
Understanding Incremental Learning of Gradient Descent: A Fine-grained Analysis of Matrix Sensing

Jikai Jin¹ Zhiyuan Li² Kaifeng Lyu³ Simon S. Du⁴ Jason D. Lee⁵

Abstract

It is believed that Gradient Descent (GD) induces an implicit bias towards good generalization in training machine learning models. This paper provides a fine-grained analysis of the dynamics of GD for the matrix sensing problem, whose goal is to recover a low-rank ground-truth matrix from near-isotropic linear measurements. It is shown that GD with small initialization behaves similarly to the greedy low-rank learning heuristics (Li et al., 2020) and follows an incremental learning procedure (Gissin et al., 2019): GD sequentially learns solutions with increasing ranks until it recovers the ground truth matrix. Compared to existing works which only analyze the first learning phase for rank-1 solutions, our result provides characterizations for the whole learning process. Moreover, besides the over-parameterized regime that many prior works focused on, our analysis of the incremental learning procedure also applies to the *under-parameterized* regime. Finally, we conduct numerical experiments to confirm our theoretical findings.

1. Introduction

Understanding the optimization and generalization properties of optimization algorithms is one of the central topics in deep learning theory (Zhang et al., 2017; Sun, 2019). It has long been a mystery why simple algorithms such as Gradient Descent (GD) or Stochastic Gradient Descent (SGD) can find global minima even for highly non-convex functions (Du et al., 2019), and why the global minima being

found can generalize well (Hardt et al., 2016).

One influential line of works provides theoretical analysis of the *implicit bias* of GD/SGD. These results typically exhibit theoretical settings where the low-loss solutions found by GD/SGD attain certain optimality conditions of a particular generalization metric, *e.g.*, the parameter norm (or the classifier margin) (Soudry et al., 2018; Gunasekar et al., 2018; Nacson et al., 2019; Lyu & Li, 2020; Ji & Telgarsky, 2020), the sharpness of local loss landscape (Blanc et al., 2020; Damian et al., 2021; Li et al., 2022a; Lyu et al., 2022).

Among these works, a line of works seek to characterize the implicit bias even when the training is away from convergence. Kalimeris et al. (2019) empirically observed that SGD learns model from simple ones, such as linear classifiers, to more complex ones. As a result, SGD always tries to fit the training data with minimal model complexity. This behavior, usually referred to as the *simplicity bias* or the *incremental learning* behavior of GD/SGD, can be a hidden mechanism of deep learning that prevents highly over-parameterized models from overfitting. In theory, Hu et al. (2020); Lyu et al. (2021); Frei et al. (2021) established that GD on two-layer nets learns linear classifiers first.

The goal of this paper is to demonstrate this simplicity bias/incremental learning in the *matrix sensing* problem, a non-convex optimization problem that arises in a wide range of real-world applications, *e.g.*, image reconstruction (Zhao et al., 2010; Peng et al., 2014), object detection (Shen & Wu, 2012; Zou et al., 2013) and array processing systems (Kalogerias & Petropulu, 2013). Moreover, this problem can serve as a standard test-bed of the implicit bias of GD/SGD in deep learning theory, since it retains many of the key phenomena in deep learning while being simpler to analyze.

Formally, the matrix sensing problem asks for recovering a ground-truth matrix $\mathbf{Z}^* \in \mathbb{R}^{d \times d}$ given m observations y_1, \dots, y_m . Each observation y_i here is resulted from a linear measurement $y_i = \langle \mathbf{A}_i, \mathbf{Z}^* \rangle$, where $\{\mathbf{A}_i\}_{1 \leq i \leq m}$ is a collection of symmetric measurement matrices. In this paper, we focus on the case where \mathbf{Z}^* is symmetric, positive semi-definite (PSD) and low-rank, *i.e.*, $\mathbf{Z}^* \succeq \mathbf{0}$ and $\text{rank}(\mathbf{Z}^*) = r_* \ll d$.

An intriguing approach to solve this problem is to use the

¹School of Mathematical Sciences, Peking University
²Department of Computer Science, Stanford University
³Department of Computer Science, Princeton University
⁴Paul G. Allen School of Computer Science and Engineering, University of Washington
⁵Department of Electrical and Computer Engineering, Princeton University. Correspondence to: Jikai Jin <jinjikai7@gmail.com>.

Proceedings of the 40th International Conference on Machine Learning, Honolulu, Hawaii, USA. PMLR 202, 2023. Copyright 2023 by the author(s).

Burer-Monteiro type decomposition $\mathbf{Z}^* = \mathbf{U}\mathbf{U}^\top$ with $\mathbf{U} \in \mathbb{R}^{d \times \hat{r}}$, and minimize the squared loss with GD:

$$\min_{\mathbf{U} \in \mathbb{R}^{d \times \hat{r}}} f(\mathbf{U}) := \frac{1}{4m} \sum_{i=1}^m (y_i - \langle \mathbf{A}_i, \mathbf{U}\mathbf{U}^\top \rangle)^2. \quad (1)$$

In the ideal case, the number of columns of \mathbf{U} , denoted as \hat{r} above, should be set to r_* . However, r_* may not be known in advance. This leads to two training regimes that are more likely to happen: the *under-parameterized* regime where $\hat{r} \leq r_*$, and the *over-parameterized* regime where $\hat{r} > r_*$.

The over-parameterized regime may lead to overfitting at first glance, but surprisingly, with small initialization, GD induces a good implicit bias towards solutions with the exact or approximate recovery of the ground truth. It was first conjectured in Gunasekar et al. (2017) that GD with small initialization finds the matrix with the minimum nuclear norm. Gunasekar et al. (2017) also proved this conjecture for a special case where all measurements are commutable. However, a series of works point out that this nuclear norm minimization view cannot capture the incremental learning behavior of GD, which, in the context of matrix sensing, refers to the phenomenon that GD tends to learn solutions with rank gradually increasing with training steps. Arora et al. (2019) exhibited this phenomenon when there is only one observation ($m = 1$). Gissin et al. (2019); Jiang et al. (2022) studied the full-observation case, where every entry of the ground truth is measured independently $f(\mathbf{U}) = \frac{1}{4d^2} \|\mathbf{Z}^* - \mathbf{U}\mathbf{U}^\top\|_{\text{F}}^2$, and GD is shown to sequentially recover singular components of the ground truth from the largest singular value to the smallest one. Li et al. (2020) provided theoretical evidence that the incremental learning behavior generally occurs for matrix sensing. They specifically provided a counterexample for Gunasekar et al. (2017)'s conjecture, where GD converges to a rank-1 solution with a very large nuclear norm. Razin & Cohen (2020) also pointed out a case where GD drives the norm to infinity while keeping the rank to be approximately 1.

Despite these progresses, theoretical understanding of the simplicity bias of GD remains limited. In fact, a vast majority of existing analyses can only show that GD is initially biased towards a rank-1 solution and cannot be generalized to higher ranks, unless additional assumptions on GD dynamics are made (Li et al., 2020, Appendix H), (Belabbas, 2020; Jacot et al., 2021; Razin et al., 2021; 2022). Recently Li et al. (2022b) shows that the implicit bias of Gunasekar et al. (2017) essentially relies on rewriting gradient flow in the space of \mathbf{U} as continuous mirror descent in the space of $\mathbf{U}\mathbf{U}^\top$, which only works a special type of reparametrized model, named ‘‘commuting parametrization’’. However, Li et al. (2022b) also shows that matrix sensing with general (non-commutable) measurements does not fall into this type.

1.1. Our Contributions

In this paper, we take a step towards understanding the generalization of GD with small initialization by firmly demonstrating the simplicity bias/incremental learning behavior in the matrix sensing setting, assuming the Restricted Isometry Property (RIP). Our main result is informally stated below. See Theorem 4.1 for the formal version.

Definition 1.1 (Best Rank- s Solution). We define the *best rank- s solution* as the unique global minimizer \mathbf{Z}_s^* of the following constrained optimization problem:

$$\begin{aligned} \min_{\mathbf{Z} \in \mathbb{R}^{d \times d}} \quad & \frac{1}{4m} \sum_{i=1}^m (y_i - \langle \mathbf{A}_i, \mathbf{Z} \rangle)^2 \\ \text{s.t.} \quad & \mathbf{Z} \succeq \mathbf{0}, \quad \text{rank}(\mathbf{Z}) \leq s. \end{aligned} \quad (2)$$

Theorem 1.2 (Informal version of Theorem 4.1). *Consider the matrix sensing problem (1) with rank- r_* ground-truth matrix \mathbf{Z}^* and measurements $\{\mathbf{A}_i\}_{i=1}^m$. Assume that the measurements satisfy the RIP condition (Definition 3.2). With small learning rate $\mu > 0$ and small initialization $\mathbf{U}_{\alpha,0} = \alpha \mathbf{U} \in \mathbb{R}^{d \times \hat{r}}$, the trajectory of $\mathbf{U}_{\alpha,t} \mathbf{U}_{\alpha,t}^\top$ during GD training enters an $o(1)$ -neighborhood of each of the best rank- s solutions in the order of $s = 1, 2, \dots, \hat{r} \wedge r_*$ when $\alpha \rightarrow 0$. Moreover, when $\hat{r} \leq r_*$, we have $\lim_{t \rightarrow \infty} \mathbf{U}_{\alpha,t} \mathbf{U}_{\alpha,t}^\top = \mathbf{Z}_{\hat{r}}^*$.*

It is shown in Li et al. (2018); Stöger & Soltanolkotabi (2021) that GD exactly recovers the ground truth under the RIP condition, but our theorem goes beyond them in a number of ways. First, in the over-parameterized regime (i.e., $\hat{r} > r_*$), it implies that GD performs *incremental learning*: learning solutions with increasing ranks until it finds the ground truth. Second, this result also shows that in the under-parameterized regime (i.e., $\hat{r} \leq r_*$), GD exhibits the same implicit bias, but finally it converges to the best low-rank solution of the matrix sensing loss. By contrast, to the best of our knowledge, only the over-parameterized setting is analyzed in existing literature.

Theorem 1.2 can also be considered as a generalization of previous results in Gissin et al. (2019); Jiang et al. (2022) which show that $\mathbf{U}_{\alpha,t} \mathbf{U}_{\alpha,t}^\top$ passes by the best low-rank solutions one by one in the full observation case of matrix sensing $f(\mathbf{U}) = \frac{1}{4d^2} \|\mathbf{Z}^* - \mathbf{U}\mathbf{U}^\top\|_{\text{F}}^2$. However, our setting has two major challenges which significantly complicate our analysis. First, since our setting only gives partial measurements, the decomposition of signal and error terms in Gissin et al. (2019); Jiang et al. (2022) cannot be applied. Instead, we adopt a different approach which is motivated by Stöger & Soltanolkotabi (2021). Second, it is well-known that the optimal rank- s solution of matrix factorization is \mathbf{X}_s (defined in Section 3), but little is known for \mathbf{Z}_s^* . In Section 4.1 we analyze the landscape of (2), establishing the uniqueness of \mathbf{Z}_s^* and local landscape properties under the

RIP condition. We find that when $U_{\alpha,t}U_{\alpha,t}^\top \approx Z_s^*$, GD follows an approximate low-rank trajectory, so that it behaves similarly to GD in the under-parameterized regime. Using our landscape results, we can finally prove [Theorem 1.2](#).

Organization. We review additional related works in [Section 2](#). In [Section 3](#), we provide an overview of necessary background and notations. We then present our main results in [Section 4](#) with proof sketch where we also state some key lemmas that are used in the proof, including [Lemma 4.3](#) and some landscape results. In [Section 5](#) we present a trajectory analysis of GD and prove [Lemma 4.3](#). Experimental results are presented in [Section 6](#) which verify our theoretical findings. Finally, in [Section 7](#), we summarize our main contributions and discuss some promising future directions. Complete proofs of all results are given in the Appendix.

2. Related work

Low-rank matrix recovery. The goal of low-rank matrix recovery is to recover an unknown low-rank matrix from a number of (possibly noisy) measurements. Examples include matrix sensing ([Recht et al., 2010](#)), matrix completion ([Candès & Recht, 2009](#); [Candès & Plan, 2010](#)) and robust PCA ([Xu et al., 2010](#); [Candès et al., 2011](#)). [Fornasier et al. \(2011\)](#); [Ngo & Saad \(2012\)](#); [Wei et al. \(2016\)](#); [Tong et al. \(2021\)](#) study efficient optimization algorithms with convergence guarantees. Interested readers can refer to [Davenport & Romberg \(2016\)](#) for an overview of this topic.

Simplicity bias/incremental learning of GD. Besides the works mentioned in the introduction, there are many other works studying the simplicity bias/incremental learning of GD on tensor factorization ([Razin et al., 2021](#); [2022](#)), deep linear networks ([Gidel et al., 2019](#)), two-layer nets with orthogonal inputs ([Boursier et al., 2022](#)).

Landscape analysis of non-convex low-rank problems. The strict saddle property ([Ge et al., 2016](#); [2015](#); [Lee et al., 2016](#)) was established for non-convex low-rank problems in a unified framework by [Ge et al. \(2017\)](#). [Tu et al. \(2016\)](#) proved a local PL property for matrix sensing with exact parameterization (i.e. the rank of parameterization and ground-truth matrix are the same). The optimization geometry of general objective function with Burer-Monteiro type factorization is studied in [Zhu et al. \(2018\)](#); [Li et al. \(2019\)](#); [Zhu et al. \(2021\)](#). We provide a comprehensive analysis in this regime for matrix factorization as well as matrix sensing that improves over their results.

3. Preliminaries

In this section, we first list the notations used in this paper, and then provide details of our theoretical setup and necessary preliminary results.

3.1. Notations

We write $\min\{a, b\}$ as $a \wedge b$ for short. For any matrix A , we use $\|A\|_F$ to denote the Frobenius norm of A , use $\|A\|$ to denote the spectral norm $\|A\|_2$, and use $\sigma_{\min}(A)$ to denote the smallest singular value of A . We use the following notation for Singular Value Decomposition (SVD):

Definition 3.1 (Singular Value Decomposition). For any matrix $A \in \mathbb{R}^{d_1 \times d_2}$ of rank r , we use $A = V_A \Sigma_A W_A^\top$ to denote a Singular Value Decomposition (SVD) of A , where $V_A \in \mathbb{R}^{d_1 \times r}$, $W_A \in \mathbb{R}^{d_2 \times r}$ satisfy $V_A^\top V_A = I$, $W_A^\top W_A = I$, and $\Sigma_A \in \mathbb{R}^{r \times r}$ is diagonal.

For the matrix sensing problem [\(1\)](#), we write the ground-truth matrix as $Z^* = X X^\top$ for some $X = [v_1, v_2, \dots, v_{r_*}] \in \mathbb{R}^{d \times r_*}$ with orthogonal columns from an orthogonal basis $\{v_i : i \in [d]\}$ of \mathbb{R}^d . We denote the singular values of X as $\sigma_1, \sigma_2, \dots, \sigma_{r_*}$, then the singular values of Z^* are $\sigma_1^2, \sigma_2^2, \dots, \sigma_{r_*}^2$. We set $\sigma_{r_*+1} := 0$ for convenience. For simplicity, we only consider the case where Z^* has distinct singular values, i.e., $\sigma_1^2 > \sigma_2^2 > \dots > \sigma_{r_*}^2 > 0$.

We use $\kappa_* := \frac{\sigma_1^2}{\min_{1 \leq s \leq r_*} \{\sigma_s^2 - \sigma_{s+1}^2\}}$ to quantify the degeneracy of the singular values of Z^* . We also use the notation $X_s = [v_1, v_2, \dots, v_s]$ for the matrix consisting of the first s columns of X and $X_s^\perp = [v_{s+1}, \dots, v_d]$. Following [Definition 3.1](#), we let $V_{X_s^\perp} = \left[\frac{v_{s+1}}{\|v_{s+1}\|}, \dots, \frac{v_d}{\|v_d\|} \right]$. Note that the best rank- s solution Z_s^* ([Definition 1.1](#)) does not equal $X_s X_s^\top$ in general.

We write the results of the measurements $\{A_i\}_{i=1}^m$ as a linear mapping $\mathcal{A} : \mathbb{R}^{d \times d} \mapsto \mathbb{R}^m$, where $[\mathcal{A}(Z)]_i = \frac{1}{\sqrt{m}} \langle A_i, Z \rangle$ for all $1 \leq i \leq m$. We use $\mathcal{A}^* : \mathbb{R}^m \rightarrow \mathbb{R}^{d \times d}$, $\mathcal{A}^*(w) = \frac{1}{\sqrt{m}} \sum_{i=1}^m w_i A_i$ to denote the adjoint operator of \mathcal{A} . Our loss function [\(1\)](#) can then be written as $f(U) = \frac{1}{4} \|\mathcal{A}(Z^* - U U^\top)\|_2^2$. The gradient is given by $\nabla f(U) = \mathcal{A}^*(y - \mathcal{A}(U U^\top)) U = \mathcal{A}^* \mathcal{A}(X X^\top - U U^\top) U$.

In this paper, we consider GD with learning rate $\mu > 0$ starting from U_0 . The update rule is

$$U_{t+1} = U_t - \mu \nabla f(U_t) = (I + \mu M_t) U_t, \quad (3)$$

where $M_t := \mathcal{A}^* \mathcal{A}(X X^\top - U_t U_t^\top)$. We specifically focus on GD with *small initialization*: letting $U_0 = \alpha \bar{U}$ for some matrix $\bar{U} \in \mathbb{R}^{d \times \hat{r}}$ with $\|\bar{U}\| = 1$, we are interested in the trajectory of GD when $\alpha \rightarrow 0$. Sometimes we write U_t as $U_{\alpha,t}$ to highlight the dependence of the trajectory on α .

3.2. Assumptions

For our theoretical analysis of the matrix sensing problem, we make the following standard assumption in the matrix sensing literature:

Definition 3.2 (Restricted Isometry Property). We say that a measurement operator \mathcal{A} satisfies the (δ, r) -RIP condition if $(1 - \delta)\|\mathbf{Z}\|_F^2 \leq \|\mathcal{A}(\mathbf{Z})\|_2^2 \leq (1 + \delta)\|\mathbf{Z}\|_F^2$ for all matrices $\mathbf{Z} \in \mathbb{R}^{d \times d}$ with $\text{rank}(\mathbf{Z}) \leq r$.

Assumption 3.3. The measurement operator \mathcal{A} satisfies the $(2r_* + 1, \delta)$ -RIP property, where $r_* = \text{rank}(\mathbf{Z}^*)$ and $\delta \leq 10^{-12} \kappa_*^{-4.5} r_*^{-1}$.

The RIP condition is the key to ensure the ground truth to be recoverable with partial observations. An important consequence of RIP is that it guarantees $\mathcal{A}^* \mathcal{A}(\mathbf{Z}) = \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{Z} \rangle \mathbf{A}_i \approx \mathbf{Z}$ when \mathbf{Z} is low-rank. This is made rigorous in the following proposition.

Proposition 3.4. (Stöger & Soltanolkotabi, 2021, Lemma 7.3) Suppose that \mathcal{A} satisfies (r, δ) -RIP with $r \geq 2$, then for all symmetric \mathbf{Z} , we have $\|(\mathcal{A}^* \mathcal{A} - \mathbf{I})\mathbf{Z}\|_2 \leq \delta \|\mathbf{Z}\|_*$, where $\|\cdot\|_*$ is the nuclear norm. Moreover, if $\text{rank}(\mathbf{Z}) \leq r - 1$, then $\|(\mathcal{A}^* \mathcal{A} - \mathbf{I})\mathbf{Z}\|_2 \leq \sqrt{r} \delta \|\mathbf{Z}\|$.

We need the following regularity condition on initialization.

Assumption 3.5. For all $1 \leq s \leq \hat{r} \wedge r_*$, $\sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \bar{\mathbf{U}}) \geq \rho$ for some constant $\rho > 0$, where $\mathbf{V}_{\mathbf{X}_s}$ is defined as Definition 3.1.

The following proposition implies that Assumption 3.5 is satisfied with high probability with a Gaussian initialization.

Proposition 3.6. Suppose that all entries of $\bar{\mathbf{U}} \in \mathbb{R}^{d \times \hat{r}}$ are independently drawn from $\mathcal{N}(0, \frac{1}{\hat{r}})$ and $\rho = \frac{\epsilon \sqrt{\hat{r} - \sqrt{\hat{r} \wedge r_* - 1}}}{\sqrt{\hat{r}}} \geq \frac{\epsilon}{2r_*}$, then $\sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \bar{\mathbf{U}}) \geq \rho$ holds for all $1 \leq s \leq \hat{r} \wedge r_*$ with probability at least $1 - \hat{r}(C\epsilon + e^{-c\hat{r}})$, where $c, C > 0$ are universal constants.

Lastly, we make the following assumption on the step size.

Assumption 3.7. The step size $\mu \leq 10^{-4} \delta \|\mathbf{X}\|^{-2}$.

3.3. Procrustes Distance

Our analysis uses the notion of Procrustes distance defined as in Goodall (1991); Tu et al. (2016).

Definition 3.8 (Procrustes Distance). The Procrustes distance between two matrices $\mathbf{U}_1, \mathbf{U}_2 \in \mathbb{R}^{d \times s}$ ($d, s > 0$) is defined as the optimal value of the classic *orthogonal Procrustes problem*:

$$\text{dist}(\mathbf{U}_1, \mathbf{U}_2) = \min_{\mathbf{R} \in \mathbb{R}^{s \times s}; \mathbf{R}^\top \mathbf{R} = \mathbf{I}} \|\mathbf{U}_1 - \mathbf{U}_2 \mathbf{R}\|_F. \quad (4)$$

We note that the Procrustes distance is well-defined because the set of $s \times s$ orthogonal matrices is compact and thus the continuous function $\|\mathbf{U}_1 - \mathbf{U}_2 \mathbf{R}\|_F$ in \mathbf{R} can attain its minimum. The Procrustes distance is a pseudometric, *i.e.*, it is symmetric and satisfies the triangle inequality.

The following lemma is borrowed from Tu et al. (2016), which connects the Procrustes distance between \mathbf{U}_1 and \mathbf{U}_2 with the distance between $\mathbf{U}_1 \mathbf{U}_1^\top$ and $\mathbf{U}_2 \mathbf{U}_2^\top$.

Lemma 3.9 (Tu et al. 2016, Lemma 5.4). For any two matrices $\mathbf{U}_1, \mathbf{U}_2 \in \mathbb{R}^{d \times r}$, we have $\|\mathbf{U}_1 \mathbf{U}_1^\top - \mathbf{U}_2 \mathbf{U}_2^\top\|_F \geq (2\sqrt{2} - 2)^{1/2} \cdot \sigma_r(\mathbf{U}_1) \cdot \text{dist}(\mathbf{U}_1, \mathbf{U}_2)$.

4. Main results

In this section, we present our main theorems and their proof sketches, following the theoretical setup in Section 3. Full proofs can be found in Appendices C to F.

Theorem 4.1. Under Assumptions 3.3, 3.5 and 3.7, consider GD (3) with initialization $\mathbf{U}_{\alpha,0} = \alpha \bar{\mathbf{U}}$ for solving the matrix sensing problem (1). There exist universal constants c, M , constant $C = C(\mathbf{X}, \bar{\mathbf{U}})$ and a sequence of time points $T_\alpha^1 < T_\alpha^2 < \dots < T_\alpha^{\hat{r} \wedge r_*}$ such that for all $1 \leq s \leq \hat{r} \wedge r_*$, the following holds when α is sufficiently small:

$$\|\mathbf{U}_{\alpha, T_\alpha^s} \mathbf{U}_{\alpha, T_\alpha^s}^\top - \mathbf{Z}_s^*\|_F \leq C \alpha^{\frac{1}{M \kappa_*^2}}, \quad (5)$$

where we recall that \mathbf{Z}_s^* is the best rank- s solution defined in Definition 1.1. Moreover, GD follows an incremental learning procedure: we have $\lim_{\alpha \rightarrow 0} \max_{1 \leq t \leq T_\alpha^s} \sigma_{s+1}(\mathbf{U}_{\alpha,t}) = 0$ for all $1 \leq s \leq \hat{r} \wedge r_*$, where $\sigma_i(\mathbf{A})$ denotes the i -th largest singular value of a matrix \mathbf{A} .

It is guaranteed that \mathbf{Z}_s^* is unique for all $1 \leq s \leq \hat{r} \wedge r_*$ under our assumptions (see Lemma 4.4). In short, Theorem 4.1 states that GD with small initialization discovers the best rank- s solution ($s = 1, 2, \dots, \hat{r} \wedge r_*$) sequentially. In particular, when $s = r_*$, the best rank- s solution is exactly the ground truth $\mathbf{X} \mathbf{X}^\top$. Hence with over-parameterization ($\hat{r} \geq r_*$), GD can discover the ground truth.

At a high level, our result characterizes the complete learning dynamics of GD and reveals an incremental learning mechanism, *i.e.*, GD starts from learning simple solutions and then gradually increases the complexity of search space until it finds the ground truth.

In the under-parameterized setting, we can further establish the following convergence result:

Theorem 4.2 (Convergence in the under-parameterized regime). Suppose that $\hat{r} \leq r^*$, then there exists a constant $\bar{\alpha} > 0$ such that when $\alpha < \bar{\alpha}$, we have $\lim_{t \rightarrow +\infty} \mathbf{U}_{\alpha,t} \mathbf{U}_{\alpha,t}^\top = \mathbf{Z}_{\hat{r}}^*$.

4.1. Key lemmas

In this section, we present some key lemmas for proving our main results. First, we can show that with small initialization, GD can get into a small neighborhood of \mathbf{Z}_s^* .

Lemma 4.3. Under Assumptions 3.3 and 3.5, there exists $\hat{T}_\alpha^s > 0$ for all $\alpha > 0$ and $1 \leq s \leq \hat{r} \wedge r_*$ such

that $\lim_{\alpha \rightarrow 0} \max_{1 \leq t \leq \hat{T}_\alpha^s} \sigma_{s+1}(\mathbf{U}_{\alpha,t}) = 0$. Furthermore, it holds that $\left\| \mathbf{U}_{\hat{T}_\alpha^s} \mathbf{U}_{\hat{T}_\alpha^s}^\top - \mathbf{Z}_s^* \right\|_F = \mathcal{O}(\kappa_*^3 r_* \delta \|\mathbf{X}\|^2)$.

The full proof can be found in [Appendix C](#). Motivated by [Stöger & Soltanolkotabi \(2021\)](#), we consider the following decomposition of \mathbf{U}_t :

$$\mathbf{U}_t = \mathbf{U}_t \mathbf{W}_t \mathbf{W}_t^\top + \mathbf{U}_t \mathbf{W}_{t,\perp} \mathbf{W}_{t,\perp}^\top, \quad (6)$$

where $\mathbf{W}_t := \mathbf{W}_{\mathbf{V}_{\mathbf{X}_s}^\top} \mathbf{U}_t \in \mathbb{R}^{\hat{r} \times s}$ is the matrix consisting of the right singular vectors of $\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t$ ([Definition 3.1](#)) and $\mathbf{W}_{t,\perp} \in \mathbb{R}^{\hat{r} \times (\hat{r}-s)}$ is any orthogonal complement of \mathbf{W}_t , i.e., $\mathbf{W}_t \mathbf{W}_t^\top + \mathbf{W}_{t,\perp} \mathbf{W}_{t,\perp}^\top = \mathbf{I}$. The dependence of $\mathbf{W}_t, \mathbf{W}_{t,\perp}$ on s is omitted but will be clear from the context.

We will refer to the term $\mathbf{U}_t \mathbf{W}_t \mathbf{W}_t^\top$ as the *parallel component* and $\mathbf{U}_t \mathbf{W}_{t,\perp} \mathbf{W}_{t,\perp}^\top$ as the *orthogonal component*. The idea is to show that the parallel component grows quickly until it gets close to the best rank- s solution at some time \hat{T}_α^s (namely $\mathbf{U}_t \mathbf{W}_t \mathbf{W}_t^\top \mathbf{U}_t^\top \approx \mathbf{Z}_s^*$ when $t = \hat{T}_\alpha^s$). Meanwhile, the orthogonal term grows exponentially slower and stays $o(1)$ before \hat{T}_α^s . See [Section 5](#) for a detailed proof sketch.

[Lemma 4.3](#) shows that $\mathbf{U}_t \mathbf{U}_t^\top$ would enter a neighborhood of \mathbf{Z}_s^* with constant radius. However, there is still a gap between [Lemma 4.3](#) and [Theorem 4.1](#), since the latter states that $\mathbf{U}_t \mathbf{U}_t^\top$ would actually get $o(1)$ -close to \mathbf{Z}_s^* .

To proceed, we define the *under-parameterized* matrix sensing loss f_s for every $1 \leq s \leq r_*$:

$$f_s(\mathbf{U}) = \frac{1}{4} \left\| \mathcal{A}(\mathbf{Z}^* - \mathbf{U} \mathbf{U}^\top) \right\|_2^2, \quad \mathbf{U} \in \mathbb{R}^{d \times s}. \quad (7)$$

While the function we are minimizing is f (defined in [\(1\)](#)) rather than f_s , [Lemma 4.3](#) suggests that for $t \leq \hat{T}_\alpha^s$, \mathbf{U}_t is always approximately rank- s , so that we use a low-rank approximation for $\mathbf{U}_{\hat{T}_\alpha^s}$ and associate the dynamics locally with the GD dynamics of f_s . We will elaborate on how this is done in [Section 4.2](#).

When $\text{dist}(\mathbf{U}_1, \mathbf{U}_2) = 0$, it can be easily shown that $f_s(\mathbf{U}_1) = f_s(\mathbf{U}_2)$ since f_s is invariant to orthogonal transformations. Moreover, we note that the global minimizer of f_s is unique up to orthogonal transformations.

Lemma 4.4. *Under [Assumption 3.3](#), if $\mathbf{U}_s^* \in \mathbb{R}^{d \times s}$ is a global minimizer of f_s , then the set of global minimizers $\arg \min f_s$ is equal to $\{\mathbf{U}_s^* \mathbf{R} : \mathbf{R} \in \mathbb{R}^{s \times s}, \mathbf{R}^\top \mathbf{R} = \mathbf{I}\}$.*

Around the global minimizers, we show that f_s satisfies the Restricted Secant Inequality (RSI) which is useful for optimization analysis.

Definition 4.5. For any $\mathbf{U} \in \mathbb{R}^{d \times s}$, we use $\Pi_s(\mathbf{U})$ to denote the set of *closest* global minimizers of f_s to \mathbf{U} , namely $\Pi_s(\mathbf{U}) = \arg \min \{\|\mathbf{U} - \mathbf{U}_s^*\|_F : \mathbf{U}_s^* \in \arg \min f_s\}$.

Lemma 4.6 (Restricted Secant Inequality). *Under [Assumption 3.3](#), if a matrix $\mathbf{U} \in \mathbb{R}^{d \times s}$ satisfies $\|\mathbf{U} - \mathbf{U}_s^*\|_F \leq 10^{-2} \kappa_*^{-1} \|\mathbf{X}\|$ for some $\mathbf{U}_s^* \in \Pi_s(\mathbf{U})$, then we have*

$$\langle \nabla f_s(\mathbf{U}), \mathbf{U} - \mathbf{U}_s^* \rangle \geq 0.1 \kappa_*^{-1} \|\mathbf{X}\|^2 \|\mathbf{U} - \mathbf{U}_s^*\|_F^2. \quad (8)$$

Remark 4.7. In general, a function $g : \mathbb{R}^n \mapsto \mathbb{R}$ satisfies the RSI condition if for some $\mu > 0$, $\langle \nabla g(x), x - \pi(x) \rangle \geq \mu \|x - \pi(x)\|^2$ holds for all x , where $\pi(x)$ is a projection of x onto $\arg \min g$. This condition can be used to prove linear convergence of GD ([Zhang & Yin, 2013](#)), but it is weaker than strong convexity and stronger than Polyak-Łojasiewicz(PL) condition ([Karimi et al., 2016](#)).

We end this subsection with the following lemma which says that all global minimizers of the f_s must be close to \mathbf{X}_s under the procrustes distance, which is used in the proof sketch of [Theorem 4.1](#) in [Section 4](#).

Lemma 4.8. *Under [Assumption 3.3](#), we have $\text{dist}(\mathbf{U}_s^*, \mathbf{X}_s) \leq 40\delta\kappa_* \|\mathbf{X}\|_F$ for any global minimizer \mathbf{U}_s^* of f_s . Moreover, $\|\mathbf{Z}_s^* - \mathbf{X}_s \mathbf{X}_s^\top\|_F \leq 160\delta\kappa_* \sqrt{r_*} \|\mathbf{X}\|^2$.*

Corollary 4.9. *Under [Assumption 3.3](#), we have $\sigma_{\min}(\mathbf{U}_s^*) \geq \frac{1}{2} \sigma_{\min}(\mathbf{X}_s) = \frac{1}{2} \sigma_s \geq \frac{1}{2} \kappa_*^{-\frac{1}{2}} \|\mathbf{X}\|$.*

The full proofs for [Lemmas 4.6](#) and [4.8](#) and [Corollary 4.9](#) can be found in [Appendix E](#).

4.2. Proof outline

Based on the key lemmas, here we provide the outlines of the proofs for our main theorems and defer the details to [Appendix F](#). We first prove [Theorem 4.2](#) which can be directly derived by combining the lemmas in [Section 4.1](#).

Proof Sketch of [Theorem 4.2](#). For any global minimizer $\mathbf{U}_{\hat{r}}^*$ of [\(1\)](#), we have

$$\begin{aligned} & \text{dist}(\mathbf{U}_{\hat{r}}^*, \mathbf{U}_{\alpha, \hat{T}_\alpha^{\hat{r}}}) \\ & \leq (2\sqrt{2} - 2)^{-1/2} \sigma_{\min}^{-1}(\mathbf{U}_{\hat{r}}^*) \left\| \mathbf{U}_{\alpha, \hat{T}_\alpha^{\hat{r}}} \mathbf{U}_{\alpha, \hat{T}_\alpha^{\hat{r}}}^\top - \mathbf{Z}_{\hat{r}}^* \right\|_F \\ & \leq \mathcal{O}(\kappa_*^{\frac{1}{2}} \|\mathbf{X}\|^{-1}) \cdot \mathcal{O}(\kappa_*^3 r_* \delta \|\mathbf{X}\|^2) \\ & = \mathcal{O}(\kappa_*^{3.5} r_* \delta \|\mathbf{X}\|), \end{aligned}$$

where the first inequality is due to [Lemma 3.9](#) and the second one is due to [Corollary 4.9](#) and [Lemma 4.3](#).

[Assumption 3.3](#) and [Lemma 4.6](#) then imply that $\mathbf{U}_{\alpha, \hat{T}_\alpha^{\hat{r}}}$ lies in the small neighborhood of the set of global minimizers of $f = f_{\hat{r}}$, in which the RSI holds. Following a standard non-convex optimization analysis ([Karimi et al., 2016](#)), we can show that GD converges linearly to $\arg \min f_{\hat{r}}$ (in the Procrustes distance), which yields the conclusion. \square

Now we turn to prove [Theorem 4.1](#). While f is not necessarily local RSI, we use a low-rank approximation for \mathbf{U}_t

and associate the dynamics in this neighborhood with the GD dynamics of f_s .

Proof sketch of Theorem 4.1. Recall that by Lemma 4.3, $U_{\alpha, \hat{T}_\alpha^s}$ is approximately rank- s . So there must exist a matrix $\bar{U}_{\alpha, 0} \in \mathbb{R}^{d \times \hat{r}}$ with $\text{rank}(\bar{U}_{\alpha, 0}) \leq s$ such that

$$\bar{U}_{\alpha, 0} \bar{U}_{\alpha, 0}^\top - U_{\alpha, T_\alpha^s} U_{\alpha, T_\alpha^s}^\top = o(1) \quad \text{as } \alpha \rightarrow 0. \quad (9)$$

Indeed, we can let $\bar{U}_{\alpha, t}$ be the parallel component of U_{α, T_α^s} because the orthogonal component stays $o(1)$ (see the discussions following (6) and Corollary 5.5 for details).

Let $\{\bar{U}_{\alpha, t}\}_{t \geq 0}$ be the trajectory of GD with step size μ , initialized at $\bar{U}_{\alpha, 0}$. Since the gradient of the objective function f is locally Lipschitz, the solution obtained by the two GD trajectories $\{\bar{U}_{\alpha, t}\}_{t \geq 0}$ and $\{U_{\alpha, \hat{T}_\alpha^s + t}\}_{t \geq 0}$ will remain $o(1)$ -close for at least constantly many steps. Indeed we can show that they will keep $o(1)$ close for some $\bar{t}_\alpha = \omega(1)$ steps, i.e., for all $t \in [0, \bar{t}_\alpha]$,

$$\bar{U}_{\alpha, t} \bar{U}_{\alpha, t}^\top - U_{\alpha, \hat{T}_\alpha^s + t} U_{\alpha, \hat{T}_\alpha^s + t}^\top = o(1). \quad (10)$$

From (3) it is evident that GD initialized at $\bar{U}_{\alpha, 0}$ actually lives in the space of matrices with $\text{rank} \leq s$. Indeed we can identify its dynamics with another GD on f_s (defined in (7)). Concretely, let $\hat{U}_{\alpha, 0} \in \mathbb{R}^{d \times s}$ be a matrix so that $\hat{U}_{\alpha, 0} \hat{U}_{\alpha, 0}^\top = \bar{U}_{\alpha, 0} \bar{U}_{\alpha, 0}^\top$, and let $\{\hat{U}_{\alpha, t}\}_{t \geq 0}$ be the trajectory of GD that optimizes f_s with step size μ starting from $\hat{U}_{\alpha, 0}$. Then we have $\hat{U}_{\alpha, t} \hat{U}_{\alpha, t}^\top = \bar{U}_{\alpha, t} \bar{U}_{\alpha, t}^\top$ for all $t \geq 0$.

We can now apply our landscape results for f_s to analyze the GD trajectory $\{\hat{U}_{\alpha, t}\}_{t \geq 0}$. By (9) and Lemma 4.3, we have $\|\hat{U}_{\alpha, 0} \hat{U}_{\alpha, 0}^\top - Z^*\|_F = \mathcal{O}(\kappa_*^3 r_* \delta \|\mathbf{X}\|^2)$, so using a similar argument as in the proof sketch of Theorem 4.2, Corollary 4.9 and Lemma 3.9 imply that the initialization $\hat{U}_{\alpha, 0}$ is within an $\mathcal{O}(\kappa_*^{3.5} r_* \delta \|\mathbf{X}\|^2)$ neighborhood of the set of global minimizers of $f_s(\mathbf{U})$. From Assumption 3.3 and Lemma 4.6 we know that $f_s(\mathbf{U})$ satisfies a local RSI condition in this neighborhood, so following standard non-convex optimization analysis (Karimi et al., 2016), we can show that $\{\hat{U}_{\alpha, t}\}_{t \geq 0}$ converges linearly to its set of global minimizers in the Procrustes distance. We need to choose a time t such that (10) remains true while this linear convergence process takes place for sufficiently many steps. This is possible since $\bar{t}_\alpha = \omega(1)$; indeed we can show that there always exists some $t = t_\alpha^s \leq \bar{t}_\alpha$ such that both $\|\hat{U}_{\alpha, t} \hat{U}_{\alpha, t}^\top - U_{\alpha, \hat{T}_\alpha^s + t} U_{\alpha, \hat{T}_\alpha^s + t}^\top\|_F$ and $\|\hat{U}_{\alpha, t} - U_s^*\|_F$ are bounded by $\mathcal{O}(\alpha^{\frac{1}{M r_*^2}})$. Hence $\|U_{\alpha, t} U_{\alpha, t}^\top - Z^*\|_F = \mathcal{O}(\alpha^{\frac{1}{M r_*^2}})$ when $t = T_\alpha^s := \hat{T}_\alpha^s + t_\alpha^s$.

For $1 \leq s < \hat{r} \wedge r_*$ and $t \leq t_\alpha^s$, since (10) holds and $\text{rank}(\hat{U}_{\alpha, t}) \leq s$, we have $\max_{1 \leq t \leq T_\alpha^s} \sigma_{s+1}(U_{\alpha, t}) \rightarrow 0$ as

$\alpha \rightarrow 0$. Finally, by Lemma 4.8 and Assumption 3.3 we have $\|Z_{s+1}^* - \mathbf{X}_{s+1} \mathbf{X}_{s+1}^\top\| = \mathcal{O}(\delta \kappa_* \sqrt{r_*}) = \mathcal{O}(\kappa_*^{-1} \|\mathbf{X}\|^2)$, so $\sigma_{s+1}(Z_{s+1}^*) \gtrsim \sigma_{s+1}^2$. Therefore, $U_{\alpha, t} U_{\alpha, t}^\top$ cannot be close to Z_{s+1}^* when $t \leq T_\alpha^s$, so we must have $T_\alpha^{s+1} > T_\alpha^s$. This completes the proof of Theorem 4.1. \square

5. Proof sketch of Lemma 4.3

In this section, we outline the proof sketch of Lemma 4.3. We divide the GD dynamics into three phases and characterize the dynamics separately. Proof details for these three phases can be found in Appendices C.1, C.2 and C.4.

5.1. The spectral phase

Starting from a small initialization, GD initially behaves similarly to power iteration since $U_{t+1} = (\mathbf{I} + \mu M_t) U_t \approx (\mathbf{I} + \mu M) U_t$, where $M := \mathcal{A}^* \mathcal{A}(\mathbf{X} \mathbf{X}^\top)$ is a symmetric matrix. Let $M = \sum_{k=1}^d \hat{\sigma}_k^2 \hat{v}_k \hat{v}_k^\top$ be the eigendecomposition of M , where $\hat{\sigma}_1 \geq \hat{\sigma}_2 \geq \dots \geq \hat{\sigma}_d \geq 0$. Using our assumption on δ (Assumption 3.3), we can show that $|\sigma_i - \hat{\sigma}_i|, 1 \leq i \leq s$ are sufficiently small so that $\hat{\sigma}_i$'s are positive and well-separated. Then we have

$$\begin{aligned} U_T &\approx (\mathbf{I} + \mu M)^T U_0 = \sum_{i=1}^d (1 + \mu \hat{\sigma}_i^2)^T \hat{v}_i \hat{v}_i^\top U_0 \\ &\approx \sum_{i=1}^s (1 + \mu \hat{\sigma}_i^2)^T \hat{v}_i \hat{v}_i^\top U_0, \end{aligned} \quad (11)$$

where the last step holds because $(1 + \mu \hat{\sigma}_s)^T \gg (1 + \mu \hat{\sigma}_{s+1})^T$. In other words, we can expect that there is an exponential separation between the parallel and orthogonal component of U_T . Formally, we can prove the following property at the end of the spectral phase:

Lemma 5.1 (Lemma C.4, simplified version). *Suppose that Assumptions 3.3, 3.5 and 3.7 hold. Then there exist positive constants $C_i = C_i(\mathbf{X}, \bar{\mathbf{U}})$, $\gamma_i = \gamma_i(\mathbf{X}, \bar{\mathbf{U}})$, $i = 2, 3$ independent of α such that $\gamma_2 < \gamma_3$ and the following inequalities hold for $t = T_\alpha^{\text{sp}} = \mathcal{O}\left(\frac{\log \alpha^{-1}}{\log(1 + \mu \|\mathbf{X}\|^2)}\right)$ when α is sufficiently small:*

$$\begin{aligned} \|U_t\| &\leq \|\mathbf{X}\|, \quad \sigma_{\min}(U_t W_t) \geq C_2 \cdot \alpha^{\gamma_2}, \\ \|U_t W_{t, \perp}\| &\leq C_3 \cdot \alpha^{\gamma_3}, \text{ and } \|V_{\mathbf{X}_s, \perp}^\top V_{U_t W_t}\| \leq 200\delta. \end{aligned}$$

5.2. The parallel improvement phase

For small α , we have $\sigma_{\min}(U_t W_t) \gg \|U_t W_{t, \perp}\|$ by the end of the spectral phase. When (11) no longer holds, we enter a new phase which we call the *parallel improvement phase*. In this phase, the ratio $\frac{\sigma_{\min}(U_t W_t)}{\|U_t W_{t, \perp}\|}$ grows exponentially in t , until the former reaches a constant scale. Formally, let $T_{\alpha, s}^{\text{pi}} =$

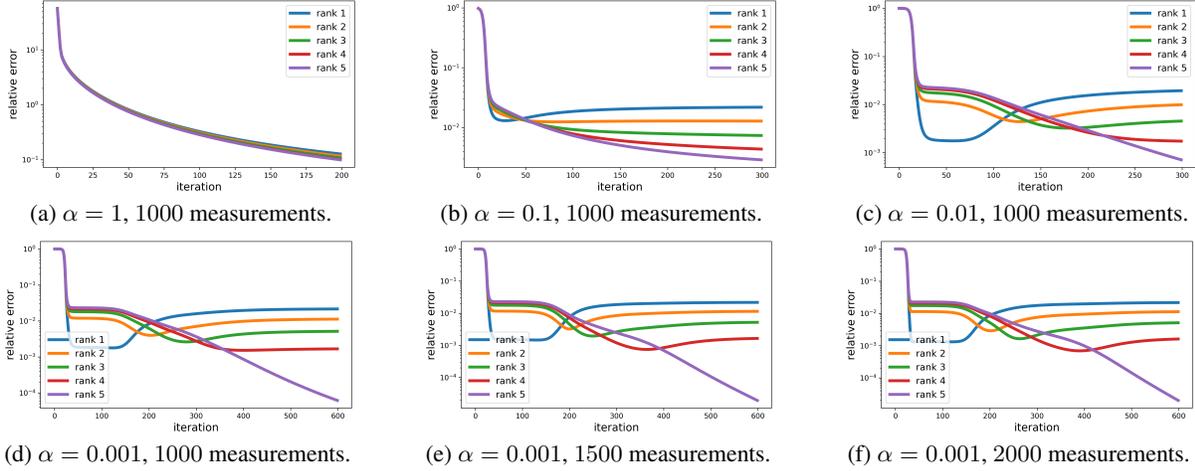


Figure 1. The evolution of relative error against the best solution of different ranks over time.

$\min \{t \geq 0 : \sigma_{\min}^2(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{\alpha,t+1}) > 0.3\kappa_*^{-1} \|\mathbf{X}\|^2\}$, then we can prove the following lemma via induction.

Lemma 5.2. *Suppose that Assumptions 3.3, 3.5 and 3.7 hold and let $c_3 = 10^4 \kappa_* r_*^{\frac{1}{2}} \delta$. Then for sufficiently small α , the following inequalities hold when $T_{\alpha,s}^{\text{SP}} \leq t < T_{\alpha,s}^{\text{PI}}$:*

$$\begin{aligned} \sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1}) &\geq \sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \mathbf{W}_t) \\ &\geq (1 + 0.5\mu(\sigma_s^2 + \sigma_{s+1}^2)) \sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t), \end{aligned} \quad (13a)$$

$$\begin{aligned} \|\mathbf{U}_{t+1} \mathbf{W}_{t+1,\perp}\| &\leq (1 + \mu(0.4\sigma_s^2 + 0.6\sigma_{s+1}^2)) \|\mathbf{U}_t \mathbf{W}_{t,\perp}\|, \end{aligned} \quad (13b)$$

$$\|\mathbf{V}_{\mathbf{X}_s,\perp}^\top \mathbf{V}_{\mathbf{U}_{t+1} \mathbf{W}_{t+1}}\| \leq c_3, \quad (13c)$$

$$\text{rank}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1}) = \text{rank}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \mathbf{W}_t) = s. \quad (13d)$$

We can immediately deduce from Lemmas 5.1 and 5.2 that the orthogonal term $\|\mathbf{U}_t \mathbf{W}_{t,\perp}\|$ remains $o(1)$ by the end of the parallel improvement phase:

Corollary 5.3 (Lemma C.11, simplified version). *Under the conditions of Lemma 5.2, when α is sufficiently small we have $\|\mathbf{U}_{T_{\alpha,s}^{\text{PI}}} \mathbf{W}_{T_{\alpha,s}^{\text{PI}}}\| \leq C_5 \cdot \alpha^{\frac{1}{4\kappa_*}}$ for some constant $C_5 = C_5(\mathbf{X}, \mathbf{U})$.*

5.3. The refinement phase

After $\sigma_{\min}(\mathbf{U}_t \mathbf{W}_t)$ grows to constant scale, we enter the *refinement phase* for which we show that $\|\mathbf{X}_s \mathbf{X}_s^\top - \mathbf{U}_t \mathbf{U}_t^\top\|_F$ keeps decreasing until it is $\mathcal{O}(\delta \kappa_*^3 r_* \|\mathbf{X}\|^2)$. Formally, let $\tau = \kappa_*^{-1} \|\mathbf{X}\|^2$ and $T_{\alpha,s}^{\text{ft}} = T_{\alpha,s}^{\text{PI}} - \frac{\log(10^{-2} \|\mathbf{X}\|^{-2} \kappa_*^{-1} c_3^{-1})}{\log(1 - \frac{1}{2} \mu \tau)} > T_{\alpha,s}^{\text{PI}}$ where c_3 is defined in Lemma 5.2, then the following lemma holds.

Lemma 5.4. *Suppose that $T_{\alpha,s}^{\text{PI}} \leq t \leq T_{\alpha,s}^{\text{ft}}$ and all the*

conditions in Lemma 5.2 hold, then we have

$$\begin{aligned} &\|\mathbf{V}_{\mathbf{X}_s}^\top (\mathbf{X} \mathbf{X}^\top - \mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top)\|_F \\ &\leq \left(1 - \frac{1}{2} \mu \tau\right) \|\mathbf{V}_{\mathbf{X}_s}^\top (\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top)\|_F \\ &\quad + 20\mu \|\mathbf{X}\|^4 (\delta + 5c_3) + 2000\mu^2 \sqrt{r_*} \|\mathbf{X}\|^6. \end{aligned}$$

Moreover, we have $\|\mathbf{U}_{t+1} \mathbf{W}_{t+1,\perp}\| \leq (1 + \sigma_s^2) \|\mathbf{U}_t \mathbf{W}_{t,\perp}\|$ and $\|\mathbf{V}_{\mathbf{X}_s,\perp} \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}\| \leq c_3$.

Using Lemma 5.4, we arrive at the following result:

Corollary 5.5. *For sufficiently small α , at $t = T_{\alpha,s}^{\text{ft}}$ we have $\|\mathbf{V}_{\mathbf{X}_s}^\top (\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top)\|_F \leq 80\delta \kappa_*^3 r_* \|\mathbf{X}\|^2$ and $\|\mathbf{U}_t \mathbf{W}_{t,\perp}\| = o(1)$ ($\alpha \rightarrow 0$).*

Concluding the proof of Lemma 4.3. At $t = T_{\alpha,s}^{\text{ft}}$, we have

$$\begin{aligned} &\|\mathbf{X}_s \mathbf{X}_s^\top - \mathbf{U}_t \mathbf{U}_t^\top\|_F \\ &\leq \|(\mathbf{X}_s \mathbf{X}_s^\top - \mathbf{U}_t \mathbf{U}_t^\top) \mathbf{V}_{\mathbf{X}_s} \mathbf{V}_{\mathbf{X}_s}^\top\|_F \end{aligned} \quad (14a)$$

$$\begin{aligned} &+ \|\mathbf{U}_t \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_s,\perp} \mathbf{V}_{\mathbf{X}_s,\perp}^\top\|_F \\ &\leq \|(\mathbf{X}_s \mathbf{X}_s^\top - \mathbf{U}_t \mathbf{U}_t^\top) \mathbf{V}_{\mathbf{X}_s} \mathbf{V}_{\mathbf{X}_s}^\top\|_F \end{aligned} \quad (14b)$$

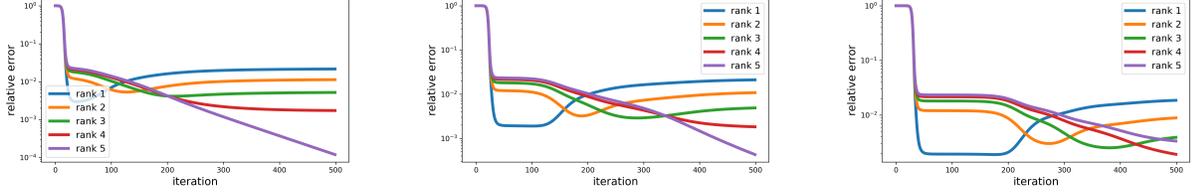
$$\begin{aligned} &+ \|\mathbf{V}_{\mathbf{X}_s,\perp}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_s,\perp}\|_F \\ &\leq \|\mathbf{V}_{\mathbf{X}_s}^\top (\mathbf{X}_s \mathbf{X}_s^\top - \mathbf{U}_t \mathbf{U}_t^\top)\|_F \\ &+ \sqrt{r_*} \|\mathbf{V}_{\mathbf{X}_s,\perp}^\top \mathbf{U}_t \mathbf{W}_t\|^2 + \sqrt{d} \|\mathbf{V}_{\mathbf{X}_s,\perp}^\top \mathbf{U}_t \mathbf{W}_{t,\perp}\|^2 \end{aligned} \quad (14c)$$

$$\begin{aligned} &\leq \|\mathbf{V}_{\mathbf{X}_s}^\top (\mathbf{X}_s \mathbf{X}_s^\top - \mathbf{U}_t \mathbf{U}_t^\top)\|_F \\ &+ 9\sqrt{r_*} \|\mathbf{X}\|^2 \|\mathbf{V}_{\mathbf{X}_s,\perp} \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}\|^2 + \sqrt{d} \|\mathbf{U}_t \mathbf{W}_{t,\perp}\|^2 \end{aligned} \quad (14d)$$

$$= \mathcal{O}(\delta \kappa_*^3 r_* \|\mathbf{X}\|^2 + \|\mathbf{X}\|^2 c_3^2 \sqrt{r_*}) + o(1) \quad (14e)$$

$$= \mathcal{O}(\delta \kappa_*^3 r_* \|\mathbf{X}\|^2), \quad (14f)$$

where (14c) uses $\|\mathbf{A}\|_F \leq \sqrt{\text{rank}(\mathbf{A})} \|\mathbf{A}\|$, (14d) uses $\|\mathbf{U}_t\| \leq 3\|\mathbf{X}\|$, (14e) follows from Lemma 5.4 and Corollary 5.5 and the last step follows from $c_3 = 10^4 \kappa_* \sqrt{r_*} \delta$ and Assumption 3.3.


 (a) $\alpha = 0.1$, 1000 measurements.

 (b) $\alpha = 0.01$, 1000 measurements.

 (c) $\alpha = 0.001$, 1000 measurements.

 Figure 2. The evolution of the loss and relative error against best solution of different ranks in the exact-parameterized case $\hat{r} = r_* = 5$.

By Lemma 4.8, the best rank- s solution is close to the matrix factorization minimizer *i.e.* $\|\mathbf{Z}_s^* - \mathbf{X}_s \mathbf{X}_s^\top\|_F = \mathcal{O}(\delta \kappa_* \sqrt{r_*} \|\mathbf{X}\|^2)$. We thus obtain that $\|\mathbf{Z}_s^* - \mathbf{U}_t \mathbf{U}_t^\top\|_F = \mathcal{O}(\delta \kappa_*^3 r_* \|\mathbf{X}\|^2)$. Finally, since $\text{rank}(\mathbf{U}_t \mathbf{W}_t) \leq s$ (recall the decomposition (6)), we have $\sigma_{s+1}(\mathbf{U}_t) \leq \|\mathbf{U}_t \mathbf{W}_{t,\perp}\| = o(1)$. The conclusion follows.

6. Experiments

In this section, we perform numerical experiments to empirically verify our theoretical findings.

Experimental setup. We consider the matrix sensing problem (1) with $d = 50$, $r_* = 5$, $\alpha \in \{1, 0.1, 0.01, 0.001\}$, $m \in \{1000, 2000, 5000\}$. We will consider different choices for \hat{r} in the experiments. The ground truth $\mathbf{Z}^* = \mathbf{X} \mathbf{X}^\top$ is generated such that the entries of \mathbf{X} are i.i.d. standard Gaussian variables. We use the same ground truth throughout our experiments.

For $i = 1, 2, \dots, m$, all entries of the measurement $\mathbf{A}_i \in \mathbb{R}^{d \times d}$ are chosen i.i.d. from the standard Gaussian $\mathcal{N}(0, 1)$. When $m \gtrsim dr_* \delta^{-2}$, this set of measurements satisfies the RIP with high probability (Recht et al., 2010, Theorem 4.2).

We solve the problem (1) via running GD for $T = 10^4$ iterations starting with small initialization with scale α . Specifically, we choose $\mathbf{U}_0 = \alpha \bar{\mathbf{U}}$ where the entries of $\bar{\mathbf{U}} \in \mathbb{R}^{d \times \hat{r}}$ are drawn i.i.d. from $\mathcal{N}(0, 1)$. We consider both the over-parameterized and the exact/under-parameterized regime. The learning rate of GD is set to be $\mu = 0.005$.

6.1. Implicit low-rank bias

In this subsection, we consider the over-parameterized setting with $r = 50$. For each iteration $t \in [T]$ and rank $s \in [r_*]$, we define the relative error $\mathcal{E}_s(t) = \frac{\|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_s \mathbf{X}_s^\top\|_F^2}{\|\mathbf{X}_s \mathbf{X}_s^\top\|_F^2}$ to measure the proximity of the GD iterates to \mathbf{X}_s . We plot the relative error in Figure 1 for different choices of α and m (which affects the measurement error δ).

Small initialization. The implicit low-rank bias of GD is evident when the initialization scale α is small. Indeed, one can observe that GD first visits a small neighborhood of \mathbf{X}_1 , spends a long period of time near it, and then moves towards

\mathbf{X}_2 . It then proceeds to learn $\mathbf{X}_3, \mathbf{X}_4, \dots$ in a similar way, until it finally fits the ground truth. This is in align with Theorem 4.1. By contrast, for large initialization we do not have this implicit bias.

The effect of measurement error. For fixed α , one can observe the relative error becomes smaller when the number of measurements increases. This is in align with Lemma 4.3 in which the bound depends on δ . In particular, for the case $s = r_*$ in the end the distance to the set of global minima goes to zero as $\alpha \rightarrow 0$.

6.2. Matrix sensing with exact parameterization

Now we study the behavior of GD in the exact parameterization regime ($r = r_*$). We fix $m = 1000$, $r = r_* = 5$ and run GD for $T = 500$ iterations. We plot the relative error in Figure 2. As predicted by Theorem 4.1, we can observe that when α is small, GD exhibits an implicit low-rank bias and takes a longer time to converge. The latter is because GD would get into a $\text{poly}(\alpha)$ -small neighborhood of the saddle point \mathbf{Z}_s and take a long time to escape the saddle. As guaranteed by Theorem 4.2, we also observe the final convergence to global minimizers for sufficiently small α .

7. Conclusion

In this paper, we study the matrix sensing problem with RIP measurements and show that GD with small initialization follows an incremental learning procedure, where GD finds near-optimal solutions with increasing ranks until it finds the ground-truth. We take a step towards understanding the optimization and generalization aspects of simple optimization methods, thereby providing insights into their success in modern applications such as deep learning (Goodfellow et al., 2016). Also, we provide a detailed landscape analysis in the under-parameterized regime, which to the best of our knowledge is the first analysis of this kind.

Although we focus on matrix sensing in this paper, it has been revealed in a line of works that the implicit regularization effect may vary for different models, including deep matrix factorization (Arora et al., 2019) and nonlinear ReLU/LeakyReLU networks (Lyu et al., 2021; Timor et al., 2022). Also, it is shown in Woodworth et al. (2020) that different initialization scales can lead to distinct inductive

bias and affect the generalization and optimization behaviors. All these results indicate that we need further studies to comprehensively understand gradient-based optimization methods from the generalization aspect.

Acknowledgements

JJ is partially supported by the elite undergraduate training program of School of Mathematical Sciences in Peking University. KL is supported by funding from NSF, ONR, Simons Foundation, DARPA and SRC. JDL acknowledges support of the ARO under MURI Award W911NF-11-1-0304, the Sloan Research Fellowship, NSF CCF 2002272, NSF IIS 2107304, NSF CIF 2212262, ONR Young Investigator Award, and NSF CAREER Award 2144994. SSD acknowledges support of NSF IIS 2110170, NSF DMS 2134106, NSF CCF 2212261, NSF IIS 2143493, NSF CCF 2019844.

References

- Arora, S., Cohen, N., Hu, W., and Luo, Y. Implicit regularization in deep matrix factorization. *Advances in Neural Information Processing Systems*, 32, 2019.
- Belabbas, M. A. On implicit regularization: Morse functions and applications to matrix factorization. *arXiv preprint arXiv:2001.04264*, 2020.
- Blanc, G., Gupta, N., Valiant, G., and Valiant, P. Implicit regularization for deep neural networks driven by an ornstein-uhlenbeck like process. In Abernethy, J. and Agarwal, S. (eds.), *Proceedings of Thirty Third Conference on Learning Theory*, volume 125 of *Proceedings of Machine Learning Research*, pp. 483–513. PMLR, 09–12 Jul 2020.
- Boursier, E., Pillaud-Vivien, L., and Flammarion, N. Gradient flow dynamics of shallow relu networks for square loss and orthogonal inputs. *arXiv preprint arXiv:2206.00939*, 2022.
- Candes, E. J. and Plan, Y. Matrix completion with noise. *Proceedings of the IEEE*, 98(6):925–936, 2010.
- Candès, E. J. and Recht, B. Exact matrix completion via convex optimization. *Foundations of Computational mathematics*, 9(6):717–772, 2009.
- Candès, E. J., Li, X., Ma, Y., and Wright, J. Robust principal component analysis? *Journal of the ACM (JACM)*, 58(3): 1–37, 2011.
- Damian, A., Ma, T., and Lee, J. D. Label noise SGD provably prefers flat global minimizers. In Ranzato, M., Beygelzimer, A., Dauphin, Y., Liang, P., and Vaughan, J. W. (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 27449–27461. Curran Associates, Inc., 2021.
- Davenport, M. A. and Romberg, J. An overview of low-rank matrix recovery from incomplete observations. *IEEE Journal of Selected Topics in Signal Processing*, 10(4): 608–622, 2016.
- Du, S., Lee, J., Li, H., Wang, L., and Zhai, X. Gradient descent finds global minima of deep neural networks. In *International Conference on Machine Learning*, pp. 1675–1685. PMLR, 2019.
- Fornasier, M., Rauhut, H., and Ward, R. Low-rank matrix recovery via iteratively reweighted least squares minimization. *SIAM Journal on Optimization*, 21(4):1614–1640, 2011.
- Frei, S., Cao, Y., and Gu, Q. Provable generalization of sgd-trained neural networks of any width in the presence of adversarial label noise. In Meila, M. and Zhang, T. (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 3427–3438. PMLR, 18–24 Jul 2021.
- Ge, R., Huang, F., Jin, C., and Yuan, Y. Escaping from saddle points—online stochastic gradient for tensor decomposition. In *Conference on learning theory*, pp. 797–842. PMLR, 2015.
- Ge, R., Lee, J. D., and Ma, T. Matrix completion has no spurious local minimum. In *Advances in Neural Information Processing Systems*, pp. 2973–2981, 2016.
- Ge, R., Jin, C., and Zheng, Y. No spurious local minima in nonconvex low rank problems: A unified geometric analysis. In *International Conference on Machine Learning*, pp. 1233–1242. PMLR, 2017.
- Gidel, G., Bach, F., and Lacoste-Julien, S. Implicit regularization of discrete gradient dynamics in linear neural networks. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- Gissin, D., Shalev-Shwartz, S., and Daniely, A. The implicit bias of depth: How incremental learning drives generalization. In *International Conference on Learning Representations*, 2019.
- Goodall, C. Procrustes methods in the statistical analysis of shape. *Journal of the Royal Statistical Society. Series B (Methodological)*, 53(2):285–339, 1991. ISSN 00359246.
- Goodfellow, I., Bengio, Y., Courville, A., and Bengio, Y. *Deep learning*, volume 1. MIT Press, 2016.

- Gunasekar, S., Woodworth, B. E., Bhojanapalli, S., Neyshabur, B., and Srebro, N. Implicit regularization in matrix factorization. *Advances in Neural Information Processing Systems*, 30, 2017.
- Gunasekar, S., Lee, J., Soudry, D., and Srebro, N. Characterizing implicit bias in terms of optimization geometry. In *International Conference on Machine Learning*, pp. 1832–1841. PMLR, 2018.
- Hardt, M., Recht, B., and Singer, Y. Train faster, generalize better: Stability of stochastic gradient descent. In *International conference on machine learning*, pp. 1225–1234. PMLR, 2016.
- Hu, W., Xiao, L., Adlam, B., and Pennington, J. The surprising simplicity of the early-time learning dynamics of neural networks. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H. (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 17116–17128. Curran Associates, Inc., 2020.
- Jacot, A., Ged, F., Gabriel, F., Şimşek, B., and Hongler, C. Deep linear networks dynamics: Low-rank biases induced by initialization scale and l2 regularization. *arXiv preprint arXiv:2106.15933*, 2021.
- Ji, Z. and Telgarsky, M. Directional convergence and alignment in deep learning. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H. (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 17176–17186. Curran Associates, Inc., 2020.
- Jiang, L., Chen, Y., and Ding, L. Algorithmic regularization in model-free overparametrized asymmetric matrix factorization. *ArXiv*, abs/2203.02839, 2022.
- Kalimeris, D., Kaplun, G., Nakkiran, P., Edelman, B., Yang, T., Barak, B., and Zhang, H. Sgd on neural networks learns functions of increasing complexity. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- Kalogerias, D. S. and Petropulu, A. P. Matrix completion in colocated mimo radar: Recoverability, bounds & theoretical guarantees. *IEEE Transactions on Signal Processing*, 62(2):309–321, 2013.
- Karimi, H., Nutini, J., and Schmidt, M. Linear convergence of gradient and proximal-gradient methods under the polyak-łojasiewicz condition. In *Joint European conference on machine learning and knowledge discovery in databases*, pp. 795–811. Springer, 2016.
- Lee, J. D., Simchowitz, M., Jordan, M. I., and Recht, B. Gradient descent only converges to minimizers. In *Conference on Learning Theory*, pp. 1246–1257, 2016.
- Li, X., Lu, J., Arora, R., Haupt, J., Liu, H., Wang, Z., and Zhao, T. Symmetry, saddle points, and global optimization landscape of nonconvex matrix factorization. *IEEE Transactions on Information Theory*, 65(6):3489–3514, 2019.
- Li, Y., Ma, T., and Zhang, H. Algorithmic regularization in over-parameterized matrix sensing and neural networks with quadratic activations. In *Conference On Learning Theory*, pp. 2–47. PMLR, 2018.
- Li, Z., Luo, Y., and Lyu, K. Towards resolving the implicit bias of gradient descent for matrix factorization: Greedy low-rank learning. In *International Conference on Learning Representations*, 2020.
- Li, Z., Wang, T., and Arora, S. What happens after SGD reaches zero loss? –a mathematical framework. In *International Conference on Learning Representations*, 2022a.
- Li, Z., Wang, T., Lee, J. D., and Arora, S. Implicit bias of gradient descent on reparametrized models: On equivalence to mirror descent. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho, K. (eds.), *Advances in Neural Information Processing Systems*, 2022b. URL https://openreview.net/forum?id=k4KHXS6_zOV.
- Lyu, K. and Li, J. Gradient descent maximizes the margin of homogeneous neural networks. In *International Conference on Learning Representations*, 2020.
- Lyu, K., Li, Z., Wang, R., and Arora, S. Gradient descent on two-layer nets: Margin maximization and simplicity bias. *Advances in Neural Information Processing Systems*, 34: 12978–12991, 2021.
- Lyu, K., Li, Z., and Arora, S. Understanding the generalization benefit of normalization layers: Sharpness reduction. *arXiv preprint arXiv:2206.07085*, 2022.
- Nacson, M. S., Gunasekar, S., Lee, J., Srebro, N., and Soudry, D. Lexicographic and depth-sensitive margins in homogeneous and non-homogeneous deep models. In *International Conference on Machine Learning*, pp. 4683–4692. PMLR, 2019.
- Ngo, T. and Saad, Y. Scaled gradients on grassmann manifolds for matrix completion. *Advances in neural information processing systems*, 25, 2012.
- Peng, Y., Suo, J., Dai, Q., and Xu, W. Reweighted low-rank matrix recovery and its application in image restoration. *IEEE transactions on cybernetics*, 44(12):2418–2430, 2014.
- Razin, N. and Cohen, N. Implicit regularization in deep learning may not be explainable by norms. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H.

- (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 21174–21187. Curran Associates, Inc., 2020.
- Razin, N., Maman, A., and Cohen, N. Implicit regularization in tensor factorization. In *International Conference on Machine Learning*, pp. 8913–8924. PMLR, 2021.
- Razin, N., Maman, A., and Cohen, N. Implicit regularization in hierarchical tensor factorization and deep convolutional neural networks. *arXiv preprint arXiv:2201.11729*, 2022.
- Recht, B., Fazel, M., and Parrilo, P. A. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM review*, 52(3):471–501, 2010.
- Rudelson, M. and Vershynin, R. Smallest singular value of a random rectangular matrix. *Communications on Pure and Applied Mathematics: A Journal Issued by the Courant Institute of Mathematical Sciences*, 62(12):1707–1739, 2009.
- Shen, X. and Wu, Y. A unified approach to salient object detection via low rank matrix recovery. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 853–860. IEEE, 2012.
- Soudry, D., Hoffer, E., Nacson, M. S., Gunasekar, S., and Srebro, N. The implicit bias of gradient descent on separable data. *Journal of Machine Learning Research*, 19(70):1–57, 2018.
- Stöger, D. and Soltanolkotabi, M. Small random initialization is akin to spectral learning: Optimization and generalization guarantees for overparameterized low-rank matrix reconstruction. *Advances in Neural Information Processing Systems*, 34, 2021.
- Sun, R. Optimization for deep learning: theory and algorithms. *arXiv preprint arXiv:1912.08957*, 2019.
- Timor, N., Vardi, G., and Shamir, O. Implicit regularization towards rank minimization in relu networks. *arXiv preprint arXiv:2201.12760*, 2022.
- Tong, T., Ma, C., and Chi, Y. Accelerating ill-conditioned low-rank matrix estimation via scaled gradient descent. *J. Mach. Learn. Res.*, 22:150–1, 2021.
- Tu, S., Boczar, R., Simchowitz, M., Soltanolkotabi, M., and Recht, B. Low-rank solutions of linear matrix equations via Procrustes flow. In *International Conference on Machine Learning*, pp. 964–973. PMLR, 2016.
- Wedin, P.-Å. Perturbation bounds in connection with singular value decomposition. *BIT Numerical Mathematics*, 12(1):99–111, 1972.
- Wei, K., Cai, J.-F., Chan, T. F., and Leung, S. Guarantees of riemannian optimization for low rank matrix recovery. *SIAM Journal on Matrix Analysis and Applications*, 37(3):1198–1222, 2016.
- Woodworth, B., Gunasekar, S., Lee, J. D., Moroshko, E., Savarese, P., Golan, I., Soudry, D., and Srebro, N. Kernel and rich regimes in overparametrized models. In *Conference on Learning Theory*, pp. 3635–3673. PMLR, 2020.
- Xu, H., Caramanis, C., and Sanghavi, S. Robust pca via outlier pursuit. *Advances in neural information processing systems*, 23, 2010.
- Zhang, C., Bengio, S., Hardt, M., Recht, B., and Vinyals, O. Understanding deep learning requires rethinking generalization. In *International Conference on Learning Representations*, 2017. URL <https://openreview.net/forum?id=Sy8gdB9xx>.
- Zhang, H. and Yin, W. Gradient methods for convex minimization: better rates under weaker conditions. *arXiv preprint arXiv:1303.4645*, 2013.
- Zhao, B., Haldar, J. P., Brinegar, C., and Liang, Z.-P. Low rank matrix recovery for real-time cardiac mri. In *2010 IEEE International Symposium on Biomedical Imaging: From nano to macro*, pp. 996–999. IEEE, 2010.
- Zhu, Z., Li, Q., Tang, G., and Wakin, M. B. Global optimality in low-rank matrix optimization. *IEEE Transactions on Signal Processing*, 66(13):3614–3628, 2018.
- Zhu, Z., Li, Q., Tang, G., and Wakin, M. B. The global optimization geometry of low-rank matrix optimization. *IEEE Transactions on Information Theory*, 67(2):1308–1331, 2021.
- Zou, W., Kpalma, K., Liu, Z., and Ronsin, J. Segmentation driven low-rank matrix recovery for saliency detection. In *24th British machine vision conference (BMVC)*, pp. 1–13, 2013.

The appendix is organized as follows: in [Appendix A](#) we present a number of results that will be used for later proof. [Appendix B](#) sketches the main idea for proving our main results. [Appendix C](#) is devoted to a rigorous proof of [Lemma 4.3](#), with some auxiliary lemmas proved in [Appendix D](#). In [Appendix E](#) we analyze the landscape of low-rank matrix sensing and prove our landscape results in [Section 4.1](#). These results are then used in [Appendix F](#) to prove [Theorems 4.1](#) and [4.2](#). Finally, [Appendix G](#) studies the landscape of rank-1 matrix sensing, which enjoys a strongly convex property, as we mentioned in [Section 4.1](#) without proof.

A. Preliminaries

In this section, we present some useful results that is needed in subsequent analysis.

A.1. The RIP condition and its properties

In this subsection, we collect a few useful properties of the RIP condition, which we recall below:

Definition A.1. We say that the measurement \mathcal{A} satisfies the (δ, r) -RIP condition if for all matrices $\mathbf{Z} \in \mathbb{R}^{d \times d}$ with $\text{rank}(\mathbf{Z}) \leq r$, we have

$$(1 - \delta)\|\mathbf{Z}\|_F^2 \leq \|\mathcal{A}(\mathbf{Z})\|_2^2 \leq (1 + \delta)\|\mathbf{Z}\|_F^2.$$

The key intuition behind RIP is that $\mathcal{A}^* \mathcal{A} \approx \mathbf{I}$, where $\mathcal{A}^* : \mathbf{v} \mapsto \frac{1}{\sqrt{m}} \sum_{i=1}^m v_i \mathbf{A}_i$ is the adjoint of \mathcal{A} . This intuition is made rigorous by the following proposition:

Proposition A.2. (*Stöger & Soltanolkotabi, 2021, Lemma 7.3*) Suppose that \mathcal{A} satisfies (r, δ) -RIP with $r \geq 2$, then for all symmetric \mathbf{Z} ,

- (1). if $\text{rank}(\mathbf{Z}) \leq r - 1$, we have $\|(\mathcal{A}^* \mathcal{A} - \mathbf{I})\mathbf{Z}\|_2 \leq \sqrt{r}\delta\|\mathbf{Z}\|$.
- (2). $\|(\mathcal{A}^* \mathcal{A} - \mathbf{I})\mathbf{Z}\|_2 \leq \delta\|\mathbf{Z}\|_*$, where $\|\cdot\|_*$ is the nuclear norm.

A.2. Matrix analysis

The following lemma is a direct corollary of [Proposition A.2](#) and will be frequently used in our proof.

Lemma A.3. Suppose that the measurement \mathcal{A} satisfies $(\delta, 2r_* + 1)$ -RIP condition, then for all matrices $\mathbf{U} \in \mathbb{R}^{d \times r}$ such that $\text{rank}(\mathbf{U}) \leq r_*$, we have

$$\|(\mathcal{A}^* \mathcal{A} - \mathbf{I})(\mathbf{X}\mathbf{X}^\top - \mathbf{U}\mathbf{U}^\top)\| \leq \delta\sqrt{r_*}(\|\mathbf{X}\|^2 + \|\mathbf{U}\|^2).$$

In our proof we will frequently make use of the Weyl's inequality for singular values:

Lemma A.4 (Weyl's inequality). Let $\mathbf{A}, \Delta \in \mathbb{R}^{d \times d}$ be two matrices, then for all $1 \leq k \leq d$, we have

$$|\sigma_k(\mathbf{A}) - \sigma_k(\mathbf{A} + \Delta)| \leq \|\Delta\|.$$

We will also need the Wedin's sin theorem for singular value decomposition:

Lemma A.5. (*Wedin, 1972, Section 3*) Define $R(\cdot)$ to be the column space of a matrix. Suppose that matrices $\mathbf{B} = \mathbf{A} + \mathbf{T}$, $\mathbf{A}_1, \mathbf{B}_1$ are the top- s components in the SVD of \mathbf{A} and \mathbf{B} respectively, and $\mathbf{A}_0 = \mathbf{A} - \mathbf{A}_1, \mathbf{B}_0 = \mathbf{B} - \mathbf{B}_1$. If $\delta = \sigma_{\min}(\mathbf{B}_1) - \sigma_{\max}(\mathbf{A}_0) > 0$, then we have

$$\|\sin \Theta(R(\mathbf{A}_1), R(\mathbf{B}_1))\| \leq \frac{\|\mathbf{T}\|}{\delta}$$

where $\Theta(\cdot, \cdot)$ denotes the angle between two subspaces.

Equipped with [Lemma A.3](#), we can have the following characterization of the eigenvalues of \mathbf{M} (recall that $\mathbf{M} = \mathcal{A}^* \mathcal{A}(\mathbf{X}\mathbf{X}^\top)$):

Lemma A.6. Let $\mathbf{M} := \mathcal{A}^* \mathcal{A}(\mathbf{X}\mathbf{X}^\top)$ and $\mathbf{M} = \sum_{k=1}^d \hat{\sigma}_k^2 \hat{\mathbf{v}}_k \hat{\mathbf{v}}_k^\top$ be the eigen-decomposition of \mathbf{M} . For $1 \leq i \leq d$ we have

$$|\sigma_i^2 - \hat{\sigma}_i^2| \leq \delta\|\mathbf{X}\|^2.$$

Proof. By Weyl's inequality we have

$$|\sigma_i^2 - \hat{\sigma}_i^2| \leq \|\mathbf{M} - \mathbf{X}\mathbf{X}^\top\| \leq \delta \|\mathbf{X}\|^2$$

as desired. \square

A.3. Optimization

Lemma A.7. *Suppose that a smooth function $f \in \mathbb{R}^m \mapsto \mathbb{R}$ with minimum value $f^* > -\infty$ satisfies the following conditions with some $\epsilon > 0$:*

- (1). $\lim_{\|\mathbf{x}\| \rightarrow +\infty} f(\mathbf{x}) = +\infty$.
- (2). *There exists an open subset $S \subset \mathbb{R}^m$ such that the set S^* of global minima of f is contained in S , and for all stationary points \mathbf{x} of f in $\mathbb{R}^m - S$, we have $f(\mathbf{x}) - f^* \geq 2\epsilon$. Moreover, we also have $f(\mathbf{x}) - f^* \geq 2\epsilon$ on ∂S .*

Then we have

$$\{\mathbf{x} \in \mathbb{R}^m : f(\mathbf{x}) - f^* \leq \epsilon\} \subset S.$$

Proof. Let \mathbf{x}^* be the minimizer of f on $\mathbb{R}^m - S$. By condition (1) we can deduce that \mathbf{x}^* always exists. Moreover, since any local minimizer of a function defined on a compact set must either be a stationary point or lie on the boundary of its domain, we can see that either $\mathbf{x}^* \in \partial S$ or $\nabla f(\mathbf{x}^*) = 0$ holds. By condition (2), either cases would imply that $f(\mathbf{x}^*) - f^* \geq 2\epsilon$, as desired. \square

Lemma A.8. *Let $\{\mathbf{x}_k\}, \{\mathbf{y}_k\} \subset \mathbb{R}^n$ be two sequences generated by $\mathbf{x}_{k+1} = \mathbf{x}_k - \mu \nabla f(\mathbf{x}_k)$ and $\mathbf{y}_{k+1} = \mathbf{y}_k - \mu \nabla f(\mathbf{y}_k)$. Suppose that $\|\mathbf{x}_k\| \leq B$ and $\|\mathbf{y}_k\| \leq B$ for all k and f is L -smooth in $\{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\| \leq B\}$, then we have*

$$\|\mathbf{x}_k - \mathbf{y}_k\| \leq (1 + \mu L)^k \|\mathbf{x}_0 - \mathbf{y}_0\|.$$

Proof. The update rule implies that

$$\begin{aligned} \|\mathbf{x}_{k+1} - \mathbf{y}_{k+1}\| &= \|\mathbf{x}_k - \mathbf{y}_k - \mu \nabla f(\mathbf{x}_k) + \mu \nabla f(\mathbf{y}_k)\| \\ &\leq \|\mathbf{x}_k - \mathbf{y}_k\| + \mu \|\nabla f(\mathbf{x}_k) - \nabla f(\mathbf{y}_k)\| \\ &\leq (1 + \mu L) \|\mathbf{x}_k - \mathbf{y}_k\| \end{aligned}$$

which yields the desired inequality. \square

A.4. Proof for Proposition 3.6

Proposition 3.6. *Suppose that all entries of $\bar{\mathbf{U}} \in \mathbb{R}^{d \times \hat{r}}$ are independently drawn from $\mathcal{N}(0, \frac{1}{\hat{r}})$ and $\rho = \epsilon \frac{\sqrt{\hat{r}} - \sqrt{\hat{r} \wedge r_* - 1}}{\sqrt{\hat{r}}} \geq \frac{\epsilon}{2r_*}$, then $\sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \bar{\mathbf{U}}) \geq \rho$ holds for all $1 \leq s \leq \hat{r} \wedge r_*$ with probability at least $1 - \hat{r}(C\epsilon + e^{-c\hat{r}})$, where $c, C > 0$ are universal constants.*

Proposition 3.6 immediately follows from the following result:

Proposition A.9. *(Restatement of Rudelson & Vershynin, 2009, Theorem 1.1) Let A be an $N \times n$ random matrix, $N \geq n$, whose elements are independent copies of a mean zero sub-gaussian random variable with unit variance. Then, for every $\epsilon > 0$, we have $\mathbb{P}\left(s_n(A) \leq \epsilon(\sqrt{N} - \sqrt{n-1})\right) \leq (C\epsilon)^{N-n+1} + e^{-cN}$ where $C, c > 0$ depend (polynomially) only on the sub-Gaussian moment.*

Now we can complete the proof of **Proposition 3.6**. Note that the entries of $U \in \mathbb{R}^{d \times \hat{r}}$ are independently drawn from $\mathcal{N}(0, \frac{1}{\hat{r}})$ and $\mathbf{V}_{\mathbf{X}_s} \in \mathbb{R}^{d \times s}$ is an orthonormal matrix. We write $\mathbf{V}_{\mathbf{X}_s}^+ \in \mathbb{R}^{d \times s \hat{r}}$ as a block diagonal matrix with \hat{r} copies of $\mathbf{V}_{\mathbf{X}_s}$ on the diagonal, and $\text{vec}(\mathbf{U}) \in \mathbb{R}^{d\hat{r}}$ be a vector formed by the concatenation of the columns of U . Then $\mathbf{V}_{\mathbf{X}_s}^+$ is still orthonormal, and $\text{vec}(\mathbf{U}) \sim \mathcal{N}(0, \frac{1}{\hat{r}}\mathbf{I})$. Since multivariate Gaussian distributions are invariant under orthonormal transformations, we deduce that $(\mathbf{V}_{\mathbf{X}_s}^+)^T \text{vec}(\mathbf{U}) \sim \mathcal{N}(0, \frac{1}{\hat{r}}\mathbf{I})$. Equivalently, the entries of $\mathbf{V}_{\mathbf{X}_s}^+ \mathbf{U}$ are i.i.d. $\mathcal{N}(0, \frac{1}{\hat{r}})$.

The matrix $\sqrt{\hat{r}}\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}$ satisfies all the conditions in [Proposition A.9](#). Thus, with probability at least $1 - (C\epsilon)^{\hat{r}-s+1} - e^{-c\hat{r}}$, we have $\sigma_{\min}(\sqrt{\hat{r}}\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}) \geq \epsilon(\sqrt{\hat{r}} - \sqrt{s-1})$, or equivalently $\sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}) \geq \epsilon \frac{\sqrt{\hat{r}} - \sqrt{s-1}}{\sqrt{\hat{r}}}$. Finally, the conclusion follows from a union bound:

$$\begin{aligned} & \mathbb{P} \left[\exists 1 \leq s \leq \hat{r} \wedge r_* \text{ s.t. } \sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}) < \frac{\epsilon}{2\hat{r}} \right] \\ & \leq \sum_{s=1}^{\hat{r} \wedge r_*} \mathbb{P} \left[\sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}) < \epsilon \frac{\sqrt{\hat{r}} - \sqrt{s-1}}{\sqrt{\hat{r}}} \right] \leq \sum_{s=1}^{\hat{r} \wedge r_*} (e^{-c\hat{r}} + (C\epsilon)^{\hat{r}-s+1}) \leq r (e^{-c\hat{r}} + C\epsilon). \end{aligned} \quad (15)$$

A.5. Procrustes Distance

Procrustes distance is introduced in [Section 3.3](#). The following characterization of the optimal \mathbf{R} in [Definition 3.8](#) is known in the literature (see e.g. [Tu et al., 2016](#), Section 5.2.1) but we provide a proof for completeness.

Lemma A.10. *Let $\mathbf{U}_1, \mathbf{U}_2 \in \mathbb{R}^{d \times r}$ where $r \leq d$. Then for any orthogonal matrix $\mathbf{R} \in \mathbb{R}^{r \times r}$ that minimizes $\|\mathbf{U}_1 - \mathbf{U}_2 \mathbf{R}\|_F$ (i.e., any orthogonal \mathbf{R} s.t. $\|\mathbf{U}_1 - \mathbf{U}_2 \mathbf{R}\|_F = \text{dist}(\mathbf{U}_1, \mathbf{U}_2)$), $\mathbf{U}_1^\top \mathbf{U}_2 \mathbf{R}$ is a symmetric positive semi-definite matrix.*

Proof. We only need to consider the case when $\mathbf{U}_2^\top \mathbf{U}_1 \neq \mathbf{0}$. Observe that

$$\begin{aligned} \|\mathbf{U}_1 - \mathbf{U}_2 \mathbf{R}\|_F^2 &= \|\mathbf{U}_1\|_F^2 + \|\mathbf{U}_2 \mathbf{R}\|_F^2 - 2 \text{tr}(\mathbf{R}^\top \mathbf{U}_2^\top \mathbf{U}_1) \\ &= \|\mathbf{U}_1\|_F^2 + \|\mathbf{U}_2\|_F^2 - 2 \text{tr}(\mathbf{R}^\top \mathbf{U}_2^\top \mathbf{U}_1). \end{aligned}$$

Let $\mathbf{A} \Sigma \mathbf{B}^\top$ be the SVD of $\mathbf{U}_2^\top \mathbf{U}_1$, where $\mathbf{A}^\top \mathbf{A} = \mathbf{I}$, $\mathbf{B}^\top \mathbf{B} = \mathbf{I}$ and $\Sigma \succ \mathbf{0}$. Then

$$\text{tr}(\mathbf{R}^\top \mathbf{U}_2^\top \mathbf{U}_1) = \text{tr}(\mathbf{B}^\top \mathbf{R}^\top \mathbf{A} \Sigma) \leq \|\mathbf{B}^\top \mathbf{R}^\top \mathbf{A}\| \text{tr}(\Sigma) = \text{tr}(\Sigma),$$

where the final step is due to orthogonality of $\mathbf{B}^\top \mathbf{R}^\top \mathbf{A} \in \mathbb{R}^{s \times s}$, and equality holds if and only if $\mathbf{B}^\top \mathbf{R}^\top \mathbf{A} = \mathbf{I}$. Let $\mathbf{C} = \mathbf{R}^\top \mathbf{A}$. Let $\mathbf{b}_i, \mathbf{c}_i \in \mathbb{R}^d$ be the i -th column of \mathbf{B} and \mathbf{C} respectively, then $\mathbf{B}^\top \mathbf{C} = \mathbf{I}$ implies that $\mathbf{b}_i^\top \mathbf{c}_i = 1$. Note that $\|\mathbf{b}_i\|_2 = \|\mathbf{c}_i\|_2 = 1$, so we must have $\mathbf{b}_i = \mathbf{c}_i$ for all i , i.e., $\mathbf{B} = \mathbf{C} = \mathbf{R}^\top \mathbf{A}$. Therefore, $\mathbf{U}_1^\top \mathbf{U}_2 \mathbf{R} = \mathbf{B} \Sigma \mathbf{A}^\top \mathbf{R} = \mathbf{B} \Sigma \mathbf{B}^\top$, which implies that $\mathbf{U}_1^\top \mathbf{U}_2 \mathbf{R}$ is symmetric and positive semi-definite. \square

B. Main idea for the proof of [Theorem 4.1](#)

In this section, we briefly introduce our main ideas for proving [Theorem 4.1](#). Motivated by [Stöger & Soltanolkotabi \(2021\)](#), we decompose the matrix \mathbf{U}_t into a parallel component and an orthogonal component. Specifically, we write

$$\mathbf{U}_t = \underbrace{\mathbf{U}_t \mathbf{W}_t \mathbf{W}_t^\top}_{\text{parallel component}} + \underbrace{\mathbf{U}_t \mathbf{W}_{t,\perp} \mathbf{W}_{t,\perp}^\top}_{\text{orthogonal component}}, \quad (16)$$

where $\mathbf{W}_t := \mathbf{W}_{\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t} \in \mathbb{R}^{\hat{r} \times s}$ is the matrix consisting of the right singular vectors of $\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t$ ([Definition 3.1](#)) and $\mathbf{W}_{t,\perp} \in \mathbb{R}^{\hat{r} \times (\hat{r}-s)}$ is an orthogonal complement of \mathbf{W}_t . Our goal is to prove that at some time t , we have $\mathbf{V}_{\mathbf{X}_s}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_s \mathbf{X}_s^\top) \approx \mathbf{0}$ and $\|\mathbf{U}_t \mathbf{W}_{t,\perp}\| \approx 0$. As we will see later, these imply that $\|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_s \mathbf{X}_s^\top\| \approx 0$. In the remaining part of this section we give a heuristic explanation for considering [\(16\)](#).

Additional Notations. Let $\mathbf{V}_{\mathbf{X}_s,\perp} \in \mathbb{R}^{d \times (d-s)}$ be an orthogonal complement of $\mathbf{V}_{\mathbf{X}_s} \in \mathbb{R}^{d \times s}$. Let $\Sigma_s = \text{diag}(\sigma_1, \dots, \sigma_s)$ and $\Sigma_{s,\perp} = \text{diag}(\sigma_{s+1}, \dots, \sigma_r, 0, \dots, 0) \in \mathbb{R}^{(d-s) \times (d-s)}$. We use $\Delta_t := (\mathcal{A}^* \mathcal{A} - \mathbf{I})(\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top)$ to denote the vector consisting of measurement errors for $\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top$.

B.1. Heuristic explanations of the decomposition

A simple and intuitive approach for showing the implicit low rank bias is to directly analyze the growth of $\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t$ versus $\mathbf{V}_{\mathbf{X}_s,\perp}^\top \mathbf{U}_t$. Ideally, the former grows faster than the latter, so that GD only learns the components in \mathbf{X}_s .

By the update rule of GD (3),

$$\begin{aligned} \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_{t+1} &= \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top [\mathbf{I} + \mu \mathcal{A}^* \mathcal{A}(\mathbf{X}\mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top)] \mathbf{U}_t \\ &= \underbrace{\mathbf{V}_{\mathbf{X}_{s,\perp}}^\top [\mathbf{I} + \mu \mathbf{X}\mathbf{X}^\top - \mu \mathbf{U}_t \mathbf{U}_t^\top]}_{=: \mathbf{G}_{t,1}} \mathbf{U}_t + \underbrace{\mu \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \Delta_t \mathbf{U}_t}_{=: \mathbf{G}_{t,2}} \\ &= \mathbf{G}_{t,1} + \mu \mathbf{G}_{t,2}. \end{aligned}$$

For the first term $\mathbf{G}_{t,1}$, we have

$$\begin{aligned} \mathbf{G}_{t,1} &= (\mathbf{I} + \mu \Sigma_{s,\perp}^2) \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t - \mu \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{U}_t \\ &= (\mathbf{I} + \mu \Sigma_{s,\perp}^2) \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t (\mathbf{I} - \mu \mathbf{U}_t \mathbf{U}_t^\top) + \mathcal{O}(\mu^2), \end{aligned}$$

where the last term $\mathcal{O}(\mu^2)$ is negligible when μ is sufficiently small. Since $\|\Sigma_{s,\perp}\| = \sigma_{s+1}$, the spectral norm of $\mathbf{G}_{t,1}$ can be bounded by

$$\begin{aligned} \|\mathbf{G}_{t,1}\| &\leq \|\mathbf{I} + \mu \Sigma_{s,\perp}^2\| \cdot \|\mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t\| \cdot \|\mathbf{I} - \mu \mathbf{U}_t \mathbf{U}_t^\top\| + \mathcal{O}(\mu^2) \\ &\leq (1 + \mu \sigma_{s+1}^2) \|\mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t\| + \mathcal{O}(\mu^2). \end{aligned}$$

However, the main difference with the full-observation case (Jiang et al., 2022) is the second term $\mathbf{G}_{t,2} := \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \Delta_t \mathbf{U}_t$. Since the measurement errors Δ_t are small but arbitrary, it is hard to compare this term with $\mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_{t+1}$. As a result, we cannot directly bound the growth of $\|\mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t\|$.

However, the aforementioned problem disappears if we turn to bound the growth of $\|\mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_{t+1} \mathbf{W}_{t,\perp}\|$. To see this, first we deduce the following by repeatedly using $\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_{t,\perp} = 0$ due to the definition of $\mathbf{W}_{t,\perp}$.

$$\begin{aligned} \mathbf{G}_{t,1} \mathbf{W}_{t,\perp} &= \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top [\mathbf{I} + \mu \mathbf{X}\mathbf{X}^\top - \mu \mathbf{U}_t \mathbf{U}_t^\top] \mathbf{U}_t \mathbf{W}_{t,\perp} \\ &= \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top (\mathbf{I} + \mu \mathbf{X}\mathbf{X}^\top) \mathbf{U}_t \mathbf{W}_{t,\perp} - \mu \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{U}_t \mathbf{W}_{t,\perp} \\ &= (\mathbf{I} + \mu \Sigma_{s,\perp}^2) \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{W}_{t,\perp} - \mu \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t (\mathbf{W}_t \mathbf{W}_t^\top + \mathbf{W}_{t,\perp} \mathbf{W}_{t,\perp}^\top) \mathbf{U}_t^\top \mathbf{U}_t \mathbf{W}_{t,\perp} \\ &= (\mathbf{I} + \mu \Sigma_{s,\perp}^2) \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{W}_{t,\perp} (\mathbf{I} - \mu \mathbf{W}_{t,\perp}^\top \mathbf{U}_t^\top \mathbf{U}_t \mathbf{W}_{t,\perp}) \\ &\quad - \mu \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{W}_t \mathbf{W}_t^\top \mathbf{U}_t^\top \mathbf{U}_t \mathbf{W}_{t,\perp} + \mathcal{O}(\mu^2), \end{aligned}$$

$$\mathbf{G}_{t,2} \mathbf{W}_{t,\perp} = \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \Delta_t \mathbf{U}_t \mathbf{W}_{t,\perp} = \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \Delta_t \mathbf{V}_{\mathbf{X}_s} \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{W}_{t,\perp},$$

So we have the following recursion:

$$\begin{aligned} \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_{t+1} \mathbf{W}_{t,\perp} &= (\mathbf{I} + \mu \Sigma_{s,\perp}^2 + \mu \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \Delta_t \mathbf{V}_{\mathbf{X}_s}) \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{W}_{t,\perp} (\mathbf{I} - \mu \mathbf{W}_{t,\perp}^\top \mathbf{U}_t^\top \mathbf{U}_t \mathbf{W}_{t,\perp}) \\ &\quad - \mu \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{W}_t \mathbf{W}_t^\top \mathbf{U}_t^\top \mathbf{U}_t \mathbf{W}_{t,\perp} + \mathcal{O}(\mu^2), \end{aligned}$$

We further note that

$$\mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_{t+1} \mathbf{W}_{t+1,\perp} = \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_{t+1} \mathbf{W}_t \mathbf{W}_t^\top \mathbf{W}_{t+1,\perp} + \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_{t+1} \mathbf{W}_{t,\perp} \mathbf{W}_{t,\perp}^\top \mathbf{W}_{t+1,\perp}, \quad (17)$$

which establishes the relationship between $\mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_{t+1} \mathbf{W}_{t,\perp}$ and $\mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_{t+1} \mathbf{W}_{t+1,\perp}$. To complete the proof we need to prove the following:

- The minimal eigenvalue of the *parallel component* $\mathbf{U}_t \mathbf{W}_t \mathbf{W}_t^\top$ grows at a linear rate with speed strictly faster than σ_{s+1} .
- The term $\left\| \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t} \right\| \ll 1$, which implies that the first term in (17) is negligible.

C. Proof of Lemma 4.3

In this section, we give the full proof of Lemma 4.3, with some additional technical lemmas left to Appendix D.

Lemma 4.3. *Under Assumptions 3.3 and 3.5, there exists $\hat{T}_\alpha^s > 0$ for all $\alpha > 0$ and $1 \leq s \leq \hat{r} \wedge r_*$ such that $\lim_{\alpha \rightarrow 0} \max_{1 \leq t \leq \hat{T}_\alpha^s} \sigma_{s+1}(\mathbf{U}_{\alpha,t}) = 0$. Furthermore, it holds that $\left\| \mathbf{U}_{\hat{T}_\alpha^s} \mathbf{U}_{\hat{T}_\alpha^s}^\top - \mathbf{Z}_s^* \right\|_F = \mathcal{O}(\kappa_*^3 r_* \delta \|\mathbf{X}\|^2)$.*

Appendices C.1 and C.2 are devoted to analyzing the spectral phase and parallel improvement phase, respectively. Appendix C.3 uses induction to characterize the low-rank GD trajectory in the parallel improvement phase. In Appendix C.4 we study the refinement phase, which allows us to derive Lemma 4.3.

C.1. The spectral phase

Starting from a small initialization $\mathbf{U}_0 = \alpha \bar{\mathbf{U}}$, $\alpha \ll 1$, we first enter the spectral phase where GD behaves similar to power iteration. As in Stöger & Soltanolkotabi (2021), we refer to this phase as the spectral phase. Specifically, we have in the spectral phase that

$$\mathbf{U}_{t+1} = (\mathbf{I} + \mu(\mathcal{A}^* \mathcal{A})(\mathbf{X}\mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top)) \mathbf{U}_t \approx (\mathbf{I} + \mu(\mathcal{A}^* \mathcal{A})(\mathbf{X}\mathbf{X}^\top)) \mathbf{U}_t.$$

The approximation holds with high accuracy as long as $\|\mathbf{U}_t\| \ll 1$. Moreover we have $\mathbf{M} := (\mathcal{A}^* \mathcal{A})(\mathbf{X}\mathbf{X}^\top) \approx \mathbf{X}\mathbf{X}^\top$ by the RIP condition; when δ is sufficiently small, we can still ensure a positive eigen-gap of \mathbf{M} . As a result, with small initialization \mathbf{U}_t would become approximately aligned with the top eigenvector u_1 of \mathbf{M} . Since $\|\mathbf{M} - \mathbf{X}\mathbf{X}^\top\| = \mathcal{O}(\delta\sqrt{r_*})$ by Proposition A.2, we have $\|u_1 - v_1\| = \mathcal{O}(\delta\sqrt{r_*})$ so that $\|\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{V}_{\mathbf{U}_t} \mathbf{w}_i\| = \mathcal{O}(\delta\sqrt{r_*})$. This proves the base case for the induction.

Formally, we define $\mathbf{M} = \mathcal{A}^* \mathcal{A}(\mathbf{X}\mathbf{X}^\top)$, $\mathbf{K}_t = (\mathbf{I} + \mu\mathbf{M})^t$ and $\mathbf{U}_t^{\text{sp}} = \mathbf{K}_t \mathbf{U}_0$. Suppose that $\mathbf{M} = \sum_{i=1}^{\text{rank}(\mathbf{M})} \hat{\sigma}_i^2 \hat{\mathbf{v}}_i \hat{\mathbf{v}}_i^\top$ is the spectral decomposition of \mathbf{M} where $\{\hat{\sigma}_i\}_{i \geq 1}$ is sorted in non-increasing order. We additionally define $\mathbf{M}_s = \sum_{i=1}^{\min\{s, \text{rank}(\mathbf{M})\}} \hat{\sigma}_i^2 \hat{\mathbf{v}}_i \hat{\mathbf{v}}_i^\top$. By Lemma A.6 and $\delta\sqrt{r_*} \leq 10^{-3} \kappa_*$ by Assumption 3.3, we have $\hat{\sigma}_s^2 \geq \sigma_s^2 - 0.01\tau$ and $\hat{\sigma}_{s+1}^2 \leq \sigma_{s+1}^2 + 0.01\tau$, where we recall that $\tau = \min_{s \in [r_*]} (\sigma_s^2 - \sigma_{s+1}^2) > 0$. Additionally, let \mathbf{L}_t be the span of the top- s left singular vectors of \mathbf{U}_t . Recall that Assumption 3.5 is made on the initialization. Let

$$t^* := \min \{i \in \mathbb{N} : \|\mathbf{U}_{i-1}^{\text{sp}} - \mathbf{U}_{i-1}\| > \|\mathbf{U}_{i-1}^{\text{sp}}\|\},$$

the following lemma bounds the error of approximating \mathbf{U}_t via \mathbf{U}_t^{sp} :

Lemma C.1. (Stöger & Soltanolkotabi, 2021, Lemma 8.1) *Suppose that \mathcal{A} satisfies the rank-1 RIP with constant δ_1 . For all integers t such that $1 \leq t \leq t^*$ it holds that*

$$\|\mathbf{E}_t\| = \|\mathbf{U}_t - \mathbf{U}_t^{\text{sp}}\| \leq 4\hat{\sigma}_1^{-2} \alpha^3 r_* (1 + \delta_1) (1 + \mu\hat{\sigma}_1^2)^{3t}. \quad (18)$$

We can derive the following lower bound on t^* from Lemma C.1.

Corollary C.2. *We have*

$$t^* \geq \frac{\log \alpha^{-1} + \frac{1}{2} \log \frac{\rho \hat{\sigma}_1^2}{4(1+\delta_1)r_*}}{\log(1 + \mu\hat{\sigma}_1^2)}.$$

Proof. By Lemma C.1 we have

$$\|\mathbf{E}_t\| \leq 4\hat{\sigma}_1^{-2} \alpha^3 r_* (1 + \delta_1) (1 + \mu\hat{\sigma}_1^2)^{3t}.$$

for all $t \leq t^*$. On the other hand, we have

$$\begin{aligned} \|\mathbf{U}_t^{\text{sp}}\| &= \alpha \left\| (\mathbf{I} + \mu\mathbf{M})^t \bar{\mathbf{U}} \right\| \\ &\geq \alpha (1 + \mu\hat{\sigma}_1^2)^t \|\hat{\mathbf{v}}_1 \hat{\mathbf{v}}_1^\top \bar{\mathbf{U}}\| \\ &\geq (1 + \mu\hat{\sigma}_1^2)^t \alpha \rho. \end{aligned}$$

Thus, it follows from $\|\mathbf{E}_{t^*}\| \geq \|\mathbf{U}_{t^*}^{\text{sp}}\|$ that

$$(1 + \mu\hat{\sigma}_1^2)^{t^*} \geq \sqrt{\frac{\rho \hat{\sigma}_1^2}{4(1 + \delta_1)r_*}} \cdot \alpha^{-1} \Rightarrow t^* \geq \frac{\log \alpha^{-1} + \frac{1}{2} \log \frac{\rho \hat{\sigma}_1^2}{4(1+\delta_1)r_*}}{\log(1 + \mu\hat{\sigma}_1^2)}$$

as desired. \square

Note that a trivial bound for the rank-1 RIP constant is $\delta_1 \leq \delta$. We can now show that for small t , GD can be viewed as approximate power iteration.

Lemma C.3. *There exists a time*

$$t = T_\alpha^{\text{SP}} := \frac{2 \log \alpha^{-1} + \log \frac{\rho \hat{\sigma}_1^2}{4r_*(1+\delta)}}{3 \log(1 + \mu \hat{\sigma}_1^2) - \log(1 + \mu \hat{\sigma}_{s+1}^2)} \leq t^*$$

such that

$$\left\| \mathbf{U}_t - \sum_{i=1}^s \alpha (1 + \mu \hat{\sigma}_i^2)^t \hat{\mathbf{v}}_i \hat{\mathbf{v}}_i^\top \bar{\mathbf{U}} \right\| \leq C_1 \cdot \alpha^\gamma$$

where $\gamma = 1 - \frac{2 \log(1 + \mu \hat{\sigma}_{s+1}^2)}{3 \log(1 + \mu \hat{\sigma}_1^2) - \log(1 + \mu \hat{\sigma}_{s+1}^2)}$ and $C_1 = C_1(\mathbf{X}, \bar{\mathbf{U}})$ is a constant that only depends on \mathbf{X} and $\bar{\mathbf{U}}$.

Proof. It's easy to check that $T_\alpha^{\text{SP}} \leq t^*$ by applying [Corollary C.2](#).

We consider the following decomposition:

$$\left\| \mathbf{U}_t - \sum_{i=1}^s \alpha (1 + \mu \hat{\sigma}_i^2)^t \hat{\mathbf{v}}_i \hat{\mathbf{v}}_i^\top \bar{\mathbf{U}} \right\| \leq \left\| \mathbf{U}_t - \mathbf{U}_t^{\text{SP}} \right\| + \left\| \mathbf{U}_t^{\text{SP}} - \sum_{i=1}^s \alpha (1 + \mu \hat{\sigma}_i^2)^t \hat{\mathbf{v}}_i \hat{\mathbf{v}}_i^\top \bar{\mathbf{U}} \right\|.$$

When $t \leq t^*$, the first term can be bounded as

$$\|\mathbf{E}_t\| \leq 4\hat{\sigma}_1^{-2} \alpha^3 r_* (1 + \delta) (1 + \mu \hat{\sigma}_1^2)^{3t}.$$

For the second term we have

$$\left\| \mathbf{U}_t^{\text{SP}} - \sum_{i=1}^s \alpha (1 + \mu \hat{\sigma}_i^2)^t \hat{\mathbf{v}}_i \hat{\mathbf{v}}_i^\top \mathbf{U} \right\| \leq \left\| \sum_{i=s+1}^{r_*} \alpha (1 + \mu \hat{\sigma}_i^2)^t \hat{\mathbf{v}}_i \hat{\mathbf{v}}_i^\top \mathbf{U} \right\| \leq \alpha (1 + \mu \hat{\sigma}_{s+1}^2)^t.$$

In particular, the definition of T_α^{SP} implies that

$$\begin{aligned} \left\| \mathbf{U}_t - \sum_{i=1}^s \alpha (1 + \mu \hat{\sigma}_i^2)^t \hat{\mathbf{v}}_i \hat{\mathbf{v}}_i^\top \mathbf{U} \right\| &\leq 2 \left(\frac{\rho \hat{\sigma}_1^2}{4r_*(1+\delta)} \right)^{\frac{1-\gamma}{2}} \alpha^\gamma \\ &\leq 2 \max \left\{ 1, \frac{\rho \hat{\sigma}_1^2}{4r_*(1+\delta)} \right\} \alpha^\gamma \\ &\leq \underbrace{\max \left\{ 2, \frac{\rho \hat{\sigma}_1^2}{r_*} \right\}}_{:=C_1} \alpha^\gamma. \end{aligned}$$

as desired. \square

We conclude this section with the following lemma, which states that initially the parallel component $\mathbf{U}_t \mathbf{W}_t$ would grow much faster than the noise term, and would become well-aligned with \mathbf{X}_s .

Lemma C.4 ([Lemma 5.1](#), formal version). *There exists positive constants $C_2 = C_2(\mathbf{X}, \bar{\mathbf{U}})$ and $C_3 = C_3(\mathbf{X}, \bar{\mathbf{U}})$ such that the following inequalities hold for $t = T_\alpha^{\text{SP}}$ when $\alpha \in \left(0, \left(\frac{\rho}{10C_1(\mathbf{X}, \bar{\mathbf{U}})} \right)^{10\kappa_*} \right)$:*

$$\|\mathbf{U}_t\| \leq \|\mathbf{X}\| \tag{19a}$$

$$\sigma_{\min}(\mathbf{U}_t \mathbf{W}_t) \geq C_2 \cdot \alpha^{1 - \frac{2 \log(1 + \mu \hat{\sigma}_s^2)}{3 \log(1 + \mu \hat{\sigma}_1^2) - \log(1 + \mu \hat{\sigma}_{s+1}^2)}} \tag{19b}$$

$$\|\mathbf{U}_t \mathbf{W}_{t,\perp}\| \leq C_3 \cdot \alpha^{1 - \frac{2 \log(1 + \mu \hat{\sigma}_{s+1}^2)}{3 \log(1 + \mu \hat{\sigma}_1^2) - \log(1 + \mu \hat{\sigma}_{s+1}^2)}} \tag{19c}$$

$$\|\mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}\| \leq 200\delta \tag{19d}$$

Proof. We prove this lemma by applying [Corollary D.3](#) to $t = T_\alpha^{\text{SP}}$ defined in the previous lemma.

The inequality (19a) can be directly verified by using [Lemma C.3](#):

$$\|\mathbf{U}_t\| \leq \alpha (1 + \mu \hat{\sigma}_1^2)^t + \alpha^\gamma \leq \left(1 + \left(\frac{C_1(\mathbf{X}, \bar{\mathbf{U}}) \hat{\sigma}_1^2}{4r_*(1+\delta)} \right)^{\frac{1}{3}} \right) \cdot \alpha^{\gamma/3} \leq \hat{C}_1(\mathbf{X}, \bar{\mathbf{U}}) \cdot \alpha^{\gamma/3} \|\mathbf{X}\|.$$

where $\hat{C}_1(\mathbf{X}, \bar{\mathbf{U}}) = 1 + \left(\frac{C_1(\mathbf{X}, \bar{\mathbf{U}}) \|\mathbf{X}\|^2}{2r_*(1+\delta)} \right)^{\frac{1}{3}}$ (the constant C_1 is defined in the previous lemma). The last inequality holds when α is sufficiently small. For the remaining inequalities, we first verify that the assumption in [Corollary D.3](#):

$$\alpha \sigma_s(\mathbf{K}_t) > 10(\alpha \sigma_{s+1}(\mathbf{K}_t) + \|\mathbf{E}_t\|). \quad (20)$$

By definition of \mathbf{K}_t , we can see that for $\alpha \leq \left(\frac{\rho}{10C_1} \right)^{10\kappa_*}$,

$$\begin{aligned} \alpha \sigma_{s+1}(\mathbf{K}_t) + \|\mathbf{E}_t\| &\leq \alpha (1 + \mu \hat{\sigma}_{s+1}^2)^t + 4\hat{\sigma}_1^{-2} \alpha^3 r_*(1+\delta) (1 + \mu \hat{\sigma}_1^2)^{3t} \\ &\leq C_1(\mathbf{X}, \bar{\mathbf{U}}) \cdot \alpha^\gamma \leq 0.1\rho \alpha^{1 - \frac{2 \log(1 + \mu \hat{\sigma}_s^2)}{3 \log(1 + \mu \hat{\sigma}_1^2) - \log(1 + \mu \hat{\sigma}_{s+1}^2)}} \\ &\leq 0.1\alpha \sigma_s(\mathbf{K}_t) \end{aligned}$$

where $\|\mathbf{E}_{T_\alpha^{\text{SP}}}\|$ is bounded in the previous lemma. Hence (20) holds. Let \mathbf{L} be the span of top- s eigenvectors of \mathbf{M} , then by [Corollary D.3](#), at $t = T_\alpha^{\text{SP}}$ we have

$$\begin{aligned} \sigma_s(\mathbf{U}_t \mathbf{W}_t) &\geq 0.4\alpha \sigma_s(\mathbf{K}_t) \sigma_{\min}(\mathbf{V}_L^\top \bar{\mathbf{U}}) \\ &\geq 0.1\alpha \rho (1 + \mu \hat{\sigma}_s^2)^t \\ &= 0.1\rho \left(\frac{\rho \hat{\sigma}_1^2}{4r_*(1+\delta)} \right)^{\frac{\log(1 + \mu \hat{\sigma}_s^2)}{3 \log(1 + \mu \hat{\sigma}_1^2) - \log(1 + \mu \hat{\sigma}_{s+1}^2)}} \alpha^{1 - \frac{2 \log(1 + \mu \hat{\sigma}_s^2)}{3 \log(1 + \mu \hat{\sigma}_1^2) - \log(1 + \mu \hat{\sigma}_{s+1}^2)}} \\ &\geq 0.1\rho \underbrace{\left(\frac{\rho \hat{\sigma}_1^2}{8r_*} \right)^{\frac{1}{10\kappa_*}}}_{:=C_2(\mathbf{X}, \bar{\mathbf{U}})} \alpha^{1 - \frac{2 \log(1 + \mu \hat{\sigma}_s^2)}{3 \log(1 + \mu \hat{\sigma}_1^2) - \log(1 + \mu \hat{\sigma}_{s+1}^2)}} \\ \|\mathbf{U}_t \mathbf{W}_{t,\perp}\| &\leq 2\alpha \sigma_{s+1}^2(\mathbf{K}_t) + \|\mathbf{E}_t\| \\ &\leq \underbrace{2C_1(\mathbf{X}, \bar{\mathbf{U}})}_{:=C_3(\mathbf{X}, \bar{\mathbf{U}})} \alpha^{1 - \frac{2 \log(1 + \mu \hat{\sigma}_s^2)}{3 \log(1 + \mu \hat{\sigma}_1^2) - \log(1 + \mu \hat{\sigma}_{s+1}^2)}} \\ \|\mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}\| &\leq 100 \left(\delta + \frac{\alpha \sigma_{s+1}(\mathbf{K}_t) + \|\mathbf{E}_t\|}{\alpha \rho \sigma_s(\mathbf{K}_t)} \right) \\ &\leq 100 \left(\delta \alpha^{\frac{2 \log(1 + \mu \hat{\sigma}_s^2) - 2 \log(1 + \mu \hat{\sigma}_{s+1}^2)}{3 \log(1 + \mu \hat{\sigma}_1^2) - \log(1 + \mu \hat{\sigma}_{s+1}^2)}} \right) \\ &\leq 200\delta. \end{aligned} \quad (21)$$

The conclusion follows. \square

C.2. The parallel improvement phase

This subsection is devoted to proving [Lemma 5.2](#) which we recall below.

Lemma 5.2. *Suppose that [Assumptions 3.3](#), [3.5](#) and [3.7](#) hold and let $c_3 = 10^4 \kappa_* r_*^{\frac{1}{2}} \delta$. Then for sufficiently small α , the following inequalities hold when $T_\alpha^{\text{SP}} \leq t < T_{\alpha, s}^{\text{P1}}$:*

$$\begin{aligned}\sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1}) &\geq \sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \mathbf{W}_t) \\ &\geq (1 + 0.5\mu(\sigma_s^2 + \sigma_{s+1}^2)) \sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t),\end{aligned}\quad (13a)$$

$$\begin{aligned}\|\mathbf{U}_{t+1} \mathbf{W}_{t+1, \perp}\| &\leq (1 + \mu(0.4\sigma_s^2 + 0.6\sigma_{s+1}^2)) \|\mathbf{U}_t \mathbf{W}_{t, \perp}\|,\end{aligned}\quad (13b)$$

$$\|\mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{V}_{\mathbf{U}_{t+1} \mathbf{W}_{t+1}}\| \leq c_3, \quad (13c)$$

$$\text{rank}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1}) = \text{rank}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \mathbf{W}_t) = s. \quad (13d)$$

C.2.1. THE PARALLEL COMPONENT

In the following we bound $\sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \mathbf{W}_t)$. We state our main result of this section in the lemma below.

Lemma C.5. *Suppose that [Assumptions 3.3](#), [3.5](#) and [3.7](#) holds, $\|\mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}\| \leq c_3 < 10^{-2}\kappa^{-1}$ and $\Delta_t = (\mathcal{A}^* \mathcal{A} - \mathbf{I})(\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top)$ satisfies $\|\Delta_t\| \leq 0.2\kappa^{-1}r_*^{-\frac{1}{2}}\|\mathbf{X}\|^2$, then we have*

$$\begin{aligned}\sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1}) &\geq \sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \mathbf{W}_t) \\ &\geq [1 + \mu(\sigma_s^2 - 5c_3\|\mathbf{X}\|^2 - 2\|\Delta_t\|) - 20\mu^2\|\mathbf{X}\|^4] (1 - \mu\sigma_{\min}^2(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t)) \sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t).\end{aligned}$$

Proof. The update rule of GD implies that

$$\begin{aligned}\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \mathbf{W}_t &= \mathbf{V}_{\mathbf{X}_s}^\top (\mathbf{I} + \mu(\mathbf{X}_s \mathbf{X}_s^\top - \mathbf{U}_t \mathbf{U}_t^\top) + \mu\Delta_t) \mathbf{U}_t \mathbf{W}_t\end{aligned}\quad (22a)$$

$$= (\mathbf{I} + \mu\Sigma_s^2) \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t - \mu \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{U}_t \mathbf{W}_t + \mu \mathbf{V}_{\mathbf{X}_s}^\top \Delta_t \mathbf{U}_t \mathbf{W}_t \quad (22b)$$

$$\begin{aligned}&= (\mathbf{I} + \mu\Sigma_s^2) \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t - \mu \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_s} \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t - \mu \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_s, \perp} \mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{U}_t \mathbf{W}_t \\ &\quad + \mu \mathbf{V}_{\mathbf{X}_s}^\top \Delta_t \mathbf{U}_t \mathbf{W}_t \\ &= (\mathbf{I} + \mu\Sigma_s^2) \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t (\mathbf{I} - \mu \mathbf{W}_t^\top \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_s} \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t) + \mu \mathbf{V}_{\mathbf{X}_s}^\top \Delta_t \mathbf{U}_t \mathbf{W}_t \\ &\quad - \mu \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_s, \perp} \mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{U}_t \mathbf{W}_t + \mu^2 \Sigma_s^2 \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t \mathbf{W}_t^\top \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_s} \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t\end{aligned}\quad (22c)$$

where (22a) follows from $\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{X} \mathbf{X}^\top = \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{X}_s \mathbf{X}_s^\top + \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{X}_{s, \perp} \mathbf{X}_{s, \perp}^\top$ and $\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{X}_{s, \perp} = 0$; (22b) follows from $\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{X}_s \mathbf{X}_s^\top = \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{V}_{\mathbf{X}_s} \Sigma_s \mathbf{V}_{\mathbf{X}_s}^\top = \Sigma_s \mathbf{V}_{\mathbf{X}_s}^\top$, and (22c) follows from $\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t = \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t \mathbf{W}_t^\top + \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_{t, \perp} \mathbf{W}_{t, \perp}^\top = \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t \mathbf{W}_t^\top$ by definition of \mathbf{W}_t and $\mathbf{W}_{t, \perp}$.

We now relate the last three terms in (22c) to $\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t$. Since $\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t$ is invertible by [Assumption 3.5](#), $\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}$, $\Sigma_{\mathbf{U}_t \mathbf{W}_t}$ and $\mathbf{W}_{\mathbf{U}_t \mathbf{W}_t}$ are also of full rank, thus we have

$$\begin{aligned}\mathbf{U}_t \mathbf{W}_t &= \mathbf{U}_t \mathbf{W}_t (\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t)^{-1} \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t \\ &= \mathbf{U}_t \mathbf{W}_t (\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t} \Sigma_{\mathbf{U}_t \mathbf{W}_t} \mathbf{W}_{\mathbf{U}_t \mathbf{W}_t}^\top)^{-1} \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t \\ &= \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t} (\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t})^{-1} \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t.\end{aligned}\quad (23)$$

Plugging (23) into the second and third terms of (22) and re-arranging, we deduce that

$$\begin{aligned}\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \mathbf{W}_t &= (\mathbf{I} + \mu(\Sigma_s^2 + \mathbf{P}_1 + \mathbf{P}_2)) \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t (\mathbf{I} - \mu \mathbf{W}_t^\top \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_s} \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t) \\ &\quad + \mu^2 (\Sigma_s^2 + \mathbf{P}_1 + \mathbf{P}_2) \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t \mathbf{W}_t^\top \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_s} \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t \\ &= \left[\mathbf{I} + \mu (\Sigma_s^2 + \mathbf{P}_1 + \mathbf{P}_2) + \mu^2 (\Sigma_s^2 + \mathbf{P}_1 + \mathbf{P}_2) \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t \mathbf{W}_t^\top \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_s} (\mathbf{I} - \mu \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t \mathbf{W}_t^\top \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_s})^{-1} \right] \cdot \\ &\quad \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t (\mathbf{I} - \mu \mathbf{W}_t^\top \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_s} \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t)\end{aligned}\quad (24)$$

where we use the equation $\mathbf{A} = (\mathbf{I} - \mu\mathbf{A}\mathbf{A}^\top)^{-1}\mathbf{A}(\mathbf{I} - \mu\mathbf{A}^\top\mathbf{A})$ with $\mathbf{A} = \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t$ (when $\mu < \frac{1}{9\|\mathbf{X}\|^2}$, $\mathbf{I} - \mu\mathbf{A}\mathbf{A}^\top$ is invertible by [Lemma D.4](#)), and

$$\begin{aligned} \mathbf{P}_1 &= \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_s, \perp} \mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t} (\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t})^{-1} \\ \mathbf{P}_2 &= \mathbf{V}_{\mathbf{X}_s}^\top \Delta_t \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t} (\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t})^{-1} \end{aligned} \quad (25)$$

By assumption we have

$$\sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}) \geq \sqrt{1 - \|\mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}\|^2} \geq \frac{1}{2},$$

so that

$$\|\mathbf{P}_1\| \leq \|\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_s, \perp}\| \cdot \|\mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}\| \cdot \|(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t})^{-1}\| \leq 5c_3 \|\mathbf{X}\|^2 \leq 0.1 \|\mathbf{X}\|^2 \quad (26)$$

and by our assumption we have

$$\|\mathbf{P}_2\| \leq \|(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t})^{-1}\| \cdot \|\Delta_t\| \leq 2\|\Delta_t\| \leq 0.2\kappa^{-1}r_*^{-\frac{1}{2}} \|\mathbf{X}\|^2. \quad (27)$$

Moreover, note that $\|\Sigma_s\|^2 = \sigma_1^2 = \|\mathbf{X}\|^2$, and since $\mu < 10^{-4}\|\mathbf{X}\|^{-2}$ by [Assumption 3.7](#), we have $\|(\mathbf{I} - \mu\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t \mathbf{W}_t^\top \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_s})^{-1}\| < 1.1$. Thus

$$\|(\Sigma_s^2 + \mathbf{P}_1 + \mathbf{P}_2) \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t \mathbf{W}_t^\top \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_s} (\mathbf{I} - \mu\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t \mathbf{W}_t^\top \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_s})^{-1}\| \leq 20\|\mathbf{X}\|^4.$$

The equation (25) implies that

$$\begin{aligned} &\sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \mathbf{W}_t) \\ &\geq \sigma_{\min} \left(\mathbf{I} + \mu (\Sigma_s^2 + \mathbf{P}_1 + \mathbf{P}_2) + (\Sigma_s^2 + \mathbf{P}_1 + \mathbf{P}_2) \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t \mathbf{W}_t^\top \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_s} (\mathbf{I} - \mu\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t \mathbf{W}_t^\top \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_s})^{-1} \right) \cdot \\ &\quad \sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t (\mathbf{I} - \mu\mathbf{W}_t^\top \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_s} \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t)) \\ &\geq (1 + \mu\sigma_{\min}^2(\Sigma_s) - \mu\|\mathbf{P}_1\| - \mu\|\mathbf{P}_2\| - 20\mu^2\|\mathbf{X}\|^4) \sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t) (1 - \mu\sigma_{\min}^2(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t)) \\ &= (1 + \mu\sigma_s^2 - \mu\|\mathbf{P}_1\| - \mu\|\mathbf{P}_2\| - 20\mu^2\|\mathbf{X}\|^4) \sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t) (1 - \mu\sigma_{\min}^2(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t)) \end{aligned}$$

Recall that \mathbf{P}_1 and \mathbf{P}_2 are bounded in (26) and (27) respectively, so we have that

$$\begin{aligned} &\sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1}) \\ &\geq \sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \mathbf{W}_t) \\ &\geq [1 + \mu(\sigma_s^2 - 5c_3\|\mathbf{X}\|^2 - 2\|\Delta_t\|^2) - 20\mu^2\|\mathbf{X}\|^4] (1 - \mu\sigma_{\min}^2(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t)) \sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t). \end{aligned}$$

The conclusion follows. \square

The corollaries below immediately follow from [Lemma C.5](#).

Corollary C.6. *Under the conditions in [Lemma C.5](#), if $\sigma_{\min}^2(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t) < 0.3\kappa^{-1}\|\mathbf{X}\|^2$, then we have*

$$\sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1}) \geq (1 + 0.5\mu(\sigma_s^2 + \sigma_{s+1}^2)) \sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t).$$

Proof. By [Lemma C.5](#) it remains to check that

$$[1 + \mu(\sigma_s^2 - 5c_3\|\mathbf{X}\|^2 - 2\|\Delta_t\|^2) - 20\mu^2\|\mathbf{X}\|^4] (1 - 0.3\mu\kappa^{-1}\|\mathbf{X}\|^2) \geq 1 + 0.5\mu(\sigma_s^2 + \sigma_{s+1}^2).$$

Indeed, recall from the conditions of [Lemma C.5](#) that $5c_3\|\mathbf{X}\|^2 \leq 0.01\kappa^{-1}\|\mathbf{X}\|^2 \leq 0.01(\sigma_s^2 - \sigma_{s+1}^2)$ and similarly $\|\Delta_t\| \leq 0.005(\sigma_s^2 - \sigma_{s+1}^2)$ and $\mu^2\|\mathbf{X}\|^4 \leq 10^{-4}\kappa^{-1}\mu\|\mathbf{X}\|^2 \leq 10^{-4}(\sigma_s^2 - \sigma_{s+1}^2)$, so that

$$\begin{aligned} &[1 + \mu(\sigma_s^2 - 5c_3\|\mathbf{X}\|^2 - 2\|\Delta_t\|^2) - 20\mu^2\|\mathbf{X}\|^4] (1 - 0.3\mu\kappa^{-1}\|\mathbf{X}\|^2) \\ &\geq (1 + \mu(0.9\sigma_s^2 + 0.1\sigma_{s+1}^2)) (1 - 0.3\mu(\sigma_s^2 - \sigma_{s+1}^2)) \\ &= 1 + \mu(0.6\sigma_s^2 + 0.4\sigma_{s+1}^2) - \mu^2\|\mathbf{X}\|^2(\sigma_s^2 - \sigma_{s+1}^2) \\ &\geq 1 + 0.5\mu(\sigma_s^2 + \sigma_{s+1}^2) \end{aligned}$$

as desired. \square

Corollary C.7. Under the conditions in [Lemma C.5](#), if

$$\sigma_{\min}^2(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t) \leq \sigma_s^2 - \mu\sigma_s^4 - 5c_3\|\mathbf{X}\|^2 - 2\|\Delta_t\|^2 - 20\mu^2\|\mathbf{X}\|^4, \quad (28)$$

then we have that $\sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1}) \geq \sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t)$.

Proof. A sufficient condition for $\sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1}) \geq \sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t)$ to hold is that

$$\begin{aligned} & [1 + \mu(\sigma_s^2 - 5c_3\|\mathbf{X}\|^2 - 2\|\Delta_t\|^2) - 20\mu^2\|\mathbf{X}\|^4] (1 - \mu\sigma_{\min}^2(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t)) \\ & \Leftrightarrow \sigma_s^2 - 5c_3\|\mathbf{X}\|^2 - 2\|\Delta_t\|^2 - 20\mu\|\mathbf{X}\|^2 - \sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t) - \mu\sigma_s^2\sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t) \geq 0. \end{aligned}$$

When (28) holds, we have

$$\begin{aligned} & \sigma_s^2 - 5c_3\|\mathbf{X}\|^2 - 2\|\Delta_t\|^2 - 20\mu\|\mathbf{X}\|^2 - \sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t) \\ & \geq \mu\sigma_s^4 \geq \mu\sigma_s^2\sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t) \end{aligned}$$

as desired. \square

C.2.2. THE ORTHOGONAL COMPONENT

In this section we turn to analyze the noise term. The main result of this section is presented in the following:

Lemma C.8. Suppose that [Assumptions 3.3](#), [3.5](#) and [3.7](#) hold, $\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \mathbf{W}_t \in \mathbb{R}^{s \times r}$ is of full rank, $\|\mathbf{V}_{\mathbf{X}_s, \perp} \mathbf{V}_{\mathbf{U}_t} \mathbf{W}_t\| \leq c_3 < 10^{-2}\kappa^{-1}$ and $\|\Delta_t\| \leq c_3\|\mathbf{X}\|^2$, then we have

$$\|\mathbf{U}_{t+1} \mathbf{W}_{t+1, \perp}\| \leq (1 + \mu\sigma_{s+1}^2 + 30\mu\|\mathbf{X}\|^2 c_3 + 0.1\mu^2\|\mathbf{X}\|^4) \|\mathbf{U}_t \mathbf{W}_{t, \perp}\|.$$

Proof. By the definition of $\mathbf{W}_{t, \perp}$, we have $\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_{t, \perp} = 0$, thus $\|\mathbf{U}_t \mathbf{W}_{t, \perp}\| = \|\mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{U}_t \mathbf{W}_{t, \perp}\|$. The latter can be decomposed as follows:

$$\mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{U}_{t+1} \mathbf{W}_{t+1, \perp} = \underbrace{\mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{U}_{t+1} \mathbf{W}_t \mathbf{W}_t^\top \mathbf{W}_{t+1, \perp}}_{=(a)} + \underbrace{\mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{U}_{t+1} \mathbf{W}_{t, \perp} \mathbf{W}_{t, \perp}^\top \mathbf{W}_{t+1, \perp}}_{=(b)}.$$

In the following, we are going to show that the term (a) is bounded by $c \cdot \mu$ where c is a small constant, while (b) grows linearly with a slow speed.

Bounding summand (a). Since

$$0 = \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \mathbf{W}_{t+1, \perp} = \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \mathbf{W}_t \mathbf{W}_t^\top \mathbf{W}_{t+1, \perp} + \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \mathbf{W}_{t, \perp} \mathbf{W}_{t, \perp}^\top \mathbf{W}_{t+1, \perp}$$

by definition, we have

$$\mathbf{W}_t^\top \mathbf{W}_{t+1, \perp} = -(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \mathbf{W}_t)^{-1} \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \mathbf{W}_{t, \perp} \mathbf{W}_{t, \perp}^\top \mathbf{W}_{t+1, \perp}. \quad (29)$$

Thus the summand (a) can be rewritten as follows:

$$\begin{aligned} & \mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{U}_{t+1} \mathbf{W}_t \mathbf{W}_t^\top \mathbf{W}_{t+1, \perp} \\ & = -\mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{U}_{t+1} \mathbf{W}_t (\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \mathbf{W}_t)^{-1} \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \mathbf{W}_{t, \perp} \mathbf{W}_{t, \perp}^\top \mathbf{W}_{t+1, \perp} \end{aligned} \quad (30a)$$

$$\begin{aligned} & = -\mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{U}_{t+1} \mathbf{W}_t (\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{V}_{\mathbf{U}_{t+1}} \mathbf{W}_t \Sigma_{\mathbf{U}_{t+1}} \mathbf{W}_t \mathbf{W}_{\mathbf{U}_{t+1}}) \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \mathbf{W}_{t, \perp} \mathbf{W}_{t, \perp}^\top \mathbf{W}_{t+1, \perp} \\ & = -\mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{V}_{\mathbf{U}_{t+1}} \mathbf{W}_t (\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{V}_{\mathbf{U}_{t+1}} \mathbf{W}_t)^{-1} \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \mathbf{W}_{t, \perp} \mathbf{W}_{t, \perp}^\top \mathbf{W}_{t+1, \perp} \end{aligned} \quad (30b)$$

$$\begin{aligned} & = -\mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{V}_{\mathbf{U}_{t+1}} \mathbf{W}_t (\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{V}_{\mathbf{U}_{t+1}} \mathbf{W}_t)^{-1} \mathbf{V}_{\mathbf{X}_s}^\top (\mathbf{I} + \mu \mathbf{A}^* \mathbf{A} (\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top)) \mathbf{U}_t \mathbf{W}_{t, \perp} \mathbf{W}_{t, \perp}^\top \mathbf{W}_{t+1, \perp} \\ & = -\mu \mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{V}_{\mathbf{U}_{t+1}} \mathbf{W}_t (\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{V}_{\mathbf{U}_{t+1}} \mathbf{W}_t)^{-1} \mathbf{V}_{\mathbf{X}_s}^\top [(\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top) + \Delta_t] \mathbf{U}_t \mathbf{W}_{t, \perp} \mathbf{W}_{t, \perp}^\top \mathbf{W}_{t+1, \perp} \end{aligned} \quad (30c)$$

$$\begin{aligned} & = \mu \mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{V}_{\mathbf{U}_{t+1}} \mathbf{W}_t (\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{V}_{\mathbf{U}_{t+1}} \mathbf{W}_t)^{-1} \mathbf{V}_{\mathbf{X}_s}^\top [\mathbf{U}_t \mathbf{U}_t^\top - \Delta_t] \mathbf{U}_t \mathbf{W}_{t, \perp} \mathbf{W}_{t, \perp}^\top \mathbf{W}_{t+1, \perp} \\ & = \mu \mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{V}_{\mathbf{U}_{t+1}} \mathbf{W}_t (\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{V}_{\mathbf{U}_{t+1}} \mathbf{W}_t)^{-1} \mathbf{M}_1 \mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{U}_t \mathbf{W}_{t, \perp} \mathbf{W}_{t, \perp}^\top \mathbf{W}_{t+1, \perp}, \end{aligned}$$

where $M_1 = \mathbf{V}_{\mathbf{X}_s}^\top [\mathbf{U}_t \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_{s,\perp}} - \Delta_t \mathbf{V}_{\mathbf{X}_{s,\perp}}]$. In (30), (30a) follows from (29), (30b) holds since $\Sigma_{\mathbf{U}_{t+1} \mathbf{W}_t} \mathbf{W}_{\mathbf{U}_{t+1} \mathbf{W}_t}^\top \in \mathbb{R}^{s \times s}$ is invertible, and in (30c) we use $\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_{t,\perp} = 0$. It follows that

$$\|(a)\| \leq \mu \|\mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{V}_{\mathbf{U}_{t+1} \mathbf{W}_t}\| \cdot \left\| (\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{V}_{\mathbf{U}_{t+1} \mathbf{W}_t})^{-1} \right\| \|\mathbf{M}_1\| \|\mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{W}_{t,\perp}\|. \quad (31)$$

By Lemma D.5 we have $\left\| \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{V}_{\mathbf{U}_{t+1} \mathbf{W}_t} \right\| \leq 0.01$, which implies that

$$\left\| (\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{V}_{\mathbf{U}_{t+1} \mathbf{W}_t})^{-1} \right\| = \sigma_{\min}^{-1}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{V}_{\mathbf{U}_{t+1} \mathbf{W}_t}) = \left(1 - \|\mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{V}_{\mathbf{U}_{t+1} \mathbf{W}_t}\|^2\right)^{-\frac{1}{2}} \leq 1.1. \quad (32)$$

Lastly, we bound M_1 as follows:

$$\begin{aligned} \|\mathbf{M}_1\| &\leq \|\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_{s,\perp}}\| + \|(\mathcal{A}^* \mathcal{A} - \mathbf{I})(\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top)\| \\ &\leq \|\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t\| \cdot \|\mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{W}_t\| + 10^{-3} \kappa^{-1} c_3 \|\mathbf{X}\|^2 \\ &\leq 10 \|\mathbf{X}\|^2 c_3. \end{aligned} \quad (33)$$

where the second inequality follows from our assumption on $\|(\mathcal{A}^* \mathcal{A} - \mathbf{I})(\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top)\|$. Combining (31), (32) and (33) yields

$$\|(a)\| \leq 20\mu \|\mathbf{X}\|^2 c_3 \|\mathbf{U}_t \mathbf{W}_{t,\perp}\|.$$

Bounding summand (b). This is the main component in the error term. We'll see that although this term can grow exponentially fast, the growth speed is slower than the minimal eigenvalue of the parallel component.

We have

$$\begin{aligned} &\mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_{t+1} \mathbf{W}_{t,\perp} \\ &= \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top [\mathbf{I} + \mu(\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top) + \mu(\mathcal{A}^* \mathcal{A} - \mathbf{I})(\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top)] \mathbf{U}_t \mathbf{W}_{t,\perp} \end{aligned} \quad (34a)$$

$$= \left(\mathbf{I} + \mu \Sigma_{s,\perp}^2 - \mu \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_{s,\perp}} + \underbrace{\mu \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \Delta_t \mathbf{V}_{\mathbf{X}_{s,\perp}}}_{=: M_2} \right) \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{W}_{t,\perp} \quad (34b)$$

$$= (\mathbf{I} + \mu \Sigma_{s,\perp}^2 - \mu \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{W}_t \mathbf{W}_t^\top \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_{s,\perp}} + \mu M_2) \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{W}_{t,\perp} (\mathbf{I} - \mu \mathbf{W}_{t,\perp}^\top \mathbf{U}_t^\top \mathbf{U}_t \mathbf{W}_{t,\perp}) \quad (34c)$$

$$+ \mu^2 (\Sigma_{s,\perp}^2 - \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{W}_t \mathbf{W}_t^\top \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_{s,\perp}} + M_2) \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{W}_{t,\perp} \mathbf{W}_{t,\perp}^\top \mathbf{U}_t^\top \mathbf{U}_t \mathbf{W}_{t,\perp} \quad (34d)$$

where we recall that $\Sigma_{s,\perp}^2 = \text{diag}(\sigma_{s+1}^2, \dots, \sigma_r^2, 0, \dots, 0) \in \mathbb{R}^{(d-s) \times (d-s)}$. In (34), (34a) follows from the update rule of GD, (34b) is obtained from $\mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{X} \mathbf{X}^\top = \Sigma_{s,\perp}^2 \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top$ and $\mathbf{U}_t \mathbf{W}_{t,\perp} = \mathbf{V}_{\mathbf{X}_s} \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_{t,\perp} + \mathbf{V}_{\mathbf{X}_{s,\perp}} \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{W}_{t,\perp} = \mathbf{V}_{\mathbf{X}_{s,\perp}} \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{W}_{t,\perp}$, and lastly in (34d) we use

$$\begin{aligned} &\mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_{s,\perp}} \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{W}_{t,\perp} \\ &= \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{W}_t \mathbf{W}_t^\top \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_{s,\perp}} \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{W}_{t,\perp} + \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{W}_{t,\perp} \mathbf{W}_{t,\perp}^\top \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_{s,\perp}} \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{W}_{t,\perp} \\ &= \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{W}_t \mathbf{W}_t^\top \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_{s,\perp}} \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{W}_{t,\perp} + \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{W}_{t,\perp} \mathbf{W}_{t,\perp}^\top \mathbf{U}_t^\top \mathbf{U}_t \mathbf{W}_{t,\perp}. \end{aligned}$$

It follows that

$$\begin{aligned} &\|\mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_{t+1} \mathbf{W}_{t,\perp}\| \\ &\leq (\|\mathbf{I} - \mu \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{W}_t \mathbf{W}_t^\top \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_{s,\perp}}\| + \mu \|\Sigma_{s,\perp}\|^2 + \mu \|M_2\|) \|\mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{W}_{t,\perp}\| (\mathbf{I} - \mu \|\mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t \mathbf{W}_{t,\perp}\|^2) \\ &\quad + \mu^2 \|\mathbf{U}_t \mathbf{W}_{t,\perp}\|^3 (\sigma_{s+1}^2 + \|\mathbf{U}_t\|^2 + 10^{-3} \kappa^{-1} c_3 \|\mathbf{X}\|^2) \\ &\leq (1 + \mu \sigma_{s+1}^2 + \mu \|\Delta_t\|) \|\mathbf{U}_t \mathbf{W}_{t,\perp}\| (1 - \mu \|\mathbf{U}_t \mathbf{W}_{t,\perp}\|^2) + 0.1 \mu^2 \|\mathbf{X}\|^4 \|\mathbf{U}_t \mathbf{W}_{t,\perp}\| \\ &\leq \|\mathbf{U}_t \mathbf{W}_{t,\perp}\| (1 + \mu \sigma_{s+1}^2 + \mu c_3 \|\mathbf{X}\|^2 + 0.1 \mu^2 \|\mathbf{X}\|^4) \end{aligned}$$

To summarize, we have

$$\|U_{t+1}W_{t+1,\perp}\| \leq (1 + \mu\sigma_{s+1}^2 + 30\mu\|X\|^2c_3 + 0.1\mu^2\|X\|^4) \|U_tW_{t,\perp}\|$$

as desired. \square

To bound the growth speed of the orthogonal component, we need to show that the quantity $\|V_{X_s,\perp}^\top V_{U_tW_t}\|$ remains small.

The following lemma serves to complete an induction step from t to $t+1$:

Lemma C.9. *Suppose $V_{X_s}^\top U_t$ is of full rank, $\|V_{X_s,\perp}^\top V_{U_tW_t}\| \leq c_3$ and $\|U_tW_{t,\perp}\| \leq \min\{\sigma_{\min}(U_tW_t), c_4\}$ with $\max\{c_3, c_4\|X\|^{-1}\} \leq 10^{-2}\kappa^{-1}$, and $\Delta_t = (\mathcal{A}^*\mathcal{A} - I)(XX^\top - U_tU_t^\top)$ satisfies $\|\Delta_t\| \leq 10^{-3}\kappa^{-1}c_3\|X\|^2$ and $\mu \leq 10^{-4}\kappa^{-1}\|X\|^{-2}c_3$, then we have $\|V_{X_s,\perp}^\top V_{U_{t+1}W_{t+1}}\| \leq c_3$.*

Proof. Let $M_t = \mathcal{A}^*\mathcal{A}(XX^\top - U_tU_t^\top)$, so the update rule of GD implies that

$$\begin{aligned} U_{t+1}W_{t+1} &= (I + \mu M_t)U_tW_{t+1} \\ &= (I + \mu M_t)(U_tW_tW_t^\top W_{t+1} + U_tW_{t,\perp}W_{t,\perp}^\top W_{t+1}) \\ &= (I + \mu M_t)(V_{U_tW_t}V_{U_tW_t}^\top U_tW_tW_t^\top W_{t+1} + U_tW_{t,\perp}W_{t,\perp}^\top W_{t+1}) \\ &= \underbrace{(I + \mu M_t)(I + P)}_{:=H} V_{U_tW_t}^\top V_{U_tW_t} U_tW_tW_t^\top W_{t+1}, \end{aligned}$$

where

$$P = U_tW_{t,\perp}W_{t,\perp}^\top W_{t+1}(V_{U_tW_t}^\top U_tW_tW_t^\top W_{t+1})^{-1}V_{U_tW_t}^\top$$

and $V_{U_tW_t}^\top U_tW_tW_t^\top W_{t+1}$ is invertible since $V_{U_tW_t}^\top U_tW_t$ is invertible by our assumption that $V_{X_s}^\top U_t$ is of full rank and $\text{rank}(U_tW_t) \geq \text{rank}(V_{X_s}^\top U_tW_t) = \text{rank}(V_{X_s}^\top U_t) = s$, and $W_t^\top W_{t+1}$ is invertible by Lemma D.7. Indeed, Lemma D.7 implies that $\sigma_{\min}(W_t^\top W_{t+1}) \geq \frac{1}{2}$ by our condition on μ .

The key observation here is that because the (square) matrix $V_{U_tW_t}^\top U_tW_tW_t^\top W_{t+1}$ is invertible, so that the column space of $U_{t+1}W_{t+1}$ is the same as that of H . Following the line of proof of Stöger & Soltanolkotabi, 2021, Lemma 9.3 (for completeness, we provide details in Lemma D.8), we deduce that

$$\begin{aligned} &\|V_{X_s,\perp}^\top V_{U_{t+1}W_{t+1}}\| = \|V_{X_s,\perp}^\top V_H W_H^\top\| \\ &\leq \left\| V_{X_s,\perp}^\top \left[\left(I + B - \frac{1}{2}V_{U_tW_t}V_{U_tW_t}^\top (B + B^\top) \right) V_{U_tW_t} - BV_{U_tW_t}V_{U_tW_t}^\top (B + B^\top) V_{U_tW_t} + D \right] \right\| \\ &\leq \left\| V_{X_s,\perp}^\top \left(I + B - \frac{1}{2}V_{U_tW_t}V_{U_tW_t}^\top (B + B^\top) \right) V_{U_tW_t} \right\| + 2\|B\|^2 + \|D\| \end{aligned} \quad (35)$$

where $B = (I + \mu M_t)(I + P) - I$ and $\|D\| \leq 100\|B\|^2$. By assumption we have

$$\begin{aligned} \|P\| &\leq \frac{\|U_tW_{t,\perp}\| \|W_{t,\perp}W_{t+1}\|}{\sigma_{\min}(U_tW_t)\sigma_{\min}(W_t^\top W_{t+1})} \\ &\leq 2\|W_{t,\perp}W_{t+1}\|, \end{aligned}$$

so that

$$\begin{aligned} &\|B - \mu(XX^\top - U_tU_t^\top)\| \\ &\leq \mu\|M_t - (XX^\top - U_tU_t^\top)\| + \|P\| + \mu\|M_t\|\|P\| \\ &\leq \mu\|\Delta_t\| + 2\|W_{t,\perp}W_{t+1}\| + 4\mu\|X\|^2\|W_{t,\perp}W_{t+1}\| \\ &\leq \mu\|\Delta_t\| + 6\|W_{t,\perp}W_{t+1}\| \\ &\leq 18\mu(10\mu\|X\|^3 + c_4)c_3\|X\| + 7\mu\|\Delta_t\| \\ &\leq 18\mu(10\mu\|X\|^3 + c_4)c_3\|X\| + 0.01\mu\kappa^{-1}c_3\|X\|^2 \end{aligned} \quad (36)$$

where we use Lemma D.7 to bound $\|W_{t,\perp}W_{t+1}\|$. Let $B_1 = \mu(XX^\top - U_tU_t^\top)$ and $R_1 = V_{X_s,\perp}^\top (I + B_1 - V_{U_tW_t}V_{U_tW_t}^\top B_1) V_{U_tW_t}$, then we have

$$\begin{aligned}
 \mathbf{R}_1 &= \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top (\mathbf{I} + \mu (\mathbf{I} - \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t} \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}^\top) (\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top)) \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t} \\
 &= (\mathbf{I} + \mu \Sigma_{s,\perp}^2) \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t} (\mathbf{I} - \mu \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}^\top \mathbf{X} \mathbf{X}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}) \\
 &\quad - \mu \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top (\mathbf{I} - \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t} \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}^\top) \mathbf{U}_t \mathbf{W}_{t,\perp} \mathbf{W}_{t,\perp}^\top \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_{s,\perp}} \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t} \\
 &\quad + \mu^2 \Sigma_{s,\perp}^2 \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t} \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}^\top \mathbf{X} \mathbf{X}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}.
 \end{aligned} \tag{37}$$

By Weyl's inequality (cf. [Lemma A.4](#)) and our assumption on c_3 ,

$$\begin{aligned}
 \sigma_{\min}(\mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}^\top \mathbf{X} \mathbf{X}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}) &\geq \sigma_{\min}(\mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}^\top \mathbf{X}_s \mathbf{X}_s^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}) - \|\mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}^\top \mathbf{X}_{s,\perp} \mathbf{X}_{s,\perp}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}\|^2 \\
 &\geq \sigma_{\min}(\mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}^\top \mathbf{X}_s \mathbf{X}_s^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}) - \sigma_{s+1}^2 \|\mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}\|^2 \\
 &\geq \sigma_s^2 \|\mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}^\top \mathbf{V}_{\mathbf{X}_s}\|^2 - \sigma_{s+1}^2 c_3^2 \\
 &= \sigma_s^2 - (\sigma_s^2 + \sigma_{s+1}^2) c_3^2 > \frac{1}{2} (\sigma_s^2 + \sigma_{s+1}^2).
 \end{aligned}$$

So we have

$$\|\mathbf{R}_1\| \leq \left(1 - \frac{\mu}{2} (\sigma_s^2 - \sigma_{s+1}^2)\right) \|\mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}\| + \mu c_3 c_4^2 + \mu^2 \|\mathbf{X}\|^4.$$

It thus follows from (35) that

$$\begin{aligned}
 &\|\mathbf{V}_{\mathbf{X}_s^\perp}^\top \mathbf{V}_{\mathbf{U}_{t+1} \mathbf{W}_{t+1}}\| \\
 &\leq \|\mathbf{R}_1\| + 2\|\mathbf{B} - \mathbf{B}_1\| + 102\|\mathbf{B}\|^2 \\
 &\leq \left(1 - \frac{\mu}{2} (\sigma_s^2 - \sigma_{s+1}^2)\right) \|\mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}\| + 40\mu c_3 c_4 \|\mathbf{X}\| + 0.02\mu \kappa^{-1} c_3 \|\mathbf{X}\|^2 + 10^3 \mu^2 \|\mathbf{X}\|^4.
 \end{aligned}$$

Since $\|\mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}\| \leq c_3$, it follows from our assumption on c_3, c_4 and μ that $\|\mathbf{V}_{\mathbf{X}_s^\perp}^\top \mathbf{V}_{\mathbf{U}_{t+1} \mathbf{W}_{t+1}}\| \leq c_3$ as well, which concludes the proof. \square

C.3. Induction

Let

$$T_{\alpha,s}^{\text{pi}} = \min \{t \geq 0 : \sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{\alpha,t+1}) > 0.3\kappa^{-1} \|\mathbf{X}\|^2\}.$$

where pi stands for the parallel improvement phase. In this section, we show that when $T_\alpha^{\text{sp}} \leq t < T_{\alpha,s}^{\text{pi}}$, the parallel component grows exponentially faster than the orthogonal component. We prove this via induction and the base case is already shown in [Lemma C.4](#).

Lemma C.10 ([Lemma 5.2](#), detailed version). *Suppose that [Assumptions 3.3, 3.5](#) and [3.7](#) hold and let $c_3 = 10^4 \kappa \sqrt{r_*} \delta, c_4 \leq 10^{-3} \kappa^{-1} \|\mathbf{X}\|$. Then the following holds for all $T_\alpha^{\text{sp}} \leq t < T_{\alpha,s}^{\text{pi}}$ as long as $\alpha \leq C_4(\mathbf{X}, \bar{\mathbf{U}}) = \left(\kappa \frac{C_2(\mathbf{X}, \bar{\mathbf{U}})^2}{C_3(\mathbf{X}, \bar{\mathbf{U}})}\right)^{-2\kappa}$ is sufficiently small:*

$$\sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1}) \geq \sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \mathbf{W}_t) \geq (1 + 0.5\mu (\sigma_s^2 + \sigma_{s+1}^2)) \sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{\alpha,t}) \tag{38a}$$

$$\|\mathbf{U}_{t+1} \mathbf{W}_{t+1,\perp}\| \leq \min \{(1 + \mu (0.4\sigma_s^2 + 0.6\sigma_{s+1}^2)) \|\mathbf{U}_t \mathbf{W}_{t,\perp}\|, c_4\} \tag{38b}$$

$$\|\mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{V}_{\mathbf{U}_{t+1} \mathbf{W}_{t+1}}\| \leq c_3. \tag{38c}$$

$$\text{rank}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1}) = \text{rank}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \mathbf{W}_t) = s. \tag{38d}$$

Proof. The base case $t = T_{\alpha,s}^{\text{pi}}$ is already proved in (19). Now suppose that the lemma holds for t , we now show that it holds for $t+1$ as well.

To begin with, we bound the term $\|\Delta_t\|$ as follows:

$$\begin{aligned}
 \|\Delta_t\| &= \|(\mathcal{A}^* \mathcal{A} - \mathbf{I})(\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top)\| \\
 &\leq \|(\mathcal{A}^* \mathcal{A} - \mathbf{I})(\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{W}_t \mathbf{W}_t^\top \mathbf{U}_t^\top)\| + \|(\mathcal{A}^* \mathcal{A} - \mathbf{I}) \mathbf{U}_t \mathbf{W}_{t,\perp} \mathbf{W}_{t,\perp}^\top \mathbf{U}_t^\top\| \\
 &\leq 10\delta \sqrt{r_*} \|\mathbf{X}\|^2 + \delta \|\mathbf{U}_t \mathbf{W}_{t,\perp} \mathbf{W}_{t,\perp}^\top \mathbf{U}_t^\top\|_* \\
 &\leq 10\delta \sqrt{r_*} \|\mathbf{X}\|^2 + \delta d \|\mathbf{U}_t \mathbf{W}_{t,\perp}\|^2
 \end{aligned} \tag{39}$$

where in the second inequality we use [Proposition A.2](#) and [Lemma A.3](#) and in the third inequality we use $\|\mathbf{A}\|_* \leq \sqrt{d}\|\mathbf{A}\|$, $\forall \mathbf{A} \in \mathbb{R}^{d \times d}$ and the induction hypotheses.

By induction hypothesis, there exists a constant $\hat{C}_4(\mathbf{X}, \bar{\mathbf{U}}) = \frac{C_2(\mathbf{X}, \bar{\mathbf{U}})}{C_3(\mathbf{X}, \bar{\mathbf{U}})}$ (see [Lemma C.4](#)) such that

$$\frac{\sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t)}{\|\mathbf{U}_t \mathbf{W}_{t,\perp}\|} \geq \frac{\sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{T_\alpha^{\text{sp}}})}{\|\mathbf{U}_{T_\alpha^{\text{sp}}} \mathbf{W}_{T_\alpha^{\text{sp}},\perp}\|} \geq \hat{C}_4 \cdot \alpha^{-\gamma_s} \quad (40)$$

where

$$\gamma_s = \frac{2(\log(1 + \mu\hat{\sigma}_s^2) - \log(1 + \mu\hat{\sigma}_{s+1}^2))}{3\log(1 + \mu\hat{\sigma}_1^2) - \log(1 + \mu\hat{\sigma}_{s+1}^2)} \geq \frac{1}{4\kappa}.$$

Since we must have $\sigma_{\min}^2(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t) \leq 0.3\kappa^{-1}\|\mathbf{X}\|^2$ by definition of $T_{\alpha,s}^{\text{pi}}$, it follows that $\|\mathbf{U}_t \mathbf{W}_{t,\perp}\|^2 \leq 10\kappa\|\mathbf{X}\|^2 \hat{C}_4^2 \alpha^{\frac{1}{2\kappa}}$, so for $\alpha \leq (\hat{C}_4^{-2}\kappa)^{-2\kappa}$, $\|\Delta_t\| \leq 11\delta\sqrt{r_*}\|\mathbf{X}\|^2$ holds.

The above inequality combined with our assumption on δ implies that the conditions on $\|\Delta_t\|$ in [Lemmas C.5](#), [C.8](#) and [C.9](#) hold. We now show that [\(38a\)](#) to [\(38d\)](#) hold for $t+1$, which completes the induction step.

First, since $t < T_{\alpha,s}^{\text{pi}}$, we have $\sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1}) \leq \kappa^{-1}\|\mathbf{X}\|^2$. Moreover, the induction hypothesis implies that $\|\mathbf{V}_{\mathbf{X}_s,\perp}^\top \mathbf{V}_{\mathbf{U}_{t-1}\mathbf{W}_{t-1}}\| \leq c_3$ and that $\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{\alpha,t}$ is of full rank. Thus the conditions of [Corollary C.6](#) are all satisfied, and we deduce that [\(38a\)](#) holds.

Second, the assumptions on c_3, c_4 and δ , combined with [Lemma C.8](#), immediately implies

$$\|\mathbf{U}_{t+1} \mathbf{W}_{t+1,\perp}\| \leq (1 + \mu(0.4\sigma_s^2 + 0.6\sigma_{s+1}^2)) \|\mathbf{U}_t \mathbf{W}_{t,\perp}\|.$$

As a result, similar to [\(40\)](#) we observe that

$$\frac{\sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1})}{\|\mathbf{U}_{t+1} \mathbf{W}_{t+1,\perp}\|} \geq \frac{\sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{T_\alpha^{\text{sp}}})}{\|\mathbf{U}_{T_\alpha^{\text{sp}}} \mathbf{W}_{T_\alpha^{\text{sp}},\perp}\|} \geq \hat{C}_4 \cdot \alpha^{-\frac{1}{4\kappa}}.$$

Since $\sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1}) \leq \|\mathbf{X}\|$, when α is sufficiently small we must have that $\|\mathbf{U}_{t+1} \mathbf{W}_{t+1,\perp}\| \leq c_4$.

Finally, [Lemma C.9](#) implies that [\(38c\)](#) is true, and [\(38d\)](#) follows from our application of [Lemma C.5](#). This concludes the proof. \square

C.4. The refinement phase and concluding the proof of [Lemma 4.3](#)

We have shown that the parallel component $\sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1})$ grows exponentially faster than the orthogonal component $\|\mathbf{U}_t \mathbf{W}_{t,\perp}\|$. In this section, we characterize the GD dynamics *after* $T_{\alpha,s}^{\text{pi}}$. We begin with the following lemma, which is straightforward from the proof of [Lemma C.10](#).

Lemma C.11 ([Corollary 5.3](#), formal version). *Under the conditions of [Lemma 5.2](#), the following inequality holds when $\alpha \leq C_4(\mathbf{X}, \bar{\mathbf{U}})$:*

$$\left\| \mathbf{U}_{T_{\alpha,s}^{\text{pi}}} \mathbf{W}_{T_{\alpha,s}^{\text{pi}},\perp} \right\| \leq C_5(\mathbf{X}, \bar{\mathbf{U}}) \cdot \alpha^{\frac{1}{4\kappa}}$$

where $C_5 = \sqrt{10\kappa}\|\mathbf{X}\| \frac{C_2(\mathbf{X}, \bar{\mathbf{U}})}{C_3(\mathbf{X}, \bar{\mathbf{U}})}$.

The following lemma states that in a certain time period after $T_{\alpha,s}^{\text{pi}}$, the parallel and orthogonal components still behave similarly to the second (parallel improvement) phase.

Lemma C.12. *Under the conditions in [Lemma C.10](#), there exists $\tilde{t}_{\alpha,s} \geq \frac{1}{\log(1+\mu\sigma_s^2)} \log\left(\frac{10^{-4}c_3\|\mathbf{X}\|^2}{\sqrt{d\kappa}C_5} \alpha^{-\frac{1}{4\kappa}}\right) = \Theta(\log \alpha^{-1})$ when $\alpha \rightarrow 0$ such that when $0 \leq t - T_{\alpha,s}^{\text{pi}} \leq \tilde{t}_{\alpha,s}$, we have*

$$\sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t) \geq \sigma_{\min}(\mathbf{U}_t \mathbf{W}_t) \geq 0.3\kappa^{-1}\|\mathbf{X}\|^2, \quad (41a)$$

$$\|\mathbf{U}_t \mathbf{W}_t\| \leq (1 + \mu(0.4\sigma_s^2 + 0.6\sigma_{s+1}^2))^{t-T_{\alpha,s}^{\text{pi}}} \left\| \mathbf{U}_{T_{\alpha,s}^{\text{pi}}} \mathbf{W}_{T_{\alpha,s}^{\text{pi}},\perp} \right\|, \quad (41b)$$

$$\|\mathbf{V}_{\mathbf{X}_s,\perp}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}\| \leq c_3. \quad (41c)$$

Proof. We choose

$$\tilde{t}_{\alpha,s} = \min \{t \geq 0 : \|\mathbf{U}_{t+1} \mathbf{W}_{t+1,\perp}\|^2 \leq c_5\} \quad (42)$$

where

$$c_5 = 10^{-4} d^{-\frac{1}{2}} \kappa^{-1} c_3 \|\mathbf{X}\|^2 \quad (43)$$

We prove (41) by induction. The proof follows the idea of Lemma C.10, except that we need to bound $\|\Delta_t\|$ in each induction step. Concretely, suppose that (41) holds at time t , then

$$\begin{aligned} \|\Delta_t\| &= \|(\mathcal{A}^* \mathcal{A} - \mathbf{I})(\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top)\| \\ &\leq \|(\mathcal{A}^* \mathcal{A} - \mathbf{I})(\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{W}_t \mathbf{W}_t^\top \mathbf{U}_t^\top)\| + \|(\mathcal{A}^* \mathcal{A} - \mathbf{I}) \mathbf{U}_t \mathbf{W}_{t,\perp} \mathbf{W}_{t,\perp}^\top \mathbf{U}_t^\top\| \\ &\leq 10\delta\sqrt{r_*} \|\mathbf{X}\|^2 + \delta \|\mathbf{U}_t \mathbf{W}_{t,\perp} \mathbf{W}_{t,\perp}^\top \mathbf{U}_t^\top\|_* \\ &\leq 10\delta\sqrt{r_*} \|\mathbf{X}\|^2 + \delta c_5 \sqrt{d} \leq 0.02\kappa^{-1} c_3 \|\mathbf{X}\|^2 \end{aligned} \quad (44)$$

where we used the definition of c_5 in the last step. As a result, we can apply the conclusion of Lemmas C.5, C.8 and C.9 which implies that (41) holds for $t+1$. Finally, combining Lemma C.11 and (41b) yields $\tilde{t}_{\alpha,s} = \Theta(\log \frac{1}{\alpha})$. \square

We now present the main result of this section:

Lemma C.13. *Suppose that $0 \leq t - T_{\alpha,s}^\pi \leq \tilde{t}_{\alpha,s}$, $\|\mathbf{V}_{\mathbf{X}_s,\perp}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}\| \leq c_3$ and the conditions in Lemma C.10 hold, then we have*

$$\begin{aligned} &\|\mathbf{V}_{\mathbf{X}_s}^\top (\mathbf{X} \mathbf{X}^\top - \mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top)\|_F \\ &\leq \left(1 - \frac{1}{2}\mu\tau\right) \|\mathbf{V}_{\mathbf{X}_s}^\top (\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top)\|_F + 20\mu \|\mathbf{X}\|^4 (\delta + 5c_3) + 2000\mu^2 \sqrt{r_*} \|\mathbf{X}\|^6. \end{aligned}$$

where we recall that $\tau = \min_{1 \leq s \leq \hat{r} \wedge r_*} (\sigma_s^2 - \sigma_{s+1}^2) > 0$.

Proof. Recall that $\mathbf{M}_t = \mathcal{A}^* \mathcal{A} (\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top)$. The update of GD implies that

$$\begin{aligned} &\mathbf{X} \mathbf{X}^\top - \mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top \\ &= \mathbf{X} \mathbf{X}^\top - (\mathbf{I} + \mu \mathbf{M}_t) \mathbf{U}_t \mathbf{U}_t^\top (\mathbf{I} + \mu \mathbf{M}_t) \\ &= \underbrace{(\mathbf{I} - \mu \mathbf{U}_t \mathbf{U}_t^\top) (\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top) (\mathbf{I} - \mu \mathbf{U}_t \mathbf{U}_t^\top)}_{=(i)} + \underbrace{\mu \Delta_t \mathbf{U}_t \mathbf{U}_t^\top}_{=(ii)} \\ &\quad + \underbrace{\mu \mathbf{U}_t \mathbf{U}_t^\top \Delta_t}_{=(iii)} + \mu^2 (\mathcal{E}_{t,1} + \mathcal{E}_{t,2}), \end{aligned}$$

where $\mathcal{E}_{t,1} = -\mathbf{U}_t \mathbf{U}_t^\top (\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top) \mathbf{U}_t \mathbf{U}_t^\top$ and $\mathcal{E}_{t,2} = -\mathbf{M}_t \mathbf{U}_t \mathbf{U}_t^\top \mathbf{M}_t$. Since $\|\mathbf{U}_t\| \leq 3\|\mathbf{X}\|$ by Lemma D.4, and $\|\mathbf{M}_t - (\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top)\| = \|\Delta_t\| \leq \|\mathbf{X}\|^2$ which is shown in (44), we have

$$\begin{aligned} \|\mathbf{V}_{\mathbf{X}_s}^\top \mathcal{E}_{t,1}\|_F &= \|\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{U}_t^\top (\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top) \mathbf{U}_t \mathbf{U}_t^\top\|_F \\ &\leq \sqrt{r_*} \|\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{U}_t^\top (\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top) \mathbf{U}_t \mathbf{U}_t^\top\|_2 \\ &\leq 10^3 \sqrt{r_*} \|\mathbf{X}\|^6 \end{aligned}$$

and

$$\begin{aligned} \|\mathbf{V}_{\mathbf{X}_s}^\top \mathcal{E}_{t,2}\|_F &= \|\mathbf{V}_{\mathbf{X}_s}^\top [(\mathcal{A}^* \mathcal{A}) (\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top)] \mathbf{U}_t \mathbf{U}_t^\top [(\mathcal{A}^* \mathcal{A}) (\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top)]\|_F \\ &\leq \sqrt{r_*} \|[(\mathcal{A}^* \mathcal{A}) (\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top)] \mathbf{U}_t \mathbf{U}_t^\top [(\mathcal{A}^* \mathcal{A}) (\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top)]\| \\ &\leq 10^3 \sqrt{r_*} \|\mathbf{X}\|^6. \end{aligned}$$

Note that we would like to bound $\|\mathbf{V}_{\mathbf{X}_s}^\top (\mathbf{X}\mathbf{X}^\top - \mathbf{U}_{t+1}\mathbf{U}_{t+1}^\top)\|_F$. We deal with the above three terms separately. For the first term, we have

$$\begin{aligned} & \|\mathbf{V}_{\mathbf{X}_s}^\top (\mathbf{I} - \mu\mathbf{U}_t\mathbf{U}_t^\top) (\mathbf{X}\mathbf{X}^\top - \mathbf{U}_t\mathbf{U}_t^\top) (\mathbf{I} - \mu\mathbf{U}_t\mathbf{U}_t^\top)\|_F \\ &= \|\mathbf{V}_{\mathbf{X}_s}^\top (\mathbf{I} - \mu\mathbf{U}_t\mathbf{U}_t^\top) \mathbf{V}_{\mathbf{X}_s} \mathbf{V}_{\mathbf{X}_s}^\top (\mathbf{X}\mathbf{X}^\top - \mathbf{U}_t\mathbf{U}_t^\top) (\mathbf{I} - \mu\mathbf{U}_t\mathbf{U}_t^\top)\|_F \\ & \quad + \|\mathbf{V}_{\mathbf{X}_s}^\top (\mathbf{I} - \mu\mathbf{U}_t\mathbf{U}_t^\top) \mathbf{V}_{\mathbf{X}_{s,\perp}} \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top (\mathbf{X}\mathbf{X}^\top - \mathbf{U}_t\mathbf{U}_t^\top) (\mathbf{I} - \mu\mathbf{U}_t\mathbf{U}_t^\top)\|_F \\ &\leq \|\mathbf{I} - \mu\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t\mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_s}\| \|\mathbf{V}_{\mathbf{X}_s}^\top (\mathbf{X}\mathbf{X}^\top - \mathbf{U}_t\mathbf{U}_t^\top)\| + \mu \|\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t\mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_{s,\perp}} \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top (\mathbf{X}\mathbf{X}^\top - \mathbf{U}_t\mathbf{U}_t^\top)\|_F \end{aligned} \quad (45a)$$

$$\leq (1 - \mu\sigma_{\min}^2(\mathbf{U}_t\mathbf{W}_t)\sigma_{\min}^2(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{V}_{\mathbf{U}_t\mathbf{W}_t})) \|\mathbf{V}_{\mathbf{X}_s}^\top (\mathbf{X}\mathbf{X}^\top - \mathbf{U}_t\mathbf{U}_t^\top)\|_F + 100\mu\|\mathbf{X}\|^4 c_3 \quad (45b)$$

$$\leq \left(1 - \frac{1}{2}\mu\tau\right) \|\mathbf{V}_{\mathbf{X}_s}^\top (\mathbf{X}\mathbf{X}^\top - \mathbf{U}_t\mathbf{U}_t^\top)\|_F + 100\mu\|\mathbf{X}\|^4 c_3, \quad (45c)$$

where in (45a) we use $\|\mathbf{I} - \mu\mathbf{U}_t\mathbf{U}_t^\top\| \leq 1$, (45b) follows from

$$\begin{aligned} \sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t\mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_s}) &= \sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t\mathbf{W}_t\mathbf{W}_t^\top \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_s}) \geq \sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t\mathbf{W}_t) \\ &\geq \sigma_{\min}^2(\mathbf{U}_t\mathbf{W}_t)\sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{V}_{\mathbf{U}_t\mathbf{W}_t}) \end{aligned}$$

and

$$\|\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t\mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_{s,\perp}}\| = \|\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t\mathbf{W}_t\mathbf{W}_t^\top \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_{s,\perp}}\| \leq \|\mathbf{U}_t\|^2 \|\mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{U}_t\mathbf{W}_t\| \leq c_3\|\mathbf{U}_t\|^2,$$

and lastly (45c) is obtained from

$$\sigma_{\min}^2(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{V}_{\mathbf{U}_t\mathbf{W}_t}) \geq 1 - \|\mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{V}_{\mathbf{U}_t\mathbf{W}_t}\|^2 \geq 1 - c_3^2.$$

For the second and the third terms, we have

$$\|\Delta_t \mathbf{U}_t \mathbf{U}_t^\top + \mathbf{U}_t \mathbf{U}_t^\top \Delta_t\| \leq 0.1\kappa c_3 \|\mathbf{X}\|^4 \quad (46)$$

where we use the estimate in (44). Combining (45) and (46) yields

$$\begin{aligned} & \|\mathbf{V}_{\mathbf{X}_s}^\top (\mathbf{X}\mathbf{X}^\top - \mathbf{U}_{t+1}\mathbf{U}_{t+1}^\top)\|_F \\ &\leq \left(1 - \frac{1}{2}\mu\tau\right) \|\mathbf{V}_{\mathbf{X}_s}^\top (\mathbf{X}\mathbf{X}^\top - \mathbf{U}_t\mathbf{U}_t^\top)\|_F + 200\mu\|\mathbf{X}\|^4 c_3 + 110\mu^2\sqrt{r_*}\|\mathbf{X}\|^6. \end{aligned}$$

□

To apply the result of Lemma 5.4, we need to verify that $\|\mathbf{V}_{\mathbf{X}_{s,\perp}} \mathbf{V}_{\mathbf{U}_t\mathbf{W}_t}\| \leq c_3$ still holds when $t \geq T_{\alpha,s}^{\text{pi}}$. In fact, this is true as long as $t - T_{\alpha,s}^{\text{pi}} \leq \mathcal{O}(\log \frac{1}{\alpha})$.

Lemma C.14. *Under the conditions in Lemma C.10, if*

$$T_{\alpha,s}^{\text{pi}} \leq t \leq T_{\alpha,s}^{\text{pi}} + \frac{\gamma_s \log \frac{c_4}{C_5(\mathbf{X}, \bar{\mathbf{U}})} \cdot \log \frac{1}{\alpha}}{\log(1 + \mu\sigma_s^2)} =: T_{\alpha,s}^{\text{re}},$$

then $\|\mathbf{U}_{t+1}\mathbf{W}_{t+1,\perp}\| \leq (1 + \mu\sigma_s^2)\|\mathbf{U}_t\mathbf{W}_{t,\perp}\|$ and $\|\mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{V}_{\mathbf{U}_t\mathbf{W}_t}\| \leq c_3$. As a consequence, we have $\|\mathbf{U}_t\mathbf{W}_{t,\perp}\| \leq (1 + \mu\sigma_s^2)^{t - T_{\alpha,s}^{\text{pi}}} C_5(\mathbf{X}, \bar{\mathbf{U}}) \cdot \alpha^{\gamma_s} \leq c_4$.

Proof. The proof is basically the same as that of Lemma C.10 and we only provide a sketch here.

We induct on t . The base case $t = T_{\alpha,s}^{\text{pi}}$ is already proved in Lemma C.10. Suppose that the lemma holds for $t - 1$ with $t < T_{\alpha,s}^{\text{re}}$, then the choice of $T_{\alpha,s}^{\text{re}}$ combined with Lemma C.8 imply that

$$\|\mathbf{U}_t\mathbf{W}_{t,\perp}\| \leq (1 + \mu\sigma_s^2)\|\mathbf{U}_{t-1}\mathbf{W}_{t-1,\perp}\| \leq (1 + \mu\sigma_s^2)^{t - T_{\alpha,s}^{\text{pi}}} \|\mathbf{U}_{T_{\alpha,s}^{\text{pi}}}\mathbf{W}_{T_{\alpha,s}^{\text{pi}}}\|.$$

Since we have $\|\mathbf{U}_{T_{\alpha,s}^{\text{pi}}}\mathbf{W}_{T_{\alpha,s}^{\text{pi}}}\| \leq C_5(\mathbf{X}, \bar{\mathbf{U}}) \cdot \alpha^{\gamma_s}$ by Lemma C.11, the choice of $T_{\alpha,s}^{\text{re}}$ implies that $\|\mathbf{U}_t\mathbf{W}_{t,\perp}\| \leq c_4$. The bound $\|\mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{V}_{\mathbf{U}_t\mathbf{W}_t}\| \leq c_3$ then follows from Lemma C.9. □

We will only use a weaker version of this lemma, namely that the bounds holds for all $T_{\alpha,s}^{\text{pi}} \leq t \leq T_{\alpha,s}^{\text{ft}}$. When α is sufficiently small, this can be directly derived from [Lemmas C.13](#) and [C.14](#) since $\tilde{t}_{\alpha,s}, T_{\alpha,s}^{\text{re}} - T_{\alpha,s}^{\text{pi}} = \Theta(\log \frac{1}{\alpha})$. Specifically, we have proven [Lemma 5.4](#) in the main text:

Lemma 5.4. *Suppose that $T_{\alpha,s}^{\text{pi}} \leq t \leq T_{\alpha,s}^{\text{ft}}$ and all the conditions in [Lemma 5.2](#) hold, then we have*

$$\begin{aligned} & \|\mathbf{V}_{\mathbf{X}_s}^\top (\mathbf{X}\mathbf{X}^\top - \mathbf{U}_{t+1}\mathbf{U}_{t+1}^\top)\|_F \\ & \leq \left(1 - \frac{1}{2}\mu\tau\right) \|\mathbf{V}_{\mathbf{X}_s}^\top (\mathbf{X}\mathbf{X}^\top - \mathbf{U}_t\mathbf{U}_t^\top)\|_F \\ & \quad + 20\mu\|\mathbf{X}\|^4(\delta + 5c_3) + 2000\mu^2\sqrt{r_*}\|\mathbf{X}\|^6. \end{aligned}$$

Moreover, we have $\|\mathbf{U}_{t+1}\mathbf{W}_{t+1,\perp}\| \leq (1 + \sigma_s^2)\|\mathbf{U}_t\mathbf{W}_{t,\perp}\|$ and $\|\mathbf{V}_{\mathbf{X}_s^\perp}\mathbf{V}_{\mathbf{U}_t\mathbf{W}_t}\| \leq c_3$.

We are now ready to present our first main result, which states that with small initialization, GD would visit the $\mathcal{O}(\delta)$ -neighborhood of the rank- s minimizer of the full observation loss i.e. $\mathbf{X}_s\mathbf{X}_s^\top$.

Theorem C.15 (Restatement of [Lemma 4.3](#)). *Under [Assumptions 3.3, 3.5](#) and [3.7](#), if the initialization scale α is sufficiently small, then for all $1 \leq s \leq \hat{r} \wedge r_*$ there exists a time $T_{\alpha,s}^{\text{ft}} \in \mathbb{Z}_+$ (where ft stands for fitting the ground-truth) such that*

$$\left\| \mathbf{X}_s\mathbf{X}_s^\top - \mathbf{U}_{T_{\alpha,s}^{\text{pi}}}\mathbf{U}_{T_{\alpha,s}^{\text{pi}}}^\top \right\|_F \leq 10^7\kappa^3r_*\|\mathbf{X}\|^2\delta.$$

Proof. First, observe that for all $t \geq 0$,

$$\begin{aligned} \|\mathbf{X}_s\mathbf{X}_s^\top - \mathbf{U}_t\mathbf{U}_t^\top\|_F & \leq \|(\mathbf{X}_s\mathbf{X}_s^\top - \mathbf{U}_t\mathbf{U}_t^\top)\mathbf{V}_{\mathbf{X}_s}\mathbf{V}_{\mathbf{X}_s}^\top\|_F + \|\mathbf{U}_t\mathbf{U}_t^\top\mathbf{V}_{\mathbf{X}_s^\perp}\mathbf{V}_{\mathbf{X}_s^\perp}^\top\|_F \\ & \leq \|(\mathbf{X}_s\mathbf{X}_s^\top - \mathbf{U}_t\mathbf{U}_t^\top)\mathbf{V}_{\mathbf{X}_s}\mathbf{V}_{\mathbf{X}_s}^\top\|_F + \|\mathbf{V}_{\mathbf{X}_s^\perp}^\top\mathbf{U}_t\mathbf{U}_t^\top\mathbf{V}_{\mathbf{X}_s^\perp}\|_F \\ & \leq \|\mathbf{V}_{\mathbf{X}_s}^\top(\mathbf{X}_s\mathbf{X}_s^\top - \mathbf{U}_t\mathbf{U}_t^\top)\|_F + \sqrt{r_*}\|\mathbf{V}_{\mathbf{X}_s^\perp}^\top\mathbf{U}_t\mathbf{W}_t\| + \sqrt{d}\|\mathbf{V}_{\mathbf{X}_s^\perp}^\top\mathbf{U}_t\mathbf{W}_{t,\perp}\|^2 \\ & \leq \|\mathbf{V}_{\mathbf{X}_s}^\top(\mathbf{X}_s\mathbf{X}_s^\top - \mathbf{U}_t\mathbf{U}_t^\top)\|_F + 9\sqrt{r_*}\|\mathbf{X}\|^2\|\mathbf{V}_{\mathbf{X}_s,\perp}\mathbf{V}_{\mathbf{U}_t\mathbf{W}_t}\|^2 + \sqrt{d}\|\mathbf{U}_t\mathbf{W}_{t,\perp}\|^2. \end{aligned} \tag{47}$$

where the last step uses [Lemma D.4](#) to bound $\|\mathbf{U}_t\|$. We set $c_3 = 10^3\kappa\sqrt{r_*}\delta$ and

$$T_{\alpha,s}^{\text{ft}} = T_{\alpha,s}^{\text{pi}} - \frac{\log(10^{-2}\|\mathbf{X}\|^{-2}\tau c_3^{-1})}{\log(1 - \frac{1}{2}\mu\tau)}, \tag{48}$$

where we recall that $\tau = \kappa^{-1}\|\mathbf{X}\|^2$, then for small α we have $T_{\alpha,s}^{\text{ft}} \leq T_{\alpha,s}^{\text{pi}} + \tilde{t}_{\alpha,s}$ (defined in [Lemma C.12](#)). Hence for $T_{\alpha,s}^{\text{pi}} \leq t < T_{\alpha,s}^{\text{ft}}$ we always have $\|\mathbf{V}_{\mathbf{X}_s,\perp}\mathbf{V}_{\mathbf{U}_t\mathbf{W}_t}\| \leq c_3$. By [Lemma C.13](#) and the choice of c_3 and δ , we have for $T_{\alpha,s}^{\text{pi}} \leq t < T_{\alpha,s}^{\text{ft}}$ that

$$\|\mathbf{V}_{\mathbf{X}_s}^\top(\mathbf{X}\mathbf{X}^\top - \mathbf{U}_{t+1}\mathbf{U}_{t+1}^\top)\|_F \leq \left(1 - \frac{1}{2}\mu\tau\right) \|\mathbf{V}_{\mathbf{X}_s}^\top(\mathbf{X}\mathbf{X}^\top - \mathbf{U}_t\mathbf{U}_t^\top)\|_F + 30\mu\|\mathbf{X}\|^4\sqrt{r_*}c_3$$

which implies that for $t = T_{\alpha,s}^{\text{ft}}$,

$$\|\mathbf{V}_{\mathbf{X}_s}^\top(\mathbf{X}\mathbf{X}^\top - \mathbf{U}_{T_s}\mathbf{U}_{T_s}^\top)\|_F \leq 80\kappa\|\mathbf{X}\|^2\sqrt{r_*}c_3.$$

Meanwhile, by [Lemma C.12](#) we have $\|\mathbf{U}_t\mathbf{W}_{t,\perp}\| \leq c_5$ (c_5 is defined in [\(43\)](#)) and $\|\mathbf{V}_{\mathbf{X}_s^\perp}^\top\mathbf{V}_{\mathbf{U}_t\mathbf{W}_t}\| \leq c_3$ at $t = T_{\alpha,s}^{\text{ft}}$. Plugging into [\(47\)](#) yields

$$\left\| \mathbf{X}_s\mathbf{X}_s^\top - \mathbf{U}_{T_{\alpha,s}^{\text{ft}}}\mathbf{U}_{T_{\alpha,s}^{\text{ft}}}^\top \right\|_F \leq 80\kappa\|\mathbf{X}\|^2\sqrt{r_*}c_3 + 9\|\mathbf{X}\|^2c_3^2\sqrt{r_*} + c_5^2\sqrt{d}.$$

By definition of c_3 and c_5 we deduce that $\left\| \mathbf{X}_s\mathbf{X}_s^\top - \mathbf{U}_{T_{\alpha,s}^{\text{ft}}}\mathbf{U}_{T_{\alpha,s}^{\text{ft}}}^\top \right\|_F \leq 10^2\tau^{-2}\|\mathbf{X}\|^6\sqrt{r_*}c_3 \leq 10^5\kappa^3r_*\|\mathbf{X}\|^2\delta$, as desired. \square

Corollary C.16. *There exists a constant*

$$C_6(\mathbf{X}, \bar{\mathbf{U}}) = C_5(\mathbf{X}, \bar{\mathbf{U}}) \cdot (1 + \mu\sigma_s^2)^{T_{\alpha,s}^{\text{ft}} - T_{\alpha,s}^{\text{pi}}} \quad (49)$$

such that

$$\max_{0 \leq t \leq T_{\alpha,s}^{\text{ft}}} \|\mathbf{U}_t \mathbf{W}_{t,\perp}\| \leq C_1 \cdot \alpha^{\frac{1}{4rc}}.$$

Proof. The case of $t \leq T_{\alpha,s}^{\text{pi}}$ directly follows from [Lemma C.11](#). For $t > T_{\alpha,s}^{\text{pi}}$, we know from [Lemma C.12](#) that

$$\|\mathbf{U}_t \mathbf{W}_{t,\perp}\| \leq \left\| \mathbf{U}_{T_{\alpha,s}^{\text{pi}}} \mathbf{W}_{T_{\alpha,s}^{\text{pi}},\perp} \right\| \cdot (1 + \mu\sigma_s^2)^{T_{\alpha,s}^{\text{ft}} - T_{\alpha,s}^{\text{pi}}}.$$

By (48), the second term is a constant independent of α , so the conclusion follows. \square

D. Auxiliary results for proving [Lemma 4.3](#)

This section contains a collection of auxiliary results that are used in the previous section.

D.1. The spectral phase

In the section, we provide auxiliary results for the analysis in the spectral phase.

Recall that $\mathbf{K}_t = (\mathbf{I} + \mu\mathbf{M})^t$ and $\mathbf{U}_t = \mathbf{U}_t^{\text{sp}} + \mathbf{E}_t = \mathbf{K}_t \mathbf{U}_0 + \mathbf{E}_t$ and $\mathbf{U}_0 = \alpha \bar{\mathbf{U}}$ with $\|\bar{\mathbf{U}}\| = 1$. Also recall that $\mathbf{M} = \sum_{i=1}^{\text{rank}(\mathbf{M})} \hat{\lambda}_i \hat{\mathbf{v}}_i \hat{\mathbf{v}}_i^\top$; we additionally define $\mathbf{M}_s = \sum_{i=1}^{\min\{s, \text{rank}(\mathbf{M})\}} \hat{\lambda}_i \hat{\mathbf{v}}_i \hat{\mathbf{v}}_i^\top$. Similarly, let L_t be the span of the top- s left singular vectors of \mathbf{U}_t . The following lemma shows that power iteration would result in large eigengap of \mathbf{U}_t .

Lemma D.1. *Let $\hat{\rho} = \sigma_{\min}(\mathbf{V}_{M_s}^\top \bar{\mathbf{U}}) > 0$, then the following three inequalities hold, given that the denominator of the third is positive.*

$$\sigma_s(\mathbf{U}_t) \geq \alpha \left(\hat{\rho} \sigma_s(\hat{\mathbf{Z}}_t) - \sigma_{s+1}(\hat{\mathbf{Z}}_t) \right) - \|\mathbf{E}_t\|, \quad (50a)$$

$$\sigma_{s+1}(\mathbf{U}_t) \leq \alpha \sigma_{s+1}(\hat{\mathbf{Z}}_t) + \|\mathbf{E}_t\|, \quad (50b)$$

$$\left\| \mathbf{V}_{M_s^\perp}^\top \mathbf{V}_{L_t} \right\| \leq \frac{\alpha \sigma_{s+1}(\hat{\mathbf{Z}}_t) + \|\mathbf{E}_t\|}{\alpha \hat{\rho} \sigma_s(\hat{\mathbf{Z}}_t) - 2(\alpha \sigma_{s+1}(\hat{\mathbf{Z}}_t) + \|\mathbf{E}_t\|)}. \quad (50c)$$

Proof. By Weyl's inequality we have

$$\begin{aligned} \sigma_{s+1}(\mathbf{U}_t) &= \sigma_{s+1}((1 + \mu\mathbf{M})^t \mathbf{U}_0) + \|\mathbf{E}_t\| \\ &= \alpha \sigma_{s+1}((1 + \mu\mathbf{M})^t \bar{\mathbf{U}}) + \|\mathbf{E}_t\| \\ &\leq \alpha \sigma_{s+1}((1 + \mu\mathbf{M}_s)^t \bar{\mathbf{U}}) + \alpha \left\| [(1 + \mu\mathbf{M})^t - (1 + \mu\mathbf{M}_s)^t] \bar{\mathbf{U}} \right\| + \|\mathbf{E}_t\| \\ &\leq \alpha(1 + \mu\hat{\lambda}_{s+1})^t + \|\mathbf{E}_t\|. \end{aligned}$$

Thus (50b) holds. Similarly,

$$\begin{aligned} \sigma_s(\mathbf{U}_t) &\geq \alpha \sigma_s(\mathbf{N}_t \mathbf{V}_{M_s} \mathbf{V}_{M_s}^\top \bar{\mathbf{U}}) - \alpha(1 + \mu\hat{\lambda}_{s+1})^t - \|\mathbf{E}_t\| \\ &\geq \alpha \sigma_s(\mathbf{N}_t \mathbf{V}_{M_s}) \sigma_{\min}(\mathbf{V}_{M_s}^\top \bar{\mathbf{U}}) - \alpha(1 + \mu\hat{\lambda}_{s+1})^t - \|\mathbf{E}_t\| \\ &\geq \alpha \hat{\rho} (1 + \mu\hat{\lambda}_s)^t - \alpha(1 + \mu\hat{\lambda}_{s+1})^t - \|\mathbf{E}_t\|. \end{aligned}$$

Finally, note that we can write

$$\alpha(1 + \mu\mathbf{M}_s)^t \bar{\mathbf{U}} = \mathbf{V}_{M_s} \underbrace{(1 + \mu\boldsymbol{\Sigma}_{M_s})^t \mathbf{V}_{M_s}^\top}_{\text{invertible}} \bar{\mathbf{U}},$$

so that the subspace spanned by the left singular vectors of $\alpha(1 + \mu\mathbf{M}_s)^t \bar{\mathbf{U}}$ coincides with the column span of \mathbf{V}_{M_s} . Since L_t is the span of top- s left singular vectors of \mathbf{U}_t , we apply Wedin's sin theorem ([Wedin, 1972](#)) and deduce (50c). \square

The next lemma relates the quantities studied in [Lemma D.1](#) with those that are needed in the induction. The proof is the same as [Stöger & Soltanolkotabi, 2021](#), Lemma 8.4, so we omit it here.

Lemma D.2. *Suppose that $\|\mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{V}_{L_t}\| \leq 0.1$ for some $t \geq 1$. Then it holds that*

$$\sigma_s(\mathbf{U}_t \mathbf{W}_t) \geq \frac{1}{2} \sigma_s(\mathbf{U}_t), \quad (51a)$$

$$\|\mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}\| \leq 10 \|\mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{V}_{L_t}\|, \quad (51b)$$

$$\|\mathbf{U}_t \mathbf{W}_{t, \perp}\| \leq 2\sigma_{s+1}(\mathbf{U}_t). \quad (51c)$$

Combining the above two lemmas, we directly obtain the following corollary:

Corollary D.3. *Suppose that $\alpha\sigma_s(\mathbf{K}_t) > 10(\alpha\sigma_{s+1}(\mathbf{K}_t) + \|\mathbf{E}_t\|)$, then we have that*

$$\begin{aligned} \sigma_s(\mathbf{U}_t \mathbf{W}_t) &\geq 0.4\alpha\sigma_{r_*}(\mathbf{K}_t) \sigma_{\min}(\mathbf{V}_L^\top \bar{\mathbf{U}}) \\ \|\mathbf{U}_t \mathbf{W}_{t, \perp}\| &\leq 2(\alpha\sigma_{s+1}(\mathbf{K}_t) + \|\mathbf{E}_t\|) \\ \|\mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}\| &\leq 100 \left(\delta + \frac{\alpha\sigma_{s+1}(\mathbf{K}_t) + \|\mathbf{E}_t\|}{\alpha\hat{\rho}_s(\mathbf{K}_t)} \right) \end{aligned} \quad (52)$$

D.2. The parallel improvement phase

In the section, we provide auxiliary results for the analysis in the parallel improvement phase.

Lemma D.4. ([Stöger & Soltanolkotabi, 2021](#), Lemma 9.4) *For sufficiently small μ and δ , suppose that $\|\mathbf{U}_t\| \leq 3\|\mathbf{X}\|$, then we also have $\|\mathbf{U}_{t+1}\| \leq 3\|\mathbf{X}\|$.*

Lemma D.5. *Under the assumptions in [Lemma C.8](#), we have*

$$\|\mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{V}_{\mathbf{U}_{t+1} \mathbf{W}_t}\| \leq 2(c_3 + 10\mu\|\mathbf{X}\|^2) \leq 0.01.$$

Proof. The proof of this lemma is essentially the same as [Stöger & Soltanolkotabi, 2021](#), Lemma B.1, and we omit it here. \square

Lemma D.6. *Under the assumptions in [Lemma C.9](#), we have*

$$\sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1}) \geq \frac{1}{2} \sigma_{\min}(\mathbf{U}_t \mathbf{W}_t).$$

Proof. We have

$$\begin{aligned} \sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1}) &\geq \sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \mathbf{W}_t) \\ &= \sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top (\mathbf{I} + \mu \mathbf{M}_t) \mathbf{U}_t \mathbf{W}_t) \\ &\geq \sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top (\mathbf{I} + \mu \mathbf{M}_t) \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}) \cdot \sigma_{\min}(\mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}^\top \mathbf{U}_t \mathbf{W}_t) \\ &\geq [\sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}) - \mu \|\mathbf{M}_t\|] \cdot \sigma_{\min}(\mathbf{U}_t \mathbf{W}_t) \\ &\geq \left(\sqrt{1 - c_3^2} - 10\mu\|\mathbf{X}\|^2 \right) \sigma_{\min}(\mathbf{U}_t \mathbf{W}_t) \geq \frac{1}{2} \sigma_{\min}(\mathbf{U}_t \mathbf{W}_t) \end{aligned}$$

where the last step follows from

$$\sigma_{\min}(\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t})^2 \geq 1 - \|\mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}\|^2 \geq 1 - c_3^2.$$

The conclusion follows. \square

Lemma D.7. *Under the assumptions in [Lemma C.9](#), we have*

$$\|\mathbf{W}_{t, \perp}^\top \mathbf{W}_{t+1}\| \leq 3\mu(10\mu\|\mathbf{X}\|^2 + c_4)c_3\|\mathbf{X}\| + \mu\|(\mathcal{A}^* \mathcal{A} - \mathbf{I})(\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top)\|.$$

Proof. The proof roughly follows [Stöger & Soltanolkotabi, 2021, Lemma B.3], but we include it here for completeness.

Since $\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} = \mathbf{V}_{t+1} \boldsymbol{\Sigma}_{t+1} \mathbf{W}_{t+1}$ and $\mathbf{V}_{t+1} \boldsymbol{\Sigma}_{t+1} \in \mathbb{R}^{s \times s}$ is invertible, we have

$$\|\mathbf{W}_{t,\perp}^\top \mathbf{W}_{t+1}\| = \left\| \mathbf{W}_{t,\perp}^\top \mathbf{U}_{t+1}^\top \mathbf{V}_{\mathbf{X}_s} (\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top \mathbf{V}_{\mathbf{X}_s})^{-\frac{1}{2}} \right\|.$$

Since

$$\begin{aligned} & \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \mathbf{W}_{t,\perp} \\ &= \mathbf{V}_{\mathbf{X}_s}^\top (\mathbf{I} + \mu \mathcal{A}^* \mathcal{A} (\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top)) \mathbf{U}_t \mathbf{W}_{t,\perp} \\ &= \mathbf{V}_{\mathbf{X}_s}^\top (\mathbf{I} + \mu (\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top)) \mathbf{U}_t \mathbf{W}_{t,\perp} + \mu \mathbf{V}_{\mathbf{X}_s}^\top \boldsymbol{\Delta}_t \mathbf{U}_t \mathbf{W}_{t,\perp} \\ &= -\mu \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{U}_t \mathbf{W}_{t,\perp} + \mu \mathbf{V}_{\mathbf{X}_s}^\top \boldsymbol{\Delta}_t \mathbf{U}_t \mathbf{W}_{t,\perp} \end{aligned} \quad (53a)$$

$$= -\mu \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t \mathbf{W}_t^\top \mathbf{U}_t^\top \mathbf{U}_t \mathbf{W}_{t,\perp} + \mu \mathbf{V}_{\mathbf{X}_s}^\top \boldsymbol{\Delta}_t \mathbf{U}_t \mathbf{W}_{t,\perp} \quad (53b)$$

$$= -\mu \underbrace{\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_t \mathbf{W}_t^\top \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_s, \perp} \mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{U}_t \mathbf{W}_{t,\perp}}_{=: \mathbf{K}_1} + \mu \underbrace{\mathbf{V}_{\mathbf{X}_s}^\top \boldsymbol{\Delta}_t \mathbf{U}_t \mathbf{W}_{t,\perp}}_{=: \mathbf{K}_2} \quad (53c)$$

where (53a) follows from $\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U}_t \mathbf{W}_{t,\perp} = \boldsymbol{\Sigma}_s \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_{t,\perp} = 0$, and in (53b) and (53c) we use $\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \mathbf{W}_{t,\perp} = 0$.

For \mathbf{K}_1 , note that

$$\begin{aligned} & \left\| (\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top \mathbf{V}_{\mathbf{X}_s})^{-\frac{1}{2}} \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_t \right\| \\ & \leq \left\| (\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top \mathbf{V}_{\mathbf{X}_s})^{-\frac{1}{2}} \mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \right\| + \mu \left\| (\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top \mathbf{V}_{\mathbf{X}_s})^{-\frac{1}{2}} \mathbf{V}_{\mathbf{X}_s}^\top \mathcal{A}^* \mathcal{A} (\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top) \mathbf{U}_t \right\| \\ & \leq 1 + 10\mu \|\mathbf{X}\|^3 \sigma_{\min}^{-1} (\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1}) \end{aligned}$$

so that

$$\left\| (\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top \mathbf{V}_{\mathbf{X}_s})^{-\frac{1}{2}} \mathbf{K}_1 \right\| \leq [1 + 10\mu \|\mathbf{X}\|^3 \sigma_{\min}^{-1} (\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1})] \|\mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{U}_t \mathbf{W}_t\| \|\mathbf{U}_t \mathbf{W}_{t,\perp}\|.$$

Plugging into (53), we deduce that

$$\begin{aligned} & \|\mathbf{W}_{t,\perp}^\top \mathbf{W}_{t+1}\| \\ & \leq 3\mu (1 + 10\mu \|\mathbf{X}\|^3 \sigma_{\min}^{-1} (\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1})) \|\mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}\| \|\mathbf{X}\| \|\mathbf{U}_t \mathbf{W}_{t,\perp}\| \\ & \quad + \mu \sigma_{\min}^{-1} (\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1}) \|\mathbf{U}_t \mathbf{W}_{t,\perp}\| \|(\mathcal{A}^* \mathcal{A} - \mathbf{I})(\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top)\| \\ & \leq 3\mu (\|\mathbf{U}_t \mathbf{W}_{t,\perp}\| + 10\mu \|\mathbf{X}\|^3) \|\mathbf{V}_{\mathbf{X}_s, \perp}^\top \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}\| \|\mathbf{X}\| \\ & \quad + \mu \sigma_{\min}^{-1} (\mathbf{V}_{\mathbf{X}_s}^\top \mathbf{U}_{t+1}) \|\mathbf{U}_t \mathbf{W}_{t,\perp}\| \|(\mathcal{A}^* \mathcal{A} - \mathbf{I})(\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top)\| \\ & \leq 3\mu (10\mu \|\mathbf{X}\|^2 + c_4) c_3 \|\mathbf{X}\| + \mu \|(\mathcal{A}^* \mathcal{A} - \mathbf{I})(\mathbf{X} \mathbf{X}^\top - \mathbf{U}_t \mathbf{U}_t^\top)\|. \end{aligned}$$

where in the last step we use Lemma D.6 and the induction hypothesis which implies that $\sigma_{\min}(\mathbf{U}_t \mathbf{W}_t) \geq \|\mathbf{U}_t \mathbf{W}_{t,\perp}\|$. \square

Lemma D.8. The matrix \mathbf{H} defined in the proof of Lemma C.9 satisfies the following:

$$\mathbf{H}(\mathbf{H}^\top \mathbf{H})^{-\frac{1}{2}} = \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t} + \mathbf{B} \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t} - \frac{1}{2} (\mathbf{I} + \mathbf{B}) \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t} \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}^\top (\mathbf{B} + \mathbf{B}^\top) \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t} - \mathbf{D},$$

where $\|\mathbf{D}\| \leq 30\|\mathbf{B}\|^2$.

Proof. By definition of \mathbf{H} we have

$$\begin{aligned} & \mathbf{H}(\mathbf{H}^\top \mathbf{H})^{-\frac{1}{2}} \\ &= (\mathbf{I} + \mu \mathbf{M})(\mathbf{I} + \mathbf{P}) \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t} (\mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}^\top (\mathbf{I} + \mathbf{P}^\top) (\mathbf{I} + \mu \mathbf{M})^2 (\mathbf{I} + \mathbf{P}) \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t})^{-\frac{1}{2}} \\ &= (\mathbf{I} + \mathbf{B}) \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t} [\mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}^\top (\mathbf{I} + \mathbf{B}^\top + \mathbf{B} + \mathbf{B}^\top \mathbf{B}) \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}]^{-\frac{1}{2}} \\ &= (\mathbf{I} + \mathbf{B}) \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t} \left[\underbrace{\mathbf{I} + \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}^\top (\mathbf{B}^\top + \mathbf{B} + \mathbf{B}^\top \mathbf{B}) \mathbf{V}_{\mathbf{U}_t \mathbf{W}_t}}_{=: \boldsymbol{\Theta}} \right]^{-\frac{1}{2}}. \end{aligned}$$

It follows from (36) and our assumptions on c_3 and c_4 that

$$\begin{aligned}\|B\| &\leq \mu \|\mathbf{X}\mathbf{X}^\top - \mathbf{U}_t\mathbf{U}_t^\top\| + 6\mu (c_3c_4\|\mathbf{X}\| + 50\|\mathbf{X}\|^2\delta) \\ &\leq 10\mu\|\mathbf{X}\|^2 + 6\mu c_3(c_4 + 1)\|\mathbf{X}\| < 0.1\end{aligned}$$

(note that this step is independent and does not rely on earlier derivations in the proof of Lemma C.9), so by Taylor's formula, we have

$$\left\|(\mathbf{I} + \Theta)^{-\frac{1}{2}} - \mathbf{I} + \frac{1}{2}\Theta\right\| \leq 3\|\Theta\|^2.$$

Hence,

$$\begin{aligned}&\left\|H(H^\top H)^{-\frac{1}{2}} - \left(V_{U_t}W_t + BV_{U_t}W_t - \frac{1}{2}(\mathbf{I} + B)V_{U_t}W_tV_{U_t}^\top(B + B^\top)V_{U_t}W_t\right)\right\| \\ &= \left\|(\mathbf{I} + B)V_{U_t}W_t \left((\mathbf{I} + \Theta)^{-\frac{1}{2}} - \mathbf{I} + \frac{1}{2}\Theta - \frac{1}{2}V_{U_t}^\top B^\top BV_{U_t}W_t\right)\right\| \\ &\leq (1 + \|B\|) \left(3\|\Theta\|^2 + \frac{1}{2}\|B\|^2\right) < 30\|B\|^2\end{aligned}$$

as desired. \square

E. Proofs for the Landscape Results in Section 4.1

In this section, we study the landscape of under-parameterized matrix sensing problem

$$f_s(\mathbf{U}) = \frac{1}{2} \|\mathcal{A}(\mathbf{U}\mathbf{U}^\top - \mathbf{X}\mathbf{X}^\top)\|_2^2, \quad \mathbf{U} \in \mathbb{R}^{d \times s}$$

Our key result in this section is Lemma 4.6, which states a local RSI condition for the matrix sensing loss. Most existing results only study the landscape of (1) in the exact- and over-parameterized case. Zhu et al. (2021) have studied the landscape of under-parameterized matrix factorization problem, but their main focus is the strict-saddle property of the loss.

E.1. Analysis of the matrix factorization loss

When the measurement satisfies the RIP condition, we can expect that the landscape of f_s looks similar to that of the (under-parameterized) matrix factorization loss:

$$F_s(\mathbf{U}) = \frac{1}{2} \|\mathbf{U}\mathbf{U}^\top - \mathbf{X}\mathbf{X}^\top\|_F^2, \quad \mathbf{U} \in \mathbb{R}^{d \times s}$$

for some $s < \hat{r}$. For this reason, we first look into the landscape of F_s before analyzing f_s .

Recall that $\mathbf{X}\mathbf{X}^\top = \sum_{i=1}^{r_*} \sigma_i^2 \mathbf{v}_i \mathbf{v}_i^\top$. The critical points of $F_s(\mathbf{U})$ is characterized by the following lemma:

Lemma E.1. $\mathbf{U} \in \mathbb{R}^{d \times s}$ is a critical point of $F_s(\mathbf{U})$ if and only if there exists an orthogonal matrix $\mathbf{R} \in \mathbb{R}^{s \times s}$, such that all columns of $\mathbf{U}\mathbf{R}$ are in $\{\sigma_i \mathbf{v}_i : 1 \leq i \leq r_*\}$.

Proof. Assume WLOG that $\mathbf{X}\mathbf{X}^\top = \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_r^2, 0, \dots, 0) =: \Sigma$. Let \mathbf{U} be a critical point of F_s , then we have that $(\mathbf{U}\mathbf{U}^\top - \mathbf{X}\mathbf{X}^\top)\mathbf{U} = 0$. Let $\mathbf{W} = \mathbf{U}\mathbf{U}^\top$, then $(\Sigma - \mathbf{W})\mathbf{W} = 0$.

Since \mathbf{W} is symmetric, so is \mathbf{W}^2 , and we obtain that $\Sigma\mathbf{W}$ is also symmetric. It is then easy to see that that if $\Sigma = \text{diag}(\lambda_1 \mathbf{I}_{m_1}, \dots, \lambda_t \mathbf{I}_{m_t})$ with $\lambda_1 > \lambda_2 > \dots > \lambda_t \geq 0$, then \mathbf{W} is also in block-diagonal form: $\mathbf{W} = \text{diag}(\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}_t)$ where $\mathbf{W}_i \in \mathbb{R}^{m_i \times m_i}$. For each $1 \leq i \leq t$, we then have the equation $(\lambda_i \mathbf{I}_{m_i} - \mathbf{W}_i)\mathbf{W}_i = 0$. Hence, there exists an orthogonal matrix \mathbf{R}_i such that $\mathbf{R}_i^\top \mathbf{W}_i \mathbf{R}_i$ is a diagonal matrix where the diagonal entries are either 0 or $\sqrt{\lambda_i} = \sigma_i$. Let $\mathbf{R} = \text{diag}(\mathbf{R}_1, \mathbf{R}_2, \dots, \mathbf{R}_t)$, then $\mathbf{R}^\top \mathbf{W} \mathbf{R}$ is diagonal and its nonzero diagonal entries form an s -subset of the multi-set $\{\sigma_i : 1 \leq i \leq r_*\}$. The conclusion follows. \square

In the case of $s = 1$, the global minimizers of F_s are $\pm \sigma_1 \mathbf{v}_1$, and we can show that F_s is locally strongly convex around these minimizers. Therefore, we can deduce that f is locally strongly-convex as well. Since our main focus is on $s > 1$, we put these details in Appendix G. When $s > 1$, $F_s(\mathbf{U})$ is not locally strongly-convex due to rotational invariance: if \mathbf{U} is a global minimizer, then so is $\mathbf{U}\mathbf{R}$ for any orthogonal matrix $\mathbf{R} \in \mathbb{R}^{s \times s}$. Instead, we establish a Restricted Secant Inequality for F_s , as shown below.

Lemma E.2. For $\mathbf{U} \in \mathbb{R}^{d \times s}$, let \mathbf{R} be an orthogonal matrix that minimizes $\|\mathbf{U} - \mathbf{X}_s \mathbf{R}\|_F$. Suppose that $\text{dist}(\mathbf{U}, \mathbf{X}_s) \leq 0.1 \|\mathbf{X}\|^{-1} \tau$ (where we recall that $\tau = \min_{s \in [r_*]} (\sigma_s^2 - \sigma_{s+1}^2)$ is the eigengap of $\mathbf{X} \mathbf{X}^\top$), then we have

$$\langle \nabla F_s(\mathbf{U}), \mathbf{U} - \mathbf{X}_s \mathbf{R} \rangle \geq 0.1 \tau \cdot \text{dist}^2(\mathbf{U}, \mathbf{X}_s).$$

Proof. Assume WLOG that $\mathbf{R} = \mathbf{I}$. Then by Lemma A.10, $\mathbf{U}^\top \mathbf{X}_s$ is symmetric and positive semi-definite. Let $\mathbf{H} = \mathbf{U} - \mathbf{X}_s$, then

$$\begin{aligned} \nabla F_s(\mathbf{U}) &= (\mathbf{U} \mathbf{U}^\top - \mathbf{X} \mathbf{X}^\top) \mathbf{U} \\ &= [(\mathbf{H} + \mathbf{X}_s)(\mathbf{H} + \mathbf{X}_s)^\top - \mathbf{X} \mathbf{X}^\top] (\mathbf{H} + \mathbf{X}_s). \end{aligned}$$

So we have

$$\begin{aligned} \langle \nabla F_s(\mathbf{U}), \mathbf{U} - \mathbf{X}_s \rangle &= \langle [(\mathbf{H} + \mathbf{X}_s)(\mathbf{H} + \mathbf{X}_s)^\top - \mathbf{X} \mathbf{X}^\top] (\mathbf{H} + \mathbf{X}_s), \mathbf{H} \rangle \\ &= \text{tr}(\mathbf{H}^\top [(\mathbf{H} + \mathbf{X}_s)(\mathbf{H} + \mathbf{X}_s)^\top - \mathbf{X} \mathbf{X}^\top] \mathbf{H} + \mathbf{H}^\top (\mathbf{H} \mathbf{H}^\top + \mathbf{H} \mathbf{X}_s^\top + \mathbf{X}_s \mathbf{H}^\top) \mathbf{X}_s) \\ &\geq -\text{tr}(\mathbf{H}^\top \mathbf{X}_{s,\perp} \mathbf{X}_{s,\perp}^\top \mathbf{H}) - 3 \|\mathbf{X}\| \|\mathbf{H}\|_F^3 + \text{tr}(\mathbf{H}^\top \mathbf{H} \mathbf{X}_s^\top \mathbf{X}_s) \end{aligned} \quad (54a)$$

$$\begin{aligned} &\geq (\sigma_s^2 - \sigma_{s+1}^2) \|\mathbf{H}\|_F^2 - 3 \|\mathbf{X}\| \|\mathbf{H}\|_F^3 \\ &\geq 0.1 \tau \|\mathbf{H}\|_F^2 \end{aligned} \quad (54b)$$

where in (54a) we use $\text{tr}(\mathbf{H}^\top \mathbf{X}_s) \geq 0$ (since $\mathbf{H}^\top \mathbf{X}_s$ is symmetric as noticed in the beginning of the proof), and (54b) is because of

$$\text{tr}(\mathbf{H}^\top \mathbf{H} \mathbf{X}_s^\top \mathbf{X}_s) \geq \sigma_{\min}(\mathbf{X}_s^\top \mathbf{X}_s) \cdot \text{tr}(\mathbf{H}^\top \mathbf{H}) = \sigma_s^2 \|\mathbf{H}\|_F^2$$

and

$$\begin{aligned} \text{tr}(\mathbf{H}^\top \mathbf{X}_{s,\perp} \mathbf{X}_{s,\perp}^\top \mathbf{H}) &= \text{tr}(\mathbf{H}^\top \mathbf{V}_{\mathbf{X}_{s,\perp}} \boldsymbol{\Sigma}_{s,\perp} \mathbf{V}_{\mathbf{X}_{s,\perp}}^\top \mathbf{H}) \\ &\leq \|\boldsymbol{\Sigma}_{s,\perp}\| \cdot \|\mathbf{H}^\top \mathbf{V}_{\mathbf{X}_{s,\perp}}\|_F^2 \leq \sigma_{s+1}^2 \|\mathbf{H}\|_F^2. \end{aligned}$$

□

Corollary E.3. Under the conditions of Lemma E.2, we have $\|\nabla F_s(\mathbf{U})\|_F \geq 0.1 \tau \text{dist}(\mathbf{U}, \mathbf{X}_s)$.

E.2. Analysis of the matrix sensing loss

The following lemma states that the minimizer of matrix sensing loss is also near-optimal for the matrix factorization loss.

Lemma E.4. Let \mathbf{Z}_s^* be a best rank- s solution as defined in Definition 1.1, then we have

$$\|\mathbf{Z}_s^* - \mathbf{X} \mathbf{X}^\top\|_F^2 \leq \|\mathbf{X}_s \mathbf{X}_s^\top - \mathbf{X} \mathbf{X}^\top\|_F^2 + 10\delta \|\mathbf{X} \mathbf{X}^\top\|_F^2.$$

Proof. By the RIP property Definition 3.2 we have

$$\begin{aligned} \|\mathbf{X} \mathbf{X}^\top - \mathbf{Z}_s^*\|_F^2 &\leq (1 - \delta)^{-1} \|\mathcal{A}(\mathbf{X} \mathbf{X}^\top - \mathbf{Z}_s^*)\|_2^2 \\ &\leq (1 - \delta)^{-1} \|\mathcal{A}(\mathbf{X} \mathbf{X}^\top - \mathbf{X}_s \mathbf{X}_s^\top)\|_2^2 \\ &\leq \frac{1 + \delta}{1 - \delta} \|\mathbf{X} \mathbf{X}^\top - \mathbf{X}_s \mathbf{X}_s^\top\|_F^2 \\ &\leq \|\mathbf{X} \mathbf{X}^\top - \mathbf{X}_s \mathbf{X}_s^\top\|_F^2 + 10\delta \|\mathbf{X} \mathbf{X}^\top\|_F^2, \end{aligned}$$

where the second inequality holds due to Definition 1.1. □

We now recall Lemma 4.8.

Lemma 4.8. Under Assumption 3.3, we have $\text{dist}(\mathbf{U}_s^*, \mathbf{X}_s) \leq 40\delta \kappa_* \|\mathbf{X}\|_F$ for any global minimizer \mathbf{U}_s^* of f_s . Moreover, $\|\mathbf{Z}_s^* - \mathbf{X}_s \mathbf{X}_s^\top\|_F \leq 160\delta \kappa_* \sqrt{r_*} \|\mathbf{X}\|^2$.

We prove the statements in this lemma separately in Lemma E.5 and Corollary E.6 below.

Lemma E.5. Suppose that [Assumption 3.3](#) holds. Let \mathbf{U}_s^* be a global minimizer of f_s , then we have

$$\text{dist}(\mathbf{U}_s^*, \mathbf{X}_s) \leq 40\delta\kappa\|\mathbf{X}\|_F.$$

Proof. Define

$$S = \{\mathbf{U} \in \mathbb{R}^{d \times s} : \text{dist}(\mathbf{U}, \mathbf{X}_s) < 0.1\kappa^{-1}\|\mathbf{X}\|\}.$$

First we can show that $\mathbf{U}_s^* \in S$. The main idea is to apply [Lemma A.7](#). Indeed, it is easy to see that

$$\lim_{\|\mathbf{U}\|_F \rightarrow +\infty} F_s(\mathbf{U}) = +\infty,$$

so the condition (1) in [Lemma A.7](#) holds. To check condition (2), we separately analyze the two cases $\mathbf{U} \in \partial S$ and $\mathbf{U} \notin S$.

Firstly, let $\mathbf{U} \in \partial S$, i.e., $\text{dist}^2(\mathbf{U}, \mathbf{X}_s) = 0.1\|\mathbf{X}\|^{-1}\tau$. Assume WLOG that $\text{dist}(\mathbf{U}, \mathbf{X}_s) = \|\mathbf{U} - \mathbf{X}_s\|_F$, then by [Lemma E.2](#) we have

$$\begin{aligned} F_s(\mathbf{U}) - F_s(\mathbf{X}_s) &= \int_0^1 t \langle \nabla F_s(t\mathbf{U} + (1-t)\mathbf{X}_s), \mathbf{U} - \mathbf{X}_s \rangle dt \\ &\geq \int_0^1 0.1\tau t^2 \|\mathbf{U} - \mathbf{X}_s\|_F^2 dt \\ &\geq 10^{-3}\|\mathbf{X}\|^{-2}\tau^3 = 10^{-3}\kappa^{-3}\|\mathbf{X}\|^2. \end{aligned}$$

Secondly, let $\mathbf{U} \notin S$ be a stationary point of f_s . Recall that all the stationary points of F_s are characterized in [Lemma E.1](#), so that for all $\mathbf{U} \notin S$ with $\nabla F_s(\mathbf{U}) = 0$, we have

$$F_s(\mathbf{U}) - F_s^* \geq 0.5(\sigma_s^4 - \sigma_{s+1}^4) \geq 0.5\tau^2.$$

On the other hand, we know from [Lemma E.4](#) that

$$F_s(\mathbf{U}_s^*) - F_s^* \leq 5\delta r_* \|\mathbf{X}\|^4. \quad (55)$$

By [Assumption 3.3](#), we have $5\delta r_* \|\mathbf{X}\|^4 < 10^{-3}\kappa^{-3}\|\mathbf{X}\|^2 < 0.5\tau^2$, so [Lemma A.7](#) implies that $\mathbf{U}_s^* \in S$.

Since $\nabla f_s(\mathbf{U}_s^*) = 0$, we have $\mathcal{A}^* \mathcal{A} (\mathbf{X}\mathbf{X}^\top - \mathbf{U}_s^* (\mathbf{U}_s^*)^\top) \mathbf{U}_s^* = 0$, so that

$$\begin{aligned} \|\nabla F_s(\mathbf{U}_s^*)\|_F &= \left\| (\mathbf{X}\mathbf{X}^\top - \mathbf{U}_s^* (\mathbf{U}_s^*)^\top) \mathbf{U}_s^* \right\|_F \\ &= \left\| (\mathcal{A}^* \mathcal{A} - \mathbf{I}) (\mathbf{X}\mathbf{X}^\top - \mathbf{U}_s^* (\mathbf{U}_s^*)^\top) \mathbf{U}_s^* \right\|_F \\ &\leq \delta \left\| \mathbf{X}\mathbf{X}^\top - \mathbf{U}_s^* (\mathbf{U}_s^*)^\top \right\|_F \|\mathbf{U}_s^*\| \\ &\leq 4\delta \|\mathbf{X}\| \cdot \|\mathbf{X}\mathbf{X}^\top\|_F. \end{aligned}$$

From $\mathbf{U}_s^* \in S$ and [Corollary E.3](#) we can deduce that

$$\text{dist}(\mathbf{U}_s^*, \mathbf{X}_s) \leq 40\delta\tau^{-1}\|\mathbf{X}\|^2\|\mathbf{X}\|_F = 40\delta\kappa\|\mathbf{X}\|_F.$$

□

Corollary E.6. Suppose that [Assumption 3.3](#) holds, then we have $\|\mathbf{Z}_s^* - \mathbf{X}_s \mathbf{X}_s^\top\|_F \leq 80\delta\kappa\sqrt{r_*}\|\mathbf{X}\|^2$ and $\sigma_{\min}((\mathbf{U}_s^*)^\top \mathbf{U}_s^*) \geq \sigma_s^2 - 80\delta\kappa\sqrt{r_*}\|\mathbf{X}\|^2$.

Proof. We assume WLOG that $\|\mathbf{U}_s^* - \mathbf{X}_s\|_F = \text{dist}(\mathbf{U}_s^*, \mathbf{X}_s)$ i.e. $\mathbf{R} = \mathbf{I}$ in [Definition 3.8](#). By [Lemma 4.8](#), we have that

$$\begin{aligned} \left\| \mathbf{U}_s^* (\mathbf{U}_s^*)^\top - \mathbf{X}_s \mathbf{X}_s^\top \right\|_F &\leq 2 \max\{\|\mathbf{U}_s^*\|, \|\mathbf{X}_s\|\} \cdot \|\mathbf{U}_s^* - \mathbf{X}_s\|_F \\ &\leq 160\delta\kappa\|\mathbf{X}\|\|\mathbf{X}\|_F \leq 160\delta\kappa\sqrt{r_*}\|\mathbf{X}\|^2. \end{aligned}$$

which proves the first inequality. Similarly, we have

$$\left\| (\mathbf{U}_s^*)^\top \mathbf{U}_s^* - \mathbf{X}_s^\top \mathbf{X}_s \right\|_F \leq 160\delta\kappa\sqrt{r_*} \|\mathbf{X}\|^2.$$

Hence $\sigma_s^2 - \sigma_{\min} \left((\mathbf{U}_s^*)^\top \mathbf{U}_s^* \right) \leq \left\| (\mathbf{U}_s^*)^\top \mathbf{U}_s^* - \mathbf{X}_s^\top \mathbf{X}_s \right\| \leq 160\delta\kappa\sqrt{r_*} \|\mathbf{X}\|^2$, as desired. \square

Corollary 4.9. *Under Assumption 3.3, we have $\sigma_{\min}(\mathbf{U}_s^*) \geq \frac{1}{2}\sigma_{\min}(\mathbf{X}_s) = \frac{1}{2}\sigma_s \geq \frac{1}{2}\kappa_*^{-\frac{1}{2}} \|\mathbf{X}\|$.*

Proof. Assumption 3.3 implies that $160\delta\kappa\sqrt{r_*} \|\mathbf{X}\|^2 \leq 0.1\kappa^{-1} \|\mathbf{X}\|^2 \leq 0.1\sigma_s^2$, so that the conclusion immediately follows from Corollary E.6. \square

Lemma E.7. *Under Assumption 3.3, suppose that $\mathbf{U}, \mathbf{U}_s^* \in \mathbb{R}^{d \times s}$ such that \mathbf{U}_s^* is a global minimizer of f_s and $\|\mathbf{U} - \mathbf{U}_s^*\|_F = \text{dist}(\mathbf{U}, \mathbf{U}_s^*) \leq 10^{-2}\kappa^{-1} \|\mathbf{X}\|$ (recall that dist is defined in Definition 3.8), then we have*

$$\langle \nabla f_s(\mathbf{U}), \mathbf{U} - \mathbf{U}_s^* \rangle \geq 0.1\kappa^{-1} \|\mathbf{X}\|^2 \|\mathbf{U} - \mathbf{U}_s^*\|_F^2.$$

Proof. By Lemma A.10, $\mathbf{U}^\top \mathbf{U}_s^*$ is symmetric and positive semi-definite. Let $\mathbf{H} = \mathbf{U} - \mathbf{U}_s^*$, then

$$\begin{aligned} \nabla f_s(\mathbf{U}) &= (\mathcal{A}^* \mathcal{A}) (\mathbf{U} \mathbf{U}^\top - \mathbf{X} \mathbf{X}^\top) \mathbf{U} \\ &= (\mathcal{A}^* \mathcal{A}) \left[(\mathbf{H} + \mathbf{U}_s^*) (\mathbf{H} + \mathbf{U}_s^*)^\top - \mathbf{X} \mathbf{X}^\top \right] (\mathbf{H} + \mathbf{U}_s^*) \\ &= \left[(\mathcal{A}^* \mathcal{A}) \left(\mathbf{H} \mathbf{H}^\top + \mathbf{U}_s^* \mathbf{H}^\top + \mathbf{H} (\mathbf{U}_s^*)^\top \right) \right] (\mathbf{H} + \mathbf{U}_s^*) - \mathcal{A}^* \mathcal{A} \left(\mathbf{X} \mathbf{X}^\top - \mathbf{U}_s^* (\mathbf{U}_s^*)^\top \right) \mathbf{H} \end{aligned}$$

where we use the first-order optimality condition

$$\mathcal{A}^* \mathcal{A} \left(\mathbf{X} \mathbf{X}^\top - \mathbf{U}_s^* (\mathbf{U}_s^*)^\top \right) \mathbf{U}_s^* = 0.$$

Since $\|\mathbf{U}_s^*\| \leq 2\|\mathbf{X}\|$ by Lemma 4.8, we may thus deduce that

$$\begin{aligned} &\left\| \nabla f_s(\mathbf{U}) - \left[(\mathbf{H} \mathbf{H}^\top + \mathbf{U}_s^* \mathbf{H}^\top + \mathbf{H} (\mathbf{U}_s^*)^\top) (\mathbf{H} + \mathbf{U}_s^*) - (\mathbf{X} \mathbf{X}^\top - \mathbf{U}_s^* (\mathbf{U}_s^*)^\top) \mathbf{H} \right] \right\|_F \\ &\leq \left\| (\mathcal{A}^* \mathcal{A} - \mathbf{I}) \left(\mathbf{H} \mathbf{H}^\top + \mathbf{U}_s^* \mathbf{H}^\top + \mathbf{H} (\mathbf{U}_s^*)^\top \right) (\mathbf{H} + \mathbf{U}_s^*) \right\|_F + \left\| (\mathcal{A}^* \mathcal{A} - \mathbf{I}) \left(\mathbf{X} \mathbf{X}^\top - \mathbf{U}_s^* (\mathbf{U}_s^*)^\top \right) \mathbf{H} \right\|_F \\ &\leq 50\delta \|\mathbf{X}\|^2 \|\mathbf{H}\|_F \end{aligned}$$

Hence

$$\begin{aligned} &\langle \nabla f_s(\mathbf{U}), \mathbf{U} - \mathbf{U}_s^* \rangle \\ &\geq \left\langle \left(\mathbf{H} \mathbf{H}^\top + \mathbf{U}_s^* \mathbf{H}^\top + \mathbf{H} (\mathbf{U}_s^*)^\top \right) (\mathbf{H} + \mathbf{U}_s^*) - \left(\mathbf{X} \mathbf{X}^\top - \mathbf{U}_s^* (\mathbf{U}_s^*)^\top \right) \mathbf{H}, \mathbf{H} \right\rangle - 50\delta \|\mathbf{X}\|^2 \|\mathbf{H}\|_F^2 \\ &\geq \text{tr} \left(\mathbf{H} (\mathbf{H} + \mathbf{U}_s^*)^\top (\mathbf{H} + \mathbf{U}_s^*) \mathbf{H}^\top + \mathbf{H}^\top \mathbf{U}_s^* \mathbf{H}^\top \mathbf{H} + \left((\mathbf{U}_s^*)^\top \mathbf{H} \right)^2 \right. \\ &\quad \left. - \mathbf{H}^\top \left(\mathbf{X} \mathbf{X}^\top - \mathbf{U}_s^* (\mathbf{U}_s^*)^\top \right) \mathbf{H} \right) - 50\delta \|\mathbf{X}\|^2 \|\mathbf{H}\|_F^2 \\ &\geq \left[\sigma_{\min} \left((\mathbf{U}_s^*)^\top \mathbf{U}_s^* \right) - \left\| \mathbf{X} \mathbf{X}^\top - \mathbf{U}_s^* (\mathbf{U}_s^*)^\top \right\| - 50\delta \|\mathbf{X}\|^2 - 3\|\mathbf{U}_s^*\| \|\mathbf{H}\| - \|\mathbf{H}\|^2 \right] \|\mathbf{H}\|_F^2. \end{aligned}$$

By Corollary E.6 we have $\sigma_{\min} \left((\mathbf{U}_s^*)^\top \mathbf{U}_s^* \right) \geq \sigma_s^2 - 80\delta\kappa \|\mathbf{X}\| \|\mathbf{X}\|_F$ and $\left\| \mathbf{X} \mathbf{X}^\top - \mathbf{U}_s^* (\mathbf{U}_s^*)^\top \right\| \leq \sigma_{s+1}^2 + 80\delta\kappa \|\mathbf{X}\|_F^2$, so that

$$\langle \nabla f_s(\mathbf{U}), \mathbf{U} - \mathbf{U}_s^* \rangle \geq (\sigma_s^2 - \sigma_{s+1}^2 - 160\delta\kappa \|\mathbf{X}\| \|\mathbf{X}\|_F - 50\delta \|\mathbf{X}\|^2 - 3\|\mathbf{U}_s^*\| \|\mathbf{H}\| - \|\mathbf{H}\|^2) \|\mathbf{H}\|_F^2.$$

When Assumption 3.3 on δ is satisfied and $\|\mathbf{H}\| \leq 10^{-2}\tau \|\mathbf{X}\|^{-1}$, the above implies that $\langle \nabla f_s(\mathbf{U}), \mathbf{U} - \mathbf{U}_s^* \rangle \geq 0.5\tau \|\mathbf{H}\|_F^2$, as desired. \square

We now prove Lemmas 4.4 and 4.6 to conclude this section. Both results follow immediately from Lemma E.7.

Lemma 4.4. Under [Assumption 3.3](#), if $\mathbf{U}_s^* \in \mathbb{R}^{d \times s}$ is a global minimizer of f_s , then the set of global minimizers $\arg \min f_s$ is equal to $\{\mathbf{U}_s^* \mathbf{R} : \mathbf{R} \in \mathbb{R}^{s \times s}, \mathbf{R}^\top \mathbf{R} = \mathbf{I}\}$.

Proof. By [Lemma E.5](#) we have that $\text{dist}(\mathbf{U}_s^*, \mathbf{X}_s^*) \leq 40\delta\kappa\|\mathbf{X}\|_F$ holds for any $\mathbf{U}_s^* \in \arg \min f_s$. Suppose now that $\mathbf{U}_s^*, \hat{\mathbf{U}}_s^* \in \arg \min f_s$ such that $\text{dist}(\mathbf{U}_s^*, \hat{\mathbf{U}}_s^*) > 0$. Then we also have that $\text{dist}(\mathbf{U}_s^*, \hat{\mathbf{U}}_s^*) \leq \text{dist}(\mathbf{U}_s^*, \mathbf{X}_s^*) + \text{dist}(\mathbf{X}_s^*, \hat{\mathbf{U}}_s^*) \leq 80\delta\kappa\|\mathbf{X}\|_F < 10^{-2}\kappa^{-1}\|\mathbf{X}\|$, where the last inequality follows from [Assumption 3.3](#). Without loss of generality, we can assume that $\|\mathbf{U}_s^* - \hat{\mathbf{U}}_s^*\|_F = \text{dist}(\mathbf{U}_s^*, \hat{\mathbf{U}}_s^*)$, so that we can apply [Lemma E.7](#) to obtain

$$\langle \nabla f_s(\hat{\mathbf{U}}_s^*), \hat{\mathbf{U}}_s^* - \mathbf{U}_s^* \rangle \geq 0.1\kappa^{-1}\|\mathbf{X}\|^2\|\hat{\mathbf{U}}_s^* - \mathbf{U}_s^*\|_F^2.$$

However, since $\hat{\mathbf{U}}_s^*$ is a global minimizer of f_s , we have $\nabla f_s(\hat{\mathbf{U}}_s^*)$ which is a contradiction. Thus the global minimizer must be unique under the procrustes distance. \square

We recall [Definition 4.5](#) which is now guaranteed to be well-defined by [Lemma 4.4](#).

Definition 4.5. For any $\mathbf{U} \in \mathbb{R}^{d \times s}$, we use $\Pi_s(\mathbf{U})$ to denote the set of *closest* global minimizers of f_s to \mathbf{U} , namely $\Pi_s(\mathbf{U}) = \arg \min\{\|\mathbf{U} - \mathbf{U}_s^*\|_F : \mathbf{U}_s^* \in \arg \min f_s\}$.

Equipped with this definition, [Lemma E.7](#) directly translates into [Lemma 4.6](#):

Lemma 4.6 (Restricted Secant Inequality). Under [Assumption 3.3](#), if a matrix $\mathbf{U} \in \mathbb{R}^{d \times s}$ satisfies $\|\mathbf{U} - \mathbf{U}_s^*\|_F \leq 10^{-2}\kappa_*^{-1}\|\mathbf{X}\|$ for some $\mathbf{U}_s^* \in \Pi_s(\mathbf{U})$, then we have

$$\langle \nabla f_s(\mathbf{U}), \mathbf{U} - \mathbf{U}_s^* \rangle \geq 0.1\kappa_*^{-1}\|\mathbf{X}\|^2\|\mathbf{U} - \mathbf{U}_s^*\|_F^2. \quad (8)$$

F. Proofs for [Theorems 4.1](#) and [4.2](#)

In this section, we prove the main theorems based on our key lemmas introduced in [Section 4.1](#).

Based on [Lemma 4.6](#), we first prove the following lemma, which shows that GD initialized near global minimizers converges linearly.

Lemma F.1. Suppose that [Assumptions 3.3](#) and [3.7](#) hold. Let $\{\hat{\mathbf{U}}_t\}_{t \geq 0}$ be a trajectory of GD that optimizes f_s with step size μ , starting with $\hat{\mathbf{U}}_0$. Also let \mathbf{U}_s^* be a global minimizer of f_s . If $\text{dist}(\hat{\mathbf{U}}_0, \mathbf{U}_s^*) \leq 10^{-2}\kappa^{-1}\|\mathbf{X}\|$, then for all $t \geq 0$,

$$\text{dist}^2(\hat{\mathbf{U}}_t, \mathbf{U}_s^*) \leq (1 - 0.05\tau\mu)^t \text{dist}^2(\hat{\mathbf{U}}_{\alpha,0}, \mathbf{U}_s^*). \quad (56)$$

Proof. We prove [\(56\)](#) by induction. It is easy to check that [\(56\)](#) holds for $t = 0$. Now we show that [\(56\)](#) holds for $t + 1$ assuming it holds for t .

Since $\text{dist}(\hat{\mathbf{U}}_0, \mathbf{U}_s^*) \leq 10^{-2}\kappa^{-1}\|\mathbf{X}\|$, we have

$$\|\hat{\mathbf{U}}_0\| \leq \|\mathbf{U}_s^*\| + 10^{-2}\kappa^{-1}\|\mathbf{X}\| \leq 2\|\mathbf{X}\|. \quad (57)$$

Let \mathbf{R} be the orthogonal matrix such that $\mathbf{U}_s^* \mathbf{R} \in \Pi(\mathbf{U}_t)$, then $\|\mathbf{U} - \mathbf{U}_s^* \mathbf{R}\|_F = \text{dist}(\mathbf{U}_t, \mathbf{U}_s^*)$. We first bound the gradient $\nabla f(\hat{\mathbf{U}}_{\alpha,t})$ as follows:

$$\begin{aligned} \|\nabla f(\hat{\mathbf{U}}_{\alpha,t})\|_F &= \left\| \mathcal{A}^* \mathcal{A} \left(\mathbf{X} \mathbf{X}^\top - \hat{\mathbf{U}}_{\alpha,t} \hat{\mathbf{U}}_{\alpha,t}^\top \right) \hat{\mathbf{U}}_{\alpha,t} \right\|_F \\ &\leq \left\| \mathcal{A}^* \mathcal{A} \left(\mathbf{X} \mathbf{X}^\top - \hat{\mathbf{U}}_{\alpha,t} \hat{\mathbf{U}}_{\alpha,t}^\top \right) \right\| \|\hat{\mathbf{U}}_{\alpha,t} - \mathbf{U}_s^*\|_F + \left\| \left(\hat{\mathbf{U}}_{\alpha,t} \hat{\mathbf{U}}_{\alpha,t}^\top - \mathbf{U}_s^* (\mathbf{U}_s^*)^\top \right) \mathbf{U}_s^* \right\|_F \\ &\leq 20\|\mathbf{X}\|^2 \|\hat{\mathbf{U}}_{\alpha,t} - \mathbf{U}_s^*\|_F \end{aligned} \quad (58)$$

where we use (57) and the RIP property (Assumption 3.3). It follows that

$$\text{dist}^2(\hat{U}_{t+1}, U_s^*) \leq \left\| \hat{U}_{t+1} - U_s^* \mathbf{R} \right\|_F^2 \quad (59a)$$

$$\begin{aligned} &= \left\| \hat{U}_t - \mu \nabla f(\hat{U}_{\alpha,t}) - U_s^* \mathbf{R} \right\|_F^2 \\ &= \left\| \hat{U}_{\alpha,t} - U_s^* \mathbf{R} \right\|_F^2 - \mu \left\langle \nabla f(\hat{U}_{\alpha,t}), \hat{U}_t - U_s^* \mathbf{R} \right\rangle + \mu^2 \left\| \nabla f(\hat{U}_{\alpha,t}) \right\|_F^2 \\ &\leq (1 - 0.1\tau\mu + 400\|\mathbf{X}\|^4\mu^2) \left\| \hat{U}_{\alpha,t} - U_s^* \mathbf{R} \right\|_F^2 \end{aligned} \quad (59b)$$

where (59a) follows from the definition of dist, and (59b) is due to Lemma 4.6 and (58). Finally, (56) follows from Assumption 3.7. \square

We then note the following proposition, which is straightforward from Lemma C.12 and Theorem C.15. In the following we use $U_{\alpha,t}$ to denote the t -th iteration of GD when initialized at $U_0 = \alpha \bar{U}$.

Proposition F.2. *Suppose that Assumptions 3.3, 3.5 and 3.7 hold and α is sufficiently small. Then there exist matrices $U_{\alpha,t}^{1r}$ for $t = -T_{\alpha,s}^{\text{ft}}, -T_{\alpha,s}^{\text{ft}} + 1, \dots, 0$ with rank $\leq s$ (where 1r stands for low rank) and a constant $C_6 = C_6(\mathbf{X}, \bar{U})$ (defined in (49)) such that*

$$\max_{T_{\alpha,s}^{\text{ft}} \leq t \leq 0} \left\| U_{\alpha,t}^{1r} - U_{\alpha, T_{\alpha,s}^{\text{ft}}+t} \right\|_F = C_6 \cdot \alpha^{\frac{1}{4\kappa}}$$

where $T_{\alpha,s}^{\text{ft}}$ is defined in Theorem C.15 and moreover

$$\left\| U_{\alpha,0}^{1r} (U_{\alpha,0}^{1r})^\top - \mathbf{Z}_s^* \right\|_F \leq 2 \times 10^5 \kappa^3 \|\mathbf{X}\|^2 r_* \delta.$$

where $\mathbf{Z}_s^* = U_s^* (U_s^*)^\top$ is the best rank- s solution as defined in Definition 1.1.

Proof. It follows from Corollary C.16 that $\max_{1 \leq t \leq T_{\alpha,s}^{\text{ft}}} \|U_t \mathbf{W}_{t,\perp}\| \leq C_6(\mathbf{X}, \bar{U}) \cdot \alpha^{\frac{1}{4\kappa}}$ (recall that C_5 is defined in Lemma C.11 and $T_{\alpha,s}^{\text{ft}}$ is defined in (48)). We choose $U_{\alpha,t}^{1r} = U_{T_{\alpha,s}^{\text{ft}}+t} \mathbf{W}_{T_{\alpha,s}^{\text{ft}}+t} \mathbf{W}_{T_{\alpha,s}^{\text{ft}}+t}^\top$, then $\text{rank}(\bar{U}_t) \leq s$ and moreover by Theorem C.15 we have $\left\| \mathbf{X}_s \mathbf{X}_s^\top - U_{\alpha,0}^{1r} (U_{\alpha,0}^{1r})^\top \right\|_F \leq 10^5 \kappa^3 \|\mathbf{X}\|^2 r_* \delta$. On the other hand, by Lemma 4.8 we have that $\left\| \mathbf{Z}_s^* - \mathbf{X}_s \mathbf{X}_s^\top \right\|_F \leq 80\delta\kappa\sqrt{r_*} \|\mathbf{X}\|^2$. Thus $\left\| U_{\alpha,0}^{1r} (U_{\alpha,0}^{1r})^\top - \mathbf{Z}_s^* \right\|_F \leq 2 \times 10^5 \kappa^3 \|\mathbf{X}\|^2 r_* \delta$ as desired. \square

Let $\hat{U}_{\alpha,0} = U_{T_{\alpha,s}^{\text{ft}}} \mathbf{W}_{T_{\alpha,s}^{\text{ft}}} \in \mathbb{R}^{d \times s}$, then it satisfies $\hat{U}_{\alpha,0} \hat{U}_{\alpha,0}^\top = U_{\alpha,0}^{1r} (U_{\alpha,0}^{1r})^\top$. The following corollary shows that $\hat{U}_{\alpha,0}$ is close to U_s^* in terms of the procrustes distance.

Corollary F.3. *We have $\text{dist}(\hat{U}_{\alpha,0}, U_s^*) \leq 3 \times 10^6 \kappa^4 r_* \|\mathbf{X}\| \delta$.*

Proof. We know from Lemma 4.8 that $\text{dist}(U_s^*, \mathbf{X}_s) \leq 40\delta\kappa \|\mathbf{X}\|_F$, so it remains to bound $\text{dist}(\hat{U}_{\alpha,0}, \mathbf{X}_s)$.

The proof idea is the same as that of Lemma 4.8, so we only provide a proof sketch here. It has been shown in the proof of Proposition F.2 that

$$F_s(\hat{U}_{\alpha,0}) := \frac{1}{2} \left\| \mathbf{X}_s \mathbf{X}_s^\top - \hat{U}_{\alpha,0} \hat{U}_{\alpha,0}^\top \right\|_F^2 \leq r_* \left\| \mathbf{X}_s \mathbf{X}_s^\top - U_{\alpha,0}^{1r} (U_{\alpha,0}^{1r})^\top \right\|_F^2 \leq 4 \times 10^{10} \kappa^6 r_*^2 \|\mathbf{X}\|^4 \delta^2 \leq 0.5\tau^2.$$

Note that F_s is the matrix factorization loss with $\mathbf{X}_s \mathbf{X}_s^\top$ being the ground-truth, so the local RSI condition (Lemma E.2) still holds. By the same reason as (55), we deduce that $\text{dist}(\hat{U}_{\alpha,0}, \mathbf{X}_s) \leq 0.1 \|\mathbf{X}\|^{-1} \tau$, i.e., $\hat{U}_{\alpha,0}$ is in the local region around \mathbf{X}_s in which the RSI condition holds. Finally, it follows from the local RSI condition that

$$\text{dist}(\hat{U}_{\alpha,0}, \mathbf{X}_s) \leq 10\tau^{-1} \left\| \nabla F_s(\hat{U}_{\alpha,0}) \right\|_F \leq 10\tau^{-1} \|\hat{U}_{\alpha,0}\| \left\| \mathbf{X}_s \mathbf{X}_s^\top - \hat{U}_{\alpha,0} \hat{U}_{\alpha,0}^\top \right\|_F \leq 3 \times 10^6 \kappa^4 r_* \|\mathbf{X}\| \delta.$$

The conclusion follows. \square

We are now ready to complete the proof of [Theorems 4.1](#) and [4.2](#).

Theorem 4.2 (Convergence in the under-parameterized regime). *Suppose that $\hat{r} \leq r^*$, then there exists a constant $\bar{\alpha} > 0$ such that when $\alpha < \bar{\alpha}$, we have $\lim_{t \rightarrow +\infty} \mathbf{U}_{\alpha,t} \mathbf{U}_{\alpha,t}^\top = \mathbf{Z}_{\hat{r}}^*$.*

Proof. When $\hat{r} \leq r_*$, the parameterization itself ensures that $\mathbf{U}_{\alpha,t}$ is low-rank, so that we can choose $\mathbf{U}_{\alpha,t}^{1r} = \mathbf{U}_{\alpha, T_{\alpha, \hat{r}}^{ft} + t}$ and $\hat{\mathbf{U}}_{\alpha,0} = \mathbf{U}_{\alpha, T_{\alpha, \hat{r}}^{ft}}$ in [Proposition F.2](#) and [Corollary F.3](#) (for $s = \hat{r}$). The proof that these choices satisfy all required conditions are identical to our proofs for these two lemmas in the general setting, and we omit them here.

Applying [Lemma F.1](#), we can thus deduce that $\lim_{t \rightarrow +\infty} \text{dist}(\hat{\mathbf{U}}_{\alpha,t}, \mathbf{U}_{\hat{r}}^*) = 0$. This means that $\lim_{t \rightarrow +\infty} \text{dist}(\mathbf{U}_{\alpha,t}, \mathbf{U}_{\hat{r}}^*) = 0$. Recall that $\mathbf{Z}_{\hat{r}}^* = \mathbf{U}_{\hat{r}}^* \mathbf{U}_{\hat{r}}^{*\top}$, so the conclusion immediately follows. \square

Theorem 4.1. *Under [Assumptions 3.3](#), [3.5](#) and [3.7](#), consider GD (3) with initialization $\mathbf{U}_{\alpha,0} = \alpha \bar{\mathbf{U}}$ for solving the matrix sensing problem (1). There exist universal constants c, M , constant $C = C(\mathbf{X}, \bar{\mathbf{U}})$ and a sequence of time points $T_\alpha^1 < T_\alpha^2 < \dots < T_\alpha^{\hat{r} \wedge r_*}$ such that for all $1 \leq s \leq \hat{r} \wedge r_*$, the following holds when α is sufficiently small:*

$$\left\| \mathbf{U}_{\alpha, T_\alpha^s} \mathbf{U}_{\alpha, T_\alpha^s}^\top - \mathbf{Z}_s^* \right\|_F \leq C \alpha^{\frac{1}{M\kappa_*^2}}, \quad (5)$$

where we recall that \mathbf{Z}_s^* is the best rank- s solution defined in [Definition 1.1](#). Moreover, GD follows an incremental learning procedure: we have $\lim_{\alpha \rightarrow 0} \max_{1 \leq t \leq T_\alpha^s} \sigma_{s+1}(\mathbf{U}_{\alpha,t}) = 0$ for all $1 \leq s \leq \hat{r} \wedge r_*$, where $\sigma_i(\mathbf{A})$ denotes the i -th largest singular value of a matrix \mathbf{A} .

Proof. Recall that $\left\| \mathbf{U}_{T_{\alpha,s}^{ft}} - \bar{\mathbf{U}}_0 \right\|_F = o(1)$ ($\alpha \rightarrow 0$) where $T_{\alpha,s}^{ft}$ is defined in [Proposition F.2](#); we omit the dependence on α to simplify notations. We also note that by the update of GD, we have $\bar{\mathbf{U}}_t \bar{\mathbf{U}}_t^\top = \hat{\mathbf{U}}_{\alpha,t} \hat{\mathbf{U}}_{\alpha,t}^\top$ for all $t \geq 0$.

By [Lemma F.1](#), we have that $\text{dist}^2(\hat{\mathbf{U}}_{\alpha,t}, \mathbf{U}_s^*) \leq (1 - 0.05\tau\mu)^t \text{dist}^2(\hat{\mathbf{U}}_{\alpha,0}, \mathbf{U}_s^*)$ and, in particular, $\left\| \hat{\mathbf{U}}_{\alpha,t} \right\| \leq 2\|\mathbf{X}\|$ for all t . Thus $\left\| \bar{\mathbf{U}}_t \right\| \leq 2\|\mathbf{X}\|$ as well. Moreover, recall that $\|\mathbf{U}_t\| \leq 3\|\mathbf{X}\|$ for all t . It's easy to see that the matrix sensing loss f is L -smooth in $\{\mathbf{U} \in \mathbb{R}^{d \times r} : \|\mathbf{U}\| \leq 3\|\mathbf{X}\|\}$ for some constant $L = \mathcal{O}(\|\mathbf{X}\|^2)$, so it follows from [Lemma A.8](#) that

$$\left\| \mathbf{U}_{T_{\alpha,s}^{ft} + t} - \bar{\mathbf{U}}_t \right\|_F \leq (1 + \mu L)^t \left\| \mathbf{U}_{T_{\alpha,s}^{ft}} - \bar{\mathbf{U}}_0 \right\|_F.$$

On the other hand, since $\text{dist}^2(\hat{\mathbf{U}}_{\alpha,t}, \mathbf{U}_s^*) \leq (1 - 0.05\tau\mu)^t \text{dist}^2(\hat{\mathbf{U}}_{\alpha,0}, \mathbf{U}_s^*)$, we can deduce that

$$\begin{aligned} \left\| \mathbf{U}_{T_{\alpha,s}^{ft} + t} \mathbf{U}_{T_{\alpha,s}^{ft} + t}^\top - \mathbf{Z}_s \right\|_F &\leq \left\| \mathbf{U}_{T_{\alpha,s}^{ft} + t} \mathbf{U}_{T_{\alpha,s}^{ft} + t}^\top - \bar{\mathbf{U}}_t \bar{\mathbf{U}}_t^\top \right\|_F + \left\| \bar{\mathbf{U}}_t \bar{\mathbf{U}}_t^\top - \mathbf{U}_s^* (\mathbf{U}_s^*)^\top \right\|_F \\ &= \left\| \mathbf{U}_{T_{\alpha,s}^{ft} + t} \mathbf{U}_{T_{\alpha,s}^{ft} + t}^\top - \bar{\mathbf{U}}_t \bar{\mathbf{U}}_t^\top \right\|_F + \left\| \hat{\mathbf{U}}_{\alpha,t} \hat{\mathbf{U}}_{\alpha,t}^\top - \mathbf{U}_s^* (\mathbf{U}_s^*)^\top \right\|_F \\ &\leq 3\|\mathbf{X}\| \left(\left\| \mathbf{U}_{T_{\alpha,s}^{ft} + t} - \bar{\mathbf{U}}_t \right\|_F + \text{dist}(\hat{\mathbf{U}}_{\alpha,t}, \mathbf{U}_s^*) \right) \\ &\leq 3\|\mathbf{X}\| \left((1 + \mu L)^t \left\| \mathbf{U}_{T_{\alpha,s}^{ft}} - \bar{\mathbf{U}}_0 \right\|_F + (1 - 0.05\tau\mu)^{\frac{t}{2}} \text{dist}^2(\hat{\mathbf{U}}_{\alpha,0}, \mathbf{U}_s^*) \right) \end{aligned}$$

Since when $\alpha \rightarrow 0$, $\left\| \mathbf{U}_{T_{\alpha,s}^{ft}} - \bar{\mathbf{U}}_0 \right\|_F = \mathcal{O}(\alpha^{\frac{1}{4\kappa}})$, it's easy to see that there exists a time $t = t_\alpha^s$ so that we have $\max_{-T_{\alpha,s}^{ft} \leq t \leq t_\alpha^s} \left\| \mathbf{U}_{T_{\alpha,s}^{ft} + t} - \bar{\mathbf{U}}_t \right\|_F = \mathcal{O}\left(\alpha^{\frac{1}{M_1\kappa^2}}\right)$ and $\left\| \mathbf{U}_{T_{\alpha,s}^{ft} + t} \mathbf{U}_{T_{\alpha,s}^{ft} + t}^\top - \mathbf{Z}_s \right\|_F = \mathcal{O}\left(\alpha^{\frac{1}{M_1\kappa^2}}\right)$ as well, where c_1 is a universal constant. Let $T_\alpha^s = T_{\alpha,s}^{ft} + t_\alpha^s$, then $\left\| \mathbf{U}_{T_\alpha^s} \mathbf{U}_{T_\alpha^s}^\top - \mathbf{Z}_s \right\|_F = o(1)$ holds. Recall that $\text{rank}(\mathbf{U}_t) \leq s$, so that $\max_{0 \leq t \leq T_\alpha^s} \sigma_{s+1}(\mathbf{U}_t) = o(1)$. Finally, for all $0 \leq s < \hat{r} \wedge r_*$, we need to show that $T_\alpha^s < T_\alpha^{s+1}$. Indeed, by [Corollary E.6](#) and the [Assumption 3.3](#) we have $\sigma_{s+1}^2(\mathbf{U}_{T_\alpha^s}) \geq \sigma_{s+1}(\mathbf{Z}_{s+1}) - o(1) \geq 0.5\sigma_{s+1}^2$, so that $T_\alpha^{s+1} > T_\alpha^s$, as desired. \square

G. The landscape of matrix sensing with rank-1 parameterization

In this section, we establish a local strong-convexity result [Lemma G.2](#) for rank-1 parameterized matrix sensing. This result is stronger than the RSI condition we established for general ranks, though the latter is sufficient for our analysis.

Lemma G.1. Define the full-observation loss with rank-1 parameterization

$$g_1(\mathbf{u}) = \frac{1}{4} \|\mathbf{u}\mathbf{u}^T - \mathbf{X}\mathbf{X}^T\|_F^2.$$

Then the global minima of g_1 are $\mathbf{u}^* = \sigma_1 \mathbf{v}_1$ and $-\mathbf{u}^*$. Moreover, suppose that $g(\mathbf{u}) - g(\mathbf{u}^*) \leq 0.5\tau_1$ where $\tau_1 = \sigma_1^2 - \sigma_2^2$ is the eigengap, then we must have

$$\|\mathbf{u} - \mathbf{u}^*\|^2 \leq 20\tau_1^{-1} (g_1(\mathbf{u}) - g_1(\mathbf{u}^*)).$$

Proof. We can assume WLOG that $\mathbf{X}\mathbf{X}^T = \text{diag}(\sigma_1^2, \dots, \sigma_{r_*}^2, 0, \dots, 0)$. Then

$$g_1(\mathbf{u}) = \frac{1}{4} \left(\|\mathbf{u}\|_2^4 - 2 \sum_{i=1}^s \sigma_i^2 \mathbf{u}_i^2 + \|\mathbf{X}^T \mathbf{X}\|_F^2 \right) \quad (60a)$$

$$\geq \frac{1}{4} (\|\mathbf{u}\|_2^4 - 2\sigma_1^2 \|\mathbf{u}\|_2^2 + \|\mathbf{X}^T \mathbf{X}\|_F^2) \quad (60b)$$

$$\geq \frac{1}{4} (\|\mathbf{X}^T \mathbf{X}\|_F^2 - \sigma_1^4) \quad (60c)$$

where equality holds if and only if $\mathbf{u}_2 = \dots = \mathbf{u}_d = 0$ and $\|\mathbf{u}\|^2 = \sigma_1^2$ i.e. $\mathbf{u} = \pm\sigma_1 \mathbf{e}_1$. Moreover, suppose that $g_1(\mathbf{u}) - g_1(\mathbf{u}^*) \leq 0.5\tau_1$, it follows from (60b) that $\tau_1 \sum_{i=2}^d \mathbf{u}_i^2 \leq 2(g_1(\mathbf{u}) - g_1(\mathbf{u}^*))$ which implies that $\sum_{i=2}^d \mathbf{u}_i^2 \leq 2\tau_1^{-1} (g_1(\mathbf{u}) - g_1(\mathbf{u}^*))$. Also (60c) yields $\|\mathbf{u}\|^2 - \sigma_1^2 \leq 4\sqrt{g_1(\mathbf{u}) - g_1(\mathbf{u}^*)}$. Assume WLOG that $\mathbf{u}_1 > 0$, then we have

$$\begin{aligned} \|\mathbf{u} - \sigma_1 \mathbf{e}_1\|^2 &\leq \sigma_1^{-2} (\mathbf{u}_1^2 - \sigma_1^2)^2 + \sum_{i=2}^d \mathbf{u}_i^2 \\ &\leq 20\tau_1^{-1} (g_1(\mathbf{u}) - g_1(\mathbf{u}^*)). \end{aligned}$$

□

Lemma G.2. Let

$$f_1(\mathbf{u}) = \frac{1}{4} \|\mathcal{A}(\mathbf{u}\mathbf{u}^T - \mathbf{X}\mathbf{X}^T)\|_2^2, \quad \mathbf{u} \in \mathbb{R}^d.$$

Suppose that $\delta \leq 10^{-3} \|\mathbf{X}\|^{-2} \tau_1$, then there exists constants a_1 and ι , such that f_1 is locally ι -strongly convex in $\mathcal{B}_1 = \mathcal{B}(\sigma_1 \mathbf{v}_1, a_1) \subset \mathbb{R}^d$. Furthermore, there is a unique global minima of f_1 inside \mathcal{B}_1 .

Proof. Recall that we defined the full observation loss $g_1(\mathbf{u}) = \frac{1}{4} \|\mathbf{u}\mathbf{u}^T - \mathbf{X}\mathbf{X}^T\|_F^2$. Let $h_1 = f_1 - g_1$, then

$$\begin{aligned} \|\nabla^2 h_1(\mathbf{u})\| &= \frac{1}{2} \|(\mathcal{A}^* \mathcal{A} - I)(\mathbf{u}\mathbf{u}^T - \mathbf{X}\mathbf{X}^T) + 2(\mathcal{A}^* \mathcal{A} - I)\mathbf{u}\mathbf{u}^T\| \\ &\leq \delta (2\|\mathbf{u}\|^2 + \|\mathbf{X}\|^2). \end{aligned}$$

When $\|\mathbf{u} - \sigma_1 \mathbf{v}_1\|^2 \leq 0.1 \min\{\sigma_1^2, \tau_1\}$ (recall $\tau_1 = \sigma_1^2 - \sigma_2^2$),

$$\sigma_{\min}(\nabla^2 g_1(\mathbf{u})) = \frac{1}{2} \sigma_{\min}(\|\mathbf{u}\|^2 I + 2\mathbf{u}\mathbf{u}^T - \mathbf{X}\mathbf{X}^T) \geq 0.4\tau_1.$$

Hence we have

$$\sigma_{\min}(\nabla^2 f_1(\mathbf{u})) \geq (\nabla^2 g_1(\mathbf{u})) - \|\nabla^2 h_1(\mathbf{u})\| \geq 0.4\tau_1 - 4\|\mathbf{X}\|^2 \delta \geq 0.2\tau_1,$$

i.e. strong-convexity holds for $a_1^2 = 0.1 \min\{\sigma_1^2, \tau_1\}$ and $\iota = 0.2\tau_1$.

Let \mathbf{u}^* be a global minima of f_1 , then we must have $\|\mathbf{u}^*\| \leq 2\|\mathbf{X}\|$ (otherwise $f_1(\mathbf{u}) > f_1(0)$). We can thus deduce that

$$\begin{aligned} g_1(\mathbf{u}^*) &\leq f_1(\mathbf{u}^*) + \frac{1}{4} |\langle \mathbf{u}\mathbf{u}^T - \mathbf{X}\mathbf{X}^T, (\mathcal{A}^* \mathcal{A} - I)(\mathbf{u}\mathbf{u}^T - \mathbf{X}\mathbf{X}^T) \rangle| \\ &\leq f_1(\mathbf{u}) + 10\delta \|\mathbf{X}\|^2 \leq g_1(\mathbf{u}) + 20\delta \|\mathbf{X}\|^2. \end{aligned}$$

It follows from Lemma G.1 and our assumption on δ that $\min\{\|\mathbf{u}^* - \sigma_1 \mathbf{v}_1\|^2, \|\mathbf{u}^* + \sigma_1 \mathbf{v}_1\|^2\} \leq \frac{1}{2} a_1^2$. Moreover, by strong convexity, there exists only one global minima in \mathcal{B}_1 , which concludes the proof. □