# APPENDIX

# A PROOF OF THEOREM 2

Our proof of Theorem 2 consists of three steps.

- First, we show that  $\mathcal{R}$  is an absorbing region for GD. Here a set is regarded as an absorbing set if the GD sequence remains within the set after its first entrance.
- Next, we show that  $\sigma_1(J_t)$  converges to zero at a linear rate, employing an SNR argument.
- Finally, we establish the linear convergence to the global minima.

Before diving deeper, we first write down the update rules for  $U_t$  and  $J_t$ . By (9), we have

$$\boldsymbol{U}_{t+1} = \boldsymbol{U}_t + \eta \boldsymbol{\Lambda}_r \boldsymbol{U}_t - \eta \boldsymbol{U}_t \boldsymbol{X}_t^\top \boldsymbol{X}_t, \tag{18}$$

$$\boldsymbol{J}_{t+1} = \boldsymbol{J}_t + \eta \boldsymbol{\Lambda}_{\text{res}} \boldsymbol{J}_t - \eta \boldsymbol{J}_t \boldsymbol{X}_t^{\top} \boldsymbol{X}_t, \tag{19}$$

where  $\Lambda_r = \operatorname{diag}(\lambda_1, \ldots, \lambda_r)$  and  $\Lambda_{\operatorname{res}} = \operatorname{diag}(\lambda_{r+1}, \ldots, \lambda_d)$ . Note that  $\Sigma_r = \operatorname{diag}(\Lambda_r, \mathbf{0})$ .

### A.1 The GD sequence remains in $\mathcal{R}$

Lemma 12 shows that  $\mathcal{R}$  is an absorbing region for GD.

**Lemma 12** Suppose  $\eta \leq \frac{\Delta^2}{36\lambda_1^3}$  and  $\mathbf{X}_t \in \mathcal{R}$ . Then  $\mathbf{X}_{t'} \in \mathcal{R}$  for all  $t' \geq t$ .

**Proof** This lemma is proved by induction. Suppose  $X_t \in \mathcal{R}$ .

- By Lemma 13 and  $\sigma_1^2(\boldsymbol{X}_t) \leq 2\lambda_1$ , we get  $\sigma_1^2(\boldsymbol{X}_{t+1}) \leq 2\lambda_1$ .
- By Lemma 14,  $\sigma_1^2(\boldsymbol{X}_t) \leq 2\lambda_1$ , and  $\sigma_1^2(\boldsymbol{J}_t) \leq \lambda_r \Delta/2$ , we get  $\sigma_1^2(\boldsymbol{J}_{t+1}) \leq \lambda_r \Delta/2$ .
- By Lemma 15 and  $X_t \in \mathcal{R}$ , we get  $\sigma_r^2(U_{t+1}) \ge \Delta/4$  and thus  $X_{t+1} \in \mathcal{R}$ .

By induction, we conclude that  $X_{t'} \in \mathcal{R}$  for all  $t' \geq t$ .

#### A.1.1 TECHNICAL LEMMAS

In this section, we summarize technical lemmas used in the proof of Lemma 12.

Lemma 13 delineates the first category of absorbing sets for GD, denoted as

$$S_1 = \{ \boldsymbol{X} \in \mathbb{R}^{d \times r} \mid \sigma_1(\boldsymbol{X}) \le a \},\$$

valid for any  $a \in [\sqrt{\lambda_1}, 1/\sqrt{3\eta}]$ .

**Lemma 13** Suppose  $\eta \leq \frac{1}{3\lambda_1}$  and  $a \in [\sqrt{\lambda_1}, 1/\sqrt{3\eta}]$ . If  $\sigma_1(\mathbf{X}_t) \leq a$ , then  $\sigma_1(\mathbf{X}_{t'}) \leq a$ ,  $\forall t' \geq t$ .

**Proof** Lemma 16 states that if  $\sigma_1(\mathbf{X}_t) \leq 1/\sqrt{3\eta}$ , then the following inequality holds

$$\sigma_1(\boldsymbol{X}_{t+1}) \leq (1 + \eta \lambda_1 - \eta \sigma_1^2(\boldsymbol{X}_t)) \cdot \sigma_1(\boldsymbol{X}_t).$$

- If  $\sqrt{\lambda_1} \leq \sigma_1(\boldsymbol{X}_t) \leq a$ , the above inequality implies that  $\sigma_1(\boldsymbol{X}_{t+1}) \leq \sigma_1(\boldsymbol{X}_t) \leq a$ .
- If  $\sigma_1(\boldsymbol{X}_t) \leq \sqrt{\lambda_1} \leq a$ , it follows that

$$\sigma_1(\boldsymbol{X}_{t+1}) \le (1 + \eta \lambda_1 - \eta \lambda_1) \sqrt{\lambda_1} \le a$$

This uses the fact that  $g_1(s) = (1 + \eta \lambda_1 - \eta s^2)s$  is increasing on  $[0, 1/\sqrt{3\eta}]$ .

By induction, we have  $\sigma_1(\mathbf{X}_{t'}) \leq a$  for all  $t' \geq t$ .

Lemma 14 demonstrates that if  $\sigma_1(\mathbf{X}_t) \leq \sqrt{2\lambda_1}$ ,  $\sigma_1^2(\mathbf{J}_t) \leq a$ , and  $a \geq \lambda_{r+1}$ , then  $\sigma_1^2(\mathbf{J}_{t+1}) \leq a$ . Combining with Lemma 13, it implies that

$$\mathcal{S}_2 = \{ oldsymbol{X} = igg( oldsymbol{U} oldsymbol{J} igg) \in \mathbb{R}^{d imes r} \mid \sigma_1(oldsymbol{X}) \leq \sqrt{2\lambda_1}, \sigma_1^2(oldsymbol{J}) \leq a \}$$

is an absorbing set for GD, provided that  $a \ge \lambda_{r+1}$  and  $\eta \le \frac{1}{12\lambda_1}$ . Here U and J are the top r rows and the (r+1)-to-d-th rows of X respectively.

**Lemma 14** Suppose  $\eta \leq \frac{1}{12\lambda_1}$ ,  $\sigma_1^2(\boldsymbol{X}_t) \leq 2\lambda_1$ , and  $a \geq \lambda_{r+1}$ . If  $\sigma_1^2(\boldsymbol{J}_t) \leq a$ , then  $\sigma_1^2(\boldsymbol{J}_{t+1}) \leq a$ .

**Proof** By Lemma 17, we have

$$\sigma_1(\boldsymbol{J}_{t+1}) \leq (1 + \eta(\lambda_{r+1} - \sigma_1^2(\boldsymbol{J}_t))) \cdot \sigma_1(\boldsymbol{J}_t).$$

- If  $\lambda_{r+1} < \sigma_1^2(\boldsymbol{J}_t) \le a$ , then it follows that  $\sigma_1^2(\boldsymbol{J}_{t+1}) \le \sigma_1^2(\boldsymbol{J}_t) \le a$ .
- If  $\sigma_1^2(\boldsymbol{J}_t) \leq \lambda_{r+1} \leq a$ , then

$$\sigma_1^2(\boldsymbol{J}_{t+1}) \le (1 + \eta(\lambda_{r+1} - \lambda_{r+1}))^2 \lambda_{r+1} \le a.$$

This uses the observation that  $g_2(s) = (1 + \eta(\lambda_{r+1} - s^2))s$  is increasing on  $[0, 1/\sqrt{3\eta}]$ .

This concludes the proof.

Lemma 15 is the last piece needed to show that region  $\mathcal{R}$  is an absorbing set for GD.

**Lemma 15** Suppose  $\eta \leq \frac{\Delta^2}{32\lambda_1^3}$ ,  $\sigma_1(\boldsymbol{X}_t) \leq \sqrt{2\lambda_1}$ , and  $\sigma_1^2(\boldsymbol{J}_t) \leq \lambda_r - \Delta/2$ . If  $\sigma_r^2(\boldsymbol{U}_t) \geq \Delta/4$ , then  $\sigma_r^2(\boldsymbol{U}_{t+1}) \geq \Delta/4$ .

**Proof** Since  $\eta \leq \frac{1}{32\lambda_1}$  and  $\sigma_1^2(J_t) \leq \lambda_r - \Delta/2$ , by Lemma 18, we have  $\sigma_1^2(I_{L-1}) \geq (1 + n\Delta - 2n\sigma_1^2(I_L)) + \sigma_1^2(I_L) - 4n^2)^3$ 

$$\begin{aligned} \sigma_r^2(\boldsymbol{U}_{t+1}) &\geq (1 + \eta\Delta - 2\eta\sigma_r^2(\boldsymbol{U}_t)) \cdot \sigma_r^2(\boldsymbol{U}_t) - 4\eta^2\lambda_1^3. \end{aligned}$$
  
Since  $g_3(s) = (1 + \eta\Delta - 2\eta s)s$  is increasing on  $(-\infty, \frac{1}{4\eta}]$  and  $\frac{\Delta}{4} \leq \sigma_r^2(\boldsymbol{U}_t) \leq 2\lambda_1 \leq \frac{1}{4\eta}$ , we have  $\sigma_r^2(\boldsymbol{U}_{t+1}) \geq (1 + \frac{\eta\Delta}{2}) \cdot \frac{\Delta}{4} - 4\eta^2\lambda_1^3 \geq \frac{\Delta}{4}, \end{aligned}$ 

where the last inequality uses  $\eta \leq \frac{\Delta^2}{32\lambda_1^3}$ .

The following lemmas give certain singular value analysis that are used in prior lemmas and subsequent analysis. Lemma 16 establishes an upper bound for  $\sigma_1(X_{t+1})$ .

**Lemma 16** If  $\sigma_1(\mathbf{X}_t) \leq 1/\sqrt{3\eta}$ , then we have

$$\sigma_1(\boldsymbol{X}_{t+1}) \le (1 + \eta \lambda_1 - \eta \sigma_1^2(\boldsymbol{X}_t)) \cdot \sigma_1(\boldsymbol{X}_t).$$

**Proof** By the singular value inequality and (9),

$$\sigma_{1}(\boldsymbol{X}_{t+1}) \leq \sigma_{1}(\boldsymbol{X}_{t}(\boldsymbol{I}_{r} - \eta \boldsymbol{X}_{t}^{\top} \boldsymbol{X}_{t})) + \eta \sigma_{1}(\boldsymbol{\Sigma} \boldsymbol{X}_{t})$$
$$\leq \sigma_{1}(\boldsymbol{X}_{t}(\boldsymbol{I}_{r} - \eta \boldsymbol{X}_{t}^{\top} \boldsymbol{X}_{t})) + \eta \lambda_{1} \sigma_{1}(\boldsymbol{X}_{t}),$$
(20)

where we use  $\sigma_1(\Sigma) = \lambda_1$ . Observe that all r singular values of  $X_t(I_r - \eta X_t^{\top} X_t)$  are given by

$$(1 - \eta \sigma_i^2(\boldsymbol{X}_t)) \cdot \sigma_i(\boldsymbol{X}_t), \ i = 1, \dots, r,$$

since  $\eta \sigma_1^2(\mathbf{X}_t) \leq 1$ . The function  $g_4(s) = (1 - \eta s^2)s$  is increasing on  $[0, 1/\sqrt{3\eta}]$ . Hence, the fact  $0 \leq \sigma_i(\mathbf{X}_t) \leq \sigma_1(\mathbf{X}_t) \leq 1/\sqrt{3\eta}$  implies that

$$\sigma_1(\boldsymbol{X}_t(\boldsymbol{I}_r - \eta \boldsymbol{X}_t^\top \boldsymbol{X}_t)) = (1 - \eta \sigma_1^2(\boldsymbol{X}_t)) \cdot \sigma_1(\boldsymbol{X}_t).$$

Substituting this equality into (20), we conclude the proof.

Lemma 17 gives an upper bound for  $\sigma_1(J_{t+1})$ .

**Lemma 17** Suppose  $\eta \leq \frac{1}{12\lambda_1}$  and  $\sigma_1(X_t) \leq \sqrt{2\lambda_1}$ , then we have

$$\sigma_1(\boldsymbol{J}_{t+1}) \leq (1 + \eta(\lambda_{r+1} - \sigma_1^2(\boldsymbol{J}_t) - \sigma_r^2(\boldsymbol{U}_t))) \cdot \sigma_1(\boldsymbol{J}_t).$$

**Proof** The update rule (19) of  $J_{t+1}$  can be decomposed as follows:

$$\boldsymbol{J}_{t+1} = \underbrace{\frac{1}{2}\boldsymbol{J}_t - \eta \boldsymbol{J}_t \boldsymbol{J}_t^\top \boldsymbol{J}_t}_{\boldsymbol{B}} + \underbrace{(\frac{1}{4}\boldsymbol{I}_{d-r} + \eta \boldsymbol{\Lambda}_{res})\boldsymbol{J}_t}_{\boldsymbol{C}} + \underbrace{\boldsymbol{J}_t(\frac{1}{4}\boldsymbol{I}_r - \eta \boldsymbol{U}_t^\top \boldsymbol{U}_t)}_{\boldsymbol{D}}.$$

By the singular value inequality,

$$\sigma_1(\boldsymbol{J}_{t+1}) \leq \sigma_1(\boldsymbol{B}) + \sigma_1(\boldsymbol{C}) + \sigma_1(\boldsymbol{D}).$$

Observe that all singular values of B are given by

$$\sigma_i(\boldsymbol{J}_t)/2 - \eta \sigma_i^3(\boldsymbol{J}_t), \quad i = 1, \dots, d - r.$$

Since  $g_5(s) = s/2 - \eta s^3$  is increasing on  $[0, 1/\sqrt{6\eta}]$ , the condition  $\sigma_i(\mathbf{J}_t) \le \sigma_1(\mathbf{J}_t) \le \sqrt{2\lambda_1} \le 1/\sqrt{6\eta}$  implies that

$$\sigma_1(\boldsymbol{B}) = \sigma_1(\boldsymbol{J}_t)/2 - \eta \sigma_1^3(\boldsymbol{J}_t).$$

For the second term C, it follows from the singular value inequality that

$$\sigma_1(\boldsymbol{C}) \leq \sigma_1(\frac{1}{4}\boldsymbol{I}_{d-r} + \eta \boldsymbol{\Lambda}_{\text{res}})\sigma_1(\boldsymbol{J}_t) \leq (1/4 + \eta \lambda_{r+1})\sigma_1(\boldsymbol{J}_t),$$

where the second inequality uses  $\eta \sigma_1(\Lambda_{\text{res}}) \leq \eta \lambda_1 \leq 1/4$ . For the third term D, since  $\eta \sigma_1^2(U_t) \leq 2\eta \lambda_1 \leq 1/4$ , we have

$$\sigma_1(\boldsymbol{D}) \le (1/4 - \eta \sigma_r^2(\boldsymbol{U}_t)) \sigma_1(\boldsymbol{J}_t)$$

Finally, we conclude the proof by combining the analysis of B, C, and D.

Lemma 18 provides an lower bound for  $\sigma_r^2(U_{t+1})$ .

**Lemma 18** Suppose  $\eta \leq \frac{1}{32\lambda_1}$  and  $\sigma_1(\mathbf{X}_t) \leq \sqrt{2\lambda_1}$ , then we have

$$\sigma_r^2(\boldsymbol{U}_{t+1}) \ge (1 + 2\eta(\lambda_r - \sigma_1^2(\boldsymbol{J}_t) - \sigma_r^2(\boldsymbol{U}_t))) \cdot \sigma_r^2(\boldsymbol{U}_t) - 4\eta^2 \lambda_1^3.$$

**Proof** Substituting the update rule (18) of  $U_{t+1}$  into  $U_{t+1}U_{t+1}^{\top}$ , we get

$$U_{t+1}U_{t+1}^{\top} = (U_t - \eta U_t X_t^{\top} X_t + \eta \Lambda_r U_t) \cdot (U_t - \eta U_t X_t^{\top} X_t + \eta \Lambda_r U_t)^{\top}$$
  
=  $B + C - \eta^2 R_1 + \eta^2 R$ 

where

$$B = U_t (\frac{1}{2} I_r - 2\eta X_t^\top X_t) U_t^\top,$$
  

$$C = (\frac{1}{\sqrt{2}} I_r + \sqrt{2} \eta \Lambda_r) U_t U_t^\top (\frac{1}{\sqrt{2}} I_r + \sqrt{2} \eta \Lambda_r),$$
  

$$R_1 = 2\Lambda_r U_t U_t^\top \Lambda_r,$$
  

$$R = (\Lambda_r U_t - U_t X_t^\top X_t) (\Lambda_r U_t - U_t X_t^\top X_t)^\top.$$

Here B is positive semi-definite (PSD) since  $2\eta\sigma_1^2(X_t) \le 4\eta\lambda_1 \le 1/2$  and  $C, R_1, R$  are all PSD. By the eigenvalue inequality and the equivalence between eigenvalues and singular values of a PSD matrix, we have

$$\sigma_r^2(\boldsymbol{U}_{t+1}) \ge \sigma_r(\boldsymbol{B}) + \sigma_r(\boldsymbol{C}) - \eta^2 \sigma_1(\boldsymbol{R}_1) + \eta^2 \sigma_r(\boldsymbol{R})$$
  
$$\ge \sigma_r(\boldsymbol{B}) + \sigma_r(\boldsymbol{C}) - \eta^2 \sigma_1(\boldsymbol{R}_1).$$
(21)

For the first term B, we decompose it into two terms:

$$\boldsymbol{B} = \underbrace{\boldsymbol{U}_t((\frac{1}{2} - 2\eta\sigma_1^2(\boldsymbol{J}_t)) \cdot \boldsymbol{I}_r - 2\eta\boldsymbol{U}_t^\top\boldsymbol{U}_t)\boldsymbol{U}_t^\top}_{\boldsymbol{B}_1} + 2\eta \cdot \underbrace{\boldsymbol{U}_t(\sigma_1^2(\boldsymbol{J}_t) \cdot \boldsymbol{I}_r - \boldsymbol{J}_t^\top\boldsymbol{J}_t)\boldsymbol{U}_t^\top}_{\boldsymbol{B}_2}$$

The inequality  $2\eta(\sigma_1^2(J_t) + \sigma_1^2(U_t)) \le 8\eta\lambda_1 \le 1/2$  implies that  $B_1$  is PSD. Since  $B_2$  is also PSD, we have  $\sigma_r(B) \ge \sigma_r(B_1)$ . To determine  $\sigma_r(B_1)$ , we write the singular values of  $B_1$  as

$$\left(\frac{1}{2}-2\eta\sigma_1^2(\boldsymbol{J}_t)\right)\cdot\sigma_i^2(\boldsymbol{U}_t)-2\eta\sigma_i^4(\boldsymbol{U}_t),\ i=1,\ldots,r.$$

Since  $1/2 - 2\eta \sigma_1^2(\boldsymbol{J}_t) \ge 1/4$ , the function  $g_6(s) = (1/2 - 2\eta \sigma_1^2(\boldsymbol{J}_t))s - 2\eta s^2$  is increasing on  $(-\infty, \frac{1}{16\eta}]$ . Then the inequality  $\sigma_i^2(\boldsymbol{U}_t) \le \sigma_1^2(\boldsymbol{U}_t) \le 2\lambda_1 \le \frac{1}{16\eta}$  implies that

$$\sigma_r(\boldsymbol{B}_1) = (\frac{1}{2} - 2\eta(\sigma_1^2(\boldsymbol{J}_t) + \sigma_r^2(\boldsymbol{U}_t))) \cdot \sigma_r^2(\boldsymbol{U}_t).$$

For the second term C, we have

$$\sigma_r(\boldsymbol{C}) \ge \sigma_r^2(\frac{1}{\sqrt{2}}\boldsymbol{I}_r + \sqrt{2}\eta\boldsymbol{\Lambda}_r)\sigma_r^2(\boldsymbol{U}_t) \ge (\frac{1}{2} + 2\eta\lambda_r)\sigma_r^2(\boldsymbol{U}_t).$$

For the third term  $\mathbf{R}_1$ , since  $\sigma_1^2(\mathbf{X}_t) \leq 2\lambda_1$ , we have

$$\sigma_1(\boldsymbol{R}_1) \le 4\lambda_1^3$$

Finally, substituting the analysis of  $B, C, R_1$  into (21) gives the desired result.

# A.2 $\sigma_1(\mathbf{J}_t)$ converges to zero linearly via an SNR argument

Lemma 19 shows that if  $X_0 \in \mathcal{R}$ , then  $\sigma_1(J_t)$  will diminish to zero at a geometric rate. A key step of the analysis is to examine the SNR  $\frac{\sigma_r^2(U_t)}{\sigma_1^2(J_t)}$ . Our analysis extends the rank-one case in Section 2 to a general rank scenario.

**Lemma 19** Suppose  $\eta \leq \Delta^2/(32\lambda_1^3)$  and  $X_0 \in \mathcal{R}$ . Then, for all  $t \geq 0$ , we have

$$\frac{\sigma_1^2(\boldsymbol{J}_{t+1})}{\sigma_r^2(\boldsymbol{U}_{t+1})} \le (1 - \eta \Delta/3) \cdot \frac{\sigma_1^2(\boldsymbol{J}_t)}{\sigma_r^2(\boldsymbol{U}_t)}.$$

Hence,  $\sigma_1^2(\boldsymbol{J}_t) \leq 8\lambda_1^2(1-\eta\Delta/3)^t/\Delta$  for all t and  $\sigma_1^2(\boldsymbol{J}_t) < \epsilon$  after

$$T_{\boldsymbol{J}}^{\epsilon} = \mathcal{O}\left(\frac{3}{\eta\Delta}\log\frac{8\lambda_1^2}{\epsilon\Delta}\right) \quad iterations.$$

**Proof** By Lemma 12, we have  $X_t \in \mathcal{R}$  for all  $t \ge 0$ . Then by Lemma 17,

$$\begin{aligned} \sigma_1^2(\boldsymbol{J}_{t+1}) &\leq (1 + 2\eta(\lambda_{r+1} - \sigma_1^2(\boldsymbol{J}_t) - \sigma_r^2(\boldsymbol{U}_t)) + 16\eta^2\lambda_1^2) \cdot \sigma_1^2(\boldsymbol{J}_t) \\ &\leq (1 - \eta\Delta/2 + 2\eta(\lambda_r - \Delta/2 - \sigma_1^2(\boldsymbol{J}_t) - \sigma_r^2(\boldsymbol{U}_t))) \cdot \sigma_1^2(\boldsymbol{J}_t) \end{aligned}$$

where the second inequality follows from  $\eta \leq \frac{\Delta}{32\lambda_1^2}$ . By Lemma 18,

$$\sigma_r^2(\boldsymbol{U}_{t+1}) \ge (1 + \eta\Delta + 2\eta(\lambda_r - \Delta/2 - \sigma_1^2(\boldsymbol{J}_t) - \sigma_r^2(\boldsymbol{U}_t))) \cdot \sigma_r^2(\boldsymbol{U}_t) - 4\eta^2 \lambda_1^3$$
  
$$\ge (1 + \eta\Delta/2 + 2\eta(\lambda_r - \Delta/2 - \sigma_1^2(\boldsymbol{J}_t) - \sigma_r^2(\boldsymbol{U}_t))) \cdot \sigma_r^2(\boldsymbol{U}_t),$$

where we use  $\sigma_r^2(U_t) \ge \Delta/4$  and  $\eta \le \frac{\Delta^2}{32\lambda_1^3}$  in the second inequality. A combination of the above two inequalities gives that

$$\frac{\sigma_1^2(\boldsymbol{J}_{t+1})}{\sigma_r^2(\boldsymbol{U}_{t+1})} \leq \frac{1 - \eta \Delta/2 + 2\eta(\lambda_r - \Delta/2 - \sigma_1^2(\boldsymbol{J}_t) - \sigma_r^2(\boldsymbol{U}_t))}{1 + \eta \Delta/2 + 2\eta(\lambda_r - \Delta/2 - \sigma_1^2(\boldsymbol{J}_t) - \sigma_r^2(\boldsymbol{U}_t))} \cdot \frac{\sigma_1^2(\boldsymbol{J}_t)}{\sigma_r^2(\boldsymbol{U}_t)}.$$

Since the function  $g_7(s) = \frac{1-\eta\Delta/2+s}{1+\eta\Delta/2+s}$  is increasing on [-1/2, 1/2], the condition  $-1/2 \leq 2\eta(\lambda_r - \Delta/2 - \sigma_1^2(\boldsymbol{J}_t) - \sigma_r^2(\boldsymbol{U}_t)) \leq 1/2$  implies that

$$\frac{\sigma_1^2(\boldsymbol{J}_{t+1})}{\sigma_r^2(\boldsymbol{U}_{t+1})} \leq \frac{3/2 - \eta \Delta/2}{3/2 + \eta \Delta/2} \cdot \frac{\sigma_1^2(\boldsymbol{J}_t)}{\sigma_r^2(\boldsymbol{U}_t)} \leq (1 - \eta \Delta/3) \cdot \frac{\sigma_1^2(\boldsymbol{J}_t)}{\sigma_r^2(\boldsymbol{U}_t)}$$

By deduction, we have

$$\sigma_1^2(\boldsymbol{J}_t) \le (1 - \eta \Delta/3)^t \cdot \sigma_r^2(\boldsymbol{U}_t) \frac{\sigma_1^2(\boldsymbol{J}_0)}{\sigma_r^2(\boldsymbol{U}_0)} \le (1 - \eta \Delta/3)^t \cdot \frac{8\lambda_1^2}{\Delta}$$

where the second inequality follows from  $\sigma_r^2(U_t) \leq 2\lambda_1$ ,  $\sigma_1^2(J_0) \leq \lambda_1$ , and  $\sigma_r^2(U_0) \geq \Delta/4$ . Therefore, for any  $\epsilon > 0$ , it takes at most  $T_J^{\epsilon} = \mathcal{O}(\frac{3}{\eta\Delta}\log\frac{8\lambda_1^2}{\epsilon\Delta})$  iterations to have  $\sigma_1^2(J_t) \leq \epsilon$ .

### A.3 FINAL CONVERGENCE

For the convergence of  $X_t X_t^{\top}$  to  $\Sigma_r$ , It remains to show that  $U_t U_t^{\top}$  converges to  $\Lambda_r$  fast, where  $\Lambda_r = \text{diag}(\lambda_1, \dots, \lambda_r)$ . Equivalently, it suffices to show that  $\sigma_1(P_t)$  converges to zero linearly, where  $P_t = \Lambda_t - U_t U_t^{\top}$ . This is established in Lemma 20.

**Lemma 20** Suppose  $\eta \leq \Delta^2/(36\lambda_1^3)$  and  $\mathbf{X}_0 \in \mathcal{R}$ . Then, for all  $t \geq 0$ , we have

$$\sigma_1(\boldsymbol{P}_{t+1}) \le \frac{100\lambda_1^2}{\eta\Delta^2} (1 - \eta\Delta/4)^{t+1}.$$

*Hence, for any*  $\epsilon > 0$ *, it takes*  $T_{\boldsymbol{P}}^{\epsilon} = \mathcal{O}\left(\frac{4}{\eta\Delta}\log\frac{100\lambda_1^2}{\eta\Delta^2\epsilon}\right)$  *iterations to reach*  $\sigma_1(\boldsymbol{P}_t) \leq \epsilon$ *.* 

**Proof** By Lemma 12,  $X_t \in \mathcal{R}$  for all  $t \ge 0$ . Using the notation of  $P_t$ , (18) can be rewritten as

$$\boldsymbol{U}_{t+1} = \boldsymbol{U}_t + \eta \boldsymbol{P}_t \boldsymbol{U}_t - \eta \boldsymbol{U}_t \boldsymbol{J}_t^{\top} \boldsymbol{J}_t.$$

By direct calculation, we have

$$\boldsymbol{P}_{t+1} = (\boldsymbol{I}_r - \eta \boldsymbol{U}_t \boldsymbol{U}_t^\top) \boldsymbol{P}_t (\boldsymbol{I}_r - \eta \boldsymbol{U}_t \boldsymbol{U}_t^\top) - \eta^2 (\boldsymbol{P}_t \boldsymbol{U}_t \boldsymbol{U}_t^\top \boldsymbol{P}_t + \boldsymbol{U}_t \boldsymbol{U}_t^\top \boldsymbol{P}_t \boldsymbol{U}_t \boldsymbol{U}_t^\top) + \boldsymbol{R}_t,$$
ere

where

$$\boldsymbol{R}_{t} = \eta(\boldsymbol{I}_{r} + \eta\boldsymbol{P}_{t})\boldsymbol{U}_{t}\boldsymbol{J}_{t}^{\top}\boldsymbol{J}_{t}\boldsymbol{U}_{t}^{\top} + \eta\boldsymbol{U}_{t}\boldsymbol{J}_{t}^{\top}\boldsymbol{J}_{t}\boldsymbol{U}_{t}^{\top}(\boldsymbol{I}_{r} + \eta\boldsymbol{P}_{t}) - \eta^{2}\boldsymbol{U}_{t}(\boldsymbol{J}_{t}^{\top}\boldsymbol{J}_{t})^{2}\boldsymbol{U}_{t}^{\top}.$$

By the singular value inequality,

$$\sigma_1(\boldsymbol{P}_{t+1}) \leq ((1 - \eta \Delta/4)^2 + 8\eta^2 \lambda_1^2) \cdot \sigma_1(\boldsymbol{P}_t) + \sigma_1(\boldsymbol{R}_t)$$
  
$$\leq (1 - \eta \Delta/4) \cdot \sigma_1(\boldsymbol{P}_t) + \sigma_1(\boldsymbol{R}_t),$$

where we use  $\Delta/4 \leq \sigma_r^2(U_t) \leq \sigma_1^2(U_t) \leq 2\lambda_1$  in the first inequality and  $\eta \leq \frac{\Delta}{36\lambda_1^2}$  in the second inequality. For the remainder term  $R_t$ , by the singular value inequality and the condition  $\eta \leq \frac{\Delta^2}{36\lambda_1^3}$ , we have

$$\sigma_1(\boldsymbol{R}_t) \le \sigma_1^2(\boldsymbol{J}_t) \le (1 - \eta \Delta/3)^t \cdot \frac{8\lambda_1^2}{\Delta},$$

where the last inequality follows from Lemma 19. Then by deduction, we have

$$\begin{aligned} \frac{\sigma_1(\boldsymbol{P}_{t+1})}{(1-\eta\Delta/4)^{t+1}} &\leq \frac{\sigma_1(\boldsymbol{P}_t)}{(1-\eta\Delta/4)^t} + \left(\frac{1-\eta\Delta/3}{1-\eta\Delta/4}\right)^t \frac{8\lambda_1^2}{(1-\eta\Delta/4)\Delta} \\ &\leq \sigma_1(\boldsymbol{P}_0) + \sum_{i=1}^t \left(\frac{1-\eta\Delta/3}{1-\eta\Delta/4}\right)^i \frac{8\lambda_1^2}{(1-\eta\Delta/4)\Delta} \\ &\leq \sigma_1(\boldsymbol{P}_0) + \frac{96\lambda_1^2}{\eta\Delta^2} \leq \frac{100\lambda_1^2}{\eta\Delta^2}, \end{aligned}$$

where the last inequality follows from  $\sigma_1(\mathbf{P}_0) \leq 2\lambda_1$ . Therefore, it takes  $T_{\mathbf{P}}^{\epsilon} = \mathcal{O}(\frac{4}{\eta\Delta}\log\frac{100\lambda_1^2}{\eta\Delta^2\epsilon})$  iterations to achieve  $\sigma_1(\mathbf{P}_t) \leq \epsilon$ .

## A.4 PROOF OF THEOREM 2

By combining Lemma 19 and Lemma 20, we can prove Theorem 2.

## Proof Observe that

$$\|\boldsymbol{\Sigma}_r - \boldsymbol{X}_t \boldsymbol{X}_t^{\top}\|_{\mathrm{F}} \le \|\boldsymbol{P}_t\|_{\mathrm{F}} + 2\|\boldsymbol{J}_t \boldsymbol{X}_t^{\top}\|_{\mathrm{F}} \le r\sigma_1(\boldsymbol{P}_t) + 2r\sqrt{2\lambda_1}\sigma_1(\boldsymbol{J}_t), \quad \forall \boldsymbol{X}_t \in \mathcal{R},$$

where we use the fact that  $\|A\|_{\rm F} \leq r\sigma_1(A)$  for any rank-*r* matrix *A*. Let

$$T^{\epsilon} = \max\left\{T_{\boldsymbol{J}}^{\epsilon^2/(32r^2\lambda_1)}, T_{\boldsymbol{P}}^{\epsilon/(2r)}\right\}.$$

Then,  $\|\boldsymbol{\Sigma}_r - \boldsymbol{X}_t \boldsymbol{X}_t^{\top}\|_{\mathrm{F}} \leq \epsilon$  for all  $t \geq T^{\epsilon}$ . Theorem 2 follows from  $T^{\epsilon} = \mathcal{O}(\frac{6}{\eta\Delta}\log\frac{200r\lambda_1^2}{\eta\Delta^2\epsilon})$ .

# **B** ANALYSIS OF LARGE INITIALIZATION

In this section, we will prove Theorem 6 as well as the results in Section 5.1. Before delving further, we first write down the update rules of  $u_{k,t}$  and  $K_{k,t}$ . Recall that  $u_{k,t}$  and  $K_{k,t}$  are the k-th and (k+1)-to-d-th rows of  $X_t$ . The update rules are given by

$$\boldsymbol{u}_{k,t+1} = \boldsymbol{u}_{k,t} + \eta \lambda_k \boldsymbol{u}_{k,t} - \eta \boldsymbol{u}_{k,t} \boldsymbol{X}_t^{\top} \boldsymbol{X}_t, \qquad (22)$$

$$\boldsymbol{K}_{k,t+1} = \boldsymbol{K}_{k,t} + \eta \boldsymbol{\Gamma}_k \boldsymbol{K}_{k,t} - \eta \boldsymbol{K}_{k,t} \boldsymbol{X}_t^{\top} \boldsymbol{X}_t,$$
(23)

where  $\Gamma_k = \text{diag}(\lambda_{k+1}, \dots, \lambda_d)$ . We also remind readers that  $u_{k,t} \in \mathbb{R}^{1 \times r}$  is a row vector. Moreover, we let  $\Pi_{u_k,t}$  denote the projection matrix associated with  $u_{k,t}$ , that is,

$$oldsymbol{\Pi}_{oldsymbol{u}_{k,t}} = oldsymbol{u}_{k,t}^{ op} (oldsymbol{u}_{k,t} oldsymbol{u}_{k,t}^{ op})^{-1} oldsymbol{u}_{k,t} \in \mathbb{R}^{r imes r}.$$

Also, we let  $G_{k,t}$  denote the first k rows of  $X_t$ .

#### B.1 PROOFS FOR SECTION 5.1: RANK-TWO MATRIX APPROXIMATION

In this section, we collect proofs related to the rank-two matrix approximation.

## B.1.1 PROOF OF LEMMA 8

**Proof** Note that  $t_{\text{init},1} \leq T_1 + T_K$ , where

$$T_1 = \min\{t \ge 0 \mid \sigma_1^2(\boldsymbol{X}_t) \le 2\lambda_1\}$$

is the first time when  $\sigma_1^2(\mathbf{X}_t)$  is smaller than  $2\lambda_1$ , and

$$T_{\boldsymbol{K}} = \min\{t \ge 0 \mid \sigma_1^2(\boldsymbol{K}_{k,t+T_1}) \le \lambda_k - \frac{3\Delta}{4}, \forall k \le r\}.$$

To prove the lemma, it suffices to analyze  $T_1$  and  $T_K$  separately.

First, we analyze  $T_1$  as follows.

- If  $\sigma_1^2(X_0) \le 2\lambda_1$ , then  $T_1 = 0$ .
- If  $2\lambda_1 < \sigma_1^2(X_0) < 1/(3\eta)$ , then by Lemma 13,  $\sigma_1^2(X_t) \le 1/(3\eta)$  for all t. Furthermore, it follows from Lemma 16 that

$$\sigma_1(\boldsymbol{X}_{t+1}) \leq (1 + \eta \lambda_1 - \eta \sigma_1^2(\boldsymbol{X}_t)) \cdot \sigma_1(\boldsymbol{X}_t) \\ \leq (1 - \eta \lambda_1) \cdot \sigma_1(\boldsymbol{X}_t), \quad \forall t < T_1,$$

where the second inequality uses  $\sigma_1^2(\mathbf{X}_t) > 2\lambda_1$  for all  $t < T_1$ . It implies that

$$\sigma_1(\boldsymbol{X}_t) \le (1 - \eta \lambda_1)^t \cdot \sigma_1(\boldsymbol{X}_0)$$

for all  $t \leq T_1$  and

$$T_1 = \mathcal{O}\left(\frac{1}{\eta\lambda_1}\log\frac{\sigma_1(\mathbf{X}_0)}{\sqrt{2\lambda_1}}\right).$$

By Lemma 13, we have  $\sigma_1^2(\mathbf{X}_t) \leq 2\lambda_1$  for all  $t \geq T_1$ .

Next, we analyze  $T_K$  and the following quantities

$$T_{\boldsymbol{K}_k} = \min\{t \ge 0 \mid \sigma_1^2(\boldsymbol{K}_{k,t+T_1}) \le \lambda_k - \frac{3\Delta}{4}\}.$$

Recall that  $K_{k,t}$  is the (k + 1)-to-d-th rows of  $X_t$ . Then by (23), we have

$$\boldsymbol{K}_{k,t+1} = \boldsymbol{K}_{k,t} + \eta \boldsymbol{\Gamma}_{k} \boldsymbol{K}_{k,t} - \eta \boldsymbol{K}_{k,t} \boldsymbol{X}_{t}^{\top} \boldsymbol{X}_{t} \\ = \underbrace{\frac{1}{2} \boldsymbol{K}_{k,t} - \eta \boldsymbol{K}_{k,t} \boldsymbol{K}_{k,t}^{\top} \boldsymbol{K}_{k,t}}_{\boldsymbol{B}} + \underbrace{(\underbrace{\frac{1}{4} \boldsymbol{I}_{d-k} + \eta \boldsymbol{\Gamma}_{k}) \boldsymbol{K}_{k,t}}_{\boldsymbol{C}} + \underbrace{\boldsymbol{K}_{k,t} (\underbrace{\frac{1}{4} \boldsymbol{I}_{k} - \eta \boldsymbol{G}_{k,t}^{\top} \boldsymbol{G}_{k,t})}_{\boldsymbol{D}},$$

where  $\Gamma_k = \text{diag}(\lambda_{k+1}, \dots, \lambda_d)$  and  $G_{k,t} \in \mathbb{R}^{k \times r}$  is the first k rows of  $X_t$ . By the singular value inequality, we obtain

$$\sigma_1(\boldsymbol{K}_{k,t+1}) \leq \sigma_1(\boldsymbol{B}) + \sigma_1(\boldsymbol{C}) + \sigma_1(\boldsymbol{D}).$$

For the first term B, similar to Lemma 17, we can show that

$$\sigma_1(\boldsymbol{B}) = \sigma_1(\boldsymbol{K}_{k,t})/2 - \eta \sigma_1^3(\boldsymbol{K}_{k,t}), \quad \forall t \ge T_1.$$

For the second term C, by the singular value inequality,

$$\sigma_1(\boldsymbol{C}) \le (rac{1}{4} + \eta \lambda_{k+1}) \cdot \sigma_1(\boldsymbol{K}_{k,t})$$

For the third term D, since  $G_{k,t}^{\top}G_{k,t}$  is PSD and  $\eta\sigma_1^2(G_{k,t}) \leq \frac{1}{4}$  for all  $t \geq T_1$ , we have

$$\sigma_1(\boldsymbol{D}) \leq \sigma_1(\boldsymbol{K}_{k,t})/4, \quad \forall t \geq T_1.$$

Combining,

$$\sigma_1(\boldsymbol{K}_{k,t+1}) \le (1 + \eta \lambda_{k+1} - \eta \sigma_1^2(\boldsymbol{K}_{k,t})) \cdot \sigma_1(\boldsymbol{K}_{k,t}), \quad \forall t \ge T_1, \quad \forall k \le r.$$
(24)

Since  $\lambda_{k+1} \leq \lambda_k - \Delta$  for  $k \leq r$ , (24) implies that

$$\sigma_1(\mathbf{K}_{k,t+T_1+1}) \le (1 - \eta \Delta/4) \cdot \sigma_1(\mathbf{K}_{k,t+T_1}), \quad \forall t < T_{\mathbf{K}_k}, \quad \forall k \le r.$$

Hence,  $\sigma_1(\mathbf{K}_{k,t+T_1}) \leq (1 - \eta \Delta/4)^t \cdot \sigma_1(\mathbf{K}_{k,T_1})$  for all  $t \leq T_{\mathbf{K}_k}$ . In particular,

$$T_{\boldsymbol{K}_{k}} = \mathcal{O}\left(\frac{2}{\eta\Delta}\log\frac{\sigma_{1}^{2}(\boldsymbol{K}_{k,T_{1}})}{\lambda_{k} - \frac{3\Delta}{4}}\right) \text{ and } T_{\boldsymbol{K}} = \mathcal{O}\left(\frac{2}{\eta\Delta}\log\frac{8\lambda_{1}}{\Delta}\right)$$

where we use  $\sigma_1^2(\mathbf{K}_{k,T_1}) \leq 2\lambda_1$  and  $\lambda_k - \frac{3\Delta}{4} \geq \frac{\Delta}{4}$ .

Finally, similar to Lemma 13 and 14, for any  $a \ge \lambda_{k+1}$ , if  $\sigma_1^2(\mathbf{K}_{k,t+T_1}) \le a$ , then  $\sigma_1^2(\mathbf{K}_{k,t'+T_1}) \le a$  for all  $t' \ge t$ . This implies that  $\sigma_1^2(\mathbf{K}_{k,t+T_1}) \le \lambda_k - \frac{3\Delta}{4}$  for all  $t \ge T_{\mathbf{K}}$  for  $k \le r$ .

#### B.1.2 PROOF OF LEMMA 9

**Proof** This lemma is a special case of Lemma 21, where we take k = 1 and  $t_{\text{init}} = t_{\text{init},1}$ . Notice that  $G_{0,t} = 0$  and  $X_t \in S$  for all  $t \ge t_{\text{init},1}$  by Lemma 8. Thus, the conditions in Lemma 21 trivially hold. Then Lemma 9 immediately follows from Lemma 21.

#### B.1.3 PROOF OF LEMMA 10

**Proof** The lemma is a special case of Lemma 22 and Lemma 23. In Lemma 22, we take k = 1 and  $t_{\text{init}} = t_{\text{init},1} + T_{u_1}$ . In Lemma 23, we take k = 1 and  $t_{\text{init}} = t_{\text{init},1} + T_{u_1} + t^*$ .

## B.1.4 PROOF OF LEMMA 11

**Proof** This lemma is a special case of Lemma 21, where we take k = 2 and  $t_{\text{init}} = t_1 + t_1^*$ .

## B.2 PROOF OF THEOREM 6

**Proof** The first property follows from Lemma 8.

To prove the remaining properties in this theorem, we will use an inductive argument. Our induction hypotheses are listed below:

$$\begin{split} \mathrm{H}(k,1) \ \ \sigma_{1}(\boldsymbol{u}_{k,t}\boldsymbol{G}_{k-1,t}^{\top}) &\leq \sqrt{\frac{\Delta}{8}}\min\{\sigma_{1}(\boldsymbol{u}_{k,t_{\mathrm{init},k}}), \sqrt{\frac{\Delta}{2}}\} \cdot (1-\eta\Delta/6)^{t-t_{\mathrm{init},k}} \text{ for all } t \geq t_{\mathrm{init},k} \\ \mathrm{H}(k,2) \ \ T_{\boldsymbol{u}_{k}} &= \mathcal{O}\left(\frac{4}{\eta\Delta}\log\frac{\Delta}{2\sigma_{1}^{2}(\boldsymbol{u}_{k,t_{\mathrm{init},k}})}\right) \text{ and } \sigma_{1}^{2}(\boldsymbol{u}_{k,t}) \geq \frac{\Delta}{2} \text{ for all } t \geq t_{\mathrm{init},k} + T_{\boldsymbol{u}_{k}}. \\ \mathrm{H}(k,3) \ \ \sigma_{1}(\boldsymbol{u}_{k,t}\boldsymbol{K}_{k,t}^{\top}) \leq (1-\eta\Delta/6)^{t-t_{k}} \text{ for all } t \geq t_{k}. \end{split}$$

Note that H(1,1) trivially holds because  $G_{0,t} = 0$ . Then we prove H(k, 1), H(k, 2), H(k, 3), H(k + 1, 1) successively until H(r, 3).

•  $\{H(j, \cdot)\}_{j < k} + H(k, 1) \to H(k, 2)$ 

This follows from Lemma 21, where we take  $t_{init} = t_{init,k}$ .

•  $\{H(j, \cdot)\}_{j < k} + H(k, 1) + H(k, 2) \rightarrow H(k, 3)$ 

This follows from Lemma 22, where we take  $t_{\text{init}} = t_{\text{init},k} + T_{\boldsymbol{u}_k}$ .

•  $\{\mathrm{H}(j,\cdot)\}_{j\leq k} \to \mathrm{H}(k+1,1)$ 

By  $\{H(j,3)\}_{j \le k}$ ,

$$\sigma_1(\boldsymbol{u}_{k+1,t}\boldsymbol{G}_{k,t}^{\top}) \leq \sum_{j \leq k} \sigma_1(\boldsymbol{u}_{j,t}\boldsymbol{K}_{j,t}^{\top}) \leq r(1 - \eta \Delta/6)^{t-t_k},$$

for all  $t \ge t_k$ . By definition of  $t_k^*$ , we have

$$r(1-\eta\Delta/6)^{t_k^*} \leq \sqrt{\frac{\Delta}{8}} \min\{\sigma_1(\boldsymbol{u}_{k,t_k+t_k^*}), \sqrt{\frac{\Delta}{2}}\}.$$

Then H(k + 1, 1) follows from the definition  $t_{\text{init},k+1} = t_k + t_k^*$ .

By induction,  $H(k, \cdot)$  holds for all  $k \leq r$ .

For all  $t \ge t_k$ , (16) follows from Lemma 23, where  $t_{init}$  is taken as  $t_k$ .

For all  $t \ge t_{\text{init},r} + T_{\boldsymbol{u}_r}$ , we have  $\sigma_1^2(\boldsymbol{u}_{k,t}) \ge \frac{\Delta}{2}$  for all  $k \le r$ . Simultaneously,

$$\sum_{j \leq r} \sigma_1(\boldsymbol{u}_{j,t}\boldsymbol{K}_{j,t}^{\top}) \leq r(1 - \eta \Delta/6)^{t - (t_{\text{init},r} + T_{\boldsymbol{u}_r} + t^*)}$$

holds for all  $t \ge t_{\text{init},r} + T_{u_r} + t^*$ . Let  $U_t$  be the first r rows of  $X_t$ . Viewing  $U_t U_t^{\top}$  as the sum of diagonal elements and off-diagonal elements, we find that

$$\sigma_r^2(\boldsymbol{U}_t) \ge \Delta/2 - r(1 - \eta\Delta/6)^{t - (t_{\text{init},r} + T_{\boldsymbol{u}_r} + t^*)}$$

for all  $t \ge t_{\text{init},r} + T_{\boldsymbol{u}_r} + t^*$ . Hence,  $\sigma_r^2(\boldsymbol{U}_t) \ge \Delta/4$  for all  $t \ge t_{\text{init},r} + T_{\boldsymbol{u}_r} + t^* + t^{\sharp}$ , where

$$t^{\sharp} = \frac{\log(\Delta/(4r))}{\log(1 - \eta \Delta/6)}$$

This implies that  $X_t \in \mathcal{R}$  for  $t \geq t_{\text{init},r} + T_{\boldsymbol{u}_r} + t^* + t^{\sharp}$ .

The sixth property is merely an application of Theorem 2.

The seventh property immediately follows from the previous six properties.

### **B.3** TECHNICAL LEMMAS

This section collects technical lemmas that are used in previous sections. Let us recall that  $u_{k,t}$  and  $K_{k,t}$  are the k-th and the (k + 1)-to-d-th rows of  $X_t$  respectively. The projection matrix associated with  $u_{k,t}$  is denoted by

$$\boldsymbol{\Pi}_{\boldsymbol{u}_k,t} = \boldsymbol{u}_{k,t}^\top (\boldsymbol{u}_{k,t} \boldsymbol{u}_{k,t}^\top)^{-1} \boldsymbol{u}_{k,t}.$$

The first k rows of  $X_t$  are denoted by  $G_{k,t}$ , and  $G_{0,t} = 0$  by definition.

#### **B.3.1** DYNAMICS

This subsection contains lemmas describing the dynamics of the GD sequence.

Lemma 21 shows that when  $\sigma_1(\boldsymbol{u}_{k,t}\boldsymbol{G}_{k-1,t}^{\top})$  is sufficiently small, the signal term  $\sigma_1^2(\boldsymbol{u}_{k,t+1})$  can rise above  $\Delta/2$  quickly. Moreover, as shown in Lemma 21, the term  $\sigma_1^2(\boldsymbol{u}_{k,t+1})$  will remain larger than  $\Delta/2$ .

**Lemma 21** Suppose  $\eta \leq \frac{1}{12\lambda_1}$ ,  $X_t \in S$ , and for some  $t_{\text{init}} \geq 0$  and  $k \leq r$ , the condition

$$\sigma_1(\boldsymbol{u}_{k,t}\boldsymbol{G}_{k-1,t}^{\top}) \leq \sqrt{\frac{\Delta}{8}} \min\{\sigma_1(\boldsymbol{u}_{k,t_{\text{init}}}), \sqrt{\frac{\Delta}{2}}\} \cdot (1 - \eta \Delta/6)^{t-t_{\text{init}}}$$

holds for all  $t \ge t_{\text{init}}$ . Then  $\sigma_1^2(\boldsymbol{u}_{k,t}) \ge \frac{\Delta}{2}$  for all  $t \ge t_{\text{init}} + T_{\boldsymbol{u}_k}$ , where

$$T_{\boldsymbol{u}_k} = \mathcal{O}\left(\frac{4}{\eta\Delta}\log\frac{\Delta}{2\sigma_1^2(\boldsymbol{u}_{k,t_{\text{init}}})}\right).$$

In addition, for all  $t \ge t_{init}$ , we have

$$\sigma_1^2(\boldsymbol{u}_{k,t+1}) \ge (1 + 2\eta\lambda_k - \eta\Delta/4 - 2\eta\sigma_1^2(\boldsymbol{u}_{k,t}) - 2\eta\sigma_1^2(\boldsymbol{K}_{k,t}\boldsymbol{\Pi}_{\boldsymbol{u}_{k,t}}))) \cdot \sigma_1^2(\boldsymbol{u}_{k,t}),$$
(25)

where  $\Pi_{u_k,t} = u_{k,t}^{\top} (u_{k,t} u_{k,t}^{\top})^{-1} u_{k,t}$  is the projection matrix associated with  $u_{k,t}$ .

**Proof** First, we show that  $\sigma_1^2(\boldsymbol{u}_{k,t}) \geq \min\{\sigma_1^2(\boldsymbol{u}_{k,t_{\text{init}}}), \frac{\Delta}{2}\}$  for all  $t \geq t_{\text{init}}$  by induction.

This is true when  $t = t_{\text{init}}$ . Now suppose  $\sigma_1^2(\boldsymbol{u}_{k,t}) \ge \min\{\sigma_1^2(\boldsymbol{u}_{k,t_{\text{init}}}), \frac{\Delta}{2}\}$  for some  $t \ge t_{\text{init}}$ . By assumption,  $\sigma_1^2(\boldsymbol{u}_{k,t}\boldsymbol{G}_{k-1,t}^{\top}) \le \frac{\Delta}{8}\min\{\sigma_1^2(\boldsymbol{u}_{k,t_{\text{init}}}), \frac{\Delta}{2}\} \le \frac{\Delta}{8}\sigma_1^2(\boldsymbol{u}_{k,t})$ . Then by Lemma 24 and  $\boldsymbol{X}_t \in \mathcal{S}$ , we have

$$\sigma_1^2(\boldsymbol{u}_{k,t+1}) \ge (1 + 2\eta\lambda_k - 2\eta\sigma_1^2(\boldsymbol{u}_{k,t}) - 2\eta\sigma_1^2(\boldsymbol{K}_{k,t}\boldsymbol{\Pi}_{\boldsymbol{u}_k,t}))) \cdot \sigma_1^2(\boldsymbol{u}_{k,t}) - \frac{\eta\Delta}{4}\sigma_1^2(\boldsymbol{u}_{k,t}) \quad (26)$$

$$\ge (1 + 5\eta\Delta/4 - 2\eta\sigma_1^2(\boldsymbol{u}_{k,t})) - \sigma_1^2(\boldsymbol{u}_{k,t}) \quad (27)$$

$$\geq (1 + 5\eta\Delta/4 - 2\eta\sigma_1^2(\boldsymbol{u}_{k,t})) \cdot \sigma_1^2(\boldsymbol{u}_{k,t}).$$
<sup>(27)</sup>

Then we consider two cases.

- If  $\sigma_1^2(\boldsymbol{u}_{k,t}) \leq \frac{5\Delta}{8}$ , then  $\sigma_1^2(\boldsymbol{u}_{k,t+1}) \geq \sigma_1^2(\boldsymbol{u}_{k,t}) \geq \min\{\sigma_1^2(\boldsymbol{u}_{k,t_{\text{init}}}), \frac{\Delta}{2}\}$ .
- If  $\sigma_1^2(\boldsymbol{u}_{k,t}) \geq \frac{5\Delta}{8}$ , then

$$\sigma_1^2(\boldsymbol{u}_{k,t+1}) \ge (1 + \frac{5\eta\Delta}{4} - \frac{5\eta\Delta}{4}) \cdot \frac{5\Delta}{8} = \frac{5\Delta}{8}$$
$$\ge \min\{\sigma_1^2(\boldsymbol{u}_{k,t_{\text{init}}}), \frac{\Delta}{2}\},$$

where the first inequality uses the fact that  $g_8(s) = (1 + \frac{5\eta\Delta}{4} - 2\eta s)s$  is increasing on  $(-\infty, 1/4\eta]$ .

In both cases, we have  $\sigma_1^2(\boldsymbol{u}_{k,t+1}) \ge \min\{\sigma_1^2(\boldsymbol{u}_{k,\text{init}}), \frac{\Delta}{2}\}$ . The claim then follows by induction. Furthermore, the above analysis shows that inequalities 26 and 27 hold for all  $t \ge t_{\text{init}}$ , which leads to the inequality 25. Let

$$T_{\boldsymbol{u}_k} = \min\{t \ge 0 \mid \sigma_1^2(\boldsymbol{u}_{k,t+t_{\text{init}}}) \ge \frac{\Delta}{2}\}.$$

Then for  $t < T_{\boldsymbol{u}_k}$ , we have  $\sigma_1^2(\boldsymbol{u}_{k,t+t_{\text{init}}}) < \frac{\Delta}{2}$  and by inequality 27,

$$\sigma_1^2(\boldsymbol{u}_{k,t+1+t_{\text{init}}}) \ge (1+\eta\Delta/4) \cdot \sigma_1^2(\boldsymbol{u}_{k,t+t_{\text{init}}}).$$

Hence, for all  $t \leq T_{\boldsymbol{u}_k}$ , we have

$$\sigma_1^2(\boldsymbol{u}_{k,t+t_{\text{init}}}) \ge (1+\eta\Delta/4)^t \cdot \sigma_1^2(\boldsymbol{u}_{k,t_{\text{init}}}),$$

and

$$T_{\boldsymbol{u}_k} = \mathcal{O}\left(\frac{4}{\eta\Delta}\log\frac{\Delta}{2\sigma_1^2(\boldsymbol{u}_{k,t_{\text{init}}})}\right).$$

Finally, by inequality 27, we have for any  $a \leq \frac{5\Delta}{8}$ , if  $\sigma_1^2(\boldsymbol{u}_{k,t}) \geq a$ , then  $\sigma_1^2(\boldsymbol{u}_{k,t+1}) \geq a$ . Thus, by induction,  $\sigma_1^2(\boldsymbol{u}_{k,t}) \geq \frac{\Delta}{2}$  for all  $t \geq t_{\text{init}} + T_{\boldsymbol{u}_k}$ .

Lemma 22 shows that when the noise terms  $\sigma_1(\boldsymbol{u}_{j,t}\boldsymbol{K}_{j,t}^{\top})$  converge linearly to zero for all j < k and the k-th signal term  $\sigma_1^2(\boldsymbol{u}_{k,t}) \geq \frac{\Delta}{2}$ , the noise term  $\sigma_1(\boldsymbol{u}_{k,t}\boldsymbol{K}_{k,t}^{\top})$  will also converge linearly to zero. The key component is to analyze the SNR  $\frac{\sigma_1^2(\boldsymbol{u}_{k,t})}{\sigma_1(\boldsymbol{u}_{k,t}\boldsymbol{K}_{k,t}^{\top})}$ .

**Lemma 22** Suppose  $\eta \leq \frac{\Delta}{100\lambda_1^2}$ ,  $X_t \in S$ , and for some  $t_{\text{init}} \geq 0$  and  $k \leq r$ , the conditions

$$\sigma_1(\boldsymbol{u}_{j,t}\boldsymbol{K}_{j,t}^{\top}) \le (1 - \eta \Delta/6)^{t - t_{\text{init}}}, \quad \forall j < k,$$
(28)

$$\sigma_1(\boldsymbol{u}_{k,t}\boldsymbol{G}_{k-1,t}^{\top}) \le \frac{\Delta}{4} (1 - \eta \Delta/6)^{t-t_{\text{init}}},$$
(29)

$$\sigma_1^2(\boldsymbol{u}_{k,t}) \ge \frac{\Delta}{2} \tag{30}$$

hold for all  $t \ge t_{init}$ . Then we have

$$\sigma_1(\boldsymbol{u}_{k,t}\boldsymbol{K}_{k,t}^{\top}) \leq (1 - \eta \Delta/6)^{t - t_{\text{init}} - t^*}$$

for all  $t \ge t_{\text{init}} + t^*$ , where

$$t^* = \log\left(\frac{\Delta^2}{8\lambda_1^3 + 144r^2\lambda_1}\right) / \log(1 - \eta\Delta/6).$$

**Proof** By condition 29, we can apply Lemma 21 to obtain

 $\sigma_1^2(\boldsymbol{u}_{k,t+1}) \ge (1 + 2\eta\lambda_k - \eta\Delta/4 - 2\eta\sigma_1^2(\boldsymbol{u}_{k,t}) - 2\eta\sigma_1^2(\boldsymbol{K}_{k,t}\boldsymbol{\Pi}_{\boldsymbol{u}_k,t}))) \cdot \sigma_1^2(\boldsymbol{u}_{k,t})$ for all  $t \ge t_{\text{init.}}$  By Lemma 25, we have

$$\sigma_{1}(\boldsymbol{u}_{k,t+1}\boldsymbol{K}_{k,t+1}^{\top})$$

$$\leq (1 + \eta\lambda_{k} + \eta\lambda_{k+1} - 2\eta\sigma_{1}^{2}(\boldsymbol{u}_{k,t}) - 2\eta\sigma_{1}^{2}(\boldsymbol{K}_{k,t}\boldsymbol{\Pi}_{\boldsymbol{u}_{k},t}) + 25\eta^{2}\lambda_{1}^{2}) \cdot \sigma_{1}(\boldsymbol{u}_{k,t}\boldsymbol{K}_{k,t}^{\top})$$

$$+ 3\eta\sigma_{1}(\boldsymbol{u}_{k,t}\boldsymbol{G}_{k-1,t}^{\top})\sigma_{1}(\boldsymbol{K}_{k,t}\boldsymbol{G}_{k-1,t}^{\top})$$

for all  $t \ge t_{\text{init}}$ . Divide both sides of the inequality by  $\sigma_1^2(\boldsymbol{u}_{k,t+1})$ . By Lemma 26 and  $\sigma_1^2(\boldsymbol{u}_{k,t+1}) \ge \frac{\Delta}{2}$ , we have

$$\frac{\sigma_1(\boldsymbol{u}_{k,t+1}\boldsymbol{K}_{k,t+1}^{\top})}{\sigma_1^2(\boldsymbol{u}_{k,t+1})} \le (1 - \eta\Delta/6) \frac{\sigma_1(\boldsymbol{u}_{k,t}\boldsymbol{K}_{k,t}^{\top})}{\sigma_1^2(\boldsymbol{u}_{k,t})} + \frac{6\eta}{\Delta} \sigma_1(\boldsymbol{u}_{k,t}\boldsymbol{G}_{k-1,t}^{\top}) \sigma_1(\boldsymbol{K}_{k,t}\boldsymbol{G}_{k-1,t}^{\top})$$
(31)

for all  $t \ge t_{\text{init}}$ . Observe that by condition 28 and definitions of  $u_{k,t}, K_{k,t}$ , and  $G_{k-1,t}$ , we have

$$\max\{\sigma_1(\boldsymbol{u}_{k,t}\boldsymbol{G}_{k-1,t}^{\top}), \sigma_1(\boldsymbol{K}_{k,t}\boldsymbol{G}_{k-1,t}^{\top})\} \le \sum_{j < k} \sigma_1(\boldsymbol{u}_{k,t}\boldsymbol{K}_{k,t}^{\top}) \le r(1 - \eta\Delta/6)^{t-t_{\text{init}}}$$
(32)

for all  $t \ge t_{\text{init}}$ . Combining (31) and (32),

$$\frac{\sigma_1(\boldsymbol{u}_{k,t+1}\boldsymbol{K}_{k,t+1}^{\top})}{\sigma_1^2(\boldsymbol{u}_{k,t+1})} \le (1 - \eta\Delta/6) \frac{\sigma_1(\boldsymbol{u}_{k,t}\boldsymbol{K}_{k,t}^{\top})}{\sigma_1^2(\boldsymbol{u}_{k,t})} + \frac{6\eta r^2}{\Delta} (1 - \eta\Delta/6)^{2(t-t_{\text{init}})}$$

for all  $t \ge t_{\text{init}}$ . Therefore, for all  $t \ge t_{\text{init}}$ ,

$$Q_{t+1} \le (1 - \eta \Delta/6) \cdot Q_t,$$

where the quantity  $Q_t$  is given by

$$\mathbf{Q}_{t} = \frac{\sigma_{1}(\boldsymbol{u}_{k,t}\boldsymbol{K}_{k,t}^{\top})}{\sigma_{1}^{2}(\boldsymbol{u}_{k,t})} + \frac{36r^{2}}{\Delta^{2}}(1 - \eta\Delta/6)^{2(t-t_{\text{init}})-1}.$$

By induction, we have

$$\frac{\sigma_1(\boldsymbol{u}_{k,t}\boldsymbol{K}_{k,t}^{\top})}{\sigma_1^2(\boldsymbol{u}_{k,t})} \le (1 - \eta\Delta/6)^{t-t_{\text{init}}} \left(\frac{\sigma_1(\boldsymbol{u}_{k,t_{\text{init}}}\boldsymbol{K}_{k,t_{\text{init}}}^{\top})}{\sigma_1^2(\boldsymbol{u}_{k,t_{\text{init}}})} + \frac{36r^2}{\Delta^2}(1 - \eta\Delta/6)^{-1}\right).$$

This implies that

$$\sigma_1(\boldsymbol{u}_{k,t}\boldsymbol{K}_{k,t}^{\top}) \leq \frac{8\lambda_1^3 + 144r^2\lambda_1}{\Delta^2} \cdot (1 - \eta\Delta/6)^{t - t_{\text{init}}},$$

where we use  $1 - \eta \Delta/6 \ge 1/2$ ,  $\sigma_1^2(\boldsymbol{X}_t) \le 2\lambda_1$ , and  $\sigma_1^2(\boldsymbol{u}_{k,t}) \ge \frac{\Delta}{2}$  for all  $t \ge t_{\text{init}}$ . By definition of  $t^*$ , we have  $(1 - \eta \Delta/6)^{t^*} \le \frac{\Delta^2}{8\lambda_1^3 + 144r^2\lambda_1}$ . Hence, for all  $t \ge t_{\text{init}} + t^*$ , we have

$$\sigma_1(\boldsymbol{u}_{k,t}\boldsymbol{K}_{k,t}^{\top}) \le (1 - \eta \Delta/6)^{t - t_{\text{init}} - t^*}$$

which concludes the proof.

Let  $p_{k,t} = \lambda_k - \sigma_1^2(\boldsymbol{u}_{k,t})$  be the error term associated with the k-th signal. Lemma 23 shows that when the noise terms  $\sigma_1(\boldsymbol{u}_{j,t}\boldsymbol{K}_{j,t}^{\top})$  converge linearly to zero for all  $j \leq k$  and the k-th signal term  $\sigma_1^2(\boldsymbol{u}_{k,t}) \geq \frac{\Delta}{2}$ , this signal term will converge fast to  $\lambda_k$ . Specifically, the error term  $|p_{k,t}|$  will converge to zero at a linear rate. The analysis is similar to Lemma 20.

**Lemma 23** Suppose  $\eta \leq \frac{\Delta}{100\lambda_t^2}$ ,  $X_t \in S$ , and for some  $t_{\text{init}} \geq 0$  and  $k \leq r$ , the conditions

$$\sigma_1(\boldsymbol{u}_{j,t}\boldsymbol{K}_{j,t}^{\top}) \le (1 - \eta \Delta/6)^{t - t_{\text{init}}}, \quad \forall j \le k,$$
(33)

$$\sigma_1^2(\boldsymbol{u}_{k,t}) \ge \frac{\Delta}{2} \tag{34}$$

hold for all  $t \ge t_{init}$ . Then for all  $t \ge t_{init}$ , we have

$$|p_{k,t}| \le (2\lambda_1 + \frac{24r}{\eta\Delta}) \cdot (1 - \eta\Delta/8)^{t - t_{\text{init}}},$$

where  $p_{k,t} = \lambda_k - \sigma_1^2(\boldsymbol{u}_{k,t})$ .

**Proof** Using the notation of  $p_{k,t}$ , (22) can be rewritten as

$$oldsymbol{u}_{k,t+1} = oldsymbol{u}_{k,t} + \eta p_{k,t} oldsymbol{u}_{k,t} - \eta oldsymbol{u}_{k,t} oldsymbol{W}_{t,t}$$

where

$$\boldsymbol{W}_t = \boldsymbol{G}_{k-1,t}^{ op} \boldsymbol{G}_{k-1,t} + \boldsymbol{K}_{k,t}^{ op} \boldsymbol{K}_{k,t}.$$

By direction calculation, we have

$$p_{k,t+1} = p_{k,t} \cdot ((1 - \eta \sigma_1^2(\boldsymbol{u}_{k,t}))^2 + \eta^2 \lambda_k \sigma_1^2(\boldsymbol{u}_{k,t})) + \operatorname{res}_t$$

where

$$\operatorname{res}_{t} = 2\eta (1 + \eta p_{k,t}) \boldsymbol{u}_{k,t} \boldsymbol{W}_{t} \boldsymbol{u}_{k,t}^{\top} - \eta^{2} \boldsymbol{u}_{k,t} \boldsymbol{W}_{t}^{2} \boldsymbol{u}_{k,t}^{\top}.$$

By the singular value inequality, for all  $t \ge t_{init}$ , we have

$$|p_{k,t+1}| \leq |p_{k,t}| \cdot ((1 - \eta \sigma_1^2(\boldsymbol{u}_{k,t}))^2 + \eta^2 \lambda_k \sigma_1^2(\boldsymbol{u}_{k,t})) + |\operatorname{res}_t| \\ \leq |p_{k,t}| \cdot ((1 - \eta \Delta/2)^2 + 2\eta^2 \lambda_1^2) + |\operatorname{res}_t| \\ \leq |p_{k,t}| \cdot (1 - \eta \Delta/2) + |\operatorname{res}_t|,$$
(35)

where the second inequality uses  $\Delta/2 \leq \sigma_1^2(\boldsymbol{u}_{k,t}) \leq 2\lambda_1$  and the third inequality use  $\eta \leq \frac{\Delta}{100\lambda_1^2}$ . Using  $\eta \leq \frac{\Delta}{100\lambda_1^2}$  and  $\sigma_1^2(\boldsymbol{X}_t) \leq 2\lambda_1$ , we have

$$|\operatorname{res}_t| \leq \sum_{j \leq k} \sigma_1(\boldsymbol{u}_{k,t}\boldsymbol{K}_{k,t}^{\top}) \leq r(1 - \eta\Delta/6)^{t - t_{\operatorname{init}}}$$

for all  $t \ge t_{\text{init}}$ . Substituting this into (35), we obtain

$$|p_{k,t+1}| \le |p_{k,t}| \cdot (1 - \eta \Delta/2) + r(1 - \eta \Delta/6)^{t - t_{\text{init}}} \le |p_{k,t}| \cdot (1 - \eta \Delta/8) + r(1 - \eta \Delta/6)^{t - t_{\text{init}}}.$$

This implies that for all  $t \ge t_{\text{init}}$ ,

$$Q_{t+1} \le Q_t + \frac{r}{1 - \eta \Delta/8} \left(\frac{1 - \eta \Delta/6}{1 - \eta \Delta/8}\right)^{t - t_{\text{init}}},$$

where

$$Q_t = \frac{|p_{k,t}|}{(1 - \eta \Delta/8)^{t - t_{\text{init}}}}.$$

By induction, for all  $t \ge t_{\text{init}}$ , we have

$$Q_t \leq Q_{t_{\text{init}}} + \frac{r}{1 - \eta \Delta/8} \sum_{i=0}^{t-1-t_{\text{init}}} \left(\frac{1 - \eta \Delta/6}{1 - \eta \Delta/8}\right)^i$$
$$\leq |p_{k,t_{\text{init}}}| + \frac{24r}{\eta \Delta}$$
$$\leq 2\lambda_1 + \frac{24r}{\eta \Delta}.$$

Hence, for all  $t \ge t_{\text{init}}$ , we have

$$|p_{k,t}| \le (2\lambda_1 + \frac{24r}{\eta\Delta}) \cdot (1 - \eta\Delta/8)^{t - t_{\text{init}}},$$

which concludes the proof.

## **B.3.2** TECHNICAL CALCULATIONS

The following lemmas provide calculations related to an SNR argument, where the SNR refers to the ratio  $\frac{\sigma_1^2(\boldsymbol{u}_{k,t})}{\sigma_1(\boldsymbol{u}_{k,t}\boldsymbol{K}_{k,t}^{\top})}$ . Recall that  $\boldsymbol{u}_{k,t}$  is the k-th row of  $\boldsymbol{X}_t$  and  $\boldsymbol{K}_{k,t}$  represents the (k+1)-to-d-th rows of  $\boldsymbol{X}_t$ . Moreover, we recall that

$$oldsymbol{\Pi}_{oldsymbol{u}_k,t} = oldsymbol{u}_{k,t}^ op (oldsymbol{u}_{k,t} oldsymbol{u}_{k,t}^ op)^{-1} oldsymbol{u}_{k,t}$$

is the projection matrix associated with  $u_{k,t}$ .  $G_{k,t}$  collects the first k rows of  $X_t$ .

Lemma 24 provides a lower bound on  $\sigma_1^2(u_{k,t+1})$  in terms of the preceding iteration.

**Lemma 24** For any k and  $t \ge 0$ , we have

$$\sigma_1^2(\boldsymbol{u}_{k,t+1}) \ge (1 + 2\eta\lambda_k - 2\eta\sigma_1^2(\boldsymbol{u}_{k,t}) - 2\eta\sigma_1^2(\boldsymbol{K}_{k,t}\boldsymbol{\Pi}_{\boldsymbol{u}_k,t})) \cdot \sigma_1^2(\boldsymbol{u}_{k,t}) - 2\eta\sigma_1^2(\boldsymbol{u}_{k,t}\boldsymbol{G}_{k-1,t}^{\top}).$$

**Proof** Substituting (22) into  $\sigma_1^2(\boldsymbol{u}_{k,t+1})$  gives that

$$\begin{aligned} \sigma_1^2(\boldsymbol{u}_{k,t+1}) &= \boldsymbol{u}_{k,t+1} \boldsymbol{u}_{k,t+1}^{\top} \\ &= \boldsymbol{u}_{k,t} (\boldsymbol{I}_r + \eta \lambda_k \boldsymbol{I}_r - \eta \boldsymbol{X}_t^{\top} \boldsymbol{X}_t)^2 \boldsymbol{u}_{k,t}^{\top} \\ &= \boldsymbol{u}_{k,t} (\boldsymbol{I}_r + 2\eta \lambda_k \boldsymbol{I}_r - 2\eta \boldsymbol{X}_t^{\top} \boldsymbol{X}_t) \boldsymbol{u}_{k,t}^{\top} + \eta^2 \boldsymbol{R}_{k,t} \\ &= \boldsymbol{u}_{k,t} (\boldsymbol{I}_r + 2\eta \lambda_k \boldsymbol{I}_r - 2\eta \sigma_1^2(\boldsymbol{u}_{k,t}) \boldsymbol{I}_r - 2\eta \sigma_1^2(\boldsymbol{K}_{k,t} \boldsymbol{\Pi}_{\boldsymbol{u}_k,t}) \boldsymbol{I}_r - 2\eta \boldsymbol{G}_{k-1,t}^{\top} \boldsymbol{G}_{k-1,t}) \boldsymbol{u}_{k,t}^{\top} \\ &+ 2\eta \boldsymbol{R}_{k,t}' + \eta^2 \boldsymbol{R}_{k,t}, \end{aligned}$$

where  $\boldsymbol{R}_{k,t}$  and  $\boldsymbol{R}_{k,t}'$  are non-negative real numbers given by

$$\begin{aligned} \boldsymbol{R}_{k,t} &= \boldsymbol{u}_{k,t} (\lambda_k \boldsymbol{I}_r - \boldsymbol{X}_t^\top \boldsymbol{X}_t)^2 \boldsymbol{u}_{k,t}^\top, \\ \boldsymbol{R}_{k,t}' &= \boldsymbol{u}_{k,t} (\sigma_1^2 (\boldsymbol{K}_{k,t} \boldsymbol{\Pi}_{u_k,t}) \boldsymbol{I}_r - \boldsymbol{\Pi}_{\boldsymbol{u}_k,t} \boldsymbol{K}_{k,t}^\top \boldsymbol{K}_{k,t} \boldsymbol{\Pi}_{\boldsymbol{u}_k,t}) \boldsymbol{u}_{k,t}^\top. \end{aligned}$$

It then follows that

$$\sigma_1^2(\boldsymbol{u}_{k,t+1}) \ge (1 + 2\eta\lambda_k - 2\eta\sigma_1^2(\boldsymbol{u}_{k,t}) - 2\eta\sigma_1^2(\boldsymbol{K}_{k,t}\boldsymbol{\Pi}_{\boldsymbol{u}_k,t})) \cdot \sigma_1^2(\boldsymbol{u}_{k,t}) - 2\eta\sigma_1^2(\boldsymbol{u}_{k,t}\boldsymbol{G}_{k-1,t}^{\top}),$$
  
which concludes the proof.

Lemma 25 provides an upper bound on  $\sigma_1(u_{k,t+1}K_{k,t+1}^{\top})$  in terms of the preceding iteration.

**Lemma 25** Suppose  $\eta \leq \frac{1}{12\lambda_1}$  and  $\sigma_1^2(\mathbf{X}_t) \leq 2\lambda_1$ . For any  $k \leq r$ , if  $\sigma_1^2(\mathbf{u}_{k,t}) > 0$ , then we have  $\sigma_1(\mathbf{u}_{k,t+1}\mathbf{K}_{k,t+1}^{\top})$ 

$$\leq (1 + \eta \lambda_k + \eta \lambda_{k+1} - 2\eta \sigma_1^2(\boldsymbol{u}_{k,t}) - 2\eta \sigma_1^2(\boldsymbol{K}_{k,t} \boldsymbol{\Pi}_{\boldsymbol{u}_k,t}) + 25\eta^2 \lambda_1^2) \cdot \sigma_1(\boldsymbol{u}_{k,t} \boldsymbol{K}_{k,t}^{\top}) \\ + 3\eta \sigma_1(\boldsymbol{u}_{k,t} \boldsymbol{G}_{k-1,t}^{\top}) \sigma_1(\boldsymbol{K}_{k,t} \boldsymbol{G}_{k-1,t}^{\top}).$$

**Proof** Substituting (22) and (23) into  $\boldsymbol{u}_{k,t+1}\boldsymbol{K}_{k,t+1}^{\top}$  gives that

$$\begin{aligned} \boldsymbol{u}_{k,t+1}\boldsymbol{K}_{k,t+1}^{\top} &= \boldsymbol{u}_{k,t}\boldsymbol{K}_{k,t}^{\top} + \eta\lambda_{k}\boldsymbol{u}_{k,t}\boldsymbol{K}_{k,t}^{\top} + \eta\boldsymbol{u}_{k,t}\boldsymbol{K}_{k,t}^{\top}\boldsymbol{\Gamma}_{k} - 2\eta\boldsymbol{u}_{k,t}\boldsymbol{X}_{t}^{\top}\boldsymbol{X}_{t}\boldsymbol{K}_{k,t}^{\top} + \eta^{2}\boldsymbol{E}, \\ &= \boldsymbol{B} + \boldsymbol{C} - 2\eta\boldsymbol{D} + \eta^{2}\boldsymbol{E}, \end{aligned}$$

where

$$B = u_{k,t} K_{k,t}^{\top} \left( \frac{1}{2} I_{d-k} - 2\eta K_{k,t} \Pi_{u_{k,t}} K_{k,t}^{\top} \right)$$

$$C = u_{k,t} K_{k,t}^{\top} \left( \frac{1}{2} I_{d-k} + \eta \lambda_k I_{d-k} - 2\eta \sigma_1^2(u_{k,t}) I_{d-k} + \eta \Gamma_k - 2\eta K_{k,t} (I_r - \Pi_{u_k,t}) K_{k,t}^{\top} \right),$$

$$D = u_{k,t} G_{k-1,t}^{\top} G_{k-1,t} K_{k,t}^{\top},$$

$$E = \lambda_k u_{k,t} K_{k,t}^{\top} \Gamma_k - u_{k,t} X_t^{\top} X_t K_t^{\top} \Gamma_k - \lambda_k u_{k,t} X_t^{\top} X_t K_{k,t}^{\top} + u_{k,t} (X_t^{\top} X_t)^2 K_{k,t}^{\top}.$$

By the singular value inequality,

$$\sigma_1(\boldsymbol{u}_{k,t+1}\boldsymbol{K}_{k,t+1}^{\top}) \leq \sigma_1(\boldsymbol{B}) + \sigma_1(\boldsymbol{C}) + 2\eta\sigma_1(\boldsymbol{D}) + \eta^2\sigma_1(\boldsymbol{E})$$

For the first term  $\boldsymbol{B}$ , observe that

$$(\boldsymbol{u}_{k,t}\boldsymbol{u}_{k,t}^{\top})^{-1/2}\boldsymbol{B} = (\boldsymbol{u}_{k,t}\boldsymbol{u}_{k,t}^{\top})^{-1/2}\boldsymbol{u}_{k,t}\boldsymbol{K}_{k,t}^{\top} \left(\frac{1}{2}\boldsymbol{I}_{d-k} - \boldsymbol{K}_{k,t}\boldsymbol{\Pi}_{\boldsymbol{u}_{k,t}}\boldsymbol{K}_{k,t}^{\top}\right)$$
$$= \left(1/2 - \sigma_{1}^{2}((\boldsymbol{u}_{k,t}\boldsymbol{u}_{k,t}^{\top})^{-1/2}\boldsymbol{u}_{k,t}\boldsymbol{K}_{k,t}^{\top}\right) \cdot (\boldsymbol{u}_{k,t}\boldsymbol{u}_{k,t}^{\top})^{-1/2}\boldsymbol{u}_{k,t}\boldsymbol{K}_{k,t}^{\top}$$
$$= \left(1/2 - \sigma_{1}^{2}(\boldsymbol{K}_{k,t}\boldsymbol{\Pi}_{\boldsymbol{u}_{k,t}})\right) \cdot (\boldsymbol{u}_{k,t}\boldsymbol{u}_{k,t}^{\top})^{-1/2}\boldsymbol{u}_{k,t}\boldsymbol{K}_{k,t}^{\top}.$$

where we use the equality  $\sigma_1(\mathbf{K}_{k,t}\mathbf{\Pi}_{\mathbf{u}_k,t}) = \sigma_1((\mathbf{u}_{k,t}\mathbf{u}_{k,t}^{\top})^{-1/2}\mathbf{u}_{k,t}\mathbf{K}_{k,t}^{\top})$ . Thus,

$$\sigma_1(\boldsymbol{B}) = (1/2 - \sigma_1^2(\boldsymbol{K}_{k,t}\boldsymbol{\Pi}_{\boldsymbol{u}_k,t})) \cdot \sigma_1(\boldsymbol{u}_{k,t}\boldsymbol{K}_{k,t}^{\top}).$$

For the second term C, by the singular value inequality,

$$\sigma_{1}(\boldsymbol{C}) \leq \sigma_{1}\left(\frac{1}{2}\boldsymbol{I}_{d-k} + \eta\lambda_{k}\boldsymbol{I}_{d-k} - 2\eta\sigma_{1}^{2}(\boldsymbol{u}_{k,t})\boldsymbol{I}_{d-r} + \eta\Gamma_{k} - 2\eta\boldsymbol{K}_{k,t}(\boldsymbol{I}_{r} - \boldsymbol{\Pi}_{\boldsymbol{u}_{k},t})\boldsymbol{K}_{k,t}^{\top}\right) \cdot \sigma_{1}(\boldsymbol{u}_{k,t}\boldsymbol{K}_{k,t}^{\top}) \leq (1/2 + \eta\lambda_{k} - 2\eta\sigma_{1}^{2}(\boldsymbol{u}_{k,t}) + \eta\lambda_{k+1}) \cdot \sigma_{1}(\boldsymbol{u}_{k,t}\boldsymbol{K}_{k,t}^{\top}).$$

For the third term D,  $\sigma_1(D) \leq \sigma_1(u_{k,t}G_{k-1,t}^{\top})\sigma_1(K_{k,t}G_{k-1,t}^{\top})$ . For the fourth term E, since  $\sigma_1^2(X_t) \leq 2\lambda_1$ , we have

$$\sigma_1(\boldsymbol{E}) \leq 25\lambda_1^2 \sigma_1(\boldsymbol{u}_{k,t}\boldsymbol{K}_{k,t}) + 8\lambda_1 \sigma_1(\boldsymbol{u}_{k,t}\boldsymbol{G}_{k-1,t}^{\top}) \sigma_1(\boldsymbol{K}_{k,t}\boldsymbol{G}_{k-1,t}^{\top}).$$

Combining, we prove the lemma.

Lemma 26 provides an upper bound on a specific ratio, which is used in the proof of Lemma 22. It serves as a new variant of the SNR argument.

**Lemma 26** Suppose  $\eta \leq \frac{\Delta}{100\lambda_1^2}$ ,  $\sigma_1^2(\mathbf{X}_t) \leq 2\lambda_1$ , and  $\lambda_{k+1} \leq \lambda_k - \Delta$ . Let

ratio := 
$$\frac{1 + \eta \lambda_k + \eta \lambda_{k+1} - 2\eta \sigma_1^2(\boldsymbol{u}_{k,t}) - 2\eta \sigma_1^2(\boldsymbol{K}_{k,t} \boldsymbol{\Pi}_{\boldsymbol{u}_k,t}) + 25\eta^2 \lambda_1^2}{1 + 2\eta \lambda_k - \eta \Delta/4 - 2\eta \sigma_1^2(\boldsymbol{u}_{k,t}) - 2\eta \sigma_1^2(\boldsymbol{K}_{k,t} \boldsymbol{\Pi}_{\boldsymbol{u}_k,t})}$$

Then ratio  $\leq 1 - \eta \Delta/6$ .

**Proof** Since  $\eta \leq \frac{\Delta}{100\lambda_1^2}$  and  $\lambda_{k+1} < \lambda_k - \Delta$ , we have

ratio 
$$\leq \frac{1 - \eta \Delta/4 + s_0}{1 + \eta \Delta/4 + s_0}$$

where

$$s_0 = 2\eta\lambda_k - \eta\Delta/2 - 2\eta\sigma_1^2(\boldsymbol{u}_{k,t}) - 2\eta\sigma_1^2(\boldsymbol{K}_{k,t}\boldsymbol{\Pi}_{k,t}) \in [-1/2, 1/2].$$

Since the function  $g_9(s) = \frac{1 - \eta \Delta/4 + s}{1 + \eta \Delta/4 + s}$  is increasing on [-1/2, 1/2], we have

ratio 
$$\leq \frac{1 - \eta \Delta/4 + 1/2}{1 + \eta \Delta/4 + 1/2} \leq 1 - \eta \Delta/6,$$

which concludes the proof.

# C PROOF OF PROPOSITION 7

**Proof** Consider X with  $\sigma_1(X) \leq \frac{1}{\sqrt{3\eta}}$ . Let  $X_t$  be the GD sequence initialized by X. By Corollary 2 of Lee et al. (2019), we know GD sequence almost surely avoids the strict saddle points. By Zhu et al. (2021), we know all the saddle points are strict and all the local minima are global minima. Therefore, we conclude that the GD sequence converges to the global minima almost surely.

Now it remains to show that Assumption 4 must hold if the GD sequence converges to the global minima. Indeed, if we suppose Assumption 4 does not hold, then the GD sequence will converge with  $\lim_{t\to\infty} \sigma_1(u_{k,t}) = 0$  for some  $k \leq r$ . This means the GD sequence converges to a saddle point, since any stationary point with some  $u_{k,t} = 0$  ( $k \leq r$ ) is a saddle point, rather than a global minimum. This leads to the contradiction.

# D ADDITIONAL EXPERIMENTS

In this section provide additional experiments to support and illustrate our theoretical results.

#### D.1 RANK-TWO MATRIX APPROXIMATION

Our first extended experiment examines rank-two matrix approximation with varying dimension dand initial magnitude  $\varpi$ . Specifically, we will choose d from the set {1000, 2000, 4000} and choose  $\varpi$  from the set {0.001, 0.5, 2}. For each d, we set  $\Sigma = \text{diag}(\boldsymbol{a}, \boldsymbol{e})$ , where  $\boldsymbol{a} \in \mathbb{R}^r$  is a decreasing arithmetic sequence starting from 1 to 0.5 and  $\boldsymbol{e} \in \mathbb{R}^{d-r}$  is an arithmetic sequence transitioning from 0.3 to zero. Let  $\boldsymbol{X}_0 = \varpi N_0$  with the entries of  $N_0$  independently drawn from  $\mathcal{N}(0, \frac{1}{d})$ . We compute the GD sequence  $\boldsymbol{X}_t$  with a step size of 0.1 and evaluate the errors  $\|\boldsymbol{\Sigma}_r - \boldsymbol{X}_t \boldsymbol{X}_t^{\mathsf{T}}\|_F$ , where  $\boldsymbol{\Sigma}_r = \text{diag}(\boldsymbol{a}, \boldsymbol{0})$  is the best rank-r approximation to  $\boldsymbol{\Sigma}$ . The error curves of GD for different settings are displayed in Figure 2.

Figure 2 demonstrate that all the error curves exhibit the similar behaviors. The only differences lie on the first stage.

- When we use a small *π* = 0.001, the error does not rapidly change at the beginning. This is because ||*X<sub>t</sub>*||<sub>F</sub> is close to zero and the error ||Σ<sub>r</sub> − *X<sub>t</sub>X<sub>t</sub>*||<sub>F</sub> is approximately ||Σ<sub>r</sub>||<sub>F</sub>. This period of time corresponds to the second property of Theorem 6.
- When we use  $\varpi = 2$ , we find the error first drops rapidly from a large value to  $\|\Sigma_r\|$ . This corresponds to the Lemma 8 and the first property in Theorem 6.
- When we use  $\varpi = 0.5$ , the first stage nearly disappears. This means that  $T_{u_1}$  in Theorem 6 is small, especially compared with the case where  $\varpi = 0.001$ .

In addition, we want to mention that if we use  $\varpi = 10$  to initialize the algorithm and keep other settings unchanged, then the GD sequence will diverge. This serves as a supplementary to the above experimental results.

# D.2 GENERAL RANK MATRIX APPROXIMATION

Our second experiment examines general rank matrix approximation, where we fix dimension d = 1000 and vary the rank  $r \operatorname{across} \{2, 6, 10\}$ . In addition, for each setting, we examine different initial magnitudes  $\varpi \in \{0.001, 0.5, 2\}$ . Our setting for  $\Sigma$  is the same as before, that is,  $\Sigma = \operatorname{diag}(a, e)$  with  $a \in \mathbb{R}^r$  and  $e^{d-r}$  being two arithmetic sequences. We initialize GD using  $x_0 = \varpi N_0$  and we compute the GD sequence and the errors  $\|\Sigma_r - X_t X_t\|_F$ . The results are displayed in Figure 3.

As the results demonstrate, the effects of  $\varpi$  is similar to the one in Section D.1. Moreover, we observe another interesting phenomenon that may need additional explanations. Figure 3 shows that the error curve for larger rank r is smoother than the one for smaller rank r. Our explanation is that for larger rank r, the differences between successive eigenvalues are smaller. Thus, it is harder to distinguish the associated eigenvectors, and all the eigenvectors may be learned together. As a result, the error curve remains decreasing along the iterations.



Figure 2: Error curves of GD, measured by  $\|\Sigma_r - X_t X_t^{\top}\|_F$ , for rank-two matrix approximation. The columns represent different initial magnitudes  $\varpi = 0.001, 0.5, 2$ . The rows represent different dimensions d = 1000, 2000, 4000.



Figure 3: Error curves of GD, measured by  $\|\Sigma_r - X_t X_t^{\top}\|_F$ , for general rank matrix approximation. The dimension *d* is set as 1000. Different rows represent different rank *r*. Different columns represent different initial magnitudes  $\varpi$ .