(\hat{\mathbf{B}}^t, \hat{\mathbf{B}}^*) \quad (48)$$

$$\leq 2\sqrt{k} \text{dist}(\hat{\mathbf{B}}^t, \hat{\mathbf{B}}^*) \quad (49)$$

where (47) holds with probability at least  $1 - e^{-110k^2 \log(rn)}$ , by combining equation (44) in (Collins et al., 2021) and (42), conditioned on  $\delta_k \leq \frac{1}{2}$  and  $\delta_d \leq \frac{1}{2}$ . Similarly, combining equation (45) and (42), conditioned on  $\delta_k \leq \frac{1}{2}$ , we have

$$\begin{aligned}\|(1-\alpha)\mathbf{w}_i^{t+1}\|_2 &\leq \left\| (1-\alpha)\hat{\mathbf{B}}^{t\top} \hat{\mathbf{B}}^* \mathbf{w}_i^* \right\|_2 + \|\mathbf{k}_i\|_2 \\ &\leq (1-\alpha)\sqrt{k} + \alpha C_s \quad (50)\end{aligned}$$

$$\leq 2\sqrt{k} \quad (51)$$

with probability at least  $1 - e^{-110k^2 \log(rn)}$ .

Next, just for simple notation, let  $\Delta_{\mathbf{s}}^t$  denote  $\mathbf{S}^* - \mathbf{S}^t$  and  $\Delta_{\mathbf{BW}}^t$  denote  $\hat{\mathbf{B}}^* \mathbf{W}^* - \hat{\mathbf{B}}^t \mathbf{W}^t$ . and in the following proof, we condition on the event

$$\mathcal{E} := \bigcap_{i=1}^{rn} \left\{ \|\mathbf{q}_i\|_2 \leq 2 \left( (1-\alpha)\sqrt{k} + \alpha C_s \right) \text{dist}(\hat{\mathbf{B}}^t, \hat{\mathbf{B}}^*) \cap \|(1-\alpha)\mathbf{w}_i^{t+1}\|_2 \leq (1-\alpha)\sqrt{k} + \alpha C_s \right\}, \quad (52)$$

which holds with probability at least  $1 - e^{-109k^2 \log(rn)}$ . Next, we consider the following matrix:

$$\begin{aligned}\frac{1}{m} \mathcal{A}^\dagger \mathcal{A}(\mathbf{Q}^{t\top}) - \mathbf{Q}^{t\top} &= \frac{1}{m} \sum_{i=1}^{rn} \sum_{j=1}^m \left\langle \mathbf{e}_i(\mathbf{x}_i^j)^\top, \mathbf{Q}^{t\top} \right\rangle \mathbf{e}_i(\mathbf{x}_i^j)^\top - \mathbf{Q}^{t\top} \\ &= \frac{1}{m} \sum_{i=1}^{rn} \sum_{j=1}^m \langle \mathbf{x}_i^j, \mathbf{q}_i \rangle \mathbf{e}_i(\mathbf{x}_i^j)^\top - \mathbf{Q}^{t\top}, \quad (53)\end{aligned}$$

further, we have

$$\frac{1}{rn} \left( \frac{1}{m} \mathcal{A}^\dagger \mathcal{A}(\mathbf{Q}^{t\top}) - \mathbf{Q}^{t\top} \right)^\top (1-\alpha)\mathbf{W}^{t+1\top} = \frac{1}{rnm} \sum_{i=1}^{rn} \sum_{j=1}^m \left( \langle \mathbf{x}_i^j, \mathbf{q}_i \rangle \mathbf{x}_i^j (1-\alpha)\mathbf{w}_i^{t+1\top} - \mathbf{q}_i (1-\alpha)\mathbf{w}_i^{t+1\top} \right). \quad (54)$$

Next, we establish similar arguments as the derivatives for Theorem 4.4.5 in (Vershynin, 2018) to bound  $\left\| \frac{1}{rnm} \sum_{i=1}^{rn} \sum_{j=1}^m \left( \langle \mathbf{x}_i^j, \mathbf{q}_i \rangle \mathbf{x}_i^j (1-\alpha)\mathbf{w}_i^{t+1\top} - \mathbf{q}_i (1-\alpha)\mathbf{w}_i^{t+1\top} \right) \right\|_2$ . let  $\mathcal{S}^{d-1}$  be the  $d$ -dimension unit sphere and  $\mathcal{S}^{k-1}$  be the  $k$ -dimension unit sphere, then let  $\mathcal{N}_d$  be the  $\frac{1}{4}$ -th net on  $\mathcal{S}^{d-1}$  and  $\mathcal{N}_k$



be the  $\frac{1}{4}$ -th net on  $\mathcal{S}^{k-1}$ , such that  $|\mathcal{N}_d| \leq 9^d$  and  $|\mathcal{N}_k| \leq 9^k$ , which exists according to Corollary 4.2.13 in (Vershynin, 2018). Using equation 4.13 in (Vershynin, 2018), we have

$$\begin{aligned}
& \left\| \frac{1}{rnm} \sum_{i=1}^{rn} \sum_{j=1}^m \left( \langle \mathbf{x}_i^j, \mathbf{q}_i \rangle \mathbf{x}_i^j (1-\alpha) \mathbf{w}_i^{t+1\top} - \mathbf{q}_i (1-\alpha) \mathbf{w}_i^{t+1\top} \right) \right\|_2 \\
& \leq 2 \max_{\mathbf{z} \in \mathcal{N}_d, \mathbf{y} \in \mathcal{N}_k} \mathbf{z}^\top \left( \sum_{i=1}^{rn} \sum_{j=1}^m \left( \frac{1}{rnm} \langle \mathbf{x}_i^j, \mathbf{q}_i \rangle \mathbf{x}_i^j (1-\alpha) \mathbf{w}_i^{t+1\top} - \frac{1}{rnm} \mathbf{q}_i (1-\alpha) \mathbf{w}_i^{t+1\top} \right) \right) \mathbf{y} \\
& = 2 \max_{\mathbf{z} \in \mathcal{N}_d, \mathbf{y} \in \mathcal{N}_k} \sum_{i=1}^{rn} \sum_{j=1}^m \left( \frac{1}{rnm} \langle \mathbf{x}_i^j, \mathbf{q}_i \rangle \langle \mathbf{z}, \mathbf{x}_i^j \rangle \langle (1-\alpha) \mathbf{w}_i^{t+1}, \mathbf{y} \rangle - \frac{1}{rnm} \langle \mathbf{z}, \mathbf{q}_i \rangle \langle (1-\alpha) \mathbf{w}_i^{t+1}, \mathbf{y} \rangle \right) \quad (55)
\end{aligned}$$

Since  $\mathbf{x}_i^j$  is  $\mathbf{I}_d$ -sub-gaussian,  $\langle \mathbf{z}, \mathbf{x}_i^j \rangle$  is sub-gaussian with norm  $\|\mathbf{z}\|_2 = c$  for any  $\mathbf{z} \in \mathcal{N}_d$ . Also  $\langle \mathbf{x}_i^j, \mathbf{q}_i \rangle$  is sub-gaussian with norm  $\|\mathbf{q}_i\|_2$ . Therefore,  $\langle \mathbf{z}, \mathbf{x}_i^j \rangle \langle \mathbf{x}_i^j, \mathbf{q}_i \rangle$  is sub-exponential with norm at most  $c \|\mathbf{q}_i\|_2$ , which indicates  $\frac{1}{rnm} \langle \mathbf{z}, \mathbf{x}_i^j \rangle \langle \mathbf{x}_i^j, \mathbf{q}_i \rangle \langle (1-\alpha) \mathbf{w}_i, \mathbf{y} \rangle$  is sub-exponential with norm at most

$$\begin{aligned}
\frac{c}{rnm} \|\mathbf{q}_i\|_2 \langle (1-\alpha) \mathbf{w}_i, \mathbf{y} \rangle & \leq \frac{c}{rnm} \|\mathbf{q}_i\|_2 \|(1-\alpha) \mathbf{w}_i\|_2 \\
& \leq \frac{c'}{rnm} \left( (1-\alpha) \sqrt{k} + \alpha C_s \right)^2 \text{dist}(\hat{\mathbf{B}}^t, \hat{\mathbf{B}}^*) \quad (56)
\end{aligned}$$

$$:= \frac{c'}{rnm} \Delta \quad (57)$$

for some absolute constant  $c'$ . Since  $\mathbb{E}[\frac{1}{rnm} \langle \mathbf{x}_i^j, \mathbf{q}_i \rangle \langle \mathbf{z}, \mathbf{x}_i^j \rangle \langle (1-\alpha) \mathbf{w}_i, \mathbf{y} \rangle - \frac{1}{rnm} \langle \mathbf{z}, \mathbf{q}_i \rangle \langle (1-\alpha) \mathbf{w}_i, \mathbf{y} \rangle] = 0$ , we have a sum of  $rnm$  independent, mean zero, sub-exponential random variables, for which we can apply Bernstein's inequality and obtain

$$\mathbb{P} \left( \sum_{i=1}^{rn} \sum_{j=1}^m \left( \frac{1}{rnm} \langle \mathbf{x}_i^j, \mathbf{q}_i \rangle \langle \mathbf{z}, \mathbf{x}_i^j \rangle \langle (1-\alpha) \mathbf{w}_i, \mathbf{y} \rangle - \frac{1}{rnm} \langle \mathbf{z}, \mathbf{q}_i \rangle \langle (1-\alpha) \mathbf{w}_i, \mathbf{y} \rangle \right) \geq s \right) \leq \exp \left( -c_2 rnm \min \left( \frac{s^2}{\Delta^2}, \frac{s}{\Delta} \right) \right). \quad (58)$$

Take union bound over all  $\mathbf{z} \in \mathcal{N}_d, \mathbf{y} \in \mathcal{N}_k$ ,

$$\mathbb{P} \left( \left\| \frac{1}{rn} \left( \frac{1}{m} \mathcal{A}^\dagger \mathcal{A} (\mathbf{Q}^{t\top}) - \mathbf{Q}^{t\top} \right) (1-\alpha) \mathbf{W}^{t+1\top} \right\|_2 \geq 2s \middle| \mathcal{E} \right) \leq 9^{d+k} \exp \left( -c_2 rnm \min \left( \frac{s^2}{\Delta^2}, \frac{s}{\Delta} \right) \right). \quad (59)$$

Let  $\frac{s}{\Delta} = \max(\varepsilon, \varepsilon^2)$  for some  $\varepsilon > 0$ , then  $\varepsilon^2 = \min \left( \frac{s^2}{\Delta^2}, \frac{s}{\Delta} \right)$ . Further, let  $\varepsilon = \sqrt{\frac{113(d+k)}{c_2 rnm}}$ , and conditioned on  $\varepsilon \leq 1$ , we obtain

$$\mathbb{P} \left( \left\| \frac{1}{rn} \left( \frac{1}{m} \mathcal{A}^\dagger \mathcal{A} (\mathbf{Q}^{t\top}) - \mathbf{Q}^{t\top} \right) \mathbf{W}^{t+1\top} \right\|_2 \geq c_4 \sqrt{\frac{d}{rnm}} \left( (1-\alpha) \sqrt{k} + \alpha C_s \right)^2 \text{dist}(\hat{\mathbf{B}}^t, \hat{\mathbf{B}}^*) \middle| \mathcal{E} \right) \leq e^{-110d} \quad (60)$$

Finally, by using  $\mathbb{P}(A) \leq \mathbb{P}(A \mid \mathcal{E}) + \mathbb{P}(\mathcal{E})$ , where

$$A := \left\{ \left\| \frac{1}{rn} \left( \frac{1}{m} \mathcal{A}^\dagger \mathcal{A}(\mathbf{Q}^{t\top}) - \mathbf{Q}^{t\top} \right) \mathbf{W}^{t+1\top} \right\|_2 \geq c_4 \sqrt{\frac{d}{rnm}} \left( (1-\alpha)\sqrt{k} + \alpha C_s \right)^2 \text{dist}(\hat{\mathbf{B}}^t, \hat{\mathbf{B}}^*) \right\}, \quad (61)$$

we complete the proof.  $\square$

### 1.3 Main Result

Recall that  $\mathbf{Q}^{t\top} = \mathbf{W}^{t+1\top} \hat{\mathbf{B}}^{t\top} - \mathbf{W}^{*\top} \hat{\mathbf{B}}^{*\top} + \hat{\mathbf{S}}^{t\top} - \hat{\mathbf{S}}^{*\top}$ , plugging this into (4), and without losing generality, we drop the subscripts of  $\mathcal{I}^t$  and obtain

$$\begin{aligned} \bar{\mathbf{B}}^{t+1} &= \hat{\mathbf{B}}^t - \frac{\eta}{rnm} \left( \mathcal{A}^\dagger \mathcal{A}(\mathbf{Q}^{t\top}) \right)^\top \mathbf{W}^{t+1\top} \\ &= \hat{\mathbf{B}}^t - \frac{\eta}{rn} \mathbf{Q}^t \mathbf{W}^{t+1\top} - \frac{\eta}{rn} \left( \frac{1}{m} \mathcal{A}^\dagger \mathcal{A}(\mathbf{Q}^{t\top}) - \mathbf{Q}^{t\top} \right)^\top \mathbf{W}^{t+1\top}. \end{aligned} \quad (62)$$

Since  $\bar{\mathbf{B}}^{t+1} = \hat{\mathbf{B}}^{t+1} \mathbf{R}^{t+1}$ , we right multiply  $(\mathbf{R}^{t+1})^{-1}$  and left multiply  $\hat{\mathbf{B}}_\perp^{*\top}$  on both sides to get

$$\hat{\mathbf{B}}_\perp^{*\top} \hat{\mathbf{B}}^{t+1} = \left( \hat{\mathbf{B}}_\perp^{*\top} \hat{\mathbf{B}}^t - \frac{\eta}{rn} \hat{\mathbf{B}}_\perp^{*\top} \mathbf{Q}^t \mathbf{W}^{t+1\top} - \frac{\eta}{rn} \hat{\mathbf{B}}_\perp^{*\top} \left( \frac{1}{m} \mathcal{A}^\dagger \mathcal{A}(\mathbf{Q}^{t\top}) - \mathbf{Q}^{t\top} \right)^\top \mathbf{W}^{t+1\top} \right) (\mathbf{R}^{t+1})^{-1}. \quad (63)$$

Then we consider the term of  $\hat{\mathbf{B}}_\perp^{*\top} \mathbf{Q}^t \mathbf{W}^{t+1\top}$ :

$$\begin{aligned} \hat{\mathbf{B}}_\perp^{*\top} \mathbf{Q}^t \mathbf{W}^{t+1\top} &= \hat{\mathbf{B}}_\perp^{*\top} \left( \hat{\mathbf{B}}^t \mathbf{W}^{t+1} - \hat{\mathbf{B}}^* \mathbf{W}^* + \hat{\mathbf{S}}^t - \hat{\mathbf{S}}^* \right) \mathbf{W}^{t+1\top} \\ &= \hat{\mathbf{B}}_\perp^{*\top} \hat{\mathbf{B}}^t \mathbf{W}^{t+1} \mathbf{W}^{t+1\top} - \hat{\mathbf{B}}_\perp^{*\top} \left( \hat{\mathbf{S}}^* - \hat{\mathbf{S}}^t \right) \mathbf{W}^{t+1\top}, \end{aligned}$$

plugging this into (63) then we reach

$$\begin{aligned} \hat{\mathbf{B}}_\perp^{*\top} \hat{\mathbf{B}}^{t+1} &= \left( \hat{\mathbf{B}}_\perp^{*\top} \hat{\mathbf{B}}^t \left( \mathbf{I}_k - \frac{\eta}{rn} \mathbf{W}^{t+1} \mathbf{W}^{t+1\top} \right) + \frac{\eta}{rn} \hat{\mathbf{B}}_\perp^{*\top} \left( \hat{\mathbf{S}}^* - \hat{\mathbf{S}}^t \right) \mathbf{W}^{t+1\top} \right. \\ &\quad \left. - \frac{\eta}{rn} \hat{\mathbf{B}}_\perp^{*\top} \left( \frac{1}{m} \mathcal{A}^\dagger \mathcal{A}(\mathbf{Q}^{t\top}) - \mathbf{Q}^{t\top} \right)^\top \mathbf{W}^{t+1\top} \right) (\mathbf{R}^{t+1})^{-1}. \end{aligned} \quad (64)$$

Therefore,

$$\begin{aligned} \text{dist}(\hat{\mathbf{B}}^{t+1}, \hat{\mathbf{B}}^*) &= \left\| \hat{\mathbf{B}}_\perp^{*\top} \hat{\mathbf{B}}^{t+1} \right\|_2 \\ &\leq \left\| \hat{\mathbf{B}}_\perp^{*\top} \hat{\mathbf{B}}^t \left( \mathbf{I}_k - \frac{\eta}{rn} (1-\alpha)^2 \mathbf{W}^{t+1} \mathbf{W}^{t+1\top} \right) \right\|_2 \left\| (\mathbf{R}^{t+1})^{-1} \right\|_2 \end{aligned}$$

$$\begin{aligned}
& + \frac{\eta}{rn} \left\| \hat{\mathbf{B}}_{\perp}^{*\top} \left( \frac{1}{m} (\mathcal{A}^\dagger \mathcal{A}(\mathbf{Q}^{t\top}) - \mathbf{Q}^{t\top}) \right)^\top (1 - \alpha) \mathbf{W}^{t+1\top} \right\|_2 \left\| (\mathbf{R}^{t+1})^{-1} \right\|_2 \\
& + \frac{\eta}{rn} \left\| \hat{\mathbf{B}}_{\perp}^{*\top} \left( \alpha \hat{\mathbf{S}}^* - \alpha \hat{\mathbf{S}}^t \right) (1 - \alpha) \mathbf{W}^{t+1\top} \right\|_2 \left\| (\mathbf{R}^{t+1})^{-1} \right\|_2. \tag{65}
\end{aligned}$$

Next, we focus on the term of  $\left\| \hat{\mathbf{B}}_{\perp}^{*\top} \hat{\mathbf{B}}^t \left( \mathbf{I}_k - \frac{\eta}{rn} \mathbf{W}^{t+1} \mathbf{W}^{t+1\top} \right) \right\|_2$ , for which we have

$$\begin{aligned}
\left\| \hat{\mathbf{B}}_{\perp}^{*\top} \hat{\mathbf{B}}^t \left( \mathbf{I}_k - \frac{\eta}{rn} (1 - \alpha)^2 \mathbf{W}^{t+1} \mathbf{W}^{t+1\top} \right) \right\|_2 & \leq \left\| \hat{\mathbf{B}}_{\perp}^{*\top} \hat{\mathbf{B}}^t \right\|_2 \left\| \mathbf{I}_k - \frac{\eta}{rn} \mathbf{W}^{t+1} (1 - \alpha) \mathbf{W}^{t+1\top} \right\|_2 \\
& \leq \text{dist}(\hat{\mathbf{B}}^t, \hat{\mathbf{B}}^*) \left\| \mathbf{I}_k - \frac{\eta}{rn} \mathbf{W}^{t+1} \mathbf{W}^{t+1\top} \right\|_2. \tag{66}
\end{aligned}$$

To bound the term of  $\left\| \mathbf{I}_k - \frac{\eta}{rn} \mathbf{W}^{t+1} \mathbf{W}^{t+1\top} \right\|_2$ , we assume that  $\frac{1}{\sqrt{rn}} \mathbf{W}^{t+1}$  has non-zero minimum singular value, defined as  $\sigma_{\min}^{t+1}$ . Then as long as  $\eta \leq (\sigma_{\min}^{t+1})^2$ , we have

$$\left\| \mathbf{I}_k - \frac{\eta}{rn} \mathbf{W}^{t+1} \mathbf{W}^{t+1\top} \right\|_2 = 1 - \eta (\sigma_{\min}^{t+1})^2. \tag{67}$$

To bound the term of  $\frac{\eta}{rn} \left\| \hat{\mathbf{B}}_{\perp}^{*\top} \left( \frac{1}{m} (\mathcal{A}^\dagger \mathcal{A}(\mathbf{Q}^{t\top}) - \mathbf{Q}^{t\top}) \right)^\top \mathbf{W}^{t+1\top} \right\|_2$ , we have

$$\begin{aligned}
\frac{\eta}{rn} \left\| \hat{\mathbf{B}}_{\perp}^{*\top} \left( \frac{1}{m} (\mathcal{A}^\dagger \mathcal{A}(\mathbf{Q}^{t\top}) - \mathbf{Q}^{t\top}) \right)^\top \mathbf{W}^{t+1\top} \right\|_2 & \leq \frac{\eta}{rn} \left\| \left( \frac{1}{m} (\mathcal{A}^\dagger \mathcal{A}(\mathbf{Q}^{t\top}) - \mathbf{Q}^{t\top}) \right)^\top \mathbf{W}^{t+1\top} \right\|_2 \\
& \leq \eta \left( \delta'_k \text{dist}(\hat{\mathbf{B}}^t, \hat{\mathbf{B}}^*) + \delta''_k \right). \tag{68}
\end{aligned}$$

Similarly,

$$\frac{\eta}{rn} \left\| \hat{\mathbf{B}}_{\perp}^{*\top} \left( \hat{\mathbf{S}}^* - \hat{\mathbf{S}}^{t+1} \right) \mathbf{W}^{t+1\top} \right\|_2 \leq \frac{\eta}{\sqrt{rn}} \left\| \hat{\mathbf{S}}^* - \hat{\mathbf{S}}^t \right\|_2 \frac{1}{\sqrt{rn}} \left\| \mathbf{W}^{t+1} \right\|_2 \leq \eta 2\sqrt{k} 6\sqrt{k} = 12\eta k, \tag{69}$$

Then, we focus on bounding  $\left\| (\mathbf{R}^{t+1})^{-1} \right\|_2$ . Just for simple notation, let  $\mathbf{U}^t := \frac{1}{m} \mathcal{A}^\dagger \mathcal{A}(\mathbf{Q}^{t\top})$ , then we have

$$\begin{aligned}
\mathbf{R}^{t+1\top} \mathbf{R}^{t+1} & = \bar{\mathbf{B}}^{t+1\top} \bar{\mathbf{B}}^{t+1} \\
& = \hat{\mathbf{B}}^{t\top} \hat{\mathbf{B}}^t - \frac{\eta}{rn} \left( \hat{\mathbf{B}}^{t\top} \mathbf{U}^{t\top} \mathbf{W}^{t+1\top} + \mathbf{W}^{t+1} \mathbf{U}^t \hat{\mathbf{B}}^t \right) + \frac{\eta^2}{(rn)^2} \mathbf{W}^{t+1} \mathbf{U}^t \mathbf{U}^{t\top} \mathbf{W}^{t+1\top} \\
& = \mathbf{I}_k - \frac{\eta}{rn} \left( \hat{\mathbf{B}}^{t\top} \mathbf{U}^{t\top} \mathbf{W}^{t+1\top} + \mathbf{W}^{t+1} \mathbf{U}^t \hat{\mathbf{B}}^t \right) + \frac{\eta^2}{(rn)^2} \mathbf{W}^{t+1} \mathbf{U}^t \mathbf{U}^{t\top} \mathbf{W}^{t+1\top}. \tag{70}
\end{aligned}$$

Using Weyl's Inequality, we reach

$$\begin{aligned}
\sigma_{\min}^2(\mathbf{R}^{t+1}) & \geq 1 - \frac{\eta}{rn} \lambda_{\max} \left( \hat{\mathbf{B}}^{t\top} \mathbf{U}^{t\top} \mathbf{W}^{t+1\top} + \mathbf{W}^{t+1} \mathbf{U}^t \hat{\mathbf{B}}^t \right) + \frac{\eta^2}{(rn)^2} \lambda_{\min} \left( \mathbf{W}^{t+1} \mathbf{U}^t \mathbf{U}^{t\top} \mathbf{W}^{t+1\top} \right) \\
& \geq 1 - \frac{\eta}{rn} \lambda_{\max} \left( \hat{\mathbf{B}}^{t\top} \mathbf{U}^{t\top} \mathbf{W}^{t+1\top} + \mathbf{W}^{t+1} \mathbf{U}^t \hat{\mathbf{B}}^t \right) \tag{71}
\end{aligned}$$

where (71) holds since  $\mathbf{W}^{t+1} \mathbf{U}^t \mathbf{U}^{t\top} \mathbf{W}^{t+1\top}$  is positive semi-definite. Further,

$$\frac{\eta}{rn} \lambda_{\max} \left( \hat{\mathbf{B}}^{t\top} \mathbf{U}^{t\top} \mathbf{W}^{t+1\top} + \mathbf{W}^{t+1} \mathbf{U}^t \hat{\mathbf{B}}^t \right)$$

$$\begin{aligned}
&= \max_{\mathbf{z}: \|\mathbf{z}\|_2=1} \frac{\eta}{rn} \left( \mathbf{z}^\top \hat{\mathbf{B}}^{t\top} \mathbf{U}^{t\top} \mathbf{W}^{t+1\top} \mathbf{z} + \mathbf{z}^\top \mathbf{W}^{t+1} \mathbf{U}^t \hat{\mathbf{B}}^t \mathbf{z} \right) \\
&= \max_{\mathbf{z}: \|\mathbf{z}\|_2=1} \frac{2\eta}{rn} \mathbf{z}^\top \mathbf{W}^{t+1} \mathbf{U}^t \hat{\mathbf{B}}^t \mathbf{z} \\
&= \max_{\mathbf{z}: \|\mathbf{z}\|_2=1} \left( \frac{2\eta}{rn} \mathbf{z}^\top \mathbf{W}^{t+1} \left( \frac{1}{m} \mathcal{A}^\dagger \mathcal{A}(\mathbf{Q}^{t\top}) - \mathbf{Q}^{t\top} \right) \hat{\mathbf{B}}^t \mathbf{z} + \frac{2\eta}{rn} \mathbf{z}^\top \mathbf{W}^{t+1} \mathbf{Q}^{t\top} \hat{\mathbf{B}}^t \mathbf{z} \right) \quad (72)
\end{aligned}$$

When considering the first term, we have

$$\max_{\mathbf{z}: \|\mathbf{z}\|_2=1} \frac{2\eta}{rn} \mathbf{z}^\top \mathbf{W}^{t+1} \left( \frac{1}{m} \mathcal{A}^\dagger \mathcal{A}(\mathbf{Q}^{t\top}) - \mathbf{Q}^{t\top} \right) \hat{\mathbf{B}}^t \mathbf{z} \leq \frac{2\eta}{rn} \left\| \mathbf{W}^{t+1} \left( \frac{1}{m} \mathcal{A}^\dagger \mathcal{A}(\mathbf{Q}^{t\top}) - \mathbf{Q}^{t\top} \right) \right\|_2 \left\| \hat{\mathbf{B}}^t \right\|_2 \leq 2\eta(\delta' + \delta'') \quad (73)$$

Then we consider the second term in (72),

$$\begin{aligned}
\max_{\mathbf{z}: \|\mathbf{z}\|_2=1} \frac{2\eta}{rn} \mathbf{z}^\top \mathbf{W}^{t+1} \mathbf{Q}^{t\top} \hat{\mathbf{B}}^t \mathbf{z} &\leq \max_{\mathbf{z}: \|\mathbf{z}\|_2=1} \frac{2\eta}{rn} \mathbf{z}^\top \left( \hat{\mathbf{B}}^{t\top} \hat{\mathbf{B}}^* \mathbf{W}^* - \mathbf{F} \right) \left( \mathbf{W}^{t+1\top} \hat{\mathbf{B}}^{t\top} - \mathbf{W}^{*\top} \hat{\mathbf{B}}^{*\top} \right) \hat{\mathbf{B}}^t \mathbf{z} \\
&\quad + \max_{\mathbf{z}: \|\mathbf{z}\|_2=1} \frac{2\eta}{rn} \mathbf{z}^\top \left( \left( \hat{\mathbf{B}}^{t\top} \hat{\mathbf{B}}^* \mathbf{W}^* - \mathbf{F} \right) \left( \hat{\mathbf{S}}^{t+1\top} - \hat{\mathbf{S}}^{*\top} \right) + \mathbf{H} \mathbf{Q}^{t\top} \right) \hat{\mathbf{B}}^t \mathbf{z} \quad (74)
\end{aligned}$$

As for the first term in (74), from equation (81) in (Collins et al., 2021) we have

$$\max_{\mathbf{z}: \|\mathbf{z}\|_2=1} \frac{2\eta}{rn} \mathbf{z}^\top \left( \hat{\mathbf{B}}^{t\top} \hat{\mathbf{B}}^* \mathbf{W}^* - \mathbf{F} \right) \left( \mathbf{W}^{t+1\top} \hat{\mathbf{B}}^{t\top} - \mathbf{W}^{*\top} \hat{\mathbf{B}}^{*\top} \right) \hat{\mathbf{B}}^t \mathbf{z} \quad (75)$$

$$\leq 4\eta \frac{\delta_k}{(1 - \delta_k)^2} \bar{\sigma}_{\max,*}^2 + 2(1 + \delta)\eta \bar{\sigma}_{\max,*} \left\| \hat{\mathbf{S}}^* - \hat{\mathbf{S}}^t \right\|_2 + 2(1 + \delta)^2 \eta \left\| \hat{\mathbf{S}}^* - \hat{\mathbf{S}}^t \right\|_2^2 \quad (76)$$

As for the second term in (74),

$$\begin{aligned}
&\frac{2\eta}{rn} \left\| \left( \left( \hat{\mathbf{B}}^{t\top} \hat{\mathbf{B}}^* \mathbf{W}^* - \mathbf{F} \right) \left( \hat{\mathbf{S}}^{t+1\top} - \hat{\mathbf{S}}^{*\top} \right) + \mathbf{H} \mathbf{Q}^{t\top} \right) \hat{\mathbf{B}}^t \right\|_2 \\
&\leq \frac{2\eta}{rn} \left\| \hat{\mathbf{B}}^{t\top} \hat{\mathbf{B}}^* \mathbf{W}^* - \mathbf{F} \right\|_2 \left\| \hat{\mathbf{S}}^{t+1\top} - \hat{\mathbf{S}}^{*\top} \right\|_2 + \frac{2\eta}{rn} \left\| \mathbf{H} \mathbf{Q}^{t\top} \right\|_2 \\
&\leq 4\eta \frac{1}{\sqrt{rn}} \left\| \mathbf{W}^* \right\|_2 \frac{1}{\sqrt{rn}} \left\| \hat{\mathbf{S}}^{t+1\top} - \hat{\mathbf{S}}^{*\top} \right\|_2 + 2\eta \frac{1}{\sqrt{rn}} \left\| \mathbf{H} \right\|_2 \frac{1}{\sqrt{rn}} \left\| \mathbf{Q} \right\|_2 \quad (77)
\end{aligned}$$

$$\leq 4\eta \bar{\sigma}_{\max,*} \left\| \hat{\mathbf{S}}^* - \hat{\mathbf{S}}^{t+1} \right\|_2 + 2\eta(1 + \delta) \left\| \hat{\mathbf{S}}^* - \hat{\mathbf{S}}^t \right\|_2 \left( 2\sqrt{k} \text{dist}(\hat{\mathbf{B}}, \hat{\mathbf{B}}^*) + (1 + \delta) \left\| \hat{\mathbf{S}}^* - \hat{\mathbf{S}}^t \right\|_2 + \left\| \hat{\mathbf{S}}^* - \hat{\mathbf{S}}^{t+1} \right\|_2 \right) \quad (78)$$

$$\leq 4\eta \bar{\sigma}_{\max,*} 2\sqrt{k} + 2\eta(1 + \delta) 2\sqrt{k} \times 8\sqrt{k} \quad (79)$$

$$= 8\eta \bar{\sigma}_{\max,*} \sqrt{k} + 32(1 + \delta)\eta k \quad (80)$$

Therefore,

$$\sigma_{\min}^2(\mathbf{R}^{t+1}) \geq 1 - 2\eta(\delta' + \delta'') - 4\eta \frac{\delta_k}{(1 - \delta_k)^2} \bar{\sigma}_{\max,*}^2 - 8\eta \bar{\sigma}_{\max,*} \sqrt{k} - 32(1 + \delta)\eta k \quad (81)$$

Finally, we have

$$\text{dist}(\hat{\mathbf{B}}^{t+1}, \hat{\mathbf{B}}^*) \leq \quad (82)$$

$$\frac{(1 - \eta\sigma_{\min}^2 + \eta\delta'_k) \text{dist}(\hat{\mathbf{B}}^t, \hat{\mathbf{B}}^*) + \eta\delta'' \|\Delta\hat{\mathbf{S}}^t\|_2 + \eta(\delta''' + 6\sqrt{k}/\sqrt{rn}) \|\Delta\hat{\mathbf{S}}^{t+1}\|_2}{\sqrt{1 - 2\eta\delta'_k \text{dist} - 4\eta \frac{\delta_k}{(1-\delta_k)^2} \bar{\sigma}_{\max,*} - 2\eta(\delta_k'' + (1+\delta)\bar{\sigma}_{\max,*}) \|\Delta\hat{\mathbf{S}}^t\|_2 - 4\eta(1+\delta)^2 \|\Delta\hat{\mathbf{S}}^t\|_2^2 - 2\eta(\delta_k''' + \frac{2\bar{\sigma}_{\max,*}}{\sqrt{rn}} \|\Delta\hat{\mathbf{S}}^{t+1}\|_2) - 4\eta\sqrt{k}(1+\delta) \|\Delta\hat{\mathbf{S}}^t\|_2 \text{dist} - 2\eta(1+\delta) \|\Delta\hat{\mathbf{S}}^t\|_2 \|\Delta\hat{\mathbf{S}}^{t+1}\|_2}} \quad (83)$$

where  $\delta_k = c \frac{k^{3/2}\sqrt{\log(rn)}}{\sqrt{m}}$ ,  $\delta'_k = c_1 k \frac{\sqrt{d}}{\sqrt{rnm}}$ ,  $\delta''_k = c_2 \frac{\sqrt{kd}}{\sqrt{rnm}}$ ,  $\delta'''_k = c_3 \frac{\sqrt{kd}}{\sqrt{rnm}}$ ,  $\delta = \frac{\delta_d}{1-\delta_k}$ ,  $\delta_d = c_4 \frac{\sqrt{d\log(rn)}}{\sqrt{m}}$

## References

- Liam Collins, Hamed Hassani, Aryan Mokhtari, and Sanjay Shakkottai. Exploiting shared representations for personalized federated learning. In *International Conference on Machine Learning*, pages 2089–2099. PMLR, 2021.
- Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.

## A Proofs