

Convergence of Gradient Descent with Small Initialization for Unregularized Matrix Completion

Jianhao Ma

JIANHAO@UMICH.EDU and Salar Fattahi

FATTAHI@UMICH.EDU

University of Michigan, Ann Arbor

Editors: Shipra Agrawal and Aaron Roth

Abstract

We study the problem of symmetric matrix completion, where the goal is to reconstruct a positive semidefinite matrix $\mathbf{X}^* \in \mathbb{R}^{d \times d}$ of rank- r , parameterized by $\mathbf{U}\mathbf{U}^\top$, from only a subset of its observed entries. We show that the vanilla gradient descent (GD) with small initialization provably converges to the ground truth \mathbf{X}^* without requiring any explicit regularization. This convergence result holds true even in the over-parameterized scenario, where the true rank r is unknown and conservatively over-estimated by a search rank $r' \gg r$. The existing results for this problem either require explicit regularization, a sufficiently accurate initial point, or exact knowledge of the true rank r .

In the over-parameterized regime where $r' \geq r$, we show that, with $\tilde{\Omega}(dr^9)$ observations, GD with an initial point $\|\mathbf{U}_0\| \leq O(\epsilon)$ converges near-linearly to an ϵ -neighborhood of \mathbf{X}^* . Consequently, smaller initial points result in increasingly accurate solutions. Surprisingly, neither the convergence rate nor the final accuracy depends on the over-parameterized search rank r' , and they are only governed by the true rank r . In the exactly-parameterized regime where $r' = r$, we further enhance this result by proving that GD converges at a faster rate to achieve an arbitrarily small accuracy $\epsilon > 0$, provided the initial point satisfies $\|\mathbf{U}_0\| = O(1/d)$. At the crux of our method lies a novel *weakly-coupled leave-one-out analysis*, which allows us to establish the global convergence of GD, extending beyond what was previously possible using the classical leave-one-out analysis.

Keywords: Matrix completion, implicit regularization, leave-one-out analysis

1. Introduction

Matrix completion is a fundamental problem in the field of machine learning, where the objective is to reconstruct a positive semidefinite (PSD) matrix of rank- r , denoted as $\mathbf{X}^* \in \mathbb{R}^{d \times d}$, from only a subset of its observed entries. The most natural approach to solve this problem involves minimizing the following mean squared error:

$$\min_{\mathbf{U} \in \mathbb{R}^{d \times r'}} f(\mathbf{U}) = \frac{1}{4p} \left\| \mathcal{P}_\Omega(\mathbf{U}\mathbf{U}^\top - \mathbf{X}^*) \right\|_F^2. \quad (\text{MC})$$

Here, p represents the probability of observing each entry in \mathbf{X}^* , Ω denotes the set of observed entries, and \mathcal{P}_Ω shows the projection operation onto the set of matrices supported by Ω . When the true rank r is unknown, it is often over-estimated by the search rank $r' \geq r$, leading to what is referred to as *over-parameterized* matrix completion.

A prominent application of matrix completion is in collaborative filtering (Gleich and Lim, 2011). Additionally, it has applications in other areas, including image reconstruction (Hu et al., 2018), fast kernel matrix approximation (Graepel, 2002; Paisley and Carin, 2010), and more recently, in teaching arithmetic to transformers (Lee et al., 2023).

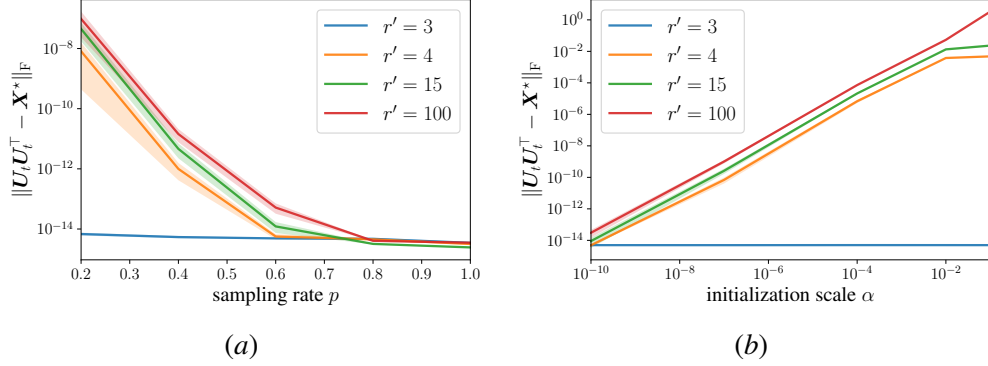


Figure 1: The performance of vanilla GD with small initialization and constant step-size on MC without any explicit regularization. In all experiments, we set the dimension $d = 100$ and the true rank $r = 3$. (a): Higher sampling rates improve the final accuracy of GD for the over-parameterized MC, but have no impact on the exactly-parameterized MC. (b): Increasing the initialization scale hampers the final accuracy of GD for the over-parameterized MC, but has no impact on the exactly-parameterized MC.

Perhaps one of the most natural approaches for solving the above optimization problem is (vanilla) gradient descent (GD): given an initial point U_0 and a fixed step-size $\eta > 0$, generate a sequence of iterates $\{U_t\}_{t=1}^T$ according to $U_{t+1} = U_t - \eta \nabla f(U_t)$. Despite its simplicity and desirable practical performance (see Figure 1 and the experiments in (Zheng and Lafferty, 2016)), the conditions under which the GD converges globally to the ground truth X^* have remained a long-standing mystery.

A line of research has been devoted to studying gradient-based algorithms with *explicit regularization* (Sun and Luo, 2016; Zheng and Lafferty, 2016; Jain et al., 2013; Ge et al., 2016). These methods typically incorporate either an $\ell_{2,\infty}$ -norm regularizer or a projection step to constrain the iterates within a set with $\ell_{2,\infty}$ -norm bounds to promote *incoherence* (see Definition 2). However, the use of $\ell_{2,\infty}$ -norm regularization or projection techniques often introduce more tuning parameters, and has been found to be unnecessary in most cases (Zheng and Lafferty, 2016; Ma et al., 2018)¹.

On the other hand, the convergence of GD without explicit regularization was initially tackled by Ma et al. (2018) in the context of symmetric matrix completion, and subsequently extended by Chen et al. (2020) to asymmetric settings. However, these studies consider a very special case of matrix completion where GD is initialized sufficiently close to the ground truth, and the rank of the ground truth r is known. In practice, however, GD converges even if it is initialized far from the ground truth and the rank is over-parameterized $r' \gg r$ (see Figure 1).

Therefore, the following question still remains open:

Why does GD with a small initialization efficiently converge to the ground truth of MC in the absence of explicit regularization, even in the general rank- r case where $r' \geq r \geq 1$?

Recently, Kim and Chung (2022) answered the above question for the special case of rank-1 symmetric matrix completion with $r' = r = 1$, but their proposed approach does not extend to the general rank- r case. It is also worth noting that the above question has been addressed for another class of matrix factorization problems satisfying a norm-preserving property called *restricted isometry property* (RIP). Problems that satisfy this property include matrix sensing (Li et al., 2018;

1. It is worth noting that explicit regularization can still offer significant advantages in scenarios with extremely limited sample sizes (Sun, 2015).

(Stöger and Soltanolkotabi, 2021; Ma and Fattahi, 2023a) and sparse recovery (Vaskevicius et al., 2019). However, a significant challenge arises with matrix completion, as it *does not* satisfy the restricted isometry property. Consequently, the existing methodologies built upon RIP are not directly applicable to matrix completion.

1.1. Summary of Contributions

In this work, we provide a complete answer to the aforementioned question. A comparison of our results with other studies on unregularized matrix completion can be found in Table 1. The key contributions of our work are as follows:

- **Convergence of GD with small initialization in over-parameterized regime:** When the rank of the ground truth r is unknown and over-parameterized by $r' \geq r \geq 1$, we prove that GD with small initialization converges to the ground truth at a near-linear rate. Surprisingly, neither the convergence rate nor the final accuracy depends on the over-parameterized search rank r' , and they are only governed by the true rank r . In particular, given an initial point that satisfies $\|U_0\| \leq O(\epsilon)$ for some $\epsilon > 0$ and a sampling rate of $p = \tilde{\Omega}(r^9 \log^6(1/\epsilon)/d)$, GD converges to ϵ -neighborhood of X^* in $O(\log^4(1/\epsilon))$ iterations. Therefore, a smaller initial point or a larger sampling rate can improve the final error of GD. The empirical observation presented in Figure 1 provides further support for this result.
- **Improved results in the exactly-parameterized regime:** We show that GD enjoys an improved convergence when the rank of the ground truth r is known and $r' = r \geq 1$. In particular, given an initial point that satisfies $\|U_0\| \leq O(1/d)$ and a sampling rate of $p = \tilde{\Omega}(r^9/d)$, GD converges to ϵ -neighborhood of X^* in $O(\log(1/\epsilon))$ iterations for any arbitrarily small $\epsilon > 0$. A key distinction from the over-parameterized setting is that the final error of GD remains unaffected by the initialization scale or the sample size, provided that they meet certain thresholds. This is also evident in Figure 1. When $r = O(1)$, the resulting sample complexity is information-theoretically optimal (modulo logarithmic factors).
- **Weakly-coupled leave-one-out analysis:** In order to establish the implicit regularization of GD for MC, a pivotal technique is a decoupling mechanism known as leave-one-out analysis, a trick rooted in probability and random matrix theory. However, the current theory based on this technique is only effective when the iterates are sufficiently close to the ground truth. At the crux of our technical analysis lies an extension of the classical leave-one-out analysis to the global setting, which we term *weakly-coupled leave-one-out analysis*. In essence, our proposed method relaxes the requirement for the initial iterates to be sufficiently close to the ground truth, making it particularly suitable for the global convergence analysis of GD.

Notations. We use bold uppercase letters X, Y to denote matrices and bold lowercase letters x, y to denote vectors. For vectors, we use $\|\cdot\|$ to denote ℓ_2 -norm, and for matrices we use $\|\cdot\|$ and $\|\cdot\|_F$ to denote operator norm and Frobenius norm, respectively. For matrix $X \in \mathbb{R}^{d_1 \times d_2}$, we denote by $X_{i,j}$ the (i, j) -th element of X , $X_{i,\cdot}$ the i -th row, and $X_{\cdot,j}$ the j -th column. The $\ell_{2,\infty}$ -norm of X , denoted as $\|X\|_{2,\infty}$, is defined as $\max_i \|X_{i,\cdot}\|$. Additionally, we define the singular values of $X \in \mathbb{R}^{d_1 \times d_2}$ as $\sigma_1(X) \geq \sigma_2(X) \geq \dots \geq \sigma_{\min\{d_1, d_2\}}(X) \geq 0$. The set of all orthogonal matrices is denoted by $\mathcal{O}_{d_1 \times d_2} := \{O \in \mathbb{R}^{d_1 \times d_2} : O^\top O = I\}$. For two matrices $X, Y \in \mathbb{R}^{d_1 \times d_2}$, we define

Algorithm	Sample complexity	Computational complexity	Global	Exact	Over-param.
(Ma et al., 2018)	$dr^3 \log^3(d)$	$\kappa^2 \log(\frac{1}{\epsilon})$	✗	✓	✗
(Chen et al., 2020)	$dr^2 \log(d)$	$\kappa^2 \log(\frac{1}{\epsilon})$	✗	✓	✗
(Kim and Chung, 2022)	$d \log^{22}(d)$	$\log(\frac{1}{\epsilon})$	✓	*	✗
Ours (Theorem 3)	$dr^9 \log^8(d)$	$\kappa^4 \log(\frac{1}{\epsilon})$	✓	✓	✗
Ours (Theorem 2)	$dr^9 \log^2(d) \log^6(\frac{1}{\epsilon})$	$\kappa^4 \log^4(\frac{1}{\epsilon})$	✓	✓	✓

Table 1: Comparisons between different algorithms for matrix completion without explicit regularization. * The result only holds for $r' = r = 1$.

their Procrustes distance as $\text{dist}(\mathbf{X}, \mathbf{Y}) = \min_{\mathbf{O} \in \mathcal{O}_{d_2 \times d_2}} \|\mathbf{X} - \mathbf{Y}\mathbf{O}\|_F$. The projection matrix onto the column space of an orthogonal matrix $\mathbf{V} \in \mathcal{O}_{d_1 \times d_2}$ is defined as $\mathcal{P}_{\mathbf{V}} := \mathbf{V}\mathbf{V}^\top$.

We use the notation $f(n) \lesssim g(n)$ or $f(n) = O(g(n))$ when a constant $C > 0$ exists such that $f(n) \leq Cg(n)$ for sufficiently large n . Conversely, $f(n) \gtrsim g(n)$ or $f(n) = \Omega(g(n))$ implies the existence of a constant $C > 0$ such that $f(n) \geq Cg(n)$ for sufficiently large n . Moreover, we use the notations $\tilde{O}(\cdot)$ and $\tilde{\Omega}(\cdot)$ to hide logarithmic dependencies on the dimension or other parameters of the problem. Additionally, we use $f(n) \asymp g(n)$ or $f(n) = \Theta(g(n))$ when $f(n) \lesssim g(n)$ and $f(n) \gtrsim g(n)$.

2. Problem Setup and Main Results

Suppose that the singular value decomposition (SVD) of \mathbf{X}^* is given by $\mathbf{X}^* = \mathbf{V}^* \mathbf{\Sigma}^* \mathbf{V}^{*\top}$, where $\mathbf{V}^* \in \mathcal{O}_{d \times r}$ and $\mathbf{\Sigma}^*$ is an $r \times r$ diagonal matrix with diagonal elements in descending order $\sigma_1^* \geq \dots \geq \sigma_r^* > 0$. We denote the condition number of \mathbf{X}^* as $\kappa = \sigma_1^* / \sigma_r^*$. Upon defining the symmetrized operator $\mathcal{R}_\Omega = \frac{1}{2p}(\mathcal{P}_\Omega + \mathcal{P}_\Omega^\top)$, the update rule for GD can be written as

$$\mathbf{U}_{t+1} = \mathbf{U}_t - \eta \nabla f(\mathbf{U}_t) = \mathbf{U}_t - \eta \mathcal{R}_\Omega(\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}^*) \mathbf{U}_t, \quad \mathbf{U}_0 = \alpha \mathbf{Z} \quad (\text{GD})$$

where $\alpha > 0$ is the initialization scale and \mathbf{Z} is the initialization matrix satisfying $\|\mathbf{Z}\| = 1$. We assume that \mathbf{Z} satisfies the following *alignment condition*.

Condition 1 (Alignment) We say the matrix $\mathbf{Z} \in \mathbb{R}^{d \times r'}$ with $\|\mathbf{Z}\| = 1$ satisfies the alignment condition if there exists a universal constant $c_0 > 0$ such that

$$\sigma_r(\mathcal{P}_{\mathbf{V}^*} \mathbf{Z}) \geq c_0. \quad (\text{alignment condition})$$

Intuitively, this condition necessitates that the initialization matrix should have a non-negligible alignment with the column space of the ground truth. The following lemma reveals that this alignment condition is satisfied for common initialization strategies with overwhelming probability.

Lemma 1 (Sufficient condition for alignment) The following statements are satisfied:

- **Gaussian initialization.** Given $0.5d \leq r' \leq d$ and $\mathbf{Z} = \mathbf{G} / \|\mathbf{G}\|$, where \mathbf{G} is a standard Gaussian matrix, *alignment condition* holds with $c_0 = 0.1$ with probability at least $1 - \exp\{-\Omega(d)\}$.

- **Orthogonal initialization.** Given $r' = d$ and $\mathbf{Z} = \mathbf{O}$ for any $\mathbf{O} \in \mathcal{O}_{d \times d}$, [alignment condition](#) is satisfied with $c_0 = 1$.
- **Spectral initialization.** Let $\mathbf{V}\Sigma\mathbf{V}^\top$ be the eigendecomposition of the best rank- r' approximation of $\mathcal{R}_\Omega(\mathbf{X}^*)$ measured in Frobenius norm. Given $r \leq r' \leq d$ and $\mathbf{Z} = \mathbf{U} / \|\mathbf{U}\|$, where $\mathbf{U} = \mathbf{V}\Sigma^{1/2}$, [alignment condition](#) holds with $c_0 = \frac{1}{2\kappa}$ with probability at least $1 - \frac{1}{d^3}$.

Before proceeding to the main theorem, we introduce two crucial conditions on the ground truth \mathbf{X}^* and the random observation set Ω .

Condition 2 (Incoherence) The rank- r PSD matrix $\mathbf{X}^* \in \mathbb{R}^{d \times d}$ with SVD $\mathbf{X}^* = \mathbf{V}^* \Sigma^* \mathbf{V}^{*\top}$ is μ -incoherent for some $\mu \geq 1$ if $\|\mathbf{V}^*\|_{2,\infty} = \sqrt{\frac{\mu}{d}} \|\mathbf{V}^*\|_F = \sqrt{\frac{\mu r}{d}}$.

Condition 3 (Random sampling model) Each entry of \mathbf{X}^* is observed independently with probability p . In other words, $\mathbb{P}((i, j) \in \Omega) = p$ independently for all $1 \leq i, j \leq d$.

[Chen \(2015\)](#) shows that the incoherence condition is necessary for the recovery of the ground truth. Intuitively, the incoherence condition with $\mu = O(1)$ entails that none of the columns of \mathbf{V}^* have significant alignment with the standard basis vectors. Such a condition ensures that the ground truth is far from being sparse. We note that successful recovery of a sparse \mathbf{X}^* is only achievable in a near-ideal scenario where $p \approx 1$ and nearly the entirety of \mathbf{X}^* is observed. On the other hand, when $\mu = O(1)$, the recovery is possible even when p scales with $\tilde{\Omega}(\text{poly}(r)/d)$ ([Chen, 2015](#)).

With the aforementioned conditions in place, we can now present our main result, which establishes the global convergence of GD on [MC](#) for the over-parameterized setting.

Theorem 2 (Convergence of GD for over-parameterized [MC](#)) Let \mathbf{X}^* be rank- r , and Conditions 2 and 3 are satisfied with a sampling rate of $p \gtrsim \frac{\kappa^6 \mu^4 r^9 \log^6(\frac{1}{\alpha}) \log^2(d)}{d}$. Consider [MC](#) with search rank $r \leq r' \leq d$. Consider the iterates of [GD](#) with the step-size $\eta \asymp \frac{\mu r}{p d \sigma_1^*}$ and the initial point $\mathbf{U}_0 = \alpha \mathbf{Z}$, where $0 < \alpha \leq \sqrt{\frac{\sigma_1^*}{d}}$ and \mathbf{Z} satisfies Condition 1. With probability at least $1 - \frac{2}{d^3}$, after $T \lesssim \frac{1}{\eta \sigma_r^*} \log\left(\frac{1}{\alpha}\right)$ iterations, we have

$$\left\| \mathbf{U}_T \mathbf{U}_T^\top - \mathbf{X}^* \right\|_F \lesssim \sqrt{\frac{\sigma_1^* \kappa^2 \mu r^2}{p}} \alpha.$$

A few observations are in order based on the above theorem.

Computational complexity. The initialization scale α governs the final accuracy of GD. Therefore, to ensure that $\left\| \mathbf{U}_T \mathbf{U}_T^\top - \mathbf{X}^* \right\|_F \leq \epsilon$, it suffices to set the initialization scale to $\alpha \lesssim \sqrt{\frac{p}{\sigma_1^* \kappa^2 \mu r^2}} \epsilon$. Moreover, assuming that $\max\{\kappa, \mu, r\} = O(1)$ and $\epsilon \leq 1/d$, this accuracy is achieved within $\tilde{O}\left(\log^4\left(\frac{1}{\epsilon}\right)\right)$ iterations, which scales only poly-logarithmically with $1/\epsilon$.

Effect of over-parameterization. The level of over-parameterization in the search rank r' does not have any impact on either the sample complexity or the convergence of GD. As a result, our results hold even if $d^2 p \ll d r'$. In such cases, [MC](#) has many global minima, some of which may not satisfy $\mathbf{U} \mathbf{U}^\top \approx \mathbf{X}^*$. This sheds light on the implicit regularization of the vanilla GD with small initialization toward low-rank solutions when applied to [MC](#).

Sample complexity. The required sample complexity is given by $d^2 p \gtrsim dr^9 \kappa^6 \mu^4 \log^6\left(\frac{1}{\alpha}\right) \log^2(d)$, which is optimal with respect to the dimension d up to a logarithmic factor. This contrasts with the direct extension of the approach by [Kim and Chung \(2022\)](#), which necessitates a sample size on the order of $d^{1+\Theta(\kappa-1)}$. Moreover, Theorem 2 highlights that the sampling rate p must scale with $\log^6\left(\frac{1}{\epsilon}\right)$ to attain an accuracy level of ϵ , which in turn leads to a mild dependency of the final error on the sampling rate. In other words, given a fixed sampling rate p , GD achieves an accuracy in the order of $\exp(-\Omega(pd))$. We suspect that the observed dependency might be an artifact of our proof technique and could potentially be relaxed with a more detailed analysis, similar to that presented in ([Stöger and Soltanolkotabi, 2021](#)). In our next theorem, we demonstrate that in the exact-parameterization regime, it is possible to relax this mild dependency, achieving a sample complexity that does not depend on the desired accuracy level or the initialization scale.

Theorem 3 (Convergence of GD for exactly-parameterized MC) *Let \mathbf{X}^* be rank- r , and Conditions 2 and 3 are satisfied with a sampling rate of $p \gtrsim \frac{\kappa^6 \mu^4 r^9 \log^8(d)}{d}$. Consider MC with search rank $r' = r$. Consider the iterates of GD with the step-size $\eta \asymp \frac{\mu r}{\sqrt{pd}\sigma_1^*}$ and the initial point $\mathbf{U}_0 = \alpha \mathbf{Z}$, where $\alpha \asymp \frac{\sigma_r^*}{\kappa^{1.5}d}$ and \mathbf{Z} satisfies Condition 1. Given any accuracy $\epsilon > 0$, with probability at least $1 - O\left(\frac{1}{d^3}\right)$ and after $T \lesssim \frac{1}{\eta\sigma_r^*} \log\left(\frac{1}{\epsilon}\right)$ iterations, we have*

$$\left\| \mathbf{U}_T \mathbf{U}_T^\top - \mathbf{X}^* \right\|_F \leq \epsilon. \quad (1)$$

We next outline the key distinctions between the two aforementioned theorems. In the exactly-parameterized regime, neither the sampling rate p nor the initialization scale α affect the final error ϵ , as long as they meet certain thresholds. In contrast, a smaller initialization scale or a larger sampling rate improves the final error in the over-parameterized regime. Moreover, the convergence rate of GD improves from $O\left(\log^4\left(\frac{1}{\epsilon}\right)\right)$ to $O\left(\log\left(\frac{1}{\epsilon}\right)\right)$. This is because the required sampling rate is smaller in the exactly-parameterized regime, allowing the algorithm to adopt a more aggressive step-size.

3. Proof Outline

In this section, we present the key ideas underpinning our proof techniques. We begin in Section 3.1 with a dynamic signal-residual decomposition. Next, in Section 3.2, we introduce the weakly-coupled leave-one-out analysis which, together with our dynamic signal-residual decomposition, completes the proof for Theorem 2. Section 3.3 explains how these techniques can be further refined to yield improved results for the exactly-parameterized regime. Throughout this section, we occasionally omit the consideration of higher-order terms involving the step-size η . We highlight that while this omission serves to streamline the presentation, our rigorous proofs in the appendix carefully account for these higher-order terms.

3.1. Dynamic Signal-residual Decomposition

We employ a dynamic projection scheme akin to that described by [Li et al. \(2018\)](#), which decomposes the iterates \mathbf{U}_t into two distinct components: a low-rank signal part, \mathbf{S}_t , and a residual part, \mathbf{E}_t . This decomposition is represented as follows:

$$\mathbf{U}_t = \mathbf{S}_t + \mathbf{E}_t, \quad \text{where} \quad \mathbf{S}_t = \mathcal{P}_{\mathbf{V}_t} \mathbf{U}_t, \text{ and } \mathbf{E}_t = \mathcal{P}_{\mathbf{V}_t}^\perp \mathbf{U}_t.$$

Upon defining $M_t = \mathcal{R}_\Omega(\mathbf{X}^* - \mathbf{U}_t \mathbf{U}_t^\top)$, the dynamic orthonormal matrix \mathbf{V}_t is recursively defined as:

$$\mathbf{V}_{t+1} = \mathbf{Z}_{t+1} \left(\mathbf{Z}_{t+1}^\top \mathbf{Z}_{t+1} \right)^{-1/2} \quad \text{where} \quad \mathbf{Z}_{t+1} = (\mathbf{I} + \eta \mathbf{M}_t) \mathbf{V}_t \text{ and } \mathbf{V}_0 = \mathbf{V}^*.$$

Upon defining the error matrix as $\Delta_t := \mathbf{X}^* - \mathbf{U}_t \mathbf{U}_t^\top$, our goal is to show that $\|\Delta_t\|_F$ decreases efficiently to $O(\alpha)$. To show the benefit of the proposed dynamic signal-residual decomposition in achieving this goal, we start by stating the one-step dynamic of the error matrix:

$$\|\Delta_{t+1}\|_F^2 = \|\Delta_t\|_F^2 - 4\eta \langle \Delta_t, \mathbf{M}_t \mathbf{U}_t \mathbf{U}_t^\top \rangle + O(\eta^2). \quad (2)$$

Therefore, to establish the convergence of $\|\Delta_t\|_F$, it suffices to provide a reasonable lower-bound for $\langle \Delta_t, \mathbf{M}_t \mathbf{U}_t \mathbf{U}_t^\top \rangle$. This can be achieved via the following descent lemma:

Lemma 4 (Descent lemma, informal) *Suppose that $\frac{\sqrt{\sigma_r^*}}{2} \leq \sigma_r(\mathbf{S}_t)$, $\|\mathbf{S}_t\| \leq 2\sqrt{\sigma_1^*}$, and $\|\mathbf{V}_t - \mathbf{V}^*\| \leq 0.1$. Then, we have*

$$\langle \Delta_t, \mathbf{M}_t \mathbf{U}_t \mathbf{U}_t^\top \rangle \geq \frac{\sigma_r^*}{15} \|\Delta_t\|_F^2 - O\left(\sqrt{\frac{\sigma_r^{*3} \mu r^2}{p}} \|\mathbf{E}_t\| + \sqrt{r} \sigma_1^* \|(\mathcal{I} - \mathcal{R}_\Omega)(\Delta_t)\| \right) \|\Delta_t\|_F.$$

By combining the descent lemma with Equation (2), we arrive at the following expression:

$$\|\Delta_{t+1}\|_F \leq \left(1 - \frac{\eta \sigma_r^*}{10} \right) \|\Delta_t\|_F + O(\eta) \left(\sqrt{\frac{\sigma_r^{*3} \mu r^2}{p}} \|\mathbf{E}_t\| + \sqrt{r} \sigma_1^* \|(\mathcal{I} - \mathcal{R}_\Omega)(\Delta_t)\| \right) + O(\eta^2). \quad (3)$$

The above inequality holds once the conditions $\frac{\sqrt{\sigma_r^*}}{2} \leq \sigma_r(\mathbf{S}_t) \leq \|\mathbf{S}_t\| \leq 2\sqrt{\sigma_1^*}$ and $\|\mathbf{V}_t - \mathbf{V}^*\| \leq 0.1$ are met. These conditions entail that during the initial phase of the algorithm, $\sigma_r(\mathbf{S}_t)$ must undergo a fast growth, whereas \mathbf{V}_t must remain close to \mathbf{V}^* . Under these conditions, GD enters a fast linear convergence phase, provided that $\|(\mathcal{I} - \mathcal{R}_\Omega)(\Delta_t)\| \ll \|\Delta_t\|_F$. In fact, we can readily establish that $\|(\mathcal{I} - \mathcal{R}_\Omega)(\Delta_t)\| \leq c \|\Delta_t\|$ for some $c > 0$. However, the challenge lies in ensuring that c remains sufficiently small so as not to negate the effect of a constant factor improvement $1 - 0.1\eta\sigma_r^*$ in Equation (3). This phase continues until GD reaches an error level controlled by $\|\mathbf{E}_t\|$. Therefore, to prove the convergence of GD, we need to establish the following properties:

- **Fast growth of \mathbf{S}_t :** Recall that $\max\{\|\mathbf{S}_0\|, \|\mathbf{E}_0\|\} = O(\alpha)$. We need to ensure efficient growth of $\sigma_r(\mathbf{S}_t)$ from $O(\alpha)$ to $\frac{\sqrt{\sigma_r^*}}{2}$, while keeping $\|\mathbf{S}_t\|$ below $2\sqrt{\sigma_1^*}$.
- **Slow growth of \mathbf{E}_t :** We need to show that while the signal term grows rapidly, the residual term \mathbf{E}_t grows at a much slower rate. Specifically, we will demonstrate that $T = O\left(\frac{1}{\eta\sigma_r^*} \log(\frac{1}{\alpha})\right)$ suffices to ensure $\frac{\sqrt{\sigma_r^*}}{2} \leq \sigma_r(\mathbf{S}_t) \leq \|\mathbf{S}_t\| \leq 2\sqrt{\sigma_1^*}$ while keeping $\|\mathbf{E}_t\| = O(\alpha)$.
- **Small values of $\|(\mathcal{I} - \mathcal{R}_\Omega)(\Delta_t)\|$ and $\|\mathbf{V}_t - \mathbf{V}^*\|$:** Equally important is maintaining control over $\|(\mathcal{I} - \mathcal{R}_\Omega)(\Delta_t)\|$ and $\|\mathbf{V}_t - \mathbf{V}^*\|$. While $\|\mathbf{V}_t - \mathbf{V}^*\| \leq 0.1$ is needed as a crucial condition for Equation (3), the value of $\|(\mathcal{I} - \mathcal{R}_\Omega)(\Delta_t)\|$ directly controls the convergence rate of GD.

To establish the above properties, we provide the one-step dynamics of \mathbf{S}_t , \mathbf{E}_t , and $(\mathcal{I} - \mathcal{R}_\Omega)(\Delta_t)$.

Lemma 5 (One-step dynamics, informal) *Under the conditions of Theorem 2, the following inequalities hold with an overwhelming probability:*

$$\|\mathbf{S}_{t+1}\| \leq \left(1 + \eta \left(\sigma_1^* - \|\mathbf{S}_t\|^2 + O(\sigma_1^* \|\mathbf{V}_t - \mathbf{V}^*\| + \|(\mathcal{I} - \mathcal{R}_\Omega)(\Delta_t)\|)\right)\right) \|\mathbf{S}_t\|, \quad (4)$$

$$\begin{aligned} \sigma_r(\mathbf{S}_{t+1}) &\geq \left(1 + \eta \left(\sigma_r^* - \sigma_r^2(\mathbf{S}_t) - O(\sigma_1^* \|\mathbf{V}_t - \mathbf{V}^*\| + \|(\mathcal{I} - \mathcal{R}_\Omega)(\Delta_t)\|)\right)\right) \sigma_r(\mathbf{S}_t) \\ &\quad + O(\eta) \cdot (\sigma_1^* \|\mathbf{V}_t - \mathbf{V}^*\| + \|(\mathcal{I} - \mathcal{R}_\Omega)(\Delta_t)\|) \|\mathbf{E}_t\|, \end{aligned} \quad (5)$$

$$\|\mathbf{E}_{t+1}\| \leq (1 + O(\eta) \cdot (\sigma_1^* \|\mathbf{V}_t - \mathbf{V}^*\| + \|(\mathcal{I} - \mathcal{R}_\Omega)(\Delta_t)\|)) \|\mathbf{E}_t\|, \quad (6)$$

$$\|(\mathcal{I} - \mathcal{R}_\Omega)(\Delta_t)\| \leq O\left(\sqrt{\frac{d}{p}} \|\Delta_t\| \left(\|\mathbf{V}^*\|_{2,\infty}^2 + \|\mathbf{V}_t\|_{2,\infty}^2\right) + \sqrt{\frac{\sigma_1^* \mu r}{p}} \|\mathbf{E}_t\|\right). \quad (7)$$

Next, we provide a high-level overview of how Lemma 5 can be used to establish the aforementioned properties. To this goal, we only focus on the initial phase of the algorithm, where both \mathbf{S}_t and \mathbf{E}_t are small. A more formal analysis for the entire trajectory is provided in Appendix C.

To use Lemma 5, it suffices to control two key quantities: $\|\mathbf{V}_t\|_{2,\infty}$ and $\|\mathbf{V}_t - \mathbf{V}^*\|$. To illustrate this, let us assume that \mathbf{V}_t inherits the incoherence of \mathbf{V}^* , that is, $\|\mathbf{V}_t\|_{2,\infty} \leq O(\sqrt{\mu r/d}) \ll 1$ for all $1 \leq t \leq T$. Then, Equation (7) suggests that $\|(\mathcal{I} - \mathcal{R}_\Omega)(\Delta_t)\| \leq c \|\Delta_t\| + O(\alpha)$, where $c = O(\mu r/\sqrt{pd}) \ll 1$, thereby ensuring the necessary control over $\|(\mathcal{I} - \mathcal{R}_\Omega)(\Delta_t)\|$.

On the other hand, the small values of $\|(\mathcal{I} - \mathcal{R}_\Omega)(\Delta_t)\|$ and $\|\mathbf{V}_t - \mathbf{V}^*\|$ play crucial roles in controlling the behavior of the signal and residual terms. To illustrate this, let us assume that at a certain point, $\|\mathbf{S}_t\| \geq 1.5\sqrt{\sigma_1^*}$. Given that we have considered $\|(\mathcal{I} - \mathcal{R}_\Omega)(\Delta_t)\|$ and $\|\mathbf{V}_t - \mathbf{V}^*\|$ to be small, Equation (4) simplifies to $\|\mathbf{S}_{t+1}\| \leq (1 - \eta\Omega(\sigma_1^*)) \|\mathbf{S}_t\| + O(\eta\sigma_1^* \|\mathbf{E}_t\|)$, effectively preventing $\|\mathbf{S}_{t+1}\|$ from further growth. With a similar reasoning, Equation (5) simplifies to $\sigma_r(\mathbf{S}_{t+1}) \geq (1 + \eta\Omega(\sigma_r^*))\sigma_r(\mathbf{S}_t)$. Here, we have leveraged the assumption that $\sigma_r^2(\mathbf{S}_t) \ll \sigma_r^*$ during the initial phase. This implies that $\sigma_r(\mathbf{S}_t)$ grows at a rate of $1 + \eta\Omega(\sigma_r^*)$. In contrast, Equation (6) implies that $\|\mathbf{E}_t\|$ grows at a rate of $1 + \eta O(\|(\mathcal{I} - \mathcal{R}_\Omega)(\Delta_t)\| + \sigma_1^* \|\mathbf{V}_t - \mathbf{V}^*\|)$, which is significantly slower than the growth rate of $\sigma_r(\mathbf{S}_t)$ because $\max\{\|(\mathcal{I} - \mathcal{R}_\Omega)(\Delta_t)\|, \sigma_1^* \|\mathbf{V}_t - \mathbf{V}^*\|\} \ll \sigma_r^*$. It is due to this discrepancy in the growth rates of $\|\mathbf{S}_t\|$ and $\|\mathbf{E}_t\|$ that GD enters the local linear convergence rate and achieves a final error of $O(\alpha)$.

3.2. Refined Leave-one-out Analysis with Weak Coupling

Indeed, it is not immediately evident why both $\|\mathbf{V}_t\|_{2,\infty}$ and $\|\mathbf{V}_t - \mathbf{V}^*\|$ would remain small. In fact, our initial intuition might suggest the opposite: recall that GD takes an aggressively large step size. Consequently, even a single GD update has the potential to disrupt the incoherence of \mathbf{V}_t . Our key contribution is to establish that such disruption does not occur, even when the iterates are arbitrarily far from the ground truth. In essence, we show that despite the gradient update $\|\eta\mathcal{R}_\Omega(\Delta_t)\mathbf{U}_t\|$ potentially having a magnitude of $\tilde{\Omega}(1)$, its impact on \mathbf{V}_t is distributed fairly evenly across its elements. As a result, it has minimal influence on $\|\mathbf{V}_t\|_{2,\infty}$ and $\|\mathbf{V}_t - \mathbf{V}^*\|$.

We start by showing that a small $\|\mathbf{V}_t\|_{2,\infty}$ implies a small $\|\mathbf{V}_t - \mathbf{V}^*\|$.

Lemma 6 (Small $\|\mathbf{V}_t\|_{2,\infty}$ implies small $\|\mathbf{V}_t - \mathbf{V}^*\|_F$, informal) Suppose that $\|\mathbf{S}_t\| \leq 2\sqrt{\sigma_1^*}$ and $\|\mathbf{V}_t - \mathbf{V}^*\| \leq \frac{1}{2\kappa}$. With an overwhelming probability, we have

$$\|\mathbf{V}_{t+1} - \mathbf{V}^*\|_F \leq \|\mathbf{V}_t - \mathbf{V}^*\|_F + O(\eta) \cdot \sigma_1^* \sqrt{\frac{dr}{p}} \left(\|\mathbf{V}^*\|_{2,\infty}^2 + \|\mathbf{V}_t\|_{2,\infty}^2 \right) + O(\eta^2).$$

Recall that, due to the incoherence of \mathbf{V}^* , we have $\|\mathbf{V}^*\|_{2,\infty} = \sqrt{\mu r/d}$. Now, suppose we can further establish that \mathbf{V}_t enjoys a similar incoherence property. In such a case, the aforementioned lemma leads to $\|\mathbf{V}_t - \mathbf{V}^*\| \leq \|\mathbf{V}_t - \mathbf{V}^*\|_F \leq O\left(T \cdot \eta \sigma_1^* \sqrt{\frac{\mu^2 r^3}{pd}}\right) + O(T \cdot \eta^2)$ for every $1 \leq t \leq T$. Given the provided bounds on T and p , this automatically establishes that $\|\mathbf{V}_t - \mathbf{V}^*\|$ remains small throughout the iterations. Therefore, it suffices to control the incoherence of $\|\mathbf{V}_t\|_{2,\infty}$.

Controlling $\|\mathbf{V}_t\|_{2,\infty}$, which necessitates estimating the ℓ_2 -norm of each row, requires a more fine-grained analysis than what is needed for the Frobenius norm. The primary challenge lies in the intricate correlations between the orthogonal matrix \mathbf{V}_t and the random observation set Ω , which preclude the straightforward application of classical concentration inequalities. To effectively decouple these correlations, we propose a technique called *weakly-coupled leave-one-out analysis*. Before introducing our proposed methodology, it is essential to grasp the core principles of the classical leave-one-out analysis.

Local leave-one-out analysis. When the search rank is exactly parameterized ($r = r'$) and the initial point is sufficiently close to the ground truth $\mathbf{U}_0 \mathbf{U}_0^\top \approx \mathbf{X}^*$, [Ma et al. \(2018\)](#) established the incoherence of the iterates via the following leave-one-out sequences $\{\mathbf{U}_t^{(l)}\}_{t=0}^T$ for each $1 \leq l \leq d$:

$$\text{dist}\left(\mathbf{U}_0, \mathbf{U}_0^{(l)}\right) \approx 0, \quad \text{and} \quad \mathbf{U}_{t+1}^{(l)} = \left(\mathbf{I} - \eta \mathcal{R}_{\Omega^{(l)}}\left(\mathbf{U}_t^{(l)} \mathbf{U}_t^{(l)\top} - \mathbf{X}^*\right)\right) \mathbf{U}_t^{(l)}, \quad (8)$$

where $\mathcal{R}_{\Omega^{(l)}}$ is the leave-one-out projection operator defined by

$$\mathcal{R}_{\Omega^{(l)}} = \frac{1}{2p} \left(\mathcal{P}_{\Omega^{(l)}} + \mathcal{P}_{\Omega^{(l)}}^\top \right), \quad \text{and} \quad [\mathcal{P}_{\Omega^{(l)}}(\mathbf{X})]_{i,j} = \begin{cases} pX_{i,j} & \text{if } i = l \text{ or } j = l, \\ X_{i,j} & \text{if } (i, j) \in \Omega, i \neq l, \text{ and } j \neq l, \\ 0 & \text{otherwise.} \end{cases}$$

The sole distinction between the projection operators $\mathcal{R}_{\Omega^{(l)}}$ and \mathcal{R}_Ω is in their l -th row and l -th column: in contrast to \mathcal{R}_Ω , the l -th row and l -th column of $\mathcal{R}_{\Omega^{(l)}}(\mathbf{X})$ are *deterministically* set to match the corresponding values of \mathbf{X} . This seemingly minor adjustment yields two important consequences: first, it ensures that $\mathbf{U}_t \approx \mathbf{U}_t^{(l)}$, and second, it guarantees that the behavior of $\mathbf{U}_t^{(l)}$ remains independent of the random measurements in the l -th row and l -th column. This decoupling technique is the key to controlling the deviation of $\|\mathbf{U}_t\|_{2,\infty}$. To formalize this intuition, let us define $\mathbf{U}^* = \mathbf{V}^* \Sigma^{*1/2}$. One can write

$$\begin{aligned} \|\mathbf{U}_t\|_{2,\infty} &= \max_{1 \leq l \leq d} \left\{ \left\| [\mathbf{U}_t \mathbf{H}_t^{(l)}]_{l,\cdot} \right\| \right\} \\ &\leq \max_{1 \leq l \leq d} \left\{ \left\| [\mathbf{U}^* - \mathbf{U}_t^{(l)} \mathbf{R}_t^{(l)}]_{l,\cdot} \right\| + \left\| [\mathbf{U}_t \mathbf{H}_t^{(l)} - \mathbf{U}_t^{(l)} \mathbf{R}_t^{(l)}]_{l,\cdot} \right\| + \left\| [\mathbf{U}^*]_{l,\cdot} \right\| \right\} \\ &\leq \underbrace{\max_{1 \leq l \leq d} \left\{ \left\| (\mathbf{U}^* - \mathbf{U}_t^{(l)} \mathbf{R}_t^{(l)})_{l,\cdot} \right\| \right\}}_{\text{leave-one-out error}} + \underbrace{\max_{1 \leq l \leq d} \left\{ \left\| \mathbf{U}_t \mathbf{H}_t^{(l)} - \mathbf{U}_t^{(l)} \mathbf{R}_t^{(l)} \right\|_F \right\}}_{\text{proximal error}} + \sqrt{\frac{\sigma_1^* \mu r}{d}}. \end{aligned} \quad (9)$$

Here $\mathbf{R}_t^{(l)}$ and $\mathbf{H}_t^{(l)}$ are orthogonal matrices defined as $\mathbf{R}_t^{(l)} = \arg \min_{\mathbf{O} \in \mathcal{O}_{r \times r}} \|\mathbf{U}_t^{(l)} \mathbf{O} - \mathbf{U}^*\|_F$ and $\mathbf{H}_t^{(l)} = \arg \min_{\mathbf{O} \in \mathcal{O}_{r \times r}} \|\mathbf{U}_t^{(l)} \mathbf{O} - \mathbf{U}_t^{(l)} \mathbf{R}_t^{(l)}\|_F$. Although it may not be immediately obvious, it can be shown that the l -th row of the matrix $\mathbf{U}^* - \mathbf{U}_t^{(l)} \mathbf{R}_t^{(l)}$ is purely *deterministic*. Consequently, it becomes possible to effectively control the leave-one-out error. To tackle the proximal error, recall that the initial point $\mathbf{U}_0 \mathbf{U}_0^\top$ is assumed to be close to \mathbf{X}^* . Within this region, the local landscape exhibits restricted strong convexity. This ensures that the true iterates \mathbf{U}_t and the leave-one-out versions $\mathbf{U}_t^{(l)}$ become increasingly close, leading to a small proximal error. By combining these two arguments, we can guarantee the incoherence of the true iterates. Furthermore, the incoherence of \mathbf{U}_t automatically implies the incoherence of \mathbf{V}_t , given that $\sqrt{\sigma_r^*} \|\mathbf{V}_t\|_{2,\infty} \lesssim \|\mathbf{U}_t\|_{2,\infty}$ when $\mathbf{U}_t \mathbf{U}_t^\top \approx \mathbf{X}^*$. For more details, we refer interested readers to the discussions in (Ma et al., 2018).

While the classical leave-one-out analysis provides precise *local* guarantees within the *exactly-parameterized* regimes, we shed light on its limitations when applied *globally* in the *over-parameterized* settings. A significant challenge arises from the discrepancy of the singular values of \mathbf{U}_t and $\mathbf{U}_t^{(l)}$: although they may remain close to the singular values of \mathbf{U}^* in the local regime, they can undergo substantial changes when positioned far from the true solution. Consequently, the original measure of proximal error based on $\text{dist}(\mathbf{U}_t, \mathbf{U}_t^{(l)})$ loses its effectiveness as a reliable metric.

Instead, recall that we only require controlling \mathbf{V}_t , which unlike \mathbf{U}_t , has *unit* singular values. This motivates us to switch to a more stable metric—the divergence between the left column spaces of \mathbf{U}_t and $\mathbf{U}_t^{(l)}$. However, an additional complication is that these left column spaces may also not align perfectly due to over-parameterization. Fortunately, by resorting to our proposed dynamic signal-residual decomposition, we can show that the iterates \mathbf{U}_t are well-approximated by the low-rank signal $\mathbf{U}_t \approx \mathbf{S}_t$. Therefore, it suffices to focus on controlling the discrepancy in the column spaces of \mathbf{S}_t and $\mathbf{S}_t^{(l)}$, i.e., $\text{dist}(\mathbf{V}_t, \mathbf{V}_t^{(l)})$. However, the new proximal error $\text{dist}(\mathbf{V}_t, \mathbf{V}_t^{(l)})$ can still grow exponentially. To explain the root cause of this exponential growth, we employ matrix Taylor expansion to derive the first-order approximations for \mathbf{V}_{t+1} and $\mathbf{V}_{t+1}^{(l)}$:

$$\mathbf{V}_{t+1} = \mathbf{V}_t + \eta \mathcal{P}_{\mathbf{V}_t}^\perp \mathbf{M}_t \mathbf{V}_t + O(\eta^2) \quad \text{and} \quad \mathbf{V}_{t+1}^{(l)} = \mathbf{V}_t^{(l)} + \eta \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \mathbf{M}_t^{(l)} \mathbf{V}_t^{(l)} + O(\eta^2), \quad (10)$$

where we define $\mathbf{M}_t^{(l)} = \mathcal{R}_{\Omega^{(l)}}(\mathbf{X}^* - \mathbf{U}_t^{(l)} \mathbf{U}_t^{(l)\top})$. To effectively control the proximal error, it is crucial to establish an upper bound for $\|\mathbf{M}_t - \mathbf{M}_t^{(l)}\|$. This distance tends to concentrate around $\|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_t^{(l)} \mathbf{U}_t^{(l)\top}\| \approx \|\mathbf{S}_t \mathbf{S}_t^\top - \mathbf{S}_t^{(l)} \mathbf{S}_t^{(l)\top}\|$ when the sampling rate p is sufficiently large. However, as previously noted, the singular values of \mathbf{S}_t and $\mathbf{S}_t^{(l)}$ may diverge. This misalignment can lead to $\|\mathbf{M}_t - \mathbf{M}_t^{(l)}\| = \Omega(\sigma_1^*)$ in the worst case. Hence, the proximal error can grow exponentially.

Weakly-coupled leave-one-out analysis. To remedy the alignment challenges identified earlier, we propose the following refined leave-one-out sequences $\{\tilde{\mathbf{V}}_t^{(l)}\}_{t=0}^T$:

$$\tilde{\mathbf{V}}_{t+1}^{(l)} = \tilde{\mathbf{Z}}_{t+1}^{(l)} \left(\tilde{\mathbf{Z}}_{t+1}^{(l)\top} \tilde{\mathbf{Z}}_{t+1}^{(l)} \right)^{-1/2} \quad \text{where} \quad \tilde{\mathbf{Z}}_{t+1}^{(l)} = \left(\mathbf{I} + \eta \tilde{\mathbf{M}}_t^{(l)} \right) \tilde{\mathbf{V}}_t^{(l)} \quad \text{and} \quad \tilde{\mathbf{V}}_0^{(l)} = \mathbf{V}^*. \quad (11)$$

In this context, $\tilde{\mathbf{M}}_t^{(l)}$ is defined as:

$$\tilde{\mathbf{M}}_t^{(l)} = \mathcal{R}_{\Omega^{(l)}} \left(\mathbf{X}^* - \tilde{\mathbf{V}}_t^{(l)} \Sigma_t \tilde{\mathbf{V}}_t^{(l)\top} \right) \quad \text{where} \quad \Sigma_t = \mathbf{V}_t^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{V}_t \in \mathbb{R}^{r \times r}. \quad (12)$$

Compared to the original $\mathbf{M}_t^{(l)}$, we replace $\Sigma_t^{(l)} = \mathbf{V}_t^{(l)\top} \mathbf{U}_t^{(l)} \mathbf{U}_t^{(l)\top} \mathbf{V}_t^{(l)}$ by $\Sigma_t = \mathbf{V}_t^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{V}_t$ in the definition of $\widetilde{\mathbf{M}}_t^{(l)}$. Our analysis indicates that this adjustment significantly improves the control over the distance $\|\mathbf{M}_t - \widetilde{\mathbf{M}}_t^{(l)}\| = o(1)$ when the sampling rate p is sufficiently large. Hence, the proximal error grows at a much slower rate.

Despite their promise, these refined leave-one-out sequences do introduce a trade-off: the statistical independence inherent in the original leave-one-out sequences is compromised due to the inclusion of Σ_t . In other words, the l -th leave-one-out sequence $\{\widetilde{\mathbf{V}}_t^{(l)}\}_{t=0}^T$ is no longer independent of the random measurements in the l -th row and l -th column. Nonetheless, we demonstrate that the resulting correlation is relatively weak, primarily because Σ_t is a comparatively small $r \times r$ matrix. To control this statistical coupling, we employ a novel *adaptive covering argument*, which can be of independent interest. This approach effectively mitigates the statistical coupling while incurring a mild increase in the required sample complexity, which remains only polynomial in r .

To formalize our arguments, we can decompose the refined leave-one-out sequences $\|\mathbf{V}_t\|_{2,\infty}$ as:

$$\begin{aligned} \|\mathbf{V}_t\|_{2,\infty} &\leq \|\mathbf{V}^* - \mathbf{V}_t\|_{2,\infty} + \|\mathbf{V}^*\|_{2,\infty} \\ &= \max_{1 \leq l \leq d} \left\{ \left\| (\mathbf{V}^* - \mathbf{V}_t)_{l,\cdot} \right\| \right\} + \|\mathbf{V}^*\|_{2,\infty} \\ &\leq \underbrace{\max_{1 \leq l \leq d} \left\{ \left\| (\mathbf{V}^* - \widetilde{\mathbf{V}}_t^{(l)})_{l,\cdot} \right\| \right\}}_{\text{refined leave-one-out error (Proposition 7)}} + \underbrace{\max_{1 \leq l \leq d} \left\{ \left\| \mathbf{V}_t - \widetilde{\mathbf{V}}_t^{(l)} \right\|_F \right\}}_{\text{refined proximal error (Proposition 8)}} + \sqrt{\frac{\mu r}{d}}. \end{aligned} \quad (13)$$

Next, we characterize the dynamic of the refined leave-one-out error.

Proposition 7 (Refined leave-one-out error) *Suppose that $p \gtrsim \frac{\log(d)}{d}$ and $\|\mathbf{V}^* - \mathbf{V}_t\| \leq \frac{1}{2\kappa}$. With probability at least $1 - \frac{1}{d^3}$, for any $1 \leq t \leq T \lesssim \frac{1}{\eta\sigma_r^*} \log\left(\frac{1}{\alpha}\right)$ and $1 \leq l \leq d$, we have*

$$\left\| (\mathbf{V}^* - \widetilde{\mathbf{V}}_{t+1}^{(l)})_{l,\cdot} \right\| \leq (1 - 0.5\eta\sigma_r^*) \left\| (\mathbf{V}^* - \widetilde{\mathbf{V}}_t^{(l)})_{l,\cdot} \right\| + O(\eta) \cdot \sigma_1^* \frac{\kappa\mu^{1.5}r^2 \log\left(\frac{1}{\alpha}\right)}{\sqrt{pd^2}}.$$

A simple inductive argument based on Proposition 7 reveals that the following inequality holds with an overwhelming probability for all $1 \leq t \leq T$:

$$\left\| (\mathbf{V}^* - \widetilde{\mathbf{V}}_t^{(l)})_{l,\cdot} \right\| \lesssim \frac{\kappa^2\mu^{1.5}r^2 \log\left(\frac{1}{\alpha}\right)}{\sqrt{pd^2}} \leq \sqrt{\frac{\mu r}{4d}}, \text{ assuming } p \gtrsim \frac{\kappa^4\mu^2r^3 \log^2\left(\frac{1}{\alpha}\right) \log(d)}{d}.$$

Next, we characterize the dynamic of the refined proximal error.

Proposition 8 (Refined proximal error) *Suppose that $p \gtrsim \frac{\kappa^6\mu^4r^9 \log^6\left(\frac{1}{\alpha}\right) \log^2(d)}{d}$, $\|\mathbf{V}^* - \mathbf{V}_t\| \leq \frac{1}{2\kappa}$, and $\|\mathbf{V}_t - \widetilde{\mathbf{V}}_t^{(l)}\|_F \leq \sqrt{\frac{\mu r}{4d}}$. With probability at least $1 - \frac{1}{d^3}$, for any $1 \leq t \leq T \lesssim \frac{1}{\eta\sigma_r^*} \log\left(\frac{1}{\alpha}\right)$ and $1 \leq l \leq d$, we have*

$$\left\| \mathbf{V}_{t+1} - \widetilde{\mathbf{V}}_{t+1}^{(l)} \right\|_F \leq \left\| \mathbf{V}_t - \widetilde{\mathbf{V}}_t^{(l)} \right\|_F + O(\eta) \cdot \sigma_1^* \sqrt{\frac{\kappa\mu^3r^{5.5} \log\left(\frac{1}{\alpha}\right) \log(d)}{\sqrt{pd} \cdot d}}.$$

The above proposition implies that

$$\|\mathbf{V}_t - \tilde{\mathbf{V}}_t^{(l)}\|_{\text{F}} \lesssim \eta \sigma_1^* \sqrt{\frac{\kappa \mu^3 r^{5.5} \log(\frac{1}{\alpha}) \log(d)}{\sqrt{pd} \cdot d}} \cdot T \lesssim \sqrt{\frac{\kappa^3 \mu^3 r^{5.5} \log^3(\frac{1}{\alpha}) \log(d)}{\sqrt{pd} \cdot d}} \leq \sqrt{\frac{\mu r}{4d}}.$$

Combining the above inequalities with the proposed decomposition in Equation (13) leads to:

$$\|\mathbf{V}_t\|_{2,\infty} \leq \sqrt{\frac{\mu r}{4d}} + \sqrt{\frac{\mu r}{4d}} + \sqrt{\frac{\mu r}{d}} \leq \sqrt{\frac{4\mu r}{d}}, \quad \text{with probability at least } 1 - \frac{2}{d^3}.$$

This establishes the incoherence of \mathbf{V}_t for all $1 \leq t \leq T$.

3.3. Improved Results for Exact Parameterization

Finally, we show that our analysis in the over-parameterized regime, combined with the following local convergence result for the exact parameterization regime by Ma et al. (2018), readily establishes the proof of Theorem 3.

Theorem 9 (Local convergence of GD (Ma et al., 2018, Theorem 2)) *Consider MC with search rank $r' = r$. Suppose that the sampling rate satisfies $p \gtrsim \frac{\mu^3 r^3 \log^3(d)}{d}$. Consider the iterates of GD with step-size $\eta \leq \frac{2}{25\kappa\sigma_1^*}$. Suppose that there exists $t_0 \geq 0$ such that \mathbf{U}_{t_0} and the leave-one-out sequences $\{\mathbf{U}_{t_0}^{(l)}\}_{l=0}^d$ defined in Equation (8) satisfy:*

$$\text{dist}(\mathbf{U}_{t_0}, \mathbf{U}^*) \leq O\left(\sqrt{\frac{\sigma_r^* \mu^3 r^3 \log(d)}{pd^2}}\right), \quad (14)$$

$$\max \left\{ \text{dist}(\mathbf{U}_{t_0}, \mathbf{U}_{t_0}^{(l)}), \text{dist}(\mathbf{U}_{t_0}^{(l)}, \mathbf{U}^*) \right\} \leq O\left(\sqrt{\frac{\sigma_r^* \mu^3 r^3 \log(d)}{pd^2}}\right), \quad \text{for all } 1 \leq l \leq d. \quad (15)$$

With probability at least $1 - O(\frac{1}{d^3})$, for all $t_0 \leq t \leq t_0 + O(d^5)$, we have

$$\|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}^*\|_{\text{F}} \leq \|\mathbf{U}_{t_0} \mathbf{U}_{t_0}^\top - \mathbf{X}^*\|_{\text{F}} (1 - 0.2\eta\sigma_r^*)^{t-t_0}.$$

To prove Theorem 3, it suffices to show that the conditions of the above theorem are met at a certain iteration $0 \leq t_0 \leq T$. This can be achieved by leveraging our result for the over-parameterized regime. In particular, upon choosing $\alpha = c \cdot \frac{\sigma_r^*}{\kappa^{1.5}d}$ for sufficiently small $c > 0$ in Theorem 2, one can show that both Conditions (14) and (15) are satisfied with an overwhelming probability after $t_0 = \tilde{O}(\frac{1}{\eta\sigma_r^*})$ iterations. From this iteration onward, Theorem 9 shows that the iterations of GD enter a local linear convergence regime, which readily establishes the final result of Theorem 3.

4. Conclusion and Future Directions

In this paper, we prove the convergence of vanilla gradient descent (GD) with small initialization for symmetric matrix completion. Existing convergence results for this problem typically require

explicit regularization or precise initializations. However, our work proves that neither condition is necessary for GD to converge. Moreover, our results also apply to the over-parameterized regime, where the rank of the true solution is unknown and over-estimated instead.

Although our required sample complexity $\tilde{O}(dr^9)$ is optimal with respect to the dimension d , it remains sub-optimal with respect to the rank r . Specifically, it exceeds the sample complexity of regularized GD, which stands at $\tilde{O}(dr^2)$ (Chen and Wainwright, 2015). We expect our analysis can be sharpened to achieve a similar sample complexity.

We anticipate that our findings will pave the way for broader results extending beyond symmetric matrix completion. In particular, our proposed weakly-coupled leave-one-out analysis relaxes several stringent conditions of classical leave-one-out analysis, making it highly applicable for the global analysis of GD. We believe that this approach, along with potential variations, holds promise for explaining the favorable performance of GD or its variants in various statistical learning problems.

Acknowledgment

We thank Hongyang R. Zhang for pointing him to the updated version of (Li et al., 2018). We thank Richard Y. Zhang and Cédric Jozs for their helpful discussions. We are grateful to the anonymous reviewers for pointing out the limitations of our sample complexity result in the overparameterized regime. This work is supported, in part, by NSF CAREER Award CCF-2337776, NSF Award DMS-2152776, and ONR Award N00014-22-1-2127.

References

- Emmanuel Abbe, Jianqing Fan, Kaizheng Wang, and Yiqiao Zhong. Entrywise eigenvector analysis of random matrices with low expected rank. *Annals of statistics*, 48(3):1452, 2020.
- Sanjeev Arora, Nadav Cohen, Noah Golowich, and Wei Hu. A convergence analysis of gradient descent for deep linear neural networks. *arXiv preprint arXiv:1810.02281*, 2018.
- Emmanuel Candes and Benjamin Recht. Exact matrix completion via convex optimization. *Communications of the ACM*, 55(6):111–119, 2012.
- Emmanuel J Candès and Terence Tao. The power of convex relaxation: Near-optimal matrix completion. *IEEE Transactions on Information Theory*, 56(5):2053–2080, 2010.
- Joshua Cape, Minh Tang, and Carey E Priebe. The two-to-infinity norm and singular subspace geometry with applications to high-dimensional statistics. 2019.
- Ji Chen and Xiaodong Li. Memory-efficient kernel pca via partial matrix sampling and nonconvex optimization: a model-free analysis of local minima. *arXiv preprint arXiv:1711.01742*, 2017.
- Ji Chen and Xiaodong Li. Model-free nonconvex matrix completion: Local minima analysis and applications in memory-efficient kernel pca. *J. Mach. Learn. Res.*, 20(142):1–39, 2019.
- Ji Chen, Dekai Liu, and Xiaodong Li. Nonconvex rectangular matrix completion via gradient descent without $\ell_{2,\infty}$ regularization. *IEEE Transactions on Information Theory*, 66(9):5806–5841, 2020.
- Yudong Chen. Incoherence-optimal matrix completion. *IEEE Transactions on Information Theory*, 61(5):2909–2923, 2015.

- Yudong Chen and Martin J Wainwright. Fast low-rank estimation by projected gradient descent: General statistical and algorithmic guarantees. *arXiv preprint arXiv:1509.03025*, 2015.
- Yuxin Chen, Yuejie Chi, Jianqing Fan, and Cong Ma. Gradient descent with random initialization: Fast global convergence for nonconvex phase retrieval. *Mathematical Programming*, 176:5–37, 2019.
- Yuejie Chi, Yue M Lu, and Yuxin Chen. Nonconvex optimization meets low-rank matrix factorization: An overview. *IEEE Transactions on Signal Processing*, 67(20):5239–5269, 2019.
- Lijun Ding and Yudong Chen. Leave-one-out approach for matrix completion: Primal and dual analysis. *IEEE Transactions on Information Theory*, 66(11):7274–7301, 2020.
- Noureddine El Karoui, Derek Bean, Peter J Bickel, Chinghway Lim, and Bin Yu. On robust regression with high-dimensional predictors. *Proceedings of the National Academy of Sciences*, 110(36):14557–14562, 2013.
- Salar Fattahi and Somayeh Sojoudi. Exact guarantees on the absence of spurious local minima for non-negative rank-1 robust principal component analysis. *Journal of machine learning research*, 2020.
- Rong Ge, Jason D Lee, and Tengyu Ma. Matrix completion has no spurious local minimum. *Advances in neural information processing systems*, 29, 2016.
- Rong Ge, Chi Jin, and Yi Zheng. No spurious local minima in nonconvex low rank problems: A unified geometric analysis. In *International Conference on Machine Learning*, pages 1233–1242. PMLR, 2017.
- Rong Ge, Yunwei Ren, Xiang Wang, and Mo Zhou. Understanding deflation process in over-parametrized tensor decomposition. *Advances in Neural Information Processing Systems*, 34:1299–1311, 2021.
- David F Gleich and Lek-heng Lim. Rank aggregation via nuclear norm minimization. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 60–68, 2011.
- Thore Graepel. Kernel matrix completion by semidefinite programming. In *International Conference on Artificial Neural Networks*, pages 694–699. Springer, 2002.
- David Gross. Recovering low-rank matrices from few coefficients in any basis. *IEEE Transactions on Information Theory*, 57(3):1548–1566, 2011.
- Suriya Gunasekar, Blake Woodworth, Srinadh Bhojanapalli, Behnam Neyshabur, and Nathan Srebro. Implicit regularization in matrix factorization. In *2018 Information Theory and Applications Workshop (ITA)*, pages 1–10. IEEE, 2018.
- Yue Hu, Xiaohan Liu, and Mathews Jacob. A generalized structured low-rank matrix completion algorithm for mr image recovery. *IEEE transactions on medical imaging*, 38(8):1841–1851, 2018.

- Prateek Jain, Praneeth Netrapalli, and Sujay Sanghavi. Low-rank matrix completion using alternating minimization. In *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*, pages 665–674, 2013.
- Daesung Kim and Hye Won Chung. Rank-1 matrix completion with gradient descent and small random initialization. *arXiv preprint arXiv:2212.09396*, 2022.
- Nayoung Lee, Kartik Sreenivasan, Jason D Lee, Kangwook Lee, and Dimitris Papailiopoulos. Teaching arithmetic to small transformers. *arXiv preprint arXiv:2307.03381*, 2023.
- Yuanzhi Li, Tengyu Ma, and Hongyang Zhang. Algorithmic regularization in over-parameterized matrix sensing and neural networks with quadratic activations. In *Conference On Learning Theory*, pages 2–47. PMLR, 2018.
- Cong Ma, Kaizheng Wang, Yuejie Chi, and Yuxin Chen. Implicit regularization in nonconvex statistical estimation: Gradient descent converges linearly for phase retrieval, matrix completion, and blind deconvolution. *arXiv preprint arXiv:1711.10467*, 2017.
- Cong Ma, Kaizheng Wang, Yuejie Chi, and Yuxin Chen. Implicit regularization in nonconvex statistical estimation: Gradient descent converges linearly for phase retrieval and matrix completion. In *International Conference on Machine Learning*, pages 3345–3354. PMLR, 2018.
- Jianhao Ma and Salar Fattahi. Blessing of depth in linear regression: Deeper models have flatter landscape around the true solution. *Advances in Neural Information Processing Systems*, 35: 34334–34346, 2022.
- Jianhao Ma and Salar Fattahi. Global convergence of sub-gradient method for robust matrix recovery: Small initialization, noisy measurements, and over-parameterization. *Journal of Machine Learning Research*, 24(96):1–84, 2023a.
- Jianhao Ma and Salar Fattahi. On the optimization landscape of burer-monteiro factorization: When do global solutions correspond to ground truth? *arXiv preprint arXiv:2302.10963*, 2023b.
- Jianhao Ma, Lingjun Guo, and Salar Fattahi. Behind the scenes of gradient descent: A trajectory analysis via basis function decomposition. *arXiv preprint arXiv:2210.00346*, 2022.
- John Paisley and Lawrence Carin. A nonparametric bayesian model for kernel matrix completion. In *2010 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 2090–2093. IEEE, 2010.
- Dominik Stöger and Mahdi Soltanolkotabi. 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.
- Ruoyu Sun. *Matrix completion via nonconvex factorization: Algorithms and theory*. PhD thesis, University of Minnesota, 2015.
- Ruoyu Sun and Zhi-Quan Luo. Guaranteed matrix completion via non-convex factorization. *IEEE Transactions on Information Theory*, 62(11):6535–6579, 2016.

- Joel A Tropp et al. An introduction to matrix concentration inequalities. *Foundations and Trends® in Machine Learning*, 8(1-2):1–230, 2015.
- Stephen Tu, Ross Boczar, Max Simchowitz, Mahdi Soltanolkotabi, and Ben Recht. Low-rank solutions of linear matrix equations via procrustes flow. In *International Conference on Machine Learning*, pages 964–973. PMLR, 2016.
- Tomas Vaskevicius, Varun Kanade, and Patrick Rebeschini. Implicit regularization for optimal sparse recovery. *Advances in Neural Information Processing Systems*, 32, 2019.
- Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- Van Vu. A simple svd algorithm for finding hidden partitions. *Combinatorics, Probability and Computing*, 27(1):124–140, 2018.
- Martin J Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge university press, 2019.
- Xiang Wang, Chenwei Wu, Jason D Lee, Tengyu Ma, and Rong Ge. Beyond lazy training for over-parameterized tensor decomposition. *Advances in Neural Information Processing Systems*, 33:21934–21944, 2020.
- Qinqing Zheng and John Lafferty. Convergence analysis for rectangular matrix completion using burer-monteiro factorization and gradient descent. *arXiv preprint arXiv:1605.07051*, 2016.

Contents

A Related Work	18
B Preliminaries	19
B.1 Outline of the Appendix	19
B.2 Additional Notations	19
B.3 Important Intermediate Lemmas	19
C Proofs for Dynamic Signal-residual Decomposition	24
C.1 Proof of Signal Dynamic	24
C.2 Proof of Residual Dynamic	28
C.3 Proof of Error Dynamic	28
D Proofs for Main Theorems	31
D.1 Proof of Theorem 2	31
D.2 Proof of Theorem 3	33
E Proofs for Incoherence Dynamic	33
E.1 Proof of Lemma 6	33
E.2 Proof of Proposition 7	35
E.3 Proof of Proposition 8	38
E.4 Proof of Proposition 20	43
F Proofs for Different Initialization Schemes	52
G Concentration Inequalities for Matrix Completion	54
H Auxiliary Lemmas	56
H.1 Concentration Inequalities	56
H.2 Matrix Norm Inequalities	56
H.3 Other Useful Inequalities	58

Appendix A. Related Work

Nonconvex matrix completion. To solve the matrix completion problem, several algorithms based on convex optimization have been developed (Candes and Recht, 2012; Candès and Tao, 2010; Gross, 2011), offering excellent theoretical guarantees. However, in high-dimensional scenarios, convex optimization techniques require significant memory and computational resources due to the iterative singular value decompositions. To overcome these limitations, researchers have shifted towards nonconvex optimization techniques using first-order methods such as GD (Sun and Luo, 2016), projected GD (Zheng and Lafferty, 2016), and alternating minimization (Jain et al., 2013). Specifically, Sun and Luo (2016) demonstrate that GD can achieve local linear convergence provided that the initialization is close to the ground truth. Subsequent studies provide the global convergence guarantees for the first-order methods by showing the benign landscape of these nonconvex optimization formulations. Specifically, they reveal that the loss landscape has no spurious local minima and all the saddle points are strict (Ge et al., 2016, 2017; Chen and Li, 2017; Fattahi and Sojoudi, 2020). Nonetheless, these advancements necessitate either an explicit $\ell_{2,\infty}$ -norm regularization or a projection step to maintain the incoherence of the iterates. Moreover, these works are only applicable in the exactly-parameterized setting Ma and Fattahi (2023b). For a more detailed exploration of matrix completion and its variants, we refer the readers to the comprehensive survey by Chi et al. (2019).

Leave-one-out analysis. Leave-one-out analysis is a powerful statistical technique designed to decouple correlations among individual entries of a stochastic process. Initially employed by El Karoui et al. (2013) to establish asymptotic sampling distributions for robust estimators in high or moderate dimensional regression, this technique has been proven invaluable across a broad spectrum of applications. For instance, Abbe et al. (2020) utilized it to control ℓ_∞ estimation errors for eigenvectors in stochastic spectral problems, enabling precise spectral clustering in community detection without the need for data cleaning or regularization. More relevantly, Ma et al. (2017) applied leave-one-out analysis to demonstrate the local linear convergence of GD for the unregularized and symmetric matrix completion. Their approach not only elucidated the convergence properties of GD in matrix completion but also paved the way for similar analyses in other low-rank recovery challenges, such as phase retrieval and blind deconvolution. Extending these insights, Chen et al. (2020) and Kim and Chung (2022) broadened the scope of this analysis to include asymmetric matrix completion and global convergence in rank-1 scenarios, respectively. Furthermore, leave-one-out analysis has facilitated advancements in Singular Value Projection (SVP) for matrix completion, as demonstrated by Ding and Chen (2020) and has been instrumental in analyzing gradient descent with random initialization for phase retrieval, as shown by Chen et al. (2019).

Implicit regularization of GD in other applications. Indeed, the conventional wisdom in statistics suggests that increasing the number of parameters beyond the true dimension without proper regularization would lead to inferior solutions due to *overfitting*. However, a growing body of works show that, for a large class of learning problems, GD leads to surprisingly good solutions, due to its *implicit regularization* property. For instance, it is known that GD recovers the true low-dimensional solutions in matrix factorization and sensing (Gunasekar et al., 2018; Li et al., 2018; Stöger and Soltanolkotabi, 2021), tensor decomposition (Wang et al., 2020; Ge et al., 2021), deep linear neural networks (Arora et al., 2018; Ma and Fattahi, 2022), and beyond (Ma et al., 2022). However, the current theory behind the success of GD in these classes of problems hinges heavily upon a norm-preserving property

of the measurements, known as the restricted isometry property (RIP), limiting its applicability in settings where RIP is not satisfied.

Appendix B. Preliminaries

B.1. Outline of the Appendix

The structure of the appendix is as follows. In the remainder of this section, we introduce additional notation. Following this, we present key intermediate lemmas (Lemmas 10 to 13) crucial for our main proofs. Section C delves into a detailed proof of the signal and residual dynamics, starting with a refined version of Lemma 5 (Proposition 14) that takes into account the incoherence of \mathbf{V}_t . Additionally, the proof of Lemma 4 is provided in this section. Moving on to Section D, we present the proofs of our main theorems. Section E presents the key novelty of our paper, focusing on establishing the incoherence of \mathbf{V}_t via weakly-coupled leave-one-out analysis. The validation of different initialization schemes, as presented in Lemma 1, is addressed in Section F. In Section G, we compile several known results on matrix completion crucial to our arguments. Lastly, Section H collects several basic lemmas, which we include for completeness.

B.2. Additional Notations

We introduce some additional notations that will be used throughout the appendix. The max-norm of \mathbf{X} , denoted as $\|\mathbf{X}\|_{\max}$, is defined as $\max_{i,j} |X_{i,j}|$. We define the operator and Frobenius norm ball as $\mathcal{B}_{\text{op}}^{d_1 \times d_2}(r) := \{\mathbf{X} \in \mathbb{R}^{d_1 \times d_2} : \|\mathbf{X}\| \leq r\}$ and $\mathcal{B}_{\text{F}}^{d_1 \times d_2}(r) := \{\mathbf{X} \in \mathbb{R}^{d_1 \times d_2} : \|\mathbf{X}\|_{\text{F}} \leq r\}$, respectively. For any matrix \mathbf{X} , we denote its SVD as $\mathbf{X} = \mathbf{L}_{\mathbf{X}} \Sigma_{\mathbf{X}} \mathbf{R}_{\mathbf{X}}^{\top}$. We denote $\mathcal{S}_{d \times d}$ as the set of all the symmetric matrices $\mathbf{X} \in \mathbb{R}^{d \times d}$. In the appendix, $\Gamma, \Gamma_1, \Gamma_2, \dots$ denote fixed universal constants, while $C, C_1, C_2, c_1, c_2, \dots$ represent universal constants whose specific values may vary depending on the context.

Throughout the appendix, our arguments are conditioned on the following good event without further explanation. We define the random observation matrix Ω as

$$\Omega_{i,j} = \begin{cases} 1 & \text{if } (i,j) \in \Omega, \\ 0 & \text{otherwise.} \end{cases} \quad (16)$$

Then, the good event can be defined as

$$E_{\text{good}} = \left\{ \left\| \frac{\Omega + \Omega^{\top}}{2p} - \mathbf{J} \right\| \leq \Gamma \sqrt{\frac{d}{p}} \right\}. \quad (17)$$

Here \mathbf{J} is the all-one matrix. According to Lemma 25, we have $\mathbb{P}(E_{\text{good}}) \geq 1 - \frac{1}{d^3}$.

B.3. Important Intermediate Lemmas

Next, we collect some useful intermediate results that will be directly used throughout our proofs. We also note that some of these intermediate results rely on concentration inequalities for matrix completion, which are thoroughly discussed in Appendix G. Before proceeding, we define the following notations

$$\Delta_t = \mathbf{X}^* - \mathbf{U}_t \mathbf{U}_t^{\top}, \quad \mathbf{M}_t = \mathcal{R}_{\Omega} \left(\mathbf{X}^* - \mathbf{U}_t \mathbf{U}_t^{\top} \right), \quad \text{and} \quad \Lambda_t = \mathbf{S}_t \mathbf{E}_t^{\top} + \mathbf{E}_t \mathbf{S}_t^{\top} + \mathbf{E}_t \mathbf{E}_t^{\top}. \quad (18)$$

Moreover, we introduce the following term

$$\mathbf{A}_t = -\eta^2 \mathbf{M}_t \mathbf{V}_t \mathbf{V}_t^\top \mathbf{M}_t \mathbf{V}_t - 0.5\eta^2 \mathbf{V}_t \mathbf{V}_t^\top \mathbf{M}_t^2 \mathbf{V}_t - 0.5\eta^3 \mathbf{M}_t \mathbf{V}_t \mathbf{V}^{\star\top} \mathbf{M}_t^2 \mathbf{V}_t + (\mathbf{I} + \eta \mathbf{M}_t) \mathbf{V}_t \mathbf{R}(\mathbf{Y}_t) \quad (19)$$

where $\mathbf{Y}_t = \mathbf{V}_t^\top (2\eta \mathbf{M}_t + \eta^2 \mathbf{M}_t^2) \mathbf{V}_t$ and $\mathbf{R}(\mathbf{X}) = \sum_{k=2}^{\infty} \frac{(-1)^k (2k)!}{4^k (k!)^2} \mathbf{X}^k$. This notion of \mathbf{A}_t will be used when controlling the higher-order terms with respect to the step-size η . We are now ready to statement our helper lemmas.

Lemma 10 (Helper lemma for \mathbf{U}_t) Suppose that $\|\mathbf{S}_t\| \leq 2\sqrt{\sigma_1^*}$ and $\|\mathbf{E}_t\| \leq \sqrt{\frac{\sigma_1^* \mu r}{d}}$. Then, we have

$$\|\mathbf{U}_t \mathbf{U}_t^\top\|_{\text{F}} \leq 8\sqrt{r} \sigma_1^*. \quad (20)$$

Proof Applying triangle inequality, we have

$$\begin{aligned} \|\mathbf{U}_t \mathbf{U}_t^\top\|_{\text{F}} &\leq \|\mathbf{S}_t \mathbf{S}_t^\top\|_{\text{F}} + 2\|\mathbf{S}_t \mathbf{E}_t^\top\|_{\text{F}} + \|\mathbf{E}_t \mathbf{E}_t^\top\|_{\text{F}} \\ &\leq \sqrt{r} \|\mathbf{S}_t\|^2 + 2\sqrt{r} \|\mathbf{S}_t\| \|\mathbf{E}_t\| + \sqrt{d} \|\mathbf{E}_t\|^2 \\ &\leq 8\sqrt{r} \sigma_1^*. \end{aligned} \quad (21)$$

Here in the last inequality, we use the assumptions $\|\mathbf{S}_t\| \leq 2\sqrt{\sigma_1^*}$ and $\|\mathbf{E}_t\| \leq \sqrt{\frac{\sigma_1^* \mu r}{d}}$ and the fact that $d \gg \mu r$. \blacksquare

Lemma 11 (Helper lemma for $\mathbf{\Lambda}_t$) Suppose that $\|\mathbf{S}_t\| \leq 2\sqrt{\sigma_1^*}$, $\|\mathbf{V}_t\|_{2,\infty} \leq 2\sqrt{\frac{\mu r}{d}}$, and $\|\mathbf{E}_t\| \leq \sqrt{\frac{\sigma_1^* \mu r}{81d}}$. Then, conditioned on E_{good} , we have

$$\begin{aligned} \|\mathbf{\Lambda}_t\| &\leq 5\sqrt{\sigma_1^*} \|\mathbf{E}_t\|, \\ \|\mathcal{R}_\Omega(\mathbf{\Lambda}_t)\| &\leq 10\Gamma \sqrt{\frac{\sigma_1^* \mu r}{p}} \|\mathbf{E}_t\|, \\ \|(\mathcal{I} - \mathcal{R}_\Omega)(\mathbf{\Lambda}_t)\| &\leq 9\Gamma \sqrt{\frac{\sigma_1^* \mu r}{p}} \|\mathbf{E}_t\|. \end{aligned} \quad (22)$$

Proof First, we can bound $\|\mathbf{\Lambda}_t\|$ as follows

$$\|\mathbf{\Lambda}_t\| \leq \|\mathbf{S}_t \mathbf{E}_t^\top\| + \|\mathbf{E}_t \mathbf{S}_t^\top\| + \|\mathbf{E}_t \mathbf{E}_t^\top\| \leq 5\sqrt{\sigma_1^*} \|\mathbf{E}_t\| \quad (23)$$

where we use the assumptions $\|\mathbf{S}_t\| \leq 2\sqrt{\sigma_1^*}$ and $\|\mathbf{E}_t\| \leq \sqrt{\frac{\sigma_1^* \mu r}{d}} \leq \sqrt{\sigma_1^*}$. Next, we control $\|(\mathcal{I} - \mathcal{R}_\Omega)(\mathbf{\Lambda}_t)\|$. To this end, we first apply triangle inequality to obtain

$$\begin{aligned} \|(\mathcal{I} - \mathcal{R}_\Omega)(\mathbf{\Lambda}_t)\| &\leq 2 \left\| (\mathcal{I} - \mathcal{R}_\Omega) \left(\mathbf{S}_t \mathbf{E}_t^\top \right) \right\| + \left\| (\mathcal{I} - \mathcal{R}_\Omega) \left(\mathbf{E}_t \mathbf{E}_t^\top \right) \right\| \\ &\stackrel{(a)}{\leq} \Gamma \sqrt{\frac{d}{p}} \|\mathbf{E}_t\|_{2,\infty} \left(\|\mathbf{E}_t\|_{2,\infty} + 2 \|\mathbf{S}_t\|_{2,\infty} \right) \\ &\stackrel{(b)}{\leq} \Gamma \sqrt{\frac{d}{p}} \|\mathbf{E}_t\| \left(\|\mathbf{E}_t\| + 8 \sqrt{\frac{\sigma_1^* \mu r}{d}} \right) \\ &\leq 9\Gamma \sqrt{\frac{\sigma_1^* \mu r}{p}} \|\mathbf{E}_t\|. \end{aligned} \tag{24}$$

Here in (a), we apply Lemma 30. In (b), we use the facts that $\|\mathbf{E}_t\|_{2,\infty} \leq \|\mathbf{E}_t\| \leq \sqrt{\frac{\sigma_1^* \mu r}{d}}$ and

$$\|\mathbf{S}_t\|_{2,\infty} = \left\| \mathbf{V}_t \mathbf{V}_t^\top \mathbf{U}_t \right\|_{2,\infty} \stackrel{\text{Lemma 30}}{\leq} \|\mathbf{V}_t\|_{2,\infty} \|\mathbf{S}_t\| \leq 2\sqrt{\frac{\mu r}{d}} \cdot 2\sqrt{\sigma_1^*} = 4\sqrt{\frac{\sigma_1^* \mu r}{d}}. \tag{25}$$

In the last inequality, we use the fact that $\|\mathbf{S}_t\| \leq 2\sqrt{\sigma_1^*}$ and $\|\mathbf{V}_t\|_{2,\infty} \leq 2\sqrt{\frac{\mu r}{d}}$. Lastly, applying triangle inequality leads to

$$\begin{aligned} \|\mathcal{R}_\Omega(\mathbf{\Lambda}_t)\| &\leq \|\mathbf{\Lambda}_t\| + \|(\mathcal{I} - \mathcal{R}_\Omega)(\mathbf{\Lambda}_t)\| \\ &\leq 5\sqrt{\sigma_1^*} \|\mathbf{E}_t\| + 9\Gamma \sqrt{\frac{\sigma_1^* \mu r}{p}} \|\mathbf{E}_t\| \\ &\leq 10\Gamma \sqrt{\frac{\sigma_1^* \mu r}{p}} \|\mathbf{E}_t\|, \end{aligned} \tag{26}$$

where the last inequality is due to $p \leq \frac{1}{25}\Gamma^2 \mu r$. ■

Lemma 12 (Helper lemma for $\mathbf{\Delta}_t$) *Under the same conditions as Lemma 11 with the additional assumption that $\|\mathbf{V}_t - \mathbf{V}^*\| \leq \Gamma_1 \frac{\kappa \mu r^{1.5} \log(\frac{1}{\alpha})}{\sqrt{pd}}$, we have*

$$\begin{aligned} \|\mathbf{\Delta}_t\| &\leq 5\sigma_1^*, \\ \|(\mathcal{I} - \mathcal{R}_\Omega)(\mathbf{\Delta}_t)\| &\leq 5\Gamma \sqrt{\frac{\mu^2 r^2}{pd}} \|\mathbf{\Delta}_t\| + 10\Gamma \sqrt{\frac{\sigma_1^* \mu r}{p}} \|\mathbf{E}_t\|, \\ \|(\mathcal{I} - \mathcal{R}_\Omega)(\mathbf{\Delta}_t)\| &\leq 21\Gamma \sigma_1^* \sqrt{\frac{\mu^2 r^2}{pd}}, \\ \|\mathbf{M}_t\| &\leq \left(1 + 5\Gamma \sqrt{\frac{\mu^2 r^2}{pd}} \right) \|\mathbf{\Delta}_t\| + 10\Gamma \sqrt{\frac{\sigma_1^* \mu r}{p}} \|\mathbf{E}_t\|, \\ \|\mathbf{M}_t\| &\leq 6\sigma_1^*, \\ \left\| \mathbf{M}_t \mathcal{P}_{\mathbf{V}_t}^\perp \right\| &\leq 2\Gamma_1 \sigma_1^* \frac{\kappa \mu r^{1.5} \log(\frac{1}{\alpha})}{\sqrt{pd}}. \end{aligned} \tag{27}$$

Proof We first control $\|\Delta_t\|$ as follows

$$\begin{aligned}
 \|\Delta_t\| &\leq \left\| \mathbf{X}^\star - \mathbf{S}_t \mathbf{S}_t^\top \right\| + \|\Lambda_t\| \\
 &\leq \max\{\|\mathbf{X}^\star\|, \|\mathbf{S}_t\|^2\} + 5\sqrt{\sigma_1^\star} \|\mathbf{E}_t\| \\
 &\leq 4\sigma_1^\star + 5\sqrt{\sigma_1^\star} \|\mathbf{E}_t\| \\
 &\leq 5\sigma_1^\star.
 \end{aligned} \tag{28}$$

Next, we control $\|(\mathcal{I} - \mathcal{R}_\Omega)(\Delta_t)\|$. To this end, applying triangle inequality leads to

$$\begin{aligned}
 \|(\mathcal{I} - \mathcal{R}_\Omega)(\Delta_t)\| &\leq \underbrace{\left\| (\mathcal{I} - \mathcal{R}_\Omega)(\mathbf{X}^\star - \mathbf{S}_t \mathbf{S}_t^\top) \right\|}_{:= (I)} + \underbrace{\|(\mathcal{I} - \mathcal{R}_\Omega)(\Lambda_t)\|}_{:= (II)}.
 \end{aligned} \tag{29}$$

For (I), applying Lemma 31, we have

$$\begin{aligned}
 (I) &\leq \Gamma \sqrt{\frac{d}{p}} \left\| \mathbf{X}^\star - \mathbf{S}_t \mathbf{S}_t^\top \right\| \left(\|\mathbf{V}^\star\|_{2,\infty}^2 + \|\mathbf{V}_t\|_{2,\infty}^2 \right) \\
 &\stackrel{(a)}{\leq} 5\Gamma \sqrt{\frac{\mu^2 r^2}{pd}} \left\| \mathbf{X}^\star - \mathbf{S}_t \mathbf{S}_t^\top \right\| \\
 &\leq 5\Gamma \sqrt{\frac{\mu^2 r^2}{pd}} (\|\Delta_t\| + \|\Lambda_t\|) \\
 &\stackrel{(b)}{\leq} 5\Gamma \sqrt{\frac{\mu^2 r^2}{pd}} \|\Delta_t\| + 25\Gamma \sqrt{\frac{\sigma_1^\star \mu^2 r^2}{pd}} \|\mathbf{E}_t\|.
 \end{aligned} \tag{30}$$

Here in (a), we use the assumption that $\|\mathbf{V}_t\|_{2,\infty} \leq 2\sqrt{\frac{\mu r}{d}}$. In (b), we use the result from Lemma 11 that $\|\Lambda_t\| \leq 5\sqrt{\sigma_1^\star} \|\mathbf{E}_t\|$. On the other hand, we know that (II) $\leq 9\Gamma \sqrt{\frac{\sigma_1^\star \mu r}{p}} \|\mathbf{E}_t\|$ due to Lemma 11. Overall, we conclude that

$$\begin{aligned}
 \|(\mathcal{I} - \mathcal{R}_\Omega)(\Delta_t)\| &\leq 5\Gamma \sqrt{\frac{\mu^2 r^2}{pd}} \|\Delta_t\| + 25\Gamma \sqrt{\frac{\sigma_1^\star \mu^2 r^2}{pd}} \|\mathbf{E}_t\| + 9\Gamma \sqrt{\frac{\sigma_1^\star \mu r}{p}} \|\mathbf{E}_t\| \\
 &\leq 5\Gamma \sqrt{\frac{\mu^2 r^2}{pd}} \|\Delta_t\| + 10\Gamma \sqrt{\frac{\sigma_1^\star \mu r}{p}} \|\mathbf{E}_t\|.
 \end{aligned} \tag{31}$$

In the final inequality, we make the assumption $d \geq 9\mu r$ without loss of generality. This assumption simplifies the presentation of the proof but does not impact the final result of the paper.

Furthermore, upon noticing that $\|\mathbf{X}^\star - \mathbf{S}_t \mathbf{S}_t^\top\| \leq \max\{\|\mathbf{X}^\star\|, \|\mathbf{S}_t\|^2\} \leq 4\sigma_1^\star$, we have

$$\begin{aligned}
 \|(\mathcal{I} - \mathcal{R}_\Omega)(\Delta_t)\| &\leq 5\Gamma \sqrt{\frac{\mu^2 r^2}{pd}} \left\| \mathbf{X}^\star - \mathbf{S}_t \mathbf{S}_t^\top \right\| + \|(\mathcal{I} - \mathcal{R}_\Omega)(\Lambda_t)\| \\
 &\leq 20\Gamma \sigma_1^\star \sqrt{\frac{\mu^2 r^2}{pd}} + 9\Gamma \sqrt{\frac{\sigma_1^\star \mu r}{p}} \|\mathbf{E}_t\| \\
 &\leq 21\Gamma \sigma_1^\star \sqrt{\frac{\mu^2 r^2}{pd}}.
 \end{aligned} \tag{32}$$

Next, combining the above two inequalities, we can control $\|\mathbf{M}_t\|$ as follows:

$$\begin{aligned}
 \|\mathbf{M}_t\| &\leq \|\mathbf{\Delta}_t\| + \|(\mathcal{I} - \mathcal{R}_\Omega)(\mathbf{\Delta}_t)\| \\
 &\leq \|\mathbf{\Delta}_t\| + 5\Gamma\sqrt{\frac{\mu^2 r^2}{pd}} \|\mathbf{\Delta}_t\| + 10\Gamma\sqrt{\frac{\sigma_1^* \mu r}{p}} \|\mathbf{E}_t\| \\
 &\leq \left(1 + 5\Gamma\sqrt{\frac{\mu^2 r^2}{pd}}\right) \|\mathbf{\Delta}_t\| + 10\Gamma\sqrt{\frac{\sigma_1^* \mu r}{p}} \|\mathbf{E}_t\|.
 \end{aligned} \tag{33}$$

Furthermore, we can also bound $\|\mathbf{M}_t\|$ as

$$\begin{aligned}
 \|\mathbf{M}_t\| &\leq \|\mathbf{\Delta}_t\| + \|(\mathcal{I} - \mathcal{R}_\Omega)(\mathbf{\Delta}_t)\| \\
 &\leq 5\sigma_1^* + 21\Gamma\sigma_1^*\sqrt{\frac{\mu^2 r^2}{pd}} \\
 &\leq 6\sigma_1^*.
 \end{aligned} \tag{34}$$

Lastly, for $\|\mathbf{M}_t \mathcal{P}_{\mathbf{V}_t}^\perp\|$, we have the following decomposition

$$\begin{aligned}
 \|\mathbf{M}_t \mathcal{P}_{\mathbf{V}_t}^\perp\| &\leq \|(\mathcal{I} - \mathcal{R}_\Omega)(\mathbf{\Delta}_t) \mathcal{P}_{\mathbf{V}_t}^\perp\| + \left\| \left(\mathbf{X}^* - \mathbf{U}_t \mathbf{U}_t^\top \right) \mathcal{P}_{\mathbf{V}_t}^\perp \right\| \\
 &\leq \|(\mathcal{I} - \mathcal{R}_\Omega)(\mathbf{\Delta}_t)\| + \left\| \mathbf{V}^* \Sigma (\mathbf{V}^* - \mathbf{V}_t)^\top \mathcal{P}_{\mathbf{V}_t}^\perp \right\| + \left\| \mathbf{U}_t \mathbf{E}_t^\top \mathcal{P}_{\mathbf{V}_t}^\perp \right\| \\
 &\leq 21\Gamma\sigma_1^*\sqrt{\frac{\mu^2 r^2}{pd}} + \sigma_1^* \|\mathbf{V}^* - \mathbf{V}_t\| + \|\mathbf{E}_t\| (\|\mathbf{S}_t\| + \|\mathbf{E}_t\|) \\
 &\leq 22\Gamma\frac{\sigma_1^* \mu r}{\sqrt{pd}} + \sigma_1^* \|\mathbf{V}^* - \mathbf{V}_t\| \\
 &\stackrel{(a)}{\leq} 22\Gamma\frac{\sigma_1^* \mu r}{\sqrt{pd}} + \Gamma_1 \sigma_1^* \frac{\kappa \mu r^{1.5} \log\left(\frac{1}{\alpha}\right)}{\sqrt{pd}} \\
 &\leq 2\Gamma_1 \sigma_1^* \frac{\kappa \mu r^{1.5} \log\left(\frac{1}{\alpha}\right)}{\sqrt{pd}}.
 \end{aligned} \tag{35}$$

Here in (a), we apply Lemma 6. ■

Lemma 13 (Helper lemma for \mathbf{A}_t) Under the same conditions as Lemma 11, we have

$$\|\mathbf{A}_t\| \leq 300\eta^2 \sigma_1^{*2}. \tag{36}$$

Proof We first use triangle inequality to bound $\|\mathbf{A}_t\|$ as follows

$$\|\mathbf{A}_t\| \leq 1.5\eta^2 \|\mathbf{M}_t\|^2 + 0.5\eta^3 \|\mathbf{M}_t\|^3 + (1 + \eta \|\mathbf{M}_t\|) \|\mathbf{R}(\mathbf{Y}_t)\|. \tag{37}$$

Next, Lemma 12 tells us that $\|\mathbf{M}_t\| \leq 6\sigma_1^*$ conditioned on E_{good} . For $\|\mathbf{R}(\mathbf{Y}_t)\|$, we first have

$$\|\mathbf{R}(\mathbf{Y}_t)\| \leq \sum_{k=2}^{\infty} \frac{(2k)!}{4^k (k!)^2} \|\mathbf{Y}_t\|^k = \frac{2 + \sqrt{1 - \|\mathbf{Y}_t\|}}{2\sqrt{1 - \|\mathbf{Y}_t\|} \left(1 + \sqrt{1 - \|\mathbf{Y}_t\|}\right)^2} \|\mathbf{Y}_t\|^2. \tag{38}$$

Note that $\|\mathbf{Y}_t\| \leq \|2\eta\mathbf{M}_t + \eta^2\mathbf{M}_t^2\| \leq 20\eta\sigma_1^*$ and the right-hand side is an increasing function of $\|\mathbf{Y}_t\|$. Therefore, we derive that

$$\|\mathbf{R}(\mathbf{Y}_t)\| \leq \frac{1}{2} \|\mathbf{Y}_t\|^2 \leq 200\eta^2\sigma_1^{*2}. \quad (39)$$

This implies that

$$\|\mathbf{A}_t\| \leq 300\eta^2\sigma_1^{*2}. \quad (40)$$

■

Appendix C. Proofs for Dynamic Signal-residual Decomposition

We first present a more precise version of the one-step dynamics of the signal and residual terms.

Proposition 14 *Suppose that $\|\mathbf{S}_t\| \leq 2\sqrt{\sigma_1^*}$, $\|\mathbf{E}_t\| \leq \sqrt{\frac{\sigma_1^*}{d}}$, $\|\mathbf{V}_t\|_{2,\infty} \leq 2\sqrt{\frac{\mu r}{d}}$ and $\|\mathbf{V}^* - \mathbf{V}_t\|_F \leq \Gamma_1 \frac{\kappa\mu r^{1.5} \log(\frac{1}{\alpha})}{\sqrt{pd}}$. Then, the following dynamics hold conditioned on E_{good} :*

$$\begin{aligned} \sigma_r(\mathbf{S}_{t+1}) &\geq (1 + 0.8\eta\sigma_r^* - \eta\sigma_r^2(\mathbf{S}_t)) \sigma_r(\mathbf{S}_t) - 6\Gamma_1\eta \frac{\sigma_1^*\kappa\mu r^{1.5} \log(\frac{1}{\alpha})}{\sqrt{pd}} \|\mathbf{E}_t\|, \\ &\hspace{15em} \text{(minimal signal dynamic)} \\ \|\mathbf{S}_{t+1}\| &\leq 2\sqrt{\sigma_1^*}, \\ &\hspace{15em} \text{(maximal signal dynamic)} \\ \|\mathbf{E}_{t+1}\| &\leq \left(1 + 2\Gamma_1\eta \frac{\sigma_1^*\kappa\mu r^{1.5} \log(\frac{1}{\alpha})}{\sqrt{pd}}\right) \|\mathbf{E}_t\|. \\ &\hspace{15em} \text{(residual dynamic)} \end{aligned}$$

Additionally, if $\sigma_r(\mathbf{S}_t) \geq \frac{\sqrt{\sigma_r^*}}{2}$, then

$$\begin{aligned} \|\mathbf{X}^* - \mathbf{U}_{t+1}\mathbf{U}_{t+1}^\top\|_F &\leq \left(1 - \frac{1}{10}\eta\sigma_r^*\right) \|\mathbf{X}^* - \mathbf{U}_t\mathbf{U}_t^\top\|_F + \Gamma_6\eta \sqrt{\frac{\sigma_1^{*3}\mu r^2}{p}} \|\mathbf{E}_t\|. \\ &\hspace{15em} \text{(error dynamic)} \end{aligned}$$

The key distinction between the above proposition and Lemma 5 lies in a finer control over the one-step dynamics, accompanied by additional assumptions on $\|\mathbf{E}_t\|$, $\|\mathbf{V}_t\|_{2,\infty}$, and $\|\mathbf{V}_t - \mathbf{V}^*\|_F$. We emphasize that the one-step dynamics in Lemma 5 are derived from the proof of this proposition. Moreover, the proof of the error dynamics (last inequality in Proposition 14) will incorporate the proof of the descent lemma (Lemma 4).

C.1. Proof of Signal Dynamic

We notice that

$$\mathbf{S}_{t+1} = \mathcal{P}_{\mathbf{V}_{t+1}}(\mathbf{I} + \eta\mathbf{M}_t)(\mathbf{S}_t + \mathbf{E}_t) = (\mathbf{I} + \eta\mathbf{M}_t)\mathbf{S}_t + \mathcal{P}_{\mathbf{V}_{t+1}}(\mathbf{I} + \eta\mathbf{M}_t)\mathbf{E}_t. \quad (41)$$

Here the second equality follows from the definition of $\mathcal{P}_{\mathbf{V}_{t+1}}$.

Maximal signal dynamic. We first provide an upper-bound for $\|\mathbf{S}_{t+1}\|$ by

$$\|\mathbf{S}_{t+1}\| \leq \underbrace{\|(\mathbf{I} + \eta \mathbf{M}_t) \mathbf{S}_t\|}_{:= (\text{I})} + \underbrace{\|\mathcal{P}_{\mathbf{V}_{t+1}}(\mathbf{I} + \eta \mathbf{M}_t) \mathbf{E}_t\|}_{:= (\text{II})}. \quad (42)$$

Next, we further control (I) by

$$\begin{aligned} (\text{I}) &\leq \|(\mathbf{I} + \eta \mathbf{\Delta}_t) \mathbf{S}_t\| + \eta \|(\mathcal{I} - \mathcal{R}_\Omega)(\mathbf{\Delta}_t)\| \|\mathbf{S}_t\| \\ &\stackrel{(a)}{\leq} \|(\mathbf{I} + \eta \mathbf{\Delta}_t) \mathbf{S}_t\| + \eta \cdot 21\Gamma\sigma_1^* \sqrt{\frac{\mu^2 r^2}{pd}} \cdot 2\sqrt{\sigma_1^*} \\ &\leq \left\| \left(\mathbf{I} + \eta \left(\mathbf{X}^* - \mathbf{S}_t \mathbf{S}_t^\top \right) \right) \mathbf{S}_t \right\| + \eta \|\mathbf{\Lambda}_t\| \|\mathbf{S}_t\| + 42\Gamma\eta \frac{\sigma_1^{*1.5} \mu r}{\sqrt{pd}} \\ &\stackrel{(b)}{\leq} \left\| \left(\mathbf{I} + \eta \left(\mathbf{X}^* - \mathbf{S}_t \mathbf{S}_t^\top \right) \right) \mathbf{S}_t \right\| + 10\eta\sigma_1^* \|\mathbf{E}_t\| + 42\Gamma\eta \frac{\sigma_1^{*1.5} \mu r}{\sqrt{pd}} \\ &\leq \underbrace{\left\| \left(\mathbf{I} + \eta \left(\mathbf{X}^* - \mathbf{S}_t \mathbf{S}_t^\top \right) \right) \mathbf{S}_t \right\|}_{:= (\text{I}_1)} + 43\Gamma\eta \frac{\sigma_1^{*1.5} \mu r}{\sqrt{pd}}. \end{aligned} \quad (43)$$

Here we apply Lemma 12 in (a) and (b). In the last inequality, we use the assumption that $\|\mathbf{E}_t\| \leq \frac{\Gamma\sqrt{\sigma_1^*} \mu r}{\sqrt{pd}}$. For (I₁), we further decompose it via triangle inequality as follows

$$\begin{aligned} (\text{I}_1) &\leq \left\| \mathcal{P}_{\mathbf{V}_t} \left(\mathbf{I} + \eta \left(\mathbf{X}^* - \mathbf{S}_t \mathbf{S}_t^\top \right) \right) \mathbf{S}_t \right\| + \left\| \mathcal{P}_{\mathbf{V}_t^\perp} \left(\mathbf{I} + \eta \left(\mathbf{X}^* - \mathbf{S}_t \mathbf{S}_t^\top \right) \right) \mathbf{S}_t \right\| \\ &= \underbrace{\left\| \left(\mathbf{I} + \eta \left(\mathbf{V}_t^\top \mathbf{V}^* \mathbf{\Sigma}^* \mathbf{V}^{*\top} \mathbf{V}_t - \bar{\mathbf{S}}_t \bar{\mathbf{S}}_t^\top \right) \right) \bar{\mathbf{S}}_t \right\|}_{:= (\text{I}_{1,1})} + \underbrace{\eta \left\| \mathcal{P}_{\mathbf{V}_t^\perp} (\mathbf{V}^* - \mathbf{V}_t) \mathbf{\Sigma}^* \mathbf{V}^{*\top} \mathbf{S}_t \right\|}_{:= (\text{I}_{1,2})}. \end{aligned} \quad (44)$$

Here we define $\bar{\mathbf{S}}_t = \mathbf{V}_t^\top \mathbf{U}_t \in \mathbb{R}^{r \times d}$. Note that $\|\bar{\mathbf{S}}_t\| = \|\mathbf{S}_t\|$. For (I_{1,1}), we have

$$\begin{aligned} (\text{I}_{1,1}) &\leq \left\| \left(\mathbf{I} - \eta \bar{\mathbf{S}}_t \bar{\mathbf{S}}_t^\top \right) \bar{\mathbf{S}}_t \right\| + \eta \left\| \mathbf{V}_t^\top \mathbf{V}^* \mathbf{\Sigma}^* \mathbf{V}^{*\top} \mathbf{V}_t \bar{\mathbf{S}}_t \right\| \\ &\stackrel{(a)}{=} \|\bar{\mathbf{S}}_t\| \left(1 - \eta \|\bar{\mathbf{S}}_t\|^2 \right) + \eta\sigma_1^* \|\bar{\mathbf{S}}_t\| \\ &= \|\mathbf{S}_t\| \left(1 + \eta\sigma_1^* - \eta \|\mathbf{S}_t\|^2 \right). \end{aligned} \quad (45)$$

where (a) follows from the fact that $\bar{\mathbf{S}}_t$ and $\bar{\mathbf{S}}_t \bar{\mathbf{S}}_t^\top$ share the same eigenvectors, and the assumption $\eta \lesssim 1/\sigma_1$. Next, we control (I_{1,2}):

$$(\text{I}_{1,2}) \leq \sigma_1^* \|\mathbf{V}^* - \mathbf{V}_t\| \|\mathbf{S}_t\| \leq 2\Gamma_1 \frac{\kappa\mu(\sigma_1^* r)^{1.5} \log\left(\frac{1}{\alpha}\right)}{\sqrt{pd}}. \quad (46)$$

Here we use Lemma 6. Therefore, we can bound (I₁) by

$$(\text{I}_1) \leq \|\mathbf{S}_t\| \left(1 + \eta\sigma_1^* - \eta \|\mathbf{S}_t\|^2 \right) + 2\Gamma_1 \eta \frac{\kappa\mu(\sigma_1^* r)^{1.5} \log\left(\frac{1}{\alpha}\right)}{\sqrt{pd}}. \quad (47)$$

This leads to

$$\begin{aligned}
 (\text{I}) &\leq \|\mathbf{S}_t\| \left(1 + \eta\sigma_1^* - \eta\|\mathbf{S}_t\|^2\right) + 2\Gamma_1\eta \frac{\kappa\mu(\sigma_1^*r)^{1.5} \log\left(\frac{1}{\alpha}\right)}{\sqrt{pd}} + 43\Gamma_1\eta \frac{\sigma_1^{*1.5}\mu r}{\sqrt{pd}} \\
 &\leq \|\mathbf{S}_t\| \left(1 + \eta\sigma_1^* - \eta\|\mathbf{S}_t\|^2\right) + 3\Gamma_1\eta \frac{\kappa\mu(\sigma_1^*r)^{1.5} \log\left(\frac{1}{\alpha}\right)}{\sqrt{pd}}.
 \end{aligned} \tag{48}$$

Next, we control (II). To this end, we first notice that

$$\mathcal{P}_{\mathbf{V}_{t+1}}(\mathbf{I} + \eta\mathbf{M}_t)\mathcal{P}_{\mathbf{V}_t}^\perp = (\mathcal{P}_{\mathbf{V}_{t+1}} - \mathcal{P}_{\mathbf{V}_t})\mathcal{P}_{\mathbf{V}_t}^\perp + \eta\mathcal{P}_{\mathbf{V}_{t+1}}\mathbf{M}_t\mathcal{P}_{\mathbf{V}_t}^\perp. \tag{49}$$

Hence, we can bound (II) by

$$\begin{aligned}
 (\text{II}) &\leq \|\mathcal{P}_{\mathbf{V}_{t+1}} - \mathcal{P}_{\mathbf{V}_t}\| \|\mathbf{E}_t\| + \eta \left\| \mathbf{M}_t \mathcal{P}_{\mathbf{V}_t}^\perp \right\| \|\mathbf{E}_t\| \\
 &\stackrel{(a)}{\leq} \left(2\|\mathbf{V}_{t+1} - \mathbf{V}_t\| + 2\Gamma_1\eta\sigma_1^* \frac{\kappa\mu r^{1.5} \log\left(\frac{1}{\alpha}\right)}{\sqrt{pd}} \right) \|\mathbf{E}_t\|.
 \end{aligned} \tag{50}$$

Here we use Lemma 43 in (a). It remains to control $\|\mathbf{V}_{t+1} - \mathbf{V}_t\|$. To this end, upon noticing that $\mathbf{V}_{t+1} - \mathbf{V}_t = \eta\mathcal{P}_{\mathbf{V}_t}^\perp\mathbf{M}_t\mathbf{V}_t + \mathbf{A}_t$, one has

$$\begin{aligned}
 \|\mathbf{V}_{t+1} - \mathbf{V}_t\| &\leq \eta \left\| \mathbf{M}_t \mathcal{P}_{\mathbf{V}_t}^\perp \right\| + \|\mathbf{A}_t\| \\
 &\stackrel{(a)}{\leq} 2\Gamma_1\eta\sigma_1^* \frac{\kappa\mu r^{1.5} \log\left(\frac{1}{\alpha}\right)}{\sqrt{pd}} + 300\eta^2\sigma_1^{*2} \\
 &\leq 3\Gamma_1\eta\sigma_1^* \frac{\kappa\mu r^{1.5} \log\left(\frac{1}{\alpha}\right)}{\sqrt{pd}}.
 \end{aligned} \tag{51}$$

Here we apply Lemma 12 and Lemma 13 in (a). Hence, we have

$$(\text{II}) \leq 5\Gamma_1\eta\sigma_1^* \frac{\kappa\mu r^{1.5} \log\left(\frac{1}{\alpha}\right)}{\sqrt{pd}} \|\mathbf{E}_t\|. \tag{52}$$

Putting everything together, we obtain that

$$\begin{aligned}
 \|\mathbf{S}_{t+1}\| &\leq \|\mathbf{S}_t\| \left(1 + \eta\sigma_1^* - \eta\|\mathbf{S}_t\|^2\right) + 3\Gamma_1\eta\sigma_1^* \frac{\kappa\mu r^{1.5} \log\left(\frac{1}{\alpha}\right)}{\sqrt{pd}} \left(\sqrt{\sigma_1^*} + 2\|\mathbf{E}_t\|\right) \\
 &\leq \|\mathbf{S}_t\| \left(1 + \eta\sigma_1^* - \eta\|\mathbf{S}_t\|^2\right) + 6\Gamma_1\eta\sigma_1^{*1.5} \frac{\kappa\mu r^{1.5} \log\left(\frac{1}{\alpha}\right)}{\sqrt{pd}}.
 \end{aligned} \tag{53}$$

Next, we consider two cases separately. First, if $\|\mathbf{S}_t\| \leq 1.5\sqrt{\sigma_1^*}$, then we simply have

$$\begin{aligned}
 \|\mathbf{S}_{t+1}\| &\leq \|\mathbf{S}_t\| \left(1 + \eta\sigma_1^* - \eta\|\mathbf{S}_t\|^2\right) + 6\Gamma_1\eta\sigma_1^{*1.5} \frac{\kappa\mu r^{1.5} \log\left(\frac{1}{\alpha}\right)}{\sqrt{pd}} \\
 &\leq 1.5\sqrt{\sigma_1^*} \cdot (1 + \eta\sigma_1^*) + 6\Gamma_1\eta\sigma_1^{*1.5} \frac{\kappa\mu r^{1.5} \log\left(\frac{1}{\alpha}\right)}{\sqrt{pd}} \\
 &\leq 2\sqrt{\sigma_1^*}.
 \end{aligned} \tag{54}$$

On the other hand, if $1.5\sqrt{\sigma_1^*} \leq \|\mathbf{S}_t\| \leq 2\sqrt{\sigma_1^*}$, then we have

$$\begin{aligned} \|\mathbf{S}_{t+1}\| &\leq \|\mathbf{S}_t\| \left(1 + \eta\sigma_1^* - \eta\|\mathbf{S}_t\|^2\right) + 6\Gamma_1\eta\sigma_1^{*1.5} \frac{\kappa\mu r^{1.5} \log\left(\frac{1}{\alpha}\right)}{\sqrt{pd}} \\ &\leq 2\sqrt{\sigma_1^*} (1 - 1.25\eta\sigma_1^*) + 6\Gamma_1\eta\sigma_1^{*1.5} \frac{\kappa\mu r^{1.5} \log\left(\frac{1}{\alpha}\right)}{\sqrt{pd}} \\ &\leq 2\sqrt{\sigma_1^*}. \end{aligned} \quad (55)$$

This completes the proof for the maximal signal dynamic.

Minimal signal dynamic. We first provide a lower-bound for $\sigma_r(\mathbf{S}_{t+1})$ as follows

$$\sigma_r(\mathbf{S}_{t+1}) \geq \underbrace{\sigma_r((\mathbf{I} + \eta\mathbf{M}_t) \mathbf{S}_t)}_{(I)} - \underbrace{\|\mathcal{P}_{\mathbf{V}_{t+1}}(\mathbf{I} + \eta\mathbf{M}_t) \mathbf{E}_t\|}_{(II)}. \quad (56)$$

For (I), applying Lemma 44, we first obtain that (I) $\geq \sigma_r(\mathbf{V}_t^\top (\mathbf{I} + \eta\mathbf{M}_t) \mathbf{S}_t)$. Next, we decompose $\mathbf{V}_t^\top (\mathbf{I} + \eta\mathbf{M}_t) \mathbf{S}_t$ as follows

$$\begin{aligned} &\mathbf{V}_t^\top (\mathbf{I} + \eta\mathbf{M}_t) \mathbf{S}_t \\ &= \underbrace{\left(\mathbf{I} + \eta\mathbf{V}_t^\top (\mathbf{M}_t + \mathbf{S}_t\mathbf{S}_t^\top) \mathbf{V}_t (\mathbf{I} - \eta\mathbf{V}_t^\top \mathbf{S}_t\mathbf{S}_t^\top \mathbf{V}_t)^{-1}\right)}_{:=\mathbf{B}_t} \underbrace{\mathbf{V}_t^\top \mathbf{S}_t (\mathbf{I} - \eta\mathbf{S}_t^\top \mathbf{S}_t)}_{:=\mathbf{C}_t}. \end{aligned} \quad (57)$$

According to Lemma 45, we have

$$\sigma_r((\mathbf{I} + \eta\mathbf{M}_t) \mathbf{S}_t) \geq \sigma_r(\mathbf{B}_t)\sigma_r(\mathbf{C}_t) = \sigma_r(\mathbf{B}_t) (1 - \eta\sigma_r^2(\mathbf{S}_t)) \sigma_r(\mathbf{S}_t), \quad (58)$$

where in the last equality, we use the fact that \mathbf{S}_t and $\mathbf{S}_t\mathbf{S}_t^\top \mathbf{S}_t$ share the same singular space and hence

$$\sigma_r(\mathbf{C}_t) = \sigma_r(\mathbf{S}_t(\mathbf{I} - \eta\mathbf{S}_t^\top \mathbf{S}_t)) = (1 - \eta\sigma_r^2(\mathbf{S}_t)) \sigma_r(\mathbf{S}_t). \quad (59)$$

Now, it suffices to provide a lower-bound for $\sigma_r(\mathbf{B}_t)$. To this end, we first notice that

$$\begin{aligned} \sigma_r(\mathbf{B}_t) &\geq 1 + \eta\sigma_r\left(\mathbf{V}_t^\top (\mathbf{M}_t + \mathbf{S}_t\mathbf{S}_t^\top) \mathbf{V}_t\right) \sigma_r\left((\mathbf{I} - \eta\mathbf{V}_t^\top \mathbf{S}_t\mathbf{S}_t^\top \mathbf{V}_t)^{-1}\right) \\ &\geq 1 + \eta\sigma_r\left(\mathbf{V}_t^\top (\mathbf{M}_t + \mathbf{S}_t\mathbf{S}_t^\top) \mathbf{V}_t\right). \end{aligned} \quad (60)$$

Here the second inequality is due to $\mathbf{I} - \eta\mathbf{V}_t^\top \mathbf{S}_t\mathbf{S}_t^\top \mathbf{V}_t \preceq \mathbf{I}$. To proceed, notice that

$$\mathbf{M}_t + \mathbf{S}_t\mathbf{S}_t^\top = \mathbf{X}^* + (\mathcal{R}_\Omega - \mathcal{I})(\boldsymbol{\Delta}_t) - \boldsymbol{\Lambda}_t. \quad (61)$$

Therefore, we have

$$\begin{aligned} \sigma_r\left(\mathbf{V}_t^\top (\mathbf{M}_t + \mathbf{S}_t\mathbf{S}_t^\top) \mathbf{V}_t\right) &\geq \sigma_r\left(\mathbf{V}_t^\top \mathbf{X}^* \mathbf{V}_t\right) - \|(\mathcal{R}_\Omega - \mathcal{I})(\boldsymbol{\Delta}_t)\| - \|\boldsymbol{\Lambda}_t\| \\ &\geq \sigma_r\left(\mathbf{V}_t^\top \mathbf{X}^* \mathbf{V}_t\right) - 21\Gamma\sigma_1^* \sqrt{\frac{\mu^2 r^2}{pd}} - 5\sqrt{\sigma_1^*} \|\mathbf{E}_t\| \\ &\geq \sigma_r\left(\mathbf{V}_t^\top \mathbf{X}^* \mathbf{V}_t\right) - 22\Gamma\sigma_1^* \sqrt{\frac{\mu^2 r^2}{pd}}. \end{aligned} \quad (62)$$

For the first term in the above inequality, we have

$$\begin{aligned}
 \sigma_r \left(\mathbf{V}_t^\top \mathbf{X}^* \mathbf{V}_t \right) &\geq \sigma_r \left(\mathbf{\Sigma}^* \mathbf{V}^{*\top} \mathbf{V}_t \right) - \sigma_1^* \|\mathbf{V}^* - \mathbf{V}_t\| \\
 &\geq \sigma_r^* - 2\sigma_1^* \|\mathbf{V}^* - \mathbf{V}_t\| \\
 &\geq \sigma_r^* - 2\Gamma_1 \sigma_1^* \frac{\kappa \mu r^{1.5} \log \left(\frac{1}{\alpha} \right)}{\sqrt{pd}} \\
 &\geq 0.95\sigma_r^*.
 \end{aligned} \tag{63}$$

Therefore, we obtain

$$\sigma_r \left(\mathbf{V}_t^\top \left(\mathbf{M}_t + \mathbf{S}_t \mathbf{S}_t^\top \right) \mathbf{V}_t \right) \geq 0.95\sigma_r^* - 22\Gamma_1 \sigma_1^* \sqrt{\frac{\mu^2 r^2}{pd}} \geq 0.9\sigma_r^*. \tag{64}$$

Combining the above arguments, we have

$$(\text{I}) \geq (1 + 0.9\eta\sigma_r^*) (1 - \eta\sigma_r^2(\mathbf{S}_t)) \sigma_r(\mathbf{S}_t) \geq (1 + 0.8\eta\sigma_r^* - \eta\sigma_r^2(\mathbf{S}_t)) \sigma_r(\mathbf{S}_t). \tag{65}$$

On the other hand, we have already derived an upper bound for (II) in the maximal signal dynamic, which is

$$(\text{II}) \leq 5\Gamma_1 \eta \sigma_1^* \frac{\kappa \mu r^{1.5} \log \left(\frac{1}{\alpha} \right)}{\sqrt{pd}} \|\mathbf{E}_t\|. \tag{66}$$

Putting everything together, we have

$$\sigma_r(\mathbf{S}_{t+1}) \geq (1 + 0.8\eta\sigma_r^* - \eta\sigma_r^2(\mathbf{S}_t)) \sigma_r(\mathbf{S}_t) - 5\Gamma_1 \eta \frac{\sigma_1^* \kappa \mu r^{1.5} \log \left(\frac{1}{\alpha} \right)}{\sqrt{pd}} \|\mathbf{E}_t\|. \tag{67}$$

C.2. Proof of Residual Dynamic

First, we can expand \mathbf{E}_{t+1} as below

$$\mathbf{E}_{t+1} = \mathcal{P}_{\mathbf{V}_{t+1}}^\perp (\mathbf{I} + \eta \mathbf{M}_t) (\mathbf{S}_t + \mathbf{E}_t) = \mathcal{P}_{\mathbf{V}_{t+1}}^\perp (\mathbf{I} + \eta \mathbf{M}_t) \mathbf{E}_t, \tag{68}$$

where in the second equality, we use the fact that $\mathcal{P}_{\mathbf{V}_{t+1}}^\perp (\mathbf{I} + \eta \mathbf{M}_t) \mathbf{S}_t = 0$. Then, by triangle inequality, we obtain

$$\|\mathbf{E}_{t+1}\| \leq \left(1 + \eta \left\| \mathbf{M}_t \mathcal{P}_{\mathbf{V}_t}^\perp \right\| \right) \|\mathbf{E}_t\| \leq \left(1 + 2\Gamma_1 \eta \frac{\sigma_1^* \kappa \mu r^{1.5} \log \left(\frac{1}{\alpha} \right)}{\sqrt{pd}} \right) \|\mathbf{E}_t\|. \tag{69}$$

Here in the last inequality, we use Lemma 12.

C.3. Proof of Error Dynamic

The core proof idea is adapted from the proof of Proposition 4.3 appeared in (Li et al., 2018). We first expand $\|\Delta_{t+1}\|_{\text{F}}^2$ as

$$\begin{aligned}
 \|\Delta_{t+1}\|_{\text{F}}^2 &= \left\| \mathbf{X}^* - (\mathbf{I} + \eta \mathbf{M}_t) \mathbf{U}_t \mathbf{U}_t^\top (\mathbf{I} + \eta \mathbf{M}_t) \right\|_{\text{F}}^2 \\
 &= \|\Delta_t\|_{\text{F}}^2 - 4\eta \underbrace{\left\langle \mathbf{X}^* - \mathbf{U}_t \mathbf{U}_t^\top, \mathbf{M}_t \mathbf{U}_t \mathbf{U}_t^\top \right\rangle}_{(\text{I})} + (\text{II}),
 \end{aligned} \tag{70}$$

where

$$\begin{aligned} (\text{II}) &= 2 \left\langle -\Delta_t + \eta M_t U_t U_t^\top, \eta^2 M_t U_t U_t^\top M_t \right\rangle \\ &\quad + \eta^2 \left\| M_t U_t U_t^\top + U_t U_t^\top M_t \right\|_F^2 + \eta^4 \left\| M_t U_t U_t^\top M_t \right\|_F^2 \end{aligned} \quad (71)$$

contains all the higher-order terms. We then provide a lower-bound for (I). To this goal, we notice that

$$\begin{aligned} (\text{I}) &= \|\Delta_t U_t\|_F^2 - \left\langle \Delta_t, (I - \mathcal{R}_\Omega) (\Delta_t) U_t U_t^\top \right\rangle \\ &\geq \|\Delta_t U_t\|_F^2 - \|\Delta_t\|_F \left\| U_t U_t^\top \right\|_F \|(I - \mathcal{R}_\Omega) (\Delta_t)\| \\ &\geq \|\Delta_t U_t\|_F^2 - \|\Delta_t\|_F \cdot 8\sqrt{r}\sigma_1^* \cdot \left(5\Gamma \sqrt{\frac{\mu^2 r^2}{pd}} \|\Delta_t\| + 10\Gamma \sqrt{\frac{\sigma_1^* \mu r}{p}} \|E_t\| \right) \\ &= \|\Delta_t U_t\|_F^2 - 40\Gamma \sigma_1^* \sqrt{\frac{\mu^2 r^3}{pd}} \|\Delta_t\|_F^2 - 80\Gamma \sqrt{\frac{\sigma_1^3 \mu r^2}{p}} \|E_t\| \|\Delta_t\|_F. \end{aligned} \quad (72)$$

Here in the last inequality, we apply Lemma 12. Next, we provide a lower-bound for $\|\Delta_t U_t\|_F^2$. To this goal, we first notice that

$$\begin{aligned} \|\Delta_t U_t\|_F &= \left\| \left(X^* - S_t S_t^\top \right) S_t - \Lambda_t S_t + \Delta_t E_t \right\|_F \\ &\geq \left\| \left(X^* - S_t S_t^\top \right) S_t \right\|_F - \|\Lambda_t\| \|S_t\|_F - \|\Delta_t\|_F \|E_t\| \\ &\geq \left\| \left(X^* - S_t S_t^\top \right) S_t \right\|_F - 10\sqrt{r}\sigma_1^* \|E_t\| - \|\Delta_t\|_F \|E_t\|. \end{aligned} \quad (73)$$

In the last inequality we use the fact that $\|\Lambda_t\| \leq 5\sqrt{\sigma_1^*} \|E_t\|$ and $\|S_t\|_F \leq \sqrt{r} \|S_t\| \leq 2\sqrt{r\sigma_1^*}$ from Lemma 11. Applying Lemma 35, we can further provide a lower-bound for $\left\| \left(X^* - S_t S_t^\top \right) S_t \right\|_F$ as

$$\left\| \left(X^* - S_t S_t^\top \right) S_t \right\|_F \geq \sigma_r(S_t) \left\| \left(X^* - S_t S_t^\top \right) V_t \right\|_F \geq \frac{\sqrt{\sigma_r^*}}{2} \left\| \left(X^* - S_t S_t^\top \right) V_t \right\|_F. \quad (74)$$

Next, we present the following intermediate lemma to control $\left\| \left(X^* - S_t S_t^\top \right) V_t \right\|_F$.

Lemma 15 *Suppose that $\|V^* - V_t\| \leq 0.1$. Then, we have*

$$\left\| \left(X^* - S_t S_t^\top \right) V_t \right\|_F^2 \geq \frac{2}{5} \left\| X^* - S_t S_t^\top \right\|_F^2. \quad (75)$$

We first use this lemma to finish the proof of the loss dynamic and defer the proof to the end of this section. Applying this lemma to Equation (74) yields

$$\begin{aligned} \left\| \left(X^* - S_t S_t^\top \right) S_t \right\|_F &\geq \sqrt{\frac{\sigma_r^*}{10}} \left\| X^* - S_t S_t^\top \right\|_F \\ &\geq \sqrt{\frac{\sigma_r^*}{10}} \|\Delta_t\|_F - \sqrt{\frac{\sigma_r^*}{10}} \|\Lambda_t\|_F \\ &\geq \sqrt{\frac{\sigma_r^*}{10}} \|\Delta_t\|_F - \sigma_1^* \sqrt{\frac{5}{2\kappa}} \|E_t\|. \end{aligned} \quad (76)$$

Combining Equation (76) and Equation (73) leads to

$$\begin{aligned}\|\Delta_t \mathbf{U}_t\|_F &\geq \sqrt{\frac{\sigma_r^*}{10}} \|\Delta_t\|_F - \sigma_1^* \sqrt{\frac{5}{2\kappa}} \|\mathbf{E}_t\| - 10\sqrt{r}\sigma_1^* \|\mathbf{E}_t\| - \|\Delta_t\|_F \|\mathbf{E}_t\| \\ &\geq \sqrt{\frac{\sigma_r^*}{10}} \|\Delta_t\|_F - 20\sqrt{r}\sigma_1^* \|\mathbf{E}_t\|.\end{aligned}\quad (77)$$

This implies that

$$\|\Delta_t \mathbf{U}_t\|_F^2 \geq \frac{\sigma_r^*}{10} \|\Delta_t\|_F^2 - 13\sqrt{\frac{r}{\kappa}} \sigma_1^{1.5} \|\Delta_t\|_F \|\mathbf{E}_t\|. \quad (78)$$

Overall, we obtain

$$\begin{aligned}(\text{I}) &\geq \frac{\sigma_r^*}{10} \|\Delta_t\|_F^2 - 13\sqrt{\frac{r}{\kappa}} \sigma_1^{1.5} \|\Delta_t\|_F \|\mathbf{E}_t\| - 40\Gamma \sigma_1^* \sqrt{\frac{\mu^2 r^3}{pd}} \|\Delta_t\|_F^2 - 80C \sqrt{\frac{\sigma_1^{*3} \mu r^2}{p}} \|\mathbf{E}_t\| \|\Delta_t\|_F \\ &\geq \frac{\sigma_r^*}{15} \|\Delta_t\|_F^2 - 81\Gamma \sqrt{\frac{\sigma_1^{*3} \mu r^2}{p}} \|\mathbf{E}_t\| \|\Delta_t\|_F.\end{aligned}\quad (79)$$

Next, we control (II). To this end, we first notice that

$$\begin{aligned}-2 \left\langle \Delta_t, \mathbf{M}_t \mathbf{U}_t \mathbf{U}_t^\top \mathbf{M}_t \right\rangle &\leq 2 \|\Delta_t\|_F \|\mathbf{M}_t\|^2 \left\| \mathbf{U}_t \mathbf{U}_t^\top \right\|_F \\ &\leq 16\sqrt{r}\sigma_1^* \|\Delta_t\|_F \left(2 \|\Delta_t\|_F + 10\Gamma \sqrt{\frac{\sigma_1^* \mu r}{p}} \|\mathbf{E}_t\| \right)^2 \\ &\leq 128\sqrt{r}\sigma_1^* \|\Delta_t\|_F \left(\|\Delta_t\|_F^2 + 25\Gamma^2 \frac{\sigma_1^* \mu r}{p} \|\mathbf{E}_t\|^2 \right).\end{aligned}\quad (80)$$

Similarly, we have

$$\begin{aligned}\left\langle \mathbf{M}_t \mathbf{U}_t \mathbf{U}_t^\top, \mathbf{M}_t \mathbf{U}_t \mathbf{U}_t^\top \mathbf{M}_t \right\rangle &\leq \|\mathbf{M}_t\|^3 \left\| \mathbf{U}_t \mathbf{U}_t^\top \right\|_F^2 \\ &\leq C_1 r \sigma_1^{*2} \left(\|\Delta_t\|_F^3 + \frac{(\sigma_1^* \mu r)^{1.5}}{p^{1.5}} \|\mathbf{E}_t\|^3 \right),\end{aligned}\quad (81)$$

$$\begin{aligned}\left\| \mathbf{M}_t \mathbf{U}_t \mathbf{U}_t^\top + \mathbf{U}_t \mathbf{U}_t^\top \mathbf{M}_t \right\|_F^2 &\leq 4 \|\mathbf{M}_t\|^2 \left\| \mathbf{U}_t \mathbf{U}_t^\top \right\|_F^2 \\ &\leq C_2 r \sigma_1^{*2} \left(\|\Delta_t\|_F^2 + \frac{\sigma_1^* \mu r}{p} \|\mathbf{E}_t\|^2 \right).\end{aligned}\quad (82)$$

$$\left\| \mathbf{M}_t \mathbf{U}_t \mathbf{U}_t^\top \mathbf{M}_t \right\|_F^2 \leq C_3 r \sigma_1^{*2} \left(\|\Delta_t\|_F^4 + \frac{\sigma_1^{*2} \mu^2 r^2}{p^2} \|\mathbf{E}_t\|^4 \right). \quad (83)$$

These inequalities lead to

$$(\text{II}) \leq C_4 \eta^2 r \sigma_1^{*2} \left(\|\Delta_t\|_F^2 + \Gamma^2 \frac{\sigma_1^* \mu r}{p} \|\mathbf{E}_t\|^2 \right). \quad (84)$$

Overall, we have

$$\begin{aligned}\|\Delta_{t+1}\|_F^2 &\leq \left(1 - \frac{1}{5}\eta\sigma_r^*\right)\|\Delta_t\|_F^2 + 324\Gamma\eta\sqrt{\frac{\sigma_1^{*3}\mu r^2}{p}}\|\mathbf{E}_t\|\|\Delta_t\|_F + C_5\eta^2r\frac{\sigma_1^{*3}\mu r}{p}\|\mathbf{E}_t\|^2 \\ &\leq \left(\left(1 - \frac{1}{10}\eta\sigma_r^*\right)\|\Delta_t\|_F + C_6\eta\sqrt{\frac{\sigma_1^{*3}\mu r^2}{p}}\|\mathbf{E}_t\|\right)^2,\end{aligned}\quad (85)$$

which implies that

$$\|\Delta_{t+1}\|_F \leq \left(1 - \frac{1}{10}\eta\sigma_r^*\right)\|\Delta_t\|_F + C_6\eta\sqrt{\frac{\sigma_1^{*3}\mu r^2}{p}}\|\mathbf{E}_t\|. \quad (86)$$

Lastly, we provide the proof of Lemma 15.

Proof of Lemma 15. First, we define $\mathbf{P} = \mathbf{V}^{*\top}\mathbf{V}_t$ and note that $\mathbf{S}_t\mathbf{S}_t^\top = \mathbf{V}_t\boldsymbol{\Sigma}_t\mathbf{V}_t^\top$. This allows us to write

$$\begin{aligned}\left\|\left(\mathbf{X}^* - \mathbf{S}_t\mathbf{S}_t^\top\right)\mathbf{V}_t\right\|_F^2 &= \|\boldsymbol{\Sigma}_t\|_F^2 + \|\boldsymbol{\Sigma}^*\mathbf{P}\|_F^2 - 2\left\langle\boldsymbol{\Sigma}_t, \mathbf{P}^\top\boldsymbol{\Sigma}^*\mathbf{P}\right\rangle, \\ \left\|\mathbf{X}^* - \mathbf{S}_t\mathbf{S}_t^\top\right\|_F^2 &= \|\boldsymbol{\Sigma}_t\|_F^2 + \|\boldsymbol{\Sigma}^*\|_F^2 - 2\left\langle\boldsymbol{\Sigma}_t, \mathbf{P}^\top\boldsymbol{\Sigma}^*\mathbf{P}\right\rangle.\end{aligned}\quad (87)$$

Substituting the above equivalent forms into Equation (75), we need to show that

$$3\|\boldsymbol{\Sigma}_t\|_F^2 + 5\|\boldsymbol{\Sigma}^*\mathbf{P}\|_F^2 \geq 2\|\boldsymbol{\Sigma}^*\|_F^2 + 6\left\langle\boldsymbol{\Sigma}_t, \mathbf{P}^\top\boldsymbol{\Sigma}^*\mathbf{P}\right\rangle. \quad (88)$$

To this end, we first apply the Cauchy-Schwartz inequality, which gives us $2\left\langle\boldsymbol{\Sigma}_t, \mathbf{P}^\top\boldsymbol{\Sigma}^*\mathbf{P}\right\rangle \leq \|\boldsymbol{\Sigma}_t\|_F^2 + \|\mathbf{P}^\top\boldsymbol{\Sigma}^*\mathbf{P}\|_F^2$. Therefore, it suffices to show that

$$5\|\boldsymbol{\Sigma}^*\mathbf{P}\|_F^2 - 2\|\boldsymbol{\Sigma}^*\|_F^2 - 3\left\|\mathbf{P}^\top\boldsymbol{\Sigma}^*\mathbf{P}\right\|_F^2 \geq 0. \quad (89)$$

This follows from

$$5\|\boldsymbol{\Sigma}^*\mathbf{P}\|_F^2 - 2\|\boldsymbol{\Sigma}^*\|_F^2 - 3\left\|\mathbf{P}^\top\boldsymbol{\Sigma}^*\mathbf{P}\right\|_F^2 = \text{tr}\left(\left(\boldsymbol{\Sigma}^*\left(\mathbf{I} - \mathbf{P}\mathbf{P}^\top\right)\boldsymbol{\Sigma}^*\right) \cdot \left(3\mathbf{P}\mathbf{P}^\top - 2\mathbf{I}\right)\right) \geq 0. \quad (90)$$

Here we use the facts that $\boldsymbol{\Sigma}^*\left(\mathbf{I} - \mathbf{P}\mathbf{P}^\top\right)\boldsymbol{\Sigma}^* \succeq 0$ since $\|\mathbf{P}\| \leq \|\mathbf{V}^*\|\|\mathbf{V}_t\| \leq 1$, and $3\mathbf{P}\mathbf{P}^\top - 2\mathbf{I} \succeq 0$ since $\sigma_r(\mathbf{P}) \geq 1 - \|\mathbf{V}^* - \mathbf{V}_t\| \geq 0.9$. This completes the proof. \blacksquare

Appendix D. Proofs for Main Theorems

In this section, we use the one-step dynamics in Proposition 14 to prove our main theorems under the conditions that $\|\mathbf{V}_t\|_{2,\infty} \leq 2\sqrt{\frac{\mu r}{d}}$ and $\|\mathbf{V}^* - \mathbf{V}_t\|_F \leq \Gamma_1 \frac{\kappa\mu r^{1.5} \log(\frac{1}{\alpha})}{\sqrt{pd}}$ for all $0 \leq t \leq T$. These two conditions will be established later in Appendix E.

D.1. Proof of Theorem 2

The proof is divided into three distinct steps.

Step 1. In the first step, we show that $\|\mathbf{S}_t\| \leq 2\sqrt{\sigma_1^*}$ and $\|\mathbf{E}_t\| \leq 2\alpha$ hold for all $0 \leq t \leq T$.

We prove this by induction. First, in the base case where $t = 0$, these two conditions are naturally met because $\|\mathbf{S}_0\| \leq \|\mathbf{U}_0\| \leq \alpha \leq 2\sqrt{\sigma_1^*}$ and $\|\mathbf{E}_0\| \leq 2\alpha$. Next, for the induction step, we assume that $\|\mathbf{S}_t\| \leq 2\sqrt{\sigma_1^*}$ and $\|\mathbf{E}_t\| \leq 2\alpha$ hold for all $0 \leq s \leq t$, with $t \leq T - 1$. Utilizing Proposition 14, we can directly derive that $\|\mathbf{S}_{t+1}\| \leq 2\sqrt{\sigma_1^*}$. Regarding $\|\mathbf{E}_{t+1}\|$, we have

$$\|\mathbf{E}_{t+1}\| \leq \left(1 + 2\Gamma_1\eta \frac{\sigma_1^* \kappa \mu r^{1.5} \log(\frac{1}{\alpha})}{\sqrt{pd}}\right)^{t+1} \stackrel{(a)}{\leq} \left(1 + 4\Gamma_1\eta \frac{\sigma_1^* \kappa \mu r^{1.5} \log(\frac{1}{\alpha})}{\sqrt{pd}} \cdot (t+1)\right) \alpha \leq 2\alpha. \quad (91)$$

Here in (a), we apply Lemma 46. This is valid since $4\Gamma_1\eta \frac{\sigma_1^* \kappa \mu r^{1.5} \log(\frac{1}{\alpha})}{\sqrt{pd}} \cdot (t+1) \leq 1$ for any $t \leq T-1 \lesssim \frac{1}{\eta\sigma_r^*} \log(\frac{1}{\alpha})$ provided that the sampling rate satisfies $p \gtrsim \frac{\kappa^4 \mu^2 r^3 \log^4(\frac{1}{\alpha})}{d}$. This completes the induction step.

Step 2. This step demonstrates that the minimal signal $\sigma_r(\mathbf{S}_t)$ grows linearly to $\frac{\sqrt{\sigma_r^*}}{2}$.

Given that we have already established $\|\mathbf{S}_t\| \leq 2\sqrt{\sigma_1^*}$ and $\|\mathbf{E}_t\| \leq 2\alpha$ for all $0 \leq t \leq T$, we can simplify the minimal signal dynamic in Proposition 14 as

$$\begin{aligned} \sigma_r(\mathbf{S}_{t+1}) &\geq (1 + 0.8\eta\sigma_r^* - \eta\sigma_r^2(\mathbf{S}_t)) \sigma_r(\mathbf{S}_t) - 12\Gamma_1\eta \frac{\sigma_1^* \kappa \mu r^{1.5} \log(\frac{1}{\alpha})}{\sqrt{pd}} \alpha \\ &\geq (1 + 0.4\eta\sigma_r^*) \sigma_r(\mathbf{S}_t) - 12\Gamma_1\eta \frac{\sigma_1^* \kappa \mu r^{1.5} \log(\frac{1}{\alpha})}{\sqrt{pd}} \alpha. \end{aligned} \quad (92)$$

This holds for any t that satisfies $\sigma_r(\mathbf{S}_t) \leq \frac{3\sqrt{\sigma_r^*}}{4}$. By applying Lemma 47, we obtain

$$\sigma_r(\mathbf{S}_t) \geq (1 + 0.4\eta\sigma_r^*)^t \left(\sigma_r(\mathbf{S}_0) - 30\Gamma_1 \frac{\kappa^2 \mu r^{1.5} \log(\frac{1}{\alpha})}{\sqrt{pd}} \alpha \right). \quad (93)$$

At initialization, it is observed that

$$\sigma_r(\mathbf{S}_0) - 30\Gamma_1 \frac{\kappa^2 \mu r^{1.5} \log(\frac{1}{\alpha})}{\sqrt{pd}} \alpha \geq c_0 \alpha - 30\Gamma_1 \frac{\kappa^2 \mu r^{1.5} \log(\frac{1}{\alpha})}{\sqrt{pd}} \alpha \geq \frac{c_0}{2} \alpha, \quad (94)$$

provided that $p \gtrsim \frac{\kappa^4 \mu^2 r^3 \log^2(\frac{1}{\alpha})}{d}$. Consequently, within $T_1 \lesssim \frac{1}{\eta\sigma_r^*} \log\left(\frac{\sigma_r^*}{\alpha}\right)$ iterations, $\sigma_r(\mathbf{S}_t)$ reaches $\frac{\sqrt{\sigma_r^*}}{2}$. It is also easy to show that $\sigma_r(\mathbf{S}_t) \geq \frac{\sqrt{\sigma_r^*}}{2}$ holds true for all $t \geq T_1$.

Step 3. This step is dedicated to demonstrating that the error $\|\mathbf{X}^* - \mathbf{U}_t \mathbf{U}_t^\top\|_F$ converges linearly to $O(\alpha)$ once $\sigma_r(\mathbf{S}_t) \geq \frac{\sqrt{\sigma_r^*}}{2}$.

Based on our arguments in **Step 2**, where we established that $\sigma_r(\mathbf{S}_t) \geq \frac{\sqrt{\sigma_r^*}}{2}$ for $t \geq T_1$, and leveraging Proposition 14 along with Lemma 47, we can derive that

$$\|\mathbf{X}^* - \mathbf{U}_t \mathbf{U}_t^\top\|_F \leq \left(1 - \frac{1}{5}\eta\sigma_r^*\right)^{t-T_1} \|\mathbf{X}^* - \mathbf{U}_{T_1} \mathbf{U}_{T_1}^\top\|_F + 3240\Gamma_1 \sqrt{\frac{\sigma_1^* \kappa^2 \mu r^2}{p}} \alpha. \quad (95)$$

Note that $\|\mathbf{X}^* - \mathbf{U}_{T_1} \mathbf{U}_{T_1}^\top\|_F \leq \|\mathbf{X}^*\|_F + \|\mathbf{U}_{T_1} \mathbf{U}_{T_1}^\top\|_F \leq 9\sqrt{r}\sigma_1^*$ according to Lemma 10. Hence, within an additional $T_2 = \frac{1}{\eta\sigma_r^*} \log\left(\frac{r\sigma_1^*}{\alpha}\right)$ iterations, the error converges to $O\left(\sqrt{\frac{\sigma_1^* \kappa^2 \mu r^2}{p}} \alpha\right)$, thus concluding the proof of Theorem 2.

D.2. Proof of Theorem 3

To prove this result, we first apply Theorem 2 to output a solution U_{t_0} and its leave-one-out versions $U_{t_0}^{(l)}$ that meet the initialization conditions in Theorem 9. Then, we apply Theorem 9 to obtain the desired result.

Establishing Condition (14). By choosing the initialization scale $\alpha = c \cdot \frac{\sigma_r^*}{\kappa^{1.5}d}$ for sufficiently small $c > 0$ and assuming the sampling rate of $p \gtrsim \frac{\kappa^2 \mu^4 r^9 \log^4(d)}{d}$, Theorem 2 guarantees that, with probability at least $1 - \frac{1}{d^2}$, the iterations of GD with step-size $\eta \asymp \frac{\mu r}{\sqrt{pd}\sigma_1^*}$ satisfy

$$\|U_{t_0} U_{t_0}^\top - X^*\|_F \leq 0.9 \Gamma_4 \sqrt{\frac{\sigma_1^* \kappa^2 \mu r^2}{p} \frac{\sigma_r^*}{\kappa^{1.5}d}}, \quad \text{for some } t_0 \lesssim \frac{1}{\eta \sigma_r^*} \log \left(\frac{\kappa^{1.5}d}{\sigma_r^*} \right). \quad (96)$$

On the other hand, Lemma 40 in the appendix implies that

$$\text{dist}(U_{t_0}, U^*) \leq \frac{1.1}{\sigma_r^*} \|U_{t_0} U_{t_0}^\top - X^*\|_F \leq \Gamma_4 \sqrt{\frac{\sigma_r^* \mu^3 r^3 \log(d)}{pd^2}},$$

which establishes Condition (14).

Establishing Condition (15). The proof of Condition (15) follows a similar logic, recognizing that the leave-one-out sequences exhibit a stronger concentration than the original iterations. Consequently, they fulfill Condition (15) within at most $t_0 \lesssim \frac{1}{\eta \sigma_r^*} \log \left(\frac{\kappa^{1.5}d}{\sigma_r^*} \right)$ iterations. Further details of this argument are omitted for brevity.

This shows that the initial conditions of Theorem 9 are satisfied after t_0 iterations. From this iteration onward, Theorem 9 shows that the iterations of GD enter a local linear convergence regime, which readily establishes the final result of Theorem 3.

Appendix E. Proofs for Incoherence Dynamic

In this section, we present our proofs for establishing the incoherence of V_t . To simplify the presentation, we omit the “ \sim ” from our notations. Therefore, $V_t^{(l)}, Z_t^{(l)}, M_t^{(l)}, \dots$ in this section refer to $\tilde{V}_t^{(l)}, \tilde{Z}_t^{(l)}, \tilde{M}_t^{(l)}, \dots$ defined in Section 3.2.

E.1. Proof of Lemma 6

We first state a finer variant of this lemma here.

Proposition 16 (Controlling $\|V_t - V^*\|_F$) *Suppose that the stepsize satisfies $\eta \asymp \frac{\mu r}{\sqrt{pd}\sigma_1^*}$. Moreover, suppose that $\|S_t\| \leq 2\sqrt{\sigma_1^*}$, $\|E_t\| \leq \sqrt{\frac{\sigma_1^*}{d}}$, $\|V_t\|_{2,\infty} \leq 2\sqrt{\frac{\mu r}{d}}$ and $\|V^* - V_t\|_F \leq \Gamma_1 \frac{\kappa \mu r^{1.5} \log(\frac{1}{\alpha})}{\sqrt{pd}}$ for all $t \leq T \lesssim \frac{1}{\eta \sigma_r^*} \log \left(\frac{1}{\alpha} \right)$. Then, conditioned on E_{good} , we have*

$$\|V_t - V^*\|_F \leq \Gamma_1 \frac{\kappa \mu r^{1.5} \log \left(\frac{1}{\alpha} \right)}{\sqrt{pd}} \quad \forall t \leq T \lesssim \frac{1}{\eta \sigma_r^*} \log \left(\frac{1}{\alpha} \right). \quad (97)$$

Proof First notice that \mathbf{V}_{t+1} can be rewritten as

$$\begin{aligned}\mathbf{V}_{t+1} &= \mathbf{Z}_{t+1} (\mathbf{Z}_{t+1}^\top \mathbf{Z}_{t+1})^{-1/2} \\ &= (\mathbf{I} + \eta \mathbf{M}_t) \mathbf{V}_t (\mathbf{V}_t^\top (\mathbf{I} + \eta \mathbf{M}_t)^2 \mathbf{V}_t)^{-1/2} \\ &= (\mathbf{I} + \eta \mathbf{M}_t) \mathbf{V}_t (\mathbf{I} + \mathbf{Y}_t)^{-1/2}\end{aligned}\tag{98}$$

where we denote $\mathbf{Y}_t = \mathbf{V}_t^\top (2\eta \mathbf{M}_t + \eta^2 \mathbf{M}_t^2) \mathbf{V}_t$. Next, we apply Taylor expansion for the matrix-valued function $f(\mathbf{X}) = (\mathbf{I} + \mathbf{X})^{-1/2}$, which states that for any \mathbf{X} satisfying $\|\mathbf{X}\| < 1$,

$$f(\mathbf{X}) = (\mathbf{I} + \mathbf{X})^{-1/2} = \mathbf{I} - \frac{1}{2} \mathbf{X} + \mathbf{R}(\mathbf{X}) \text{ where } \mathbf{R}(\mathbf{X}) = \sum_{k=2}^{\infty} \frac{(-1)^k (2k)!}{4^k (k!)^2} \mathbf{X}^k. \tag{99}$$

Then, upon setting $\mathbf{X} = \mathbf{Y}_t$ in the above equation and plugging it into Equation (98) and rearranging the subterms, we have

$$\mathbf{V}_{t+1} = (\mathbf{I} + \eta \mathbf{M}_t) \mathbf{V}_t \left(\mathbf{I} - \frac{1}{2} \mathbf{V}_t^\top (2\eta \mathbf{M}_t + \eta^2 \mathbf{M}_t^2) \mathbf{V}_t + \mathbf{R}(\mathbf{Y}_t) \right) = (\mathbf{I} + \eta \mathcal{P}_{\mathbf{V}_t}^\perp \mathbf{M}_t) \mathbf{V}_t + \mathbf{A}_t, \tag{100}$$

where

$$\mathbf{A}_t = -\eta^2 \mathbf{M}_t \mathbf{V}_t \mathbf{V}_t^\top \mathbf{M}_t \mathbf{V}_t - 0.5 \eta^2 \mathbf{V}_t \mathbf{V}_t^\top \mathbf{M}_t^2 \mathbf{V}_t - 0.5 \eta^3 \mathbf{M}_t \mathbf{V}_t \mathbf{V}_t^* \mathbf{V}_t^\top \mathbf{M}_t^2 \mathbf{V}_t + (\mathbf{I} + \eta \mathbf{M}_t) \mathbf{V}_t \mathbf{R}(\mathbf{Y}_t) \tag{101}$$

contains all the higher-order terms. Next, according to triangle inequality, we can provide an upper bound for $\|\mathbf{V}^* - \mathbf{V}_{t+1}\|_F$ as follows

$$\|\mathbf{V}^* - \mathbf{V}_{t+1}\|_F \leq \underbrace{\left\| \mathbf{V}^* - (\mathbf{I} + \eta \mathcal{P}_{\mathbf{V}_t}^\perp \mathbf{M}_t) \mathbf{V}_t \right\|_F}_{:= (\text{I})} + \|\mathbf{A}_t\|_F. \tag{102}$$

We first control the leading term (I). To this goal, we apply triangle inequality and obtain

$$\begin{aligned}(\text{I}) &\leq \left\| \mathbf{V}^* - (\mathbf{I} + \eta \mathcal{P}_{\mathbf{V}_t}^\perp (\mathbf{X}^* - \mathbf{U}_t \mathbf{U}_t^\top)) \mathbf{V}_t \right\|_F + \eta \left\| (\mathcal{I} - \mathcal{R}_\Omega) (\mathbf{X}^* - \mathbf{U}_t \mathbf{U}_t^\top) \mathbf{V}_t \right\|_F \\ &\leq \underbrace{\left\| \mathbf{V}^* - (\mathbf{I} + \eta \mathcal{P}_{\mathbf{V}_t}^\perp \Delta_t) \mathbf{V}_t \right\|_F}_{:= (\text{I}_1)} + \eta \sqrt{r} \underbrace{\left\| (\mathcal{I} - \mathcal{R}_\Omega) (\Delta_t) \right\|}_{:= (\text{I}_2)}.\end{aligned}\tag{103}$$

To control (I₁), we further decompose it as

$$\begin{aligned}(\text{I}_1) &\leq \left\| \mathbf{V}^* - (\mathbf{I} + \eta \mathcal{P}_{\mathbf{V}_t}^\perp \mathbf{V}^* \Sigma^* \mathbf{V}^{*\top}) \mathbf{V}_t \right\|_F + \eta \left\| \mathcal{P}_{\mathbf{V}_t}^\perp \mathbf{U}_t \mathbf{U}_t^\top \mathbf{V}_t \right\|_F \\ &\stackrel{(a)}{\leq} \left\| \mathbf{V}^* - \mathbf{V}_t - \eta \mathcal{P}_{\mathbf{V}_t}^\perp (\mathbf{V}^* - \mathbf{V}_t) \Sigma^* \mathbf{V}^{*\top} \mathbf{V}_t \right\|_F + \eta \sqrt{r} \left\| \mathbf{E}_t \mathbf{U}_t^\top \mathbf{V}_t \right\| \\ &\stackrel{(b)}{\leq} \underbrace{\left\| \mathbf{V}^* - \mathbf{V}_t - \eta \mathcal{P}_{\mathbf{V}_t}^\perp (\mathbf{V}^* - \mathbf{V}_t) \Sigma^* \mathbf{V}^{*\top} \mathbf{V}_t \right\|_F}_{:= (\text{I}_{1,1})} + 2\eta \sqrt{r \sigma_1^*} \|\mathbf{E}_t\|.\end{aligned}\tag{104}$$

Here in (a), we use the fact that $\mathcal{P}_{\mathbf{V}_t}^\perp \mathbf{V}_t = 0$ and the definition $\mathbf{E}_t = \mathcal{P}_{\mathbf{V}_t}^\perp \mathbf{U}_t$. In (b), we use the assumption that $\|\mathbf{S}_t\| \leq 2\sqrt{\sigma_1^*}$. Next, according to the orthogonality of \mathbf{V}_t and \mathbf{V}_t^\perp , we can upper-bound $(\mathbf{I}_{1,1})$ as follows

$$\begin{aligned}
 (\mathbf{I}_{1,1})^2 &= \|\mathcal{P}_{\mathbf{V}_t}(\mathbf{V}^* - \mathbf{V}_t)\|_F^2 + \left\| \mathcal{P}_{\mathbf{V}_t}^\perp(\mathbf{V}^* - \mathbf{V}_t) \left(\mathbf{I} - \eta \boldsymbol{\Sigma}^* \mathbf{V}^{*\top} \mathbf{V}_t \right) \right\|_F^2 \\
 &\leq \|\mathcal{P}_{\mathbf{V}_t}(\mathbf{V}^* - \mathbf{V}_t)\|_F^2 + \left\| \mathcal{P}_{\mathbf{V}_t}^\perp(\mathbf{V}^* - \mathbf{V}_t) \right\|_F^2 \left\| \mathbf{I} - \eta \boldsymbol{\Sigma}^* \mathbf{V}^{*\top} \mathbf{V}_t \right\|^2 \\
 &\stackrel{(a)}{\leq} \|\mathcal{P}_{\mathbf{V}_t}(\mathbf{V}^* - \mathbf{V}_t)\|_F^2 + \left\| \mathcal{P}_{\mathbf{V}_t}^\perp(\mathbf{V}^* - \mathbf{V}_t) \right\|_F^2 \\
 &= \|\mathbf{V}^* - \mathbf{V}_t\|_F^2.
 \end{aligned} \tag{105}$$

Here (a) is due to the fact that $\|\mathbf{I} - \eta \boldsymbol{\Sigma}^* \mathbf{V}^{*\top} \mathbf{V}_t\| \leq \|\mathbf{I} - \eta \boldsymbol{\Sigma}^*\| + \eta \sigma_1^* \|\mathbf{V}^* - \mathbf{V}_t\| \leq 1 - \eta(\sigma_r^* - \sigma_1^* \|\mathbf{V}^* - \mathbf{V}_t\|) \leq 1$ since we assume $\|\mathbf{V}^* - \mathbf{V}_t\| \leq \frac{1}{2\kappa}$. Therefore, we derive that

$$(\mathbf{I}_1) \leq \|\mathbf{V}^* - \mathbf{V}_t\|_F + 2\eta\sqrt{r\sigma_1^*} \|\mathbf{E}_t\|. \tag{106}$$

On the other hand, Lemma 12 tells us that, conditioned on E_{good} , we have $(\mathbf{I}_2) \leq 21\Gamma\sigma_1^* \sqrt{\frac{\mu^2 r^2}{pd}}$. Therefore, we can conclude that

$$(\mathbf{I}) \leq \|\mathbf{V}^* - \mathbf{V}_t\|_F + 2\eta\sqrt{r\sigma_1^*} \|\mathbf{E}_t\| + \eta\sqrt{r} \cdot 21\Gamma\sigma_1^* \sqrt{\frac{\mu^2 r^2}{pd}} \leq \|\mathbf{V}^* - \mathbf{V}_t\|_F + 22\Gamma\eta \frac{\sigma_1^* \mu r^{1.5}}{\sqrt{pd}}. \tag{107}$$

Next, according to Lemma 13, we can control $\|\mathbf{A}_t\|_F$ as

$$\|\mathbf{A}_t\|_F \leq \sqrt{r} \|\mathbf{A}_t\| \leq 300\sqrt{r}\eta^2 \sigma_1^{*2}. \tag{108}$$

Putting everything together, we have

$$\begin{aligned}
 \|\mathbf{V}^* - \mathbf{V}_{t+1}\|_F &\leq \|\mathbf{V}^* - \mathbf{V}_t\|_F + 22\Gamma\eta \frac{\sigma_1^* \mu r^{1.5}}{\sqrt{pd}} + 300\sqrt{r}\eta^2 \sigma_1^{*2} \\
 &\leq \|\mathbf{V}^* - \mathbf{V}_t\|_F + 23\Gamma\eta \frac{\sigma_1^* \mu r^{1.5}}{\sqrt{pd}},
 \end{aligned} \tag{109}$$

provided that $\eta \asymp \frac{\mu r}{\sqrt{pd}\sigma_1^*}$. This completes the proof. \blacksquare

E.2. Proof of Proposition 7

We restate the proposition here for clarity.

Proposition 17 (Dynamic of $\left\|(\mathbf{V}^* - \mathbf{V}_t^{(l)})_{l,\cdot}\right\|$) *Under the same conditions as Proposition 16, for all $t \leq T \lesssim \frac{1}{\eta\sigma_r^*} \log\left(\frac{1}{\alpha}\right)$, we have*

$$\left\|(\mathbf{V}^* - \mathbf{V}_{t+1}^{(l)})_{l,\cdot}\right\| \leq (1 - 0.5\eta\sigma_r^*) \left\|(\mathbf{V}^* - \mathbf{V}_t^{(l)})_{l,\cdot}\right\| + \Gamma_2\eta\sigma_1^* \frac{\kappa\mu^{1.5}r^2 \log\left(\frac{1}{\alpha}\right)}{\sqrt{pd^2}}. \tag{110}$$

Proof Similar to Equation (100), we can express $\mathbf{V}_{t+1}^{(l)}$ as

$$\mathbf{V}_{t+1}^{(l)} = \left(\mathbf{I} + \eta \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \mathbf{M}_t^{(l)} \right) \mathbf{V}_t^{(l)} + \mathbf{A}_t^{(l)}, \quad (111)$$

where $\mathbf{M}_t^{(l)} = \mathcal{R}_{\Omega^{(l)}} \left(\mathbf{X}^* - \mathbf{V}_t^{(l)} \boldsymbol{\Sigma}_t \mathbf{V}_t^{(l)\top} \right)$ and $\mathbf{A}_t^{(l)}$ is defined as

$$\begin{aligned} \mathbf{A}_t^{(l)} = & -\eta^2 \mathbf{M}_t^{(l)} \mathbf{V}_t^{(l)} \mathbf{V}_t^{(l)\top} \mathbf{M}_t^{(l)} \mathbf{V}_t^{(l)} - 0.5\eta^2 \mathbf{V}_t^{(l)} \mathbf{V}_t^{(l)\top} \mathbf{M}_t^{(l)2} \mathbf{V}_t^{(l)} \\ & - 0.5\eta^3 \mathbf{M}_t^{(l)} \mathbf{V}_t^{(l)} \mathbf{V}_t^{(l)\top} \mathbf{M}_t^{(l)2} \mathbf{V}_t^{(l)} + (\mathbf{I} + \eta \mathbf{M}_t) \mathbf{V}_t^{(l)} \mathbf{R} \left(\mathbf{Y}_t^{(l)} \right) \end{aligned} \quad (112)$$

containing all the higher-order terms. Applying triangle inequality yields

$$\begin{aligned} \left\| \left(\mathbf{V}^* - \mathbf{V}_{t+1}^{(l)} \right)_{l,\cdot} \right\| & \leq \left\| \left(\mathbf{V}^* - \left(\mathbf{I} + \eta \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \mathbf{M}_t^{(l)} \right) \mathbf{V}_t^{(l)} \right)_{l,\cdot} \right\| + \left\| \left(\mathbf{A}_t^{(l)} \right)_{l,\cdot} \right\| \\ & = \underbrace{\left\| \left(\mathbf{V}^* - \left(\mathbf{I} + \eta \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \boldsymbol{\Xi}_t^{(l)} \right) \mathbf{V}_t^{(l)} \right)_{l,\cdot} \right\|}_{:= (\text{I})} + \left\| \left(\mathbf{A}_t^{(l)} \right)_{l,\cdot} \right\|. \end{aligned} \quad (113)$$

Here in the last equality, we use the fact that $\mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \mathbf{M}_t^{(l)} = \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \boldsymbol{\Xi}_t^{(l)}$ where $\boldsymbol{\Xi}_t^{(l)} = \mathbf{M}_t^{(l)} - \mathbf{V}_t^{(l)} \boldsymbol{\Sigma}^* \mathbf{V}_t^{(l)\top} + \mathbf{V}_t^{(l)} \boldsymbol{\Sigma}_t \mathbf{V}_t^{(l)\top}$. Upon noticing that $\mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp = \mathbf{I} - \mathcal{P}_{\mathbf{V}_t^{(l)}}$, we further decompose (I) as follows

$$(\text{I}) \leq \underbrace{\left\| \left(\mathbf{V}^* - \left(\mathbf{I} + \eta \boldsymbol{\Xi}_t^{(l)} \right) \mathbf{V}_t^{(l)} \right)_{l,\cdot} \right\|}_{:= (\text{I}_1)} + \underbrace{\eta \left\| \left(\mathcal{P}_{\mathbf{V}_t^{(l)}} \boldsymbol{\Xi}_t^{(l)} \mathbf{V}_t^{(l)} \right)_{l,\cdot} \right\|}_{:= (\text{I}_2)}. \quad (114)$$

To control (I₁), notice that the l -th row of $\mathbf{M}_t^{(l)}$ is equal to the l -th row of $\mathbf{X}^* - \mathbf{V}_t^{(l)} \boldsymbol{\Sigma}_t \mathbf{V}_t^{(l)\top}$ due to our choice of $\mathcal{R}_{\Omega^{(l)}}$. Therefore, we have

$$\begin{aligned} (\text{I}_1) & = \left\| \left(\mathbf{V}^* - \left(\mathbf{I} + \eta \left(\mathbf{X}^* - \mathbf{V}_t^{(l)} \boldsymbol{\Sigma}_t \mathbf{V}_t^{(l)\top} - \mathbf{V}_t^{(l)} \boldsymbol{\Sigma}^* \mathbf{V}_t^{(l)\top} - \mathbf{V}_t^{(l)} \boldsymbol{\Sigma}_t \mathbf{V}_t^{(l)\top} \right) \right) \mathbf{V}_t^{(l)} \right)_{l,\cdot} \right\| \\ & = \left\| \left(\mathbf{V}^* - \left(\mathbf{I} + \eta \left(\left(\mathbf{V}^* - \mathbf{V}_t^{(l)} \right) \boldsymbol{\Sigma}^* \mathbf{V}_t^{(l)\top} \right) \right) \mathbf{V}_t^{(l)} \right)_{l,\cdot} \right\| \\ & = \left\| \left(\left(\mathbf{V}^* - \mathbf{V}_t^{(l)} \right) \left(\mathbf{I} - \eta \boldsymbol{\Sigma}^* \mathbf{V}_t^{(l)\top} \mathbf{V}_t^{(l)} \right) \right)_{l,\cdot} \right\| \\ & \leq \left\| \left(\mathbf{V}^* - \mathbf{V}_t^{(l)} \right)_{l,\cdot} \right\| \left\| \mathbf{I} - \eta \boldsymbol{\Sigma}^* \mathbf{V}_t^{(l)\top} \mathbf{V}_t^{(l)} \right\| \\ & \stackrel{(a)}{\leq} (1 - 0.5\eta\sigma_r^*) \left\| \left(\mathbf{V}^* - \mathbf{V}_t^{(l)} \right)_{l,\cdot} \right\|. \end{aligned} \quad (115)$$

Here in (a), we use the fact that $\left\| \mathbf{I} - \eta \boldsymbol{\Sigma}^* \mathbf{V}_t^{(l)\top} \mathbf{V}_t^{(l)} \right\| \leq \left\| \mathbf{I} - \eta \boldsymbol{\Sigma}^* \right\| + \eta \sigma_1^* \left\| \mathbf{V}^* - \mathbf{V}_t^{(l)} \right\| \leq 1 - \eta(\sigma_r^* - \sigma_1^*) \left\| \mathbf{V}^* - \mathbf{V}_t^{(l)} \right\| \leq 1 - 0.5\eta\sigma_r^*$ since we have $\left\| \mathbf{V}^* - \mathbf{V}_t^{(l)} \right\| \leq \frac{1}{2\kappa}$ according to the following proposition.

Lemma 18 (Frobenius norm control) *Under the same conditions as Proposition 16, for any $1 \leq t \leq T = \frac{100}{\eta\sigma_1^*} \log\left(\frac{1}{\alpha}\right)$, for all $1 \leq l \leq d$, we have*

$$\left\| \mathbf{V}^* - \mathbf{V}_t^{(l)} \right\|_{\text{F}} \leq \Gamma_1 \frac{\kappa\mu r^{1.5} \log\left(\frac{1}{\alpha}\right)}{\sqrt{pd}}. \quad (116)$$

The proof of the above lemma is the same as that of Proposition 16 and hence omitted here. Next, for (I₂), we first have

$$(I_2) = \left\| \left(\mathbf{V}_t^{(l)} \mathbf{V}_t^{(l)\top} \boldsymbol{\Xi}_t^{(l)} \mathbf{V}_t^{(l)} \right)_{l,\cdot} \right\| \leq \left\| \left(\mathbf{V}_t^{(l)} \right)_{l,\cdot} \right\| \left\| \mathbf{V}_t^{(l)\top} \boldsymbol{\Xi}_t^{(l)} \mathbf{V}_t^{(l)} \right\| \leq \left\| \left(\mathbf{V}_t^{(l)} \right)_{l,\cdot} \right\| \left\| \boldsymbol{\Xi}_t^{(l)} \right\|. \quad (117)$$

For the first part, we have

$$\left\| \left(\mathbf{V}_t^{(l)} \right)_{l,\cdot} \right\| \leq \left\| \mathbf{V}_{l,\cdot}^* \right\| + \left\| \left(\mathbf{V}^* - \mathbf{V}_t^{(l)} \right)_{l,\cdot} \right\| \leq \sqrt{\frac{\mu r}{d}} + \left\| \left(\mathbf{V}^* - \mathbf{V}_t^{(l)} \right)_{l,\cdot} \right\| \leq 2\sqrt{\frac{\mu r}{d}}. \quad (118)$$

Here we use the assumption that $\left\| \left(\mathbf{V}^* - \mathbf{V}_t^{(l)} \right)_{l,\cdot} \right\| \leq \sqrt{\frac{\mu r}{d}}$. For the second part, we have

$$\begin{aligned} \left\| \boldsymbol{\Xi}_t^{(l)} \right\| &= \left\| \left(\mathbf{V}^* - \mathbf{V}_t^{(l)} \right) \boldsymbol{\Sigma}^* \mathbf{V}^{*\top} - (\mathcal{I} - \mathcal{R}_{\Omega^{(l)}}) \left(\mathbf{X}^* - \mathbf{V}_t^{(l)} \boldsymbol{\Sigma}_t \mathbf{V}_t^{(l)\top} \right) \right\| \\ &\leq \left\| \left(\mathbf{V}^* - \mathbf{V}_t^{(l)} \right) \boldsymbol{\Sigma}^* \mathbf{V}^{*\top} \right\| + \left\| (\mathcal{I} - \mathcal{R}_{\Omega^{(l)}}) \left(\mathbf{X}^* - \mathbf{V}_t^{(l)} \boldsymbol{\Sigma}_t \mathbf{V}_t^{(l)\top} \right) \right\| \\ &\stackrel{(a)}{\leq} \sigma_1^* \left\| \mathbf{V}^* - \mathbf{V}_t^{(l)} \right\| + 21\Gamma \frac{\sigma_1^* \mu r}{\sqrt{pd}} \\ &\stackrel{(b)}{\leq} \sigma_1^* \cdot \Gamma_1 \frac{\kappa\mu r^{1.5} \log\left(\frac{1}{\alpha}\right)}{\sqrt{pd}} + 21\Gamma \frac{\sigma_1^* \mu r}{\sqrt{pd}} \\ &\leq 2\Gamma_1 \sigma_1^* \frac{\kappa\mu r^{1.5} \log\left(\frac{1}{\alpha}\right)}{\sqrt{pd}}. \end{aligned} \quad (119)$$

Here in (a), we apply Lemma 12. In (b), we apply Lemma 18. Invoking Equations (118) and (119) in Equation (117), we obtain that

$$(I_2) \leq 2\sqrt{\frac{\mu r}{d}} \cdot 2\Gamma_1 \sigma_1^* \frac{\kappa\mu r^{1.5} \log\left(\frac{1}{\alpha}\right)}{\sqrt{pd}} = 4\Gamma_1 \sigma_1^* \frac{\kappa\mu^{1.5} r^2 \log\left(\frac{1}{\alpha}\right)}{\sqrt{pd^2}}. \quad (120)$$

Next, we control the higher-order term $\left\| \left(\mathbf{A}_t^{(l)} \right)_{l,\cdot} \right\|$. To this end, we first notice that the l -th row of $\mathbf{A}_t^{(l)}$ is equal to the l -th row of

$$\begin{aligned} & -\eta^2 \boldsymbol{\Delta}_t^{(l)} \mathcal{P}_{\mathbf{V}_t^{(l)}} \mathbf{V}_t^{(l)\top} \mathbf{M}_t^{(l)} \mathbf{V}_t^{(l)} - 0.5\eta^2 \mathcal{P}_{\mathbf{V}_t^{(l)}} \mathbf{M}_t^{(l)2} \mathbf{V}_t^{(l)} - 0.5\eta^3 \boldsymbol{\Delta}_t^{(l)} \mathcal{P}_{\mathbf{V}_t^{(l)}} \mathbf{M}_t^{(l)2} \mathbf{V}_t^{(l)} \\ & + \left(\mathbf{I} + \eta \boldsymbol{\Delta}_t^{(l)} \right) \mathbf{V}_t^{(l)} \mathbf{R} \left(\mathbf{Y}_t^{(l)} \right). \end{aligned} \quad (121)$$

Therefore, we can upper-bound its operator norm by

$$\begin{aligned} \left\| \left(\mathbf{A}_t^{(l)} \right)_{l,\cdot} \right\| &\leq \left\| \left(\boldsymbol{\Delta}_t^{(l)} \right)_{l,\cdot} \right\| \left(\eta^2 \left\| \mathbf{M}_t^{(l)} \right\| + 0.5\eta^3 \left\| \mathbf{M}_t^{(l)} \right\|^2 + \eta \left\| \mathbf{R} \left(\mathbf{Y}_t^{(l)} \right) \right\| \right) \\ &\quad + \left\| \left(\mathbf{V}_t^{(l)} \right)_{l,\cdot} \right\| \left(0.5\eta^2 \left\| \mathbf{M}_t^{(l)} \right\|^2 + \left\| \mathbf{R} \left(\mathbf{Y}_t^{(l)} \right) \right\| \right). \end{aligned} \quad (122)$$

Next, we notice that $\left\| \left(\mathbf{V}_t^{(l)} \right)_{l,\cdot} \right\| \leq 2\sqrt{\frac{\mu r}{d}}$ and

$$\begin{aligned} \left\| \left(\Delta_t^{(l)} \right)_{l,\cdot} \right\| &\leq \left\| \left(\mathbf{X}^* - \mathbf{V}_t^{(l)} \Sigma_t \mathbf{V}_t^{(l)\top} \right)_{l,\cdot} \right\| \\ &\leq \sigma_1^* \sqrt{\frac{\mu r}{d}} + \left\| \left(\mathbf{V}_t^{(l)} \right)_{l,\cdot} \right\| \|\Sigma_t\| \\ &\leq 5\sigma_1^* \sqrt{\frac{\mu r}{d}}. \end{aligned} \quad (123)$$

On the other hand, according to Lemma 13, we have $\left\| \mathbf{M}_t^{(l)} \right\| \leq 6\sqrt{\sigma_1^*}$ and $\left\| \mathbf{R} \left(\mathbf{Y}_t^{(l)} \right) \right\| \leq 200\eta^2 \sigma_1^{*2}$. Therefore, we derive that

$$\left\| \left(\mathbf{A}_t^{(l)} \right)_{l,\cdot} \right\| \leq 470\eta^2 \sigma_1^{*2} \sqrt{\frac{\mu r}{d}}. \quad (124)$$

Overall, we obtain

$$\begin{aligned} \left\| \left(\mathbf{V}^* - \mathbf{V}_{t+1}^{(l)} \right)_{l,\cdot} \right\| &\leq (1 - 0.5\eta\sigma_r^*) \left\| \left(\mathbf{V}^* - \mathbf{V}_t^{(l)} \right)_{l,\cdot} \right\| + 4\Gamma_1\eta\sigma_1^* \frac{\kappa\mu^{1.5}r^2 \log\left(\frac{1}{\alpha}\right)}{\sqrt{pd^2}} + 470\eta^2 \sigma_1^{*2} \sqrt{\frac{\mu r}{d}} \\ &\leq (1 - 0.5\eta\sigma_r^*) \left\| \left(\mathbf{V}^* - \mathbf{V}_t^{(l)} \right)_{l,\cdot} \right\| + 8\Gamma_1\eta\sigma_1^* \frac{\kappa\mu^{1.5}r^2 \log\left(\frac{1}{\alpha}\right)}{\sqrt{pd^2}}, \end{aligned} \quad (125)$$

provided that $\eta \asymp \frac{\mu r}{\sqrt{pd}\sigma_1^*}$. This completes the proof of Proposition 17. \blacksquare

E.3. Proof of Proposition 8

We restate the proposition here for clarity.

Proposition 19 (One-step dynamic of $\|\mathbf{V}_t - \mathbf{V}_t^{(l)}\|_F$) Suppose that the sampling rate satisfies $p \gtrsim \frac{\kappa^6 \mu^4 r^9 \log^4(\frac{1}{\alpha}) \log^2(d)}{d}$. Suppose that $\|\mathbf{V}^* - \mathbf{V}_t\| \leq \frac{1}{2\kappa}$ and $\|\mathbf{V}_t - \mathbf{V}_t^{(l)}\|_F \leq \sqrt{\frac{\mu r}{4d}}$. With probability at least $1 - \frac{1}{d^3}$, for any $1 \leq t \leq T \lesssim \frac{1}{\eta\sigma_r^*} \log\left(\frac{1}{\alpha}\right)$, we have

$$\left\| \mathbf{V}_{t+1} - \mathbf{V}_{t+1}^{(l)} \right\|_F \leq \left\| \mathbf{V}_t - \mathbf{V}_t^{(l)} \right\|_F + \Gamma\eta\sigma_1^* \sqrt{\frac{\kappa\mu^3 r^{5.5} \log\left(\frac{1}{\alpha}\right) \log(d)}{\sqrt{pd} \cdot d}}. \quad (126)$$

Proof Note that

$$\mathbf{V}_{t+1} = \left(\mathbf{I} + \eta \mathcal{P}_{\mathbf{V}_t}^\perp \mathbf{M}_t \right) \mathbf{V}_t + \mathbf{A}_t \quad \text{and} \quad \mathbf{V}_{t+1}^{(l)} = \left(\mathbf{I} + \eta \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \mathbf{M}_t^{(l)} \right) \mathbf{V}_t^{(l)} + \mathbf{A}_t^{(l)}. \quad (127)$$

Hence, we can expand $\left\| \mathbf{V}_{t+1} - \mathbf{V}_{t+1}^{(l)} \right\|_F^2$ as

$$\begin{aligned} \left\| \mathbf{V}_{t+1} - \mathbf{V}_{t+1}^{(l)} \right\|_F^2 &= \left\| \mathbf{V}_t - \mathbf{V}_t^{(l)} \right\|_F^2 + 2\eta \underbrace{\left\langle \mathbf{V}_t - \mathbf{V}_t^{(l)}, \mathcal{P}_{\mathbf{V}_t}^\perp \mathbf{M}_t \mathbf{V}_t - \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \mathbf{M}_t^{(l)} \mathbf{V}_t^{(l)} \right\rangle}_{:=\text{(I)}} + \text{(II)}, \end{aligned} \quad (128)$$

where

$$\begin{aligned}
 (\text{II}) &= \eta^2 \left\| \mathcal{P}_{\mathbf{V}_t}^\perp \mathbf{M}_t \mathbf{V}_t - \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \mathbf{M}_t^{(l)} \mathbf{V}_t^{(l)} \right\|_F^2 + \left\| \mathbf{A}_t - \mathbf{A}_t^{(l)} \right\|_F^2 \\
 &\quad + 2 \left\langle \left(\mathbf{I} + \eta \mathcal{P}_{\mathbf{V}_t}^\perp \mathbf{M}_t \right) \mathbf{V}_t - \left(\mathbf{I} + \eta \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \mathbf{M}_t^{(l)} \right) \mathbf{V}_t^{(l)}, \mathbf{A}_t - \mathbf{A}_t^{(l)} \right\rangle
 \end{aligned} \tag{129}$$

contains all the higher-order terms.

We first control (I). Notice that $\mathcal{P}_{\mathbf{V}_t}^\perp \mathbf{M}_t = \mathcal{P}_{\mathbf{V}_t}^\perp \boldsymbol{\Xi}_t$ and $\mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \mathbf{M}_t^{(l)} = \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \boldsymbol{\Xi}_t^{(l)}$. Here we define $\boldsymbol{\Xi}_t = \mathbf{M}_t - \mathbf{V}_t \boldsymbol{\Sigma}^* \mathbf{V}_t^{\star\top} + \mathbf{V}_t \boldsymbol{\Sigma}_t \mathbf{V}_t^\top$ and $\boldsymbol{\Xi}_t^{(l)} = \mathbf{M}_t^{(l)} - \mathbf{V}_t^{(l)} \boldsymbol{\Sigma}^* \mathbf{V}_t^{\star\top} + \mathbf{V}_t^{(l)} \boldsymbol{\Sigma}_t \mathbf{V}_t^{(l)\top}$, respectively. Therefore, we can decompose (I) as follows

$$\begin{aligned}
 (\text{I}) &= \left\langle \mathbf{V}_t - \mathbf{V}_t^{(l)}, \mathcal{P}_{\mathbf{V}_t}^\perp \boldsymbol{\Xi}_t \mathbf{V}_t - \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \boldsymbol{\Xi}_t^{(l)} \mathbf{V}_t^{(l)} \right\rangle \\
 &= \underbrace{\left\langle \mathbf{V}_t - \mathbf{V}_t^{(l)}, \left(\mathcal{P}_{\mathbf{V}_t}^\perp - \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \right) \boldsymbol{\Xi}_t^{(l)} \mathbf{V}_t^{(l)} \right\rangle}_{:= (\text{I}_1)} + \underbrace{\left\langle \mathbf{V}_t - \mathbf{V}_t^{(l)}, \mathcal{P}_{\mathbf{V}_t}^\perp \boldsymbol{\Xi}_t \left(\mathbf{V}_t - \mathbf{V}_t^{(l)} \right) \right\rangle}_{:= (\text{I}_2)} \\
 &\quad + \underbrace{\left\langle \mathbf{V}_t - \mathbf{V}_t^{(l)}, \mathcal{P}_{\mathbf{V}_t}^\perp \left(\boldsymbol{\Xi}_t - \boldsymbol{\Xi}_t^{(l)} \right) \mathbf{V}_t^{(l)} \right\rangle}_{:= (\text{I}_3)}.
 \end{aligned} \tag{130}$$

Next, we provide upper bounds for these terms separately. For (I₁), applying Cauchy-Schwarz inequality leads to

$$\begin{aligned}
 (\text{I}_1) &\leq \left\| \mathbf{V}_t - \mathbf{V}_t^{(l)} \right\|_F \left\| \left(\mathcal{P}_{\mathbf{V}_t}^\perp - \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \right) \boldsymbol{\Xi}_t^{(l)} \mathbf{V}_t^{(l)} \right\|_F \\
 &\leq \left\| \mathbf{V}_t - \mathbf{V}_t^{(l)} \right\|_F \left\| \mathcal{P}_{\mathbf{V}_t}^\perp - \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \right\|_F \left\| \boldsymbol{\Xi}_t^{(l)} \right\| \\
 &\stackrel{(a)}{\leq} 2 \left\| \mathbf{V}_t - \mathbf{V}_t^{(l)} \right\|_F^2 \left\| \boldsymbol{\Xi}_t^{(l)} \right\| \\
 &\stackrel{(b)}{\leq} 4 \left\| \mathbf{V}_t - \mathbf{V}_t^{(l)} \right\|_F^2 \cdot \Gamma_1 \sigma_1^* \frac{\kappa \mu r^{1.5} \log \left(\frac{1}{\alpha} \right)}{\sqrt{pd}}.
 \end{aligned} \tag{131}$$

Here in (a), we apply Lemma 43. In (b), we use the result from Equation (119). Similarly, (I₂), we have

$$(\text{I}_2) \leq \left\| \mathbf{V}_t - \mathbf{V}_t^{(l)} \right\|_F^2 \left\| \boldsymbol{\Xi}_t \right\| \leq 2 \left\| \mathbf{V}_t - \mathbf{V}_t^{(l)} \right\|_F^2 \cdot \Gamma_1 \sigma_1^* \frac{\kappa \mu r^{1.5} \log \left(\frac{1}{\alpha} \right)}{\sqrt{pd}}. \tag{132}$$

Next, for (I_3) , we further decompose it as

$$\begin{aligned}
 (I_3) &= \left\langle \mathbf{V}_t - \mathbf{V}_t^{(l)}, \mathcal{P}_{\mathbf{V}_t}^\perp \left(\mathbf{M}_t - \mathbf{M}_t^{(l)} + \mathbf{V}_t^{(l)} \boldsymbol{\Sigma}^* \mathbf{V}^{\star\top} - \mathbf{V}_t^{(l)} \boldsymbol{\Sigma}_t \mathbf{V}_t^{(l)\top} \right) \mathbf{V}_t^{(l)} \right\rangle \\
 &= \underbrace{\left\langle \mathbf{V}_t - \mathbf{V}_t^{(l)}, \mathcal{P}_{\mathbf{V}_t}^\perp (\mathcal{R}_\Omega - \mathcal{R}_{\Omega^{(l)}}) \left(\mathbf{X}^* - \mathbf{V}_t^{(l)} \boldsymbol{\Sigma}_t \mathbf{V}_t^{(l)\top} \right) \mathbf{V}_t^{(l)} \right\rangle}_{:= (I_{3,1})} \\
 &\quad + \underbrace{\left\langle \mathbf{V}_t - \mathbf{V}_t^{(l)}, \mathcal{P}_{\mathbf{V}_t}^\perp (\mathcal{I} - \mathcal{R}_\Omega) \left(\mathbf{V}_t \boldsymbol{\Sigma}_t \left(\mathbf{V}_t - \mathbf{V}_t^{(l)} \right)^\top \right) \mathbf{V}_t^{(l)} \right\rangle}_{:= (I_{3,2})} \\
 &\quad + \underbrace{\left\langle \mathbf{V}_t - \mathbf{V}_t^{(l)}, \mathcal{P}_{\mathbf{V}_t}^\perp (\mathcal{I} - \mathcal{R}_\Omega) \left(\left(\mathbf{V}_t - \mathbf{V}_t^{(l)} \right) \boldsymbol{\Sigma}_t \mathbf{V}_t^{(l)\top} \right) \mathbf{V}_t^{(l)} \right\rangle}_{:= (I_{3,3})} \\
 &\quad + \underbrace{\left\langle \mathbf{V}_t - \mathbf{V}_t^{(l)}, \mathcal{P}_{\mathbf{V}_t}^\perp \mathbf{V}_t^{(l)} \boldsymbol{\Sigma}^* \mathbf{V}^{\star\top} \mathbf{V}_t^{(l)} \right\rangle}_{:= (I_{3,4})}.
 \end{aligned} \tag{133}$$

We control these terms separately. First, we present the following key proposition, the proof of which is deferred to the end of this section.

Proposition 20 *For all $0 \leq t \leq T$, we have*

$$\left\| (\mathcal{R}_\Omega - \mathcal{R}_{\Omega^{(l)}}) \left(\mathbf{X}^* - \mathbf{V}_t^{(l)} \boldsymbol{\Sigma}_t \mathbf{V}_t^{(l)\top} \right) \mathbf{V}_t^{(l)} \right\|_F \leq \Gamma_3 \sigma_1^* \sqrt{\frac{\kappa \mu^3 r^{5.5} \log\left(\frac{1}{\alpha}\right) \log(d)}{\sqrt{pd} \cdot d}}. \tag{134}$$

For $(I_{3,1})$, we have

$$\begin{aligned}
 (I_{3,1}) &\leq \left\| \mathbf{V}_t - \mathbf{V}_t^{(l)} \right\|_F \left\| (\mathcal{R}_\Omega - \mathcal{R}_{\Omega^{(l)}}) \left(\mathbf{X}^* - \mathbf{V}_t^{(l)} \boldsymbol{\Sigma}_t \mathbf{V}_t^{(l)\top} \right) \mathbf{V}_t^{(l)} \right\|_F \\
 &\leq \Gamma_3 \sigma_1^* \sqrt{\frac{\kappa \mu^3 r^{5.5} \log\left(\frac{1}{\alpha}\right) \log(d)}{\sqrt{pd} \cdot d}} \left\| \mathbf{V}_t - \mathbf{V}_t^{(l)} \right\|_F.
 \end{aligned} \tag{135}$$

Here in the last inequality, we use Proposition 20. Next, we apply Lemma 29 to control $(I_{3,2})$. Specifically, upon setting $\mathbf{A} = \mathbf{V}_t$, $\mathbf{B} = \mathcal{P}_{\mathbf{V}_t}^\perp (\mathbf{V}_t - \mathbf{V}_t^{(l)})$, $\mathbf{C} = (\mathbf{V}_t - \mathbf{V}_t^{(l)}) \boldsymbol{\Sigma}_t$, and $\mathbf{D} = \mathbf{V}_t^{(l)}$ in Lemma 29, with probability at least $1 - d^{-3}$, one has

$$\begin{aligned}
 (I_{3,2}) &= \left\langle (\mathcal{I} - \mathcal{R}_\Omega) (\mathbf{A} \mathbf{C}^\top), \mathbf{B} \mathbf{D}^\top \right\rangle \\
 &\leq \Gamma \sqrt{\frac{d}{p}} \|\mathbf{A}\|_{2,\infty} \|\mathbf{B}\|_F \cdot \|\mathbf{C}\|_F \|\mathbf{D}\|_{2,\infty} \\
 &= \Gamma \sqrt{\frac{d}{p}} \|\mathbf{V}_t\|_{2,\infty} \left\| \mathbf{V}_t^{(l)} \right\|_{2,\infty} \left\| \mathcal{P}_{\mathbf{V}_t}^\perp (\mathbf{V}_t - \mathbf{V}_t^{(l)}) \right\|_F \left\| (\mathbf{V}_t - \mathbf{V}_t^{(l)}) \boldsymbol{\Sigma}_t \right\|_F \\
 &\leq 16 \Gamma \sigma_1^* \sqrt{\frac{\mu^2 r^2}{pd}} \left\| \mathbf{V}_t - \mathbf{V}_t^{(l)} \right\|_F^2.
 \end{aligned} \tag{136}$$

Similarly, we apply Lemma 29 to control $(I_{3,3})$. Upon setting $\mathbf{A} = (\mathbf{V}_t - \mathbf{V}_t^{(l)}) \boldsymbol{\Sigma}_t$, $\mathbf{B} = \mathcal{P}_{\mathbf{V}_t}^\perp (\mathbf{V}_t - \mathbf{V}_t^{(l)})$, $\mathbf{C} = \mathbf{V}_t$, and $\mathbf{D} = \mathbf{V}_t^{(l)}$ in Lemma 29, with probability at least $1 - d^{-3}$, one has

$$\begin{aligned} (I_{3,3}) &\leq \Gamma \sqrt{\frac{d}{p}} \|\mathbf{V}_t\|_F \|\mathbf{V}_t^{(l)}\|_{2,\infty} \|\mathcal{P}_{\mathbf{V}_t}^\perp (\mathbf{V}_t - \mathbf{V}_t^{(l)})\|_F \|(\mathbf{V}_t - \mathbf{V}_t^{(l)}) \boldsymbol{\Sigma}_t\|_{2,\infty} \\ &\leq 8\Gamma\sigma_1^* \sqrt{\frac{\mu r^2}{p}} \|\mathbf{V}_t - \mathbf{V}_t^{(l)}\|_F \|\mathbf{V}_t - \mathbf{V}_t^{(l)}\|_{2,\infty}. \end{aligned} \quad (137)$$

To control $\|\mathbf{V}_t - \mathbf{V}_t^{(l)}\|_{2,\infty}$, we apply Lemma 39:

$$\begin{aligned} \|\mathbf{V}_t - \mathbf{V}_t^{(l)}\|_{2,\infty} &\leq \|\mathbf{V}_t - \mathbf{V}_t^{(l)}\|_F \|\mathbf{L}_{\mathbf{V}_t - \mathbf{V}_t^{(l)}}\|_{2,\infty} \\ &\leq \|\mathbf{V}_t - \mathbf{V}_t^{(l)}\|_F \sqrt{\|\mathbf{L}_{\mathbf{V}_t}\|_{2,\infty}^2 + \|\mathbf{L}_{\mathbf{V}_t^{(l)}}\|_{2,\infty}^2} \\ &\leq 3\sqrt{\frac{\mu r}{d}} \|\mathbf{V}_t - \mathbf{V}_t^{(l)}\|_F. \end{aligned} \quad (138)$$

Hence, we have

$$(I_{3,3}) \leq 24\Gamma\sigma_1^* \sqrt{\frac{\mu^2 r^3}{pd}} \|\mathbf{V}_t - \mathbf{V}_t^{(l)}\|_F^2. \quad (139)$$

Lastly, for $(I_{3,4})$, we notice that

$$\begin{aligned} (I_{3,4}) &= -\left\langle \mathbf{V}_t - \mathbf{V}_t^{(l)}, \mathcal{P}_{\mathbf{V}_t}^\perp (\mathbf{V}_t - \mathbf{V}_t^{(l)}) \boldsymbol{\Sigma}^* \mathbf{V}^{*\top} \mathbf{V}_t^{(l)} \right\rangle \\ &\leq -\left\langle \mathbf{V}_t - \mathbf{V}_t^{(l)}, \mathcal{P}_{\mathbf{V}_t}^\perp (\mathbf{V}_t - \mathbf{V}_t^{(l)}) \boldsymbol{\Sigma}^* \right\rangle + \sigma_1^* \|\mathbf{V}^* - \mathbf{V}_t^{(l)}\|_F \|\mathbf{V}_t - \mathbf{V}_t^{(l)}\|_F^2 \\ &\stackrel{(a)}{\leq} \sigma_1^* \|\mathbf{V}^* - \mathbf{V}_t^{(l)}\|_F \|\mathbf{V}_t - \mathbf{V}_t^{(l)}\|_F^2 \\ &\stackrel{(b)}{\leq} \Gamma_1 \sigma_1^* \frac{\kappa \mu r^{1.5} \log(\frac{1}{\alpha})}{\sqrt{pd}} \|\mathbf{V}_t - \mathbf{V}_t^{(l)}\|_F^2. \end{aligned} \quad (140)$$

Here in (a), we use the fact that

$$\left\langle \mathbf{V}_t - \mathbf{V}_t^{(l)}, \mathcal{P}_{\mathbf{V}_t}^\perp (\mathbf{V}_t - \mathbf{V}_t^{(l)}) \boldsymbol{\Sigma}^* \right\rangle = \left\| \mathcal{P}_{\mathbf{V}_t}^\perp (\mathbf{V}_t - \mathbf{V}_t^{(l)}) \boldsymbol{\Sigma}^{*1/2} \right\|_F^2 \geq 0. \quad (141)$$

In (b), we apply Proposition 16. Therefore, we obtain

$$(I_3) \leq C_1 \sigma_1^* \frac{\kappa \mu r^{1.5} \log(\frac{1}{\alpha})}{\sqrt{pd}} \|\mathbf{V}_t - \mathbf{V}_t^{(l)}\|_F^2 + \Gamma_3 \sigma_1^* \sqrt{\frac{\kappa \mu^3 r^{5.5} \log(\frac{1}{\alpha}) \log(d)}{\sqrt{pd} \cdot d}} \|\mathbf{V}_t - \mathbf{V}_t^{(l)}\|_F, \quad (142)$$

which implies that

$$(I) \leq C_2 \sigma_1^* \frac{\kappa \mu r^{1.5} \log(\frac{1}{\alpha})}{\sqrt{pd}} \|\mathbf{V}_t - \mathbf{V}_t^{(l)}\|_F^2 + \Gamma_3 \sigma_1^* \sqrt{\frac{\kappa \mu^3 r^{5.5} \log(\frac{1}{\alpha}) \log(d)}{\sqrt{pd} \cdot d}} \|\mathbf{V}_t - \mathbf{V}_t^{(l)}\|_F. \quad (143)$$

Next, we move on to controlling (II). First, by the basic inequality $2ab \leq a^2 + b^2$, we bound (II) as

$$(II) \leq 2\eta^2 \left\| \mathcal{P}_{\mathbf{V}_t}^\perp \mathbf{M}_t \mathbf{V}_t - \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \mathbf{M}_t^{(l)} \mathbf{V}_t^{(l)} \right\|_F^2 + 2 \left\| \mathbf{A}_t - \mathbf{A}_t^{(l)} \right\|_F^2 + 2 \left\langle \mathbf{V}_t - \mathbf{V}_t^{(l)}, \mathbf{A}_t - \mathbf{A}_t^{(l)} \right\rangle. \quad (144)$$

Next, we provide the control over $\left\| \mathcal{P}_{\mathbf{V}_t}^\perp \mathbf{M}_t \mathbf{V}_t - \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \mathbf{M}_t^{(l)} \mathbf{V}_t^{(l)} \right\|_F$. The remaining two terms can be controlled in a similar fashion, so we omit their analysis. We first apply triangle inequality to obtain

$$\begin{aligned} \left\| \mathcal{P}_{\mathbf{V}_t}^\perp \mathbf{M}_t \mathbf{V}_t - \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \mathbf{M}_t^{(l)} \mathbf{V}_t^{(l)} \right\|_F &\leq 2 \left\| \mathbf{M}_t \right\| \left\| \mathbf{V}_t - \mathbf{V}_t^{(l)} \right\|_F + \left\| \left(\mathbf{M}_t - \mathbf{M}_t^{(l)} \right) \mathbf{V}_t^{(l)} \right\|_F \\ &\leq 12\sigma_1^* \left\| \mathbf{V}_t - \mathbf{V}_t^{(l)} \right\|_F + \left\| \left(\mathbf{M}_t - \mathbf{M}_t^{(l)} \right) \mathbf{V}_t^{(l)} \right\|_F. \end{aligned} \quad (145)$$

For the second term, we further decompose it as

$$\begin{aligned} \left\| \left(\mathbf{M}_t - \mathbf{M}_t^{(l)} \right) \mathbf{V}_t^{(l)} \right\|_F &\leq \left\| (\mathcal{R}_\Omega - \mathcal{R}_{\Omega^{(l)}}) \left(\mathbf{X}^* - \mathbf{V}_t^{(l)} \boldsymbol{\Sigma}_t \mathbf{V}_t^{(l)\top} \right) \mathbf{V}_t^{(l)} \right\|_F \\ &\quad + \left\| (\mathcal{I} - \mathcal{R}_\Omega) \left(\mathbf{V}_t \boldsymbol{\Sigma}_t \left(\mathbf{V}_t - \mathbf{V}_t^{(l)} \right)^\top \right) \mathbf{V}_t \right\|_F \\ &\quad + \left\| (\mathcal{I} - \mathcal{R}_\Omega) \left(\left(\mathbf{V}_t - \mathbf{V}_t^{(l)} \right) \boldsymbol{\Sigma}_t \mathbf{V}_t^\top \right) \mathbf{V}_t \right\|_F \\ &\quad + 4\sigma_1^* \left\| \mathbf{V}_t - \mathbf{V}_t^{(l)} \right\|_F. \end{aligned} \quad (146)$$

The first term can be controlled by Proposition 20. For the second term, we have

$$\begin{aligned} \left\| (\mathcal{I} - \mathcal{R}_\Omega) \left(\mathbf{V}_t \boldsymbol{\Sigma}_t \left(\mathbf{V}_t - \mathbf{V}_t^{(l)} \right)^\top \right) \mathbf{V}_t \right\|_F &= \sup_{\|\mathbf{Z}\|_F \leq 1} \left\langle (\mathcal{I} - \mathcal{R}_\Omega) \left(\mathbf{V}_t \boldsymbol{\Sigma}_t \left(\mathbf{V}_t - \mathbf{V}_t^{(l)} \right)^\top \right), \mathbf{Z} \mathbf{V}_t \right\rangle \\ &\stackrel{(a)}{\leq} \Gamma \sqrt{\frac{d}{p}} \left\| \mathbf{V}_t \right\|_{2,\infty}^2 \left\| \left(\mathbf{V}_t - \mathbf{V}_t^{(l)} \right) \boldsymbol{\Sigma}_t \right\|_F \\ &\leq 8\Gamma\sigma_1^* \sqrt{\frac{\mu^2 r^2}{pd}} \left\| \mathbf{V}_t - \mathbf{V}_t^{(l)} \right\|_F. \end{aligned} \quad (147)$$

Here we apply Corollary 29 in (a). Similarly, we can also control the third term as follows

$$\begin{aligned} \left\| (\mathcal{I} - \mathcal{R}_\Omega) \left(\left(\mathbf{V}_t - \mathbf{V}_t^{(l)} \right) \boldsymbol{\Sigma}_t \mathbf{V}_t^\top \right) \mathbf{V}_t \right\|_F &= \sup_{\|\mathbf{Z}\|_F \leq 1} \left\langle (\mathcal{I} - \mathcal{R}_\Omega) \left(\left(\mathbf{V}_t - \mathbf{V}_t^{(l)} \right) \boldsymbol{\Sigma}_t \mathbf{V}_t^\top \right), \mathbf{Z} \mathbf{V}_t \right\rangle \\ &\leq 8\Gamma \sqrt{\frac{\mu r^2}{p}} \left\| \mathbf{V}_t - \mathbf{V}_t^{(l)} \right\|_{2,\infty} \\ &\leq 24\Gamma\sigma_1^* \sqrt{\frac{\mu^2 r^3}{pd}} \left\| \mathbf{V}_t - \mathbf{V}_t^{(l)} \right\|_F. \end{aligned} \quad (148)$$

Combining the above inequalities leads to

$$\left\| \mathcal{P}_{\mathbf{V}_t}^\perp \mathbf{M}_t \mathbf{V}_t - \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \mathbf{M}_t^{(l)} \mathbf{V}_t^{(l)} \right\|_F \leq C\sigma_1^* \left\| \mathbf{V}_t - \mathbf{V}_t^{(l)} \right\|_F. \quad (149)$$

Similarly, we can derive that $\|\mathbf{A}_t - \mathbf{A}_t^{(l)}\|_F \leq C\sigma_1^* \|\mathbf{V}_t - \mathbf{V}_t^{(l)}\|_F + C\sigma_1^* \sqrt{\frac{\kappa\mu^3 r^{5.5} \log(\frac{1}{\alpha}) \log(d)}{\sqrt{pd} \cdot d}}$. Therefore, we have

$$(II) \leq C_3 \eta^2 \left(\sigma_1^* \|\mathbf{V}_t - \mathbf{V}_t^{(l)}\|_F^2 + \sigma_1^{*2} \frac{\kappa\mu^3 r^{5.5} \log(\frac{1}{\alpha}) \log(d)}{\sqrt{pd} \cdot d} \right). \quad (150)$$

Putting everything together, we obtain

$$\begin{aligned} \|\mathbf{V}_{t+1} - \mathbf{V}_{t+1}^{(l)}\|_F &\leq \left(1 + \Gamma_1 \eta \sigma_1^* \frac{\kappa\mu r^{1.5} \log(\frac{1}{\alpha})}{\sqrt{pd}} \right) \|\mathbf{V}_t - \mathbf{V}_t^{(l)}\|_F + \Gamma_3 \eta \sigma_1^* \sqrt{\frac{\kappa\mu^3 r^{5.5} \log(\frac{1}{\alpha}) \log(d)}{\sqrt{pd} \cdot d}} \\ &\leq \|\mathbf{V}_t - \mathbf{V}_t^{(l)}\|_F + 2\Gamma_3 \eta \sigma_1^* \sqrt{\frac{\kappa\mu^3 r^{5.5} \log(\frac{1}{\alpha}) \log(d)}{\sqrt{pd} \cdot d}}. \end{aligned} \quad (151)$$

■

E.4. Proof of Proposition 20

To prove this proposition, we propose a novel argument based on *adaptive ϵ -nets*. Upon fixing an ϵ -net \mathcal{N}_ϵ with respect to Frobenius norm for the operator norm ball $\mathcal{B}_{\text{op}}^{r \times r}(4\sigma_1^*)$ with $\epsilon = \frac{c}{d}^2$, we first construct a series of adaptive ϵ -nets $\{\mathcal{V}_{\epsilon,t}^{(l)}\}_{t=0}^T$ in the following recursive manner.

$$\begin{aligned} \mathcal{V}_{\epsilon,t+1}^{(l)} &= \left\{ \mathbf{Y}_{\epsilon,t+1}^{(l)} : \|\mathbf{Y}_{\epsilon,t+1}^{(l)}\|_{2,\infty} \leq 2\sqrt{\frac{\mu r}{d}}, \text{ and } \mathbf{Y}_{\epsilon,t+1}^{(l)} = \mathbf{Z}_{\epsilon,t+1}^{(l)} \left(\mathbf{Z}_{\epsilon,t+1}^{(l)\top} \mathbf{Z}_{\epsilon,t+1}^{(l)} \right)^{-1/2} \text{ where} \right. \\ &\quad \left. \mathbf{Z}_{\epsilon,t+1}^{(l)} = \left(\mathbf{I} + \eta \mathcal{R}_{\Omega^{(l)}} \left(\mathbf{X}^* - \mathbf{V}_{\epsilon,t}^{(l)} \boldsymbol{\Sigma}_{\epsilon,t} \mathbf{V}_{\epsilon,t}^{(l)\top} \right) \right) \mathbf{V}_{\epsilon,t}^{(l)}, \forall \mathbf{V}_{\epsilon,t}^{(l)} \in \mathcal{V}_{\epsilon,t}^{(l)}, \boldsymbol{\Sigma}_{\epsilon,t} \in \mathcal{N}_\epsilon \right\}. \end{aligned} \quad (152)$$

2. We call a set \mathcal{N}_ϵ an ϵ -net with respect to a norm $\|\cdot\|$ for a given set S if for any element $s \in S$, there exists an element $s' \in \mathcal{N}_\epsilon$ such that $\|s - s'\| \leq \epsilon$.

Moreover, we set $\mathcal{V}_{\epsilon,0}^{(l)} = \{\mathbf{V}^*\}$. Then, we denote the best approximation of $\mathbf{V}_t^{(l)}$ in $\mathcal{V}_{\epsilon,t}^{(l)}$ as $\mathbf{V}_{\epsilon,t}^{(l)} = \arg \min_{\mathbf{V} \in \mathcal{V}_{\epsilon,t}^{(l)}} \|\mathbf{V}_t^{(l)} - \mathbf{V}\|_F$. Hence, we have the following decomposition

$$\begin{aligned}
 & \left\| (\mathcal{R}_\Omega - \mathcal{R}_{\Omega^{(l)}}) \left(\mathbf{X}^* - \mathbf{V}_t^{(l)} \boldsymbol{\Sigma}_t \mathbf{V}_t^{(l)\top} \right) \mathbf{V}_t^{(l)} \right\|_F \\
 & \leq \underbrace{\left\| (\mathcal{R}_\Omega - \mathcal{R}_{\Omega^{(l)}}) \left(\mathbf{X}^* - \mathbf{V}_{\epsilon,t}^{(l)} \boldsymbol{\Sigma}_{\epsilon,t} \mathbf{V}_{\epsilon,t}^{(l)\top} \right) \mathbf{V}_{\epsilon,t}^{(l)} \right\|_F}_{\text{(I)}} \\
 & \quad + \underbrace{\left\| (\mathcal{R}_\Omega - \mathcal{R}_{\Omega^{(l)}}) \left(\mathbf{X}^* - \mathbf{V}_t^{(l)} \boldsymbol{\Sigma}_t \mathbf{V}_t^{(l)\top} \right) \left(\mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)} \right) \right\|_F}_{\text{(II)}} \\
 & \quad + \underbrace{\left\| (\mathcal{R}_\Omega - \mathcal{R}_{\Omega^{(l)}}) \left(\mathbf{V}_t^{(l)} \boldsymbol{\Sigma}_t \left(\mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)} \right)^\top \right) \mathbf{V}_{\epsilon,t}^{(l)} \right\|_F}_{\text{(III)}} \\
 & \quad + \underbrace{\left\| (\mathcal{R}_\Omega - \mathcal{R}_{\Omega^{(l)}}) \left(\mathbf{V}_t^{(l)} (\boldsymbol{\Sigma}_t - \boldsymbol{\Sigma}_{\epsilon,t}) \mathbf{V}_{\epsilon,t}^{(l)\top} \right) \mathbf{V}_{\epsilon,t}^{(l)} \right\|_F}_{\text{(IV)}} \\
 & \quad + \underbrace{\left\| (\mathcal{R}_\Omega - \mathcal{R}_{\Omega^{(l)}}) \left(\left(\mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)} \right) \boldsymbol{\Sigma}_{\epsilon,t} \mathbf{V}_{\epsilon,t}^{(l)\top} \right) \mathbf{V}_{\epsilon,t}^{(l)} \right\|_F}_{\text{(V)}}.
 \end{aligned} \tag{153}$$

To proceed, we require the following key lemmas.

Lemma 21 *For any $1 \leq i \leq j \leq d$, let $p_{i,j}$ be a Bernoulli random variable taking the value $p_{i,j} = 1$ if and only if $(i, j) \in \Omega$. Let $r_{i,j} = \frac{1}{2p}(p_{i,j} + p_{j,i})$. The following statements hold:*

- **Independent case:** Suppose the sampling rate $p \gtrsim \frac{\mu r}{d} \log\left(\frac{4r}{\delta}\right)$. Suppose that $\mathbf{X} \in \mathcal{S}_{d \times d}$ and $\mathbf{V} \in \mathcal{O}_{d \times r}$ with $\|\mathbf{V}\|_{2,\infty} \leq 2\sqrt{\frac{\mu r}{d}}$ are independent of $r_{l,1}, \dots, r_{l,d}$. Then, with probability at least $1 - \delta$, we have

$$\left\| (\mathcal{R}_\Omega - \mathcal{R}_{\Omega^{(l)}}) (\mathbf{X}) \mathbf{V} \right\|_F^2 \leq \frac{32\mu r \log(4r/\delta)}{p} \|\mathbf{X}\|_{\max}^2. \tag{154}$$

- **General case:** Suppose the sampling rate $p \gtrsim \frac{\log(d)}{d}$. For arbitrary $\mathbf{X} \in \mathbb{R}^{d \times d}$ and $\mathbf{V} \in \mathbb{R}^{d \times r}$, with probability at least $1 - \frac{1}{d^3}$, we have

$$\left\| (\mathcal{R}_\Omega - \mathcal{R}_{\Omega^{(l)}}) (\mathbf{X}) \mathbf{V} \right\|_F^2 \leq \frac{2d}{p} \|\mathbf{X}\|_{\max}^2 \|\mathbf{V}\|_F^2. \tag{155}$$

Lemma 22 *Suppose the sampling rate satisfies $p \gtrsim \frac{\mu r}{d} \log\left(\frac{4r}{\delta}\right)$. For any $0 \leq t \leq T = \frac{100}{\eta \sigma_r^*} \log\left(\frac{1}{\alpha}\right)$, $1 \leq l \leq d$, with probability at least $1 - \frac{1}{d^3}$, we have*

$$\left\| (\mathcal{R}_\Omega - \mathcal{R}_{\Omega^{(l)}}) \left(\mathbf{X}^* - \mathbf{V}_{\epsilon,t}^{(l)} \boldsymbol{\Sigma}_{\epsilon,t} \mathbf{V}_{\epsilon,t}^{(l)\top} \right) \mathbf{V}_{\epsilon,t}^{(l)} \right\|_F^2 \leq \Gamma \frac{\kappa \mu^3 r^{5.5} \log\left(\frac{1}{\alpha}\right) \log\left(\frac{r \sigma_1^*}{\epsilon}\right)}{\sqrt{p d^3}}. \tag{156}$$

Lemma 23 *Under the same conditions as Lemma 22, we have*

$$\left\| \mathbf{V}_{t+1}^{(l)} - \mathbf{V}_{\epsilon, t+1}^{(l)} \right\|_F \leq \left(1 + \Gamma \eta \sigma_1^* \frac{\kappa \mu r^{1.5} \log\left(\frac{1}{\alpha}\right)}{\sqrt{pd}} \right) \left\| \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon, t}^{(l)} \right\|_F + C_3 \eta \frac{\sigma_1^* \mu r}{\sqrt{pd}} \epsilon. \quad (157)$$

The proof of these lemmas is deferred to the end of this section. Now we use Lemma 23 to control the dynamic of $\left\| \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon, t}^{(l)} \right\|_F$. Note that $\left\| \mathbf{V}_0^{(l)} - \mathbf{V}_{\epsilon, 0}^{(l)} \right\|_F = 0$. Hence, applying Lemma 47, we have

$$\begin{aligned} \left\| \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon, t}^{(l)} \right\|_F &\leq C \left(\left(1 + \Gamma \eta \sigma_1^* \frac{\kappa \mu r^{1.5} \log\left(\frac{1}{\alpha}\right)}{\sqrt{pd}} \right)^t - 1 \right) \frac{\eta \frac{\sigma_1^* \mu r}{\sqrt{pd}} \epsilon}{\eta \sigma_1^* \frac{\kappa \mu r^{1.5} \log\left(\frac{1}{\alpha}\right)}{\sqrt{pd}}} \\ &\leq C \eta \frac{\sigma_1^* \mu r}{\sqrt{pd}} \epsilon \cdot t. \end{aligned} \quad (158)$$

Now we are ready to control (I) to (V) separately. First, Lemma 22 directly implies that

$$(I) \leq \Gamma \frac{\sqrt{\kappa \mu^3 r^{5.5} \log\left(\frac{1}{\alpha}\right) \log\left(\frac{r \sigma_1^*}{\epsilon}\right)}}{\sqrt[4]{pd^3}}. \quad (159)$$

On the other hand, applying Lemma 21 to (II) leads to

$$(II) \leq \sqrt{\frac{2d}{p}} \left\| \mathbf{X}^* - \mathbf{V}_t^{(l)} \boldsymbol{\Sigma}_t \mathbf{V}_t^{(l)\top} \right\|_{\max} \left\| \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon, t}^{(l)} \right\|_F \leq \Gamma \sqrt{\frac{\mu^2 r^2}{pd}} \left\| \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon, t}^{(l)} \right\|_F. \quad (160)$$

Similarly, we have

$$\begin{aligned} (III) &\leq \sqrt{\frac{2d}{p}} \left\| \mathbf{V}_t^{(l)} \boldsymbol{\Sigma}_t \left(\mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon, t}^{(l)} \right)^\top \right\|_{\max} \left\| \mathbf{V}_{\epsilon, t}^{(l)} \right\|_F \leq \Gamma \sqrt{\frac{\mu r^2}{p}} \left\| \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon, t}^{(l)} \right\|_F, \\ (IV) &\leq \sqrt{\frac{2d}{p}} \left\| \mathbf{V}_t^{(l)} (\boldsymbol{\Sigma}_t - \boldsymbol{\Sigma}_{\epsilon, t}) \mathbf{V}_{\epsilon, t}^{(l)\top} \right\|_{\max} \left\| \mathbf{V}_{\epsilon, t}^{(l)} \right\|_F \leq \Gamma \sqrt{\frac{\mu^2 r^3}{pd}} \epsilon, \\ (V) &\leq \sqrt{\frac{2d}{p}} \left\| \left(\mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon, t}^{(l)} \right) \boldsymbol{\Sigma}_{\epsilon, t} \mathbf{V}_{\epsilon, t}^{(l)\top} \right\|_{\max} \left\| \mathbf{V}_{\epsilon, t}^{(l)} \right\|_F \leq \Gamma \sqrt{\frac{\mu r^2}{p}} \left\| \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon, t}^{(l)} \right\|_F. \end{aligned} \quad (161)$$

Overall, we derive that

$$\begin{aligned} &\left\| (\mathcal{R}_\Omega - \mathcal{R}_{\Omega^{(l)}}) \left(\mathbf{X}^* - \mathbf{V}_t^{(l)} \boldsymbol{\Sigma}_t \mathbf{V}_t^{(l)\top} \right) \mathbf{V}_t^{(l)} \right\|_F \\ &\leq \Gamma \frac{\sqrt{\kappa \mu^3 r^{5.5} \log\left(\frac{1}{\alpha}\right) \log\left(\frac{r \sigma_1^*}{\epsilon}\right)}}{\sqrt[4]{pd^3}} + \Gamma \sqrt{\frac{\mu r^2}{p}} \left\| \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon, t}^{(l)} \right\|_F \\ &\leq \Gamma \frac{\sqrt{\kappa \mu^3 r^{5.5} \log\left(\frac{1}{\alpha}\right) \log\left(\frac{r \sigma_1^*}{\epsilon}\right)}}{\sqrt[4]{pd^3}} + \Gamma \sqrt{\frac{\mu r^2}{p}} \cdot \eta \frac{\sigma_1^* \mu r}{\sqrt{pd}} \epsilon \cdot t \\ &\leq 2\Gamma \sigma_1^* \frac{\sqrt{\kappa \mu^3 r^{5.5} \log\left(\frac{1}{\alpha}\right) \log(d)}}{\sqrt[4]{pd^3}}. \end{aligned} \quad (162)$$

The last inequality follows from the fact that $\epsilon = \frac{c}{d}$. This completes the proof of Proposition 20. ■

Next, we proceed to present the proofs of Lemma 21, Lemma 22, and Lemma 23.

Proof of Lemma 21. We first expand $\|(\mathcal{R}_\Omega - \mathcal{R}_{\Omega^{(l)}})(\mathbf{X}) \mathbf{V}\|_F^2$ as follows

$$\|(\mathcal{R}_\Omega - \mathcal{R}_{\Omega^{(l)}})(\mathbf{X}) \mathbf{V}\|_F^2 = \sum_{j=1}^r \left(\sum_{k=1}^d (r_{l,k} - 1) X_{l,k} V_{k,j} \right)^2 + \sum_{i \neq l} (r_{i,l} - 1)^2 X_{i,l}^2 \|\mathbf{V}_{l,\cdot}\|^2. \quad (163)$$

Next, we prove these two cases separately.

Independent case. First, we control $\sum_{k=1}^d (r_{l,k} - 1) X_{l,k} V_{k,j}$ for all $1 \leq j \leq r$ via Bernstein's inequality (Lemma 32). To this goal, for a fixed $1 \leq j \leq r$, upon defining $Z_k = (r_{l,k} - 1) X_{l,k} V_{k,j}$ with $\mathbb{E}[Z_k] = 0$, we have

$$\begin{aligned} M &:= \max_k |Z_k| \leq \frac{1}{p} \|\mathbf{X}\|_{\max} \|\mathbf{V}\|_{2,\infty}, \\ \nu^2 &:= \sum_{k=1}^d \text{Var}[Z_k^2] \leq \frac{1}{p} \sum_{k=1}^d X_{l,k}^2 V_{k,j}^2 \leq \frac{1}{p} \|\mathbf{X}\|_{\max}^2 \sum_{k=1}^d V_{k,j}^2 = \frac{1}{p} \|\mathbf{X}\|_{\max}^2. \end{aligned} \quad (164)$$

Here in the last equality we use the fact that $\sum_{k=1}^d V_{k,j}^2 = 1$ since $\mathbf{V} \in \mathcal{O}_{d \times r}$. Therefore, due to Bernstein's inequality, with probability at least $1 - \frac{\delta}{2r}$, one has

$$\begin{aligned} \left| \sum_{k=1}^d Z_k \right| &\leq 2\nu \sqrt{\log \left(\frac{4r}{\delta} \right)} + \frac{4}{3} M \log \left(\frac{4r}{\delta} \right) \\ &\leq \frac{2 \|\mathbf{X}\|_{\max}}{\sqrt{p}} \sqrt{\log \left(\frac{4r}{\delta} \right)} + \frac{4 \|\mathbf{X}\|_{\max} \|\mathbf{V}\|_{2,\infty}}{3p} \log \left(\frac{4r}{\delta} \right) \\ &\leq \frac{4 \|\mathbf{X}\|_{\max}}{\sqrt{p}} \sqrt{\log \left(\frac{4r}{\delta} \right)}, \end{aligned} \quad (165)$$

where we use the assumption $p \gtrsim \frac{\mu r}{d} \log \left(\frac{4r}{\delta} \right)$. Hence, via a union bound, we know that with probability at least $1 - \frac{\delta}{2}$, we have

$$\sum_{j=1}^r \left(\sum_{k=1}^d (r_{l,k} - 1) X_{l,k} V_{k,j} \right)^2 \leq \frac{16r \|\mathbf{X}\|_{\max}^2 \log \left(\frac{4r}{\delta} \right)}{p}. \quad (166)$$

Next, we can control $\sum_{i \neq l} (r_{i,l} - 1)^2 X_{i,l}^2 \|\mathbf{V}_{l,\cdot}\|^2$ as

$$\sum_{i \neq l} (r_{i,l} - 1)^2 X_{i,l}^2 \|\mathbf{V}_{l,\cdot}\|^2 \leq \sum_{i=1}^d (r_{i,l} - 1)^2 \|\mathbf{X}\|_{\max}^2 \|\mathbf{V}\|_{2,\infty}^2. \quad (167)$$

Then, we apply Bernstein's inequality to control $\sum_{k=1}^d (r_{l,k} - 1)^2$. To this end, notice that

$$\begin{aligned}\mathbb{E} \left[(r_{l,k} - 1)^2 \right] &= \frac{1-p}{2p} \leq \frac{1}{2p}, \\ M &:= \max_k (r_{l,k} - 1)^2 \leq \left(\frac{1}{p} - 1 \right)^2 \leq \frac{1}{p^2}, \\ \nu^2 &:= \sum_{k=1}^d \text{Var} \left[(r_{l,k} - 1)^2 \right] \leq \frac{d}{p^3}.\end{aligned}\tag{168}$$

Therefore, with probability at least $1 - \frac{\delta}{2}$, we have

$$\sum_{k=1}^d (r_{l,k} - 1)^2 \leq \frac{d}{2p} + \frac{4}{3} \frac{1}{p^2} \log \left(\frac{2}{\delta} \right) + 2 \sqrt{\frac{d}{p^3}} \sqrt{\log \left(\frac{2}{\delta} \right)} \leq \frac{d}{p}\tag{169}$$

since we set $p \gtrsim \frac{\mu r}{d} \log \left(\frac{4r}{\delta} \right)$. This implies that with probability at least $1 - \frac{\delta}{2}$, we have

$$\sum_{i \neq l} (r_{i,l} - 1)^2 X_{i,l}^2 \|\mathbf{V}_{l,\cdot}\|^2 \leq \frac{d}{p} \|\mathbf{X}\|_{\max}^2 \|\mathbf{V}\|_{2,\infty}^2.\tag{170}$$

Finally, taking a union bound, we conclude that with probability at least $1 - \delta$, we have

$$\begin{aligned}\|(\mathcal{R}_\Omega - \mathcal{R}_{\Omega^{(l)}})(\mathbf{X}) \mathbf{V}\|_{\text{F}}^2 &\leq \frac{16r \|\mathbf{X}\|_{\max}^2 \log(4r/\delta)}{p} + \frac{d}{p} \|\mathbf{X}\|_{\max}^2 \|\mathbf{V}\|_{2,\infty}^2 \\ &\leq \frac{32\mu r \log(4r/\delta)}{p} \|\mathbf{X}\|_{\max}^2.\end{aligned}\tag{171}$$

General case. First, we apply Cauchy-Schwarz inequality to obtain

$$\begin{aligned}\sum_{j=1}^r \left(\sum_{k=1}^d (r_{l,k} - 1) X_{l,k} V_{k,j} \right)^2 &\leq \sum_{j=1}^r \left(\sum_{k=1}^d (r_{l,k} - 1)^2 \right) \left(\sum_{k=1}^d X_{l,k}^2 V_{k,j}^2 \right) \\ &\leq \left(\sum_{k=1}^d (r_{l,k} - 1)^2 \right) \cdot \sum_{j=1}^r \|\mathbf{X}\|_{\max}^2 \sum_{k=1}^d V_{k,j}^2 \\ &= \left(\sum_{k=1}^d (r_{l,k} - 1)^2 \right) \cdot \|\mathbf{X}\|_{\max}^2 \|\mathbf{V}\|_{\text{F}}^2.\end{aligned}\tag{172}$$

Next, we can control $\sum_{i \neq l} (r_{i,l} - 1)^2 X_{i,l}^2 \|\mathbf{V}_{l,\cdot}\|^2$ as follows

$$\begin{aligned}\sum_{i \neq l} (r_{i,l} - 1)^2 X_{i,l}^2 \|\mathbf{V}_{l,\cdot}\|^2 &\leq \sum_{i=1}^d (r_{i,l} - 1)^2 \|\mathbf{X}\|_{\max}^2 \|\mathbf{V}\|_{\text{F}}^2 \\ &= \left(\sum_{k=1}^d (r_{l,k} - 1)^2 \right) \|\mathbf{X}\|_{\max}^2 \|\mathbf{V}\|_{\text{F}}^2.\end{aligned}\tag{173}$$

Here we use the fact that $r_{i,j} = r_{j,i}, \forall i, j \in [d]$. Next, we apply Bernstein's inequality to control $\sum_{k=1}^d (r_{l,k} - 1)^2$. To this end, notice that

$$\begin{aligned} \mathbb{E} \left[(r_{l,k} - 1)^2 \right] &= \frac{1-p}{2p} \leq \frac{1}{2p}, \\ M &:= \max_k (r_{l,k} - 1)^2 \leq \left(\frac{1}{p} - 1 \right)^2 \leq \frac{1}{p^2}, \\ \nu^2 &:= \sum_{k=1}^d \text{Var} \left[(r_{l,k} - 1)^2 \right] \leq \frac{d}{p^3}. \end{aligned} \quad (174)$$

Therefore, with probability at least $1 - \delta$, we have

$$\sum_{k=1}^d (r_{l,k} - 1)^2 \leq \frac{d}{2p} + \frac{4}{3} \frac{1}{p^2} \log \left(\frac{2}{\delta} \right) + 2 \sqrt{\frac{d}{p^3}} \sqrt{\log \left(\frac{2}{\delta} \right)}. \quad (175)$$

Specifically, upon setting $\delta = \frac{1}{d^3}$, we obtain that with probability at least $1 - \frac{1}{d^3}$, one has

$$\sum_{k=1}^d (r_{l,k} - 1)^2 \leq \frac{d}{p}, \quad (176)$$

since we set $p \gtrsim \frac{\log(d)}{d}$. Overall, we have

$$\|(\mathcal{R}_\Omega - \mathcal{R}_{\Omega^{(l)}})(\mathbf{X}) \mathbf{V}\|_{\text{F}}^2 \leq 2 \left(\sum_{k=1}^d (r_{l,k} - 1)^2 \right) \|\mathbf{X}\|_{\text{max}}^2 \|\mathbf{V}\|_{\text{F}}^2 \leq \frac{2d}{p} \|\mathbf{X}\|_{\text{max}}^2 \|\mathbf{V}\|_{\text{F}}^2, \quad (177)$$

with probability at least $1 - \frac{1}{d^3}$. This completes the proof. \blacksquare

Proof of Lemma 22. First, for fixed $\Sigma_{\epsilon,t} \in \mathcal{N}_\epsilon$ and $\mathbf{V}_{\epsilon,t}^{(l)} \in \mathcal{V}_{\epsilon,t}^{(l)}$, Lemma 21 implies that with probability at least $1 - \delta$, we have

$$\left\| (\mathcal{R}_\Omega - \mathcal{R}_{\Omega^{(l)}}) \left(\mathbf{X}^* - \mathbf{V}_{\epsilon,t}^{(l)} \Sigma_{\epsilon,t} \mathbf{V}_{\epsilon,t}^{(l)\top} \right) \mathbf{V}_{\epsilon,t}^{(l)} \right\|_{\text{F}}^2 \leq \frac{32\mu r \log(4r/\delta)}{p} \left\| \mathbf{X}^* - \mathbf{V}_{\epsilon,t}^{(l)} \Sigma_{\epsilon,t} \mathbf{V}_{\epsilon,t}^{(l)\top} \right\|_{\text{max}}^2. \quad (178)$$

Note that

$$\begin{aligned} \left\| \mathbf{X}^* - \mathbf{V}_{\epsilon,t}^{(l)} \Sigma_{\epsilon,t} \mathbf{V}_{\epsilon,t}^{(l)\top} \right\|_{\text{max}} &\leq \|\mathbf{X}^*\|_{\text{max}} + \left\| \mathbf{V}_{\epsilon,t}^{(l)} \Sigma_{\epsilon,t} \mathbf{V}_{\epsilon,t}^{(l)\top} \right\|_{\text{max}} \\ &\stackrel{(a)}{\leq} \|\Sigma^*\| \|\mathbf{V}^*\|_{2,\infty}^2 + \|\Sigma_{\epsilon,t}\| \left\| \mathbf{V}_{\epsilon,t}^{(l)} \right\|_{2,\infty}^2 \\ &\leq 9\sigma_1^* \frac{\mu r}{d}. \end{aligned} \quad (179)$$

Therefore, with probability at least $1 - \delta$, we have

$$\left\| (\mathcal{R}_\Omega - \mathcal{R}_{\Omega^{(l)}}) \left(\mathbf{X}^* - \mathbf{V}_{\epsilon,t}^{(l)} \Sigma_{\epsilon,t} \mathbf{V}_{\epsilon,t}^{(l)\top} \right) \mathbf{V}_{\epsilon,t}^{(l)} \right\|_{\text{F}}^2 \leq \frac{288\mu^3 r^3 \log(4r/\delta)}{pd^2}. \quad (180)$$

Lastly, we apply the union bound to finalize the desired result. To this end, note that $\mathcal{B}_{\text{op}}^{r \times r}(4\sigma_1^*) \subset \mathcal{B}_{\text{F}}^{r \times r}(4\sqrt{r}\sigma_1^*)$. Hence, according to Lemma 42, we know that $|\mathcal{N}_\epsilon| \leq \left(\frac{12\sqrt{r}\sigma_1^*}{\epsilon}\right)^{r^2}$ which implies that $|\mathcal{V}_{\epsilon,t}^{(l)}| \leq \left(\frac{6\sqrt{r}\sigma_1^*}{\epsilon}\right)^{r^2t}$. Therefore, the total cardinality of $\bigcup_{l=1}^d \bigcup_{t=0}^T \mathcal{V}_{\epsilon,t}^{(l)}$ is upper-bounded by

$$\left| \bigcup_{l=1}^d \bigcup_{t=0}^T \mathcal{V}_{\epsilon,t}^{(l)} \right| \leq d \cdot \sum_{t=0}^T \left(\frac{12\sqrt{r}\sigma_1^*}{\epsilon} \right)^{r^2t} \leq 2d \left(\frac{12\sqrt{r}\sigma_1^*}{\epsilon} \right)^{r^2T}. \quad (181)$$

Hence, once we set $\delta = \frac{1}{2d^4} \left(\frac{6\sqrt{r}\sigma_1^*}{\epsilon} \right)^{-r^2T}$, we obtain that with probability at least $1 - d^{-3}$, for any $0 \leq t \leq T, 1 \leq l \leq d$ and $\mathbf{V}_{\epsilon,t}^{(l)} \in \mathcal{V}_{\epsilon,t}^{(l)}$,

$$\begin{aligned} \left\| (\mathcal{R}_\Omega - \mathcal{R}_{\Omega^{(l)}}) \left(\mathbf{X}^* - \mathbf{V}_{\epsilon,t}^{(l)} \boldsymbol{\Sigma}_{\epsilon,t} \mathbf{V}_{\epsilon,t}^{(l)\top} \right) \mathbf{V}_{\epsilon,t}^{(l)} \right\|_{\text{F}}^2 &\leq \Gamma \frac{\mu^3 r^3}{pd^2} \left(\log(d) + r^2 T \left(\frac{r\sigma_1^*}{\epsilon} \right) \right) \\ &\leq \Gamma \frac{\kappa \mu^3 r^{5.5} \log\left(\frac{1}{\alpha}\right) \log\left(\frac{r\sigma_1^*}{\epsilon}\right)}{\sqrt{pd^3}}. \end{aligned} \quad (182)$$

Here we use the fact that $T \leq \frac{100}{\eta\sigma_r^*} \log(1/\alpha)$ and $\eta = \Theta\left(\frac{\mu}{\sigma_1^*} \sqrt{\frac{r}{pd}}\right)$. This completes the proof. \blacksquare

Proof of Lemma 23. The derivation of $\left\| \mathbf{V}_{t+1}^{(l)} - \mathbf{V}_{\epsilon,t+1}^{(l)} \right\|_{\text{F}}$ is nearly the same as that of $\left\| \mathbf{V}_{t+1} - \mathbf{V}_{\epsilon,t+1}^{(l)} \right\|_{\text{F}}$. First, note that

$$\mathbf{V}_{\epsilon,t+1}^{(l)} = \left(\mathbf{I} + \eta \mathcal{P}_{\mathbf{V}_{\epsilon,t}^{(l)}}^\perp \mathbf{M}_{\epsilon,t}^{(l)} \right) \mathbf{V}_{\epsilon,t}^{(l)} + \mathbf{A}_{\epsilon,t}^{(l)}, \quad (183)$$

where $\mathbf{M}_{\epsilon,t}^{(l)}$ and $\mathbf{A}_{\epsilon,t}^{(l)}$ are defined similar to $\mathbf{M}_t^{(l)}$ and $\mathbf{A}_t^{(l)}$. Hence, we can expand $\left\| \mathbf{V}_{t+1}^{(l)} - \mathbf{V}_{\epsilon,t+1}^{(l)} \right\|_{\text{F}}^2$ as

$$\left\| \mathbf{V}_{t+1}^{(l)} - \mathbf{V}_{\epsilon,t+1}^{(l)} \right\|_{\text{F}}^2 = \left\| \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)} \right\|_{\text{F}}^2 + 2\eta \underbrace{\left\langle \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)}, \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \mathbf{M}_t^{(l)} \mathbf{V}_t^{(l)} - \mathcal{P}_{\mathbf{V}_{\epsilon,t}^{(l)}}^\perp \mathbf{M}_{\epsilon,t}^{(l)} \mathbf{V}_{\epsilon,t}^{(l)} \right\rangle}_{:=\text{(I)}} + \text{(II)}. \quad (184)$$

where

$$\begin{aligned} \text{(II)} &= \eta^2 \left\| \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \mathbf{M}_t^{(l)} \mathbf{V}_t^{(l)} - \mathcal{P}_{\mathbf{V}_{\epsilon,t}^{(l)}}^\perp \mathbf{M}_{\epsilon,t}^{(l)} \mathbf{V}_{\epsilon,t}^{(l)} \right\|_{\text{F}}^2 + \left\| \mathbf{A}_t^{(l)} - \mathbf{A}_{\epsilon,t}^{(l)} \right\|_{\text{F}}^2 \\ &\quad + 2 \left\langle \left(\mathbf{I} + \eta \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \mathbf{M}_t^{(l)} \right) \mathbf{V}_t^{(l)} - \left(\mathbf{I} + \eta \mathcal{P}_{\mathbf{V}_{\epsilon,t}^{(l)}}^\perp \mathbf{M}_{\epsilon,t}^{(l)} \right) \mathbf{V}_{\epsilon,t}^{(l)}, \mathbf{A}_t^{(l)} - \mathbf{A}_{\epsilon,t}^{(l)} \right\rangle \end{aligned} \quad (185)$$

contains all the higher-order terms. Next, we further decompose (I) as follows

$$\begin{aligned}
 (\text{I}) &= \left\langle \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)}, \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \boldsymbol{\Xi}_t^{(l)} \mathbf{V}_t^{(l)} - \mathcal{P}_{\mathbf{V}_{\epsilon,t}^{(l)}}^\perp \boldsymbol{\Xi}_{\epsilon,t}^{(l)} \mathbf{V}_{\epsilon,t}^{(l)} \right\rangle \\
 &= \underbrace{\left\langle \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)}, \left(\mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp - \mathcal{P}_{\mathbf{V}_{\epsilon,t}^{(l)}}^\perp \right) \boldsymbol{\Xi}_{\epsilon,t}^{(l)} \mathbf{V}_{\epsilon,t}^{(l)} \right\rangle}_{:= (\text{I}_1)} + \underbrace{\left\langle \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)}, \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \boldsymbol{\Xi}_t^{(l)} \left(\mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)} \right) \right\rangle}_{:= (\text{I}_2)} \\
 &\quad + \underbrace{\left\langle \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)}, \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \left(\boldsymbol{\Xi}_t^{(l)} - \boldsymbol{\Xi}_{\epsilon,t}^{(l)} \right) \mathbf{V}_{\epsilon,t}^{(l)} \right\rangle}_{:= (\text{I}_3)}.
 \end{aligned} \tag{186}$$

Here we define $\boldsymbol{\Xi}_t^{(l)} = \mathbf{M}_t^{(l)} - \mathbf{V}_t^{(l)} \boldsymbol{\Sigma}^* \mathbf{V}^{*\top} + \mathbf{V}_t^{(l)} \boldsymbol{\Sigma}_t \mathbf{V}_t^{(l)\top}$ and $\boldsymbol{\Xi}_{\epsilon,t}^{(l)} = \mathbf{M}_{\epsilon,t}^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)} \boldsymbol{\Sigma}^* \mathbf{V}^{*\top} + \mathbf{V}_{\epsilon,t}^{(l)} \boldsymbol{\Sigma}_{\epsilon,t} \mathbf{V}_{\epsilon,t}^{(l)\top}$. We control (I₁) as follows

$$(\text{I}_1) \leq \left\| \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)} \right\|_{\text{F}} \left\| \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp - \mathcal{P}_{\mathbf{V}_{\epsilon,t}^{(l)}}^\perp \right\|_{\text{F}} \left\| \boldsymbol{\Xi}_{\epsilon,t}^{(l)} \right\| \leq 2 \left\| \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)} \right\|_{\text{F}}^2 \left\| \boldsymbol{\Xi}_{\epsilon,t}^{(l)} \right\|. \tag{187}$$

Here we apply Lemma 43 in the second inequality. For $\left\| \boldsymbol{\Xi}_{\epsilon,t}^{(l)} \right\|$, following the same analysis as in Equation (119), we have

$$\begin{aligned}
 \left\| \boldsymbol{\Xi}_{\epsilon,t}^{(l)} \right\| &\leq \sigma_1^* \left\| \mathbf{V}^* - \mathbf{V}_{\epsilon,t}^{(l)} \right\| + 10\Gamma \frac{\sigma_1^* \mu r}{\sqrt{pd}} \\
 &\leq \sigma_1^* \left(\left\| \mathbf{V}^* - \mathbf{V}_t^{(l)} \right\| + \left\| \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)} \right\| \right) + 10\Gamma \frac{\sigma_1^* \mu r}{\sqrt{pd}} \\
 &\leq \sigma_1^* \left\| \mathbf{V}^* - \mathbf{V}_t^{(l)} \right\| + 20\Gamma \frac{\sigma_1^* \mu r}{\sqrt{pd}}.
 \end{aligned} \tag{188}$$

Here in the last inequality, we use the fact that $\left\| \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)} \right\| \leq \epsilon \leq 10\Gamma \frac{\mu r}{\sqrt{pd}}$. Plugging this inequality into Equation (187), we obtain

$$(\text{I}_1) \leq 2\sigma_1^* \left\| \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)} \right\|_{\text{F}}^2 \left(\sigma_1^* \left\| \mathbf{V}^* - \mathbf{V}_t^{(l)} \right\| + 20\Gamma \frac{\sigma_1^* \mu r}{\sqrt{pd}} \right). \tag{189}$$

Similar to our derivation for (I_1) , it follows that $(I_2) \leq \sigma_1^* \left\| \mathbf{V}^* - \mathbf{V}_t^{(l)} \right\|_F^2 \left(\sigma_1^* \left\| \mathbf{V}^* - \mathbf{V}_t^{(l)} \right\| + 10\Gamma \frac{\sigma_1^* \mu r}{\sqrt{pd}} \right)$. Lastly, to bound (I_3) , we further decompose it as

$$\begin{aligned}
 (I_3) &= \left\langle \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)}, \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \left(\mathbf{M}_t^{(l)} - \mathbf{M}_{\epsilon,t}^{(l)} + \mathbf{V}_{\epsilon,t}^{(l)} \boldsymbol{\Sigma}^* \mathbf{V}^{*\top} - \mathbf{V}_{\epsilon,t}^{(l)} \boldsymbol{\Sigma}_t^{(l)} \mathbf{V}_{\epsilon,t}^{(l)\top} \right) \mathbf{V}_{\epsilon,t}^{(l)} \right\rangle \\
 &= \underbrace{\left\langle \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)}, \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp (\mathcal{I} - \mathcal{R}_{\Omega^{(l)}}) \left(\mathbf{V}_{\epsilon,t}^{(l)} \boldsymbol{\Sigma}_{\epsilon,t}^{(l)} \left(\mathbf{V}_{\epsilon,t}^{(l)} - \mathbf{V}_t^{(l)} \right)^\top \right) \mathbf{V}_{\epsilon,t}^{(l)} \right\rangle}_{:= (I_{3,1})} \\
 &\quad + \underbrace{\left\langle \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)}, \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp (\mathcal{I} - \mathcal{R}_{\Omega^{(l)}}) \left(\left(\mathbf{V}_{\epsilon,t}^{(l)} - \mathbf{V}_t^{(l)} \right) \boldsymbol{\Sigma}_{\epsilon,t}^{(l)} \mathbf{V}_t^{(l)\top} \right) \mathbf{V}_{\epsilon,t}^{(l)} \right\rangle}_{:= (I_{3,2})} \\
 &\quad + \underbrace{\left\langle \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)}, \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp (\mathcal{I} - \mathcal{R}_{\Omega^{(l)}}) \left(\mathbf{V}_t^{(l)} (\boldsymbol{\Sigma}_{\epsilon,t} - \boldsymbol{\Sigma}_t) \mathbf{V}_t^{(l)\top} \right) \mathbf{V}_{\epsilon,t}^{(l)} \right\rangle}_{:= (I_{3,3})} \\
 &\quad + \underbrace{\left\langle \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)}, \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \mathbf{V}_{\epsilon,t}^{(l)} \boldsymbol{\Sigma}^* \mathbf{V}^{*\top} \mathbf{V}_{\epsilon,t}^{(l)} \right\rangle}_{:= (I_{3,4})}.
 \end{aligned} \tag{190}$$

Next, we control these terms separately. First, we apply Lemma 29 to control $(I_{3,1})$. Specifically, upon setting $\mathbf{A} = \mathbf{V}_{\epsilon,t}^{(l)}$, $\mathbf{B} = \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \left(\mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)} \right)$, $\mathbf{C} = \left(\mathbf{V}_{\epsilon,t}^{(l)} - \mathbf{V}_t^{(l)} \right) \boldsymbol{\Sigma}_{\epsilon,t}$, $\mathbf{D} = \mathbf{V}_{\epsilon,t}^{(l)}$ in Lemma 29, with probability at least $1 - d^{-3}$, we have

$$\begin{aligned}
 (I_{3,1}) &= \left\langle (\mathcal{I} - \mathcal{R}_{\Omega}) \left(\mathbf{A} \mathbf{C}^\top \right), \mathbf{B} \mathbf{D}^\top \right\rangle \\
 &\leq \Gamma \sqrt{\frac{d}{p}} \|\mathbf{A}\|_{2,\infty} \|\mathbf{B}\|_F \cdot \|\mathbf{C}\|_F \|\mathbf{D}\|_{2,\infty} \\
 &= \Gamma \sqrt{\frac{d}{p}} \left\| \mathbf{V}_{\epsilon,t}^{(l)} \right\|_{2,\infty}^2 \left\| \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \left(\mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)} \right) \right\|_F \left\| \left(\mathbf{V}_{\epsilon,t}^{(l)} - \mathbf{V}_t^{(l)} \right) \boldsymbol{\Sigma}_{\epsilon,t} \right\|_F \\
 &\leq 16\Gamma \sigma_1^* \sqrt{\frac{\mu^2 r^2}{pd}} \left\| \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)} \right\|_F^2.
 \end{aligned} \tag{191}$$

Here in the last inequality, we use the fact that $\|\boldsymbol{\Sigma}_{\epsilon,t}\| \leq 4\sigma_1^*$ and $\left\| \mathbf{V}_{\epsilon,t}^{(l)} \right\|_{2,\infty} \leq 2\sqrt{\frac{\mu r}{d}}$. In a manner akin to our derivation for $(I_{3,1})$, we can also show that

$$\begin{aligned}
 (I_{3,2}) &\leq 16\Gamma \sigma_1^* \sqrt{\frac{\mu^2 r^2}{pd}} \left\| \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)} \right\|_F^2 \\
 (I_{3,3}) &\leq 4\Gamma \sigma_1^* \sqrt{\frac{\mu^2 r^2}{pd}} \epsilon \left\| \mathbf{V}_{\epsilon,t}^{(l)} - \mathbf{V}_t^{(l)} \right\|_F.
 \end{aligned} \tag{192}$$

Lastly, for $(I_{3,4})$, we observe that

$$\begin{aligned}
(I_{3,4}) &= - \left\langle \mathbf{V}_{\epsilon,t}^{(l)}, \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \mathbf{V}_{\epsilon,t}^{(l)} \boldsymbol{\Sigma}^* \mathbf{V}^{\star\top} \mathbf{V}_{\epsilon,t}^{(l)} \right\rangle \\
&\leq - \left\langle \mathbf{V}_{\epsilon,t}^{(l)}, \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \mathbf{V}_{\epsilon,t}^{(l)} \boldsymbol{\Sigma}^* \right\rangle + \sigma_1^* \left\| \mathbf{V}^* - \mathbf{V}_{\epsilon,t}^{(l)} \right\|_F \left\| \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)} \right\|_F^2 \\
&\stackrel{(a)}{\leq} \sigma_1^* \left\| \mathbf{V}^* - \mathbf{V}_{\epsilon,t}^{(l)} \right\|_F \left\| \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)} \right\|_F^2 \\
&\leq \sigma_1^* \left(\left\| \mathbf{V}^* - \mathbf{V}_t^{(l)} \right\|_F + \left\| \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)} \right\|_F \right) \left\| \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)} \right\|_F^2 \\
&\stackrel{(b)}{\leq} \sigma_1^* \left(\Gamma_1 \frac{\kappa \mu r^{1.5} \log(\frac{1}{\alpha})}{\sqrt{pd}} + \left\| \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)} \right\|_F \right) \left\| \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)} \right\|_F^2.
\end{aligned} \tag{193}$$

Here in (a), we use the fact that $\left\langle \mathbf{V}_{\epsilon,t}^{(l)}, \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \mathbf{V}_{\epsilon,t}^{(l)} \boldsymbol{\Sigma}^* \right\rangle = \left\| \mathcal{P}_{\mathbf{V}_t^{(l)}}^\perp \mathbf{V}_{\epsilon,t}^{(l)} \boldsymbol{\Sigma}^{*1/2} \right\|_F^2 \geq 0$. In (b), we apply Proposition 16.

Putting everything together, we obtain that

$$(I_3) \leq \left\| \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)} \right\|_F \left(32\Gamma\sigma_1^* \sqrt{\frac{\mu^2 r^2}{pd}} \left\| \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)} \right\|_F + 4\Gamma\sigma_1^* \sqrt{\frac{\mu^2 r^2}{pd}} \epsilon \right), \tag{194}$$

which in turn leads to

$$(I) \leq \Gamma\sigma_1^* \frac{\kappa \mu r^{1.5} \log(\frac{1}{\alpha})}{\sqrt{pd}} \left\| \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)} \right\|_F^2 + \left\| \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)} \right\|_F \cdot 4\Gamma\sigma_1^* \sqrt{\frac{\mu^2 r^2}{pd}} \epsilon. \tag{195}$$

Therefore, our final bound is established as

$$\left\| \mathbf{V}_{t+1}^{(l)} - \mathbf{V}_{\epsilon,t+1}^{(l)} \right\|_F \leq \left(1 + \Gamma\eta\sigma_1^* \frac{\kappa \mu r^{1.5} \log(\frac{1}{\alpha})}{\sqrt{pd}} \right) \left\| \mathbf{V}_t^{(l)} - \mathbf{V}_{\epsilon,t}^{(l)} \right\|_F + C_3\eta \frac{\sigma_1^* \mu r}{\sqrt{pd}} \epsilon. \tag{196}$$

This completes the proof. \blacksquare

Appendix F. Proofs for Different Initialization Schemes

Gaussian initialization. Suppose the search rank $0.5d \leq r' \leq d$. Let $\mathbf{Z} = \mathbf{G}/\|\mathbf{G}\|$ where $\mathbf{G} \in \mathbb{R}^{d \times r'}$ is a standard Gaussian matrix. Before proceeding, we first introduce the following lemma, which characterizes the concentration of the largest and smallest singular values of a standard Gaussian matrix.

Lemma 24 (Adapted from Theorem 6.1 in (Wainwright, 2019)) Suppose that $\mathbf{G} \in \mathbb{R}^{d_1 \times d_2}$ is a standard Gaussian matrix where $d_1 \geq d_2$. Then, for any $\delta > 0$, we have

$$\begin{aligned}
\mathbb{P} \left(\|\mathbf{G}\| \geq (2 + \delta) \sqrt{d_1} \right) &\leq \exp \left\{ -\frac{d_1 \delta^2}{2} \right\}, \\
\mathbb{P} \left(\sigma_{\min}(\mathbf{G}) \geq (1 - \delta) \sqrt{d_1} - \sqrt{d_2} \right) &\leq \exp \left\{ -\frac{d_1 \delta^2}{2} \right\}.
\end{aligned} \tag{197}$$

Notice that $\sigma_r(\mathcal{P}_{V^*} \mathbf{G}) = \sigma_{\min}(\mathbf{V}^{*\top} \mathbf{G})$ where $\mathbf{V}^{*\top} \mathbf{G} \in \mathbb{R}^{r \times r'}$ is another standard Gaussian matrix. Hence, upon setting $\delta = \frac{1}{2}$ and noting that $r' \geq \frac{d}{2}$, Lemma 24 implies that with probability at least $1 - \exp\{d/16\}$, we have

$$\sigma_r(\mathcal{P}_{V^*} \mathbf{G}) \geq \frac{1}{2} \sqrt{r'} - \sqrt{r} \geq \frac{1}{4} \sqrt{d}. \quad (198)$$

On the other hand, upon setting $\delta = \frac{1}{2}$, Lemma 24 implies that with probability at least $1 - \exp\{d/8\}$, we have

$$\mathbf{G} \leq \frac{5}{2} \sqrt{d}. \quad (199)$$

Via a union bound, we know that with probability at least $1 - 2 \exp\{d/16\}$, we have

$$\sigma_r(\mathcal{P}_{V^*} \mathbf{Z}) = \frac{\sigma_r(\mathcal{P}_{V^*} \mathbf{G})}{\|\mathbf{G}\|} \geq \frac{\sqrt{d}/4}{5\sqrt{d}/2} = \frac{1}{10}. \quad (200)$$

Hence, with probability at least $1 - 2 \exp\{d/16\}$, Condition 1 is satisfied with $c_0 = 0.1$.

Orthogonal initialization. Suppose the search rank satisfies $r' = d$. Upon choosing $\mathbf{Z} = \mathbf{O}$ for some $\mathbf{O} \in \mathcal{O}_{d \times d}$, we have

$$\sigma_r(\mathcal{P}_{V^*} \mathbf{O}) = \sigma_r(\mathbf{V}^* \mathbf{V}^{*\top}) = 1. \quad (201)$$

Hence, Condition 1 is satisfied with $c_0 = 1$.

Spectral initialization. Let $\mathbf{V} \mathbf{\Sigma} \mathbf{V}^\top$ be the eigendecomposition of the best rank- r' approximation of $\mathcal{R}_\Omega(\mathbf{X}^*)$ measured in Frobenius norm. Suppose $r \leq r' \leq d$ and $\mathbf{Z} = \mathbf{U} / \|\mathbf{U}\|$, where $\mathbf{U} = \mathbf{V} \mathbf{\Sigma}^{1/2}$. Corollary 30 tells us that, with probability at least $1 - \frac{1}{d^3}$, we have

$$\|\mathcal{R}_\Omega(\mathbf{X}^*) - \mathbf{X}^*\| \leq \Gamma \sigma_1^* \sqrt{\frac{\mu^2 r^2}{pd}}. \quad (202)$$

Conditioned on this event, we have

$$\|\mathbf{U}\|^2 = \|\mathcal{R}_\Omega(\mathbf{X}^*)\| \leq \|\mathbf{X}^*\| + \|\mathcal{R}_\Omega(\mathbf{X}^*) - \mathbf{X}^*\| \leq 2\sigma_1^*. \quad (203)$$

On the other hand, by Weyl's inequality, we have

$$\begin{aligned} \sigma_r^2(\mathcal{P}_{V^*} \mathbf{U}) &= \sigma_r(\mathbf{V}^{*\top} \mathbf{U} \mathbf{U}^\top \mathbf{V}^*) \\ &\geq \sigma_r(\mathbf{V}^{*\top} \mathcal{R}_\Omega(\mathbf{X}^*) \mathbf{V}^*) - \sigma_{r+1}(\mathcal{R}_\Omega(\mathbf{X}^*)) \\ &\geq \sigma_r(\mathbf{V}^{*\top} \mathbf{X}^* \mathbf{V}^*) - \|\mathcal{R}_\Omega(\mathbf{X}^*) - \mathbf{X}^*\| - \sigma_{r+1}(\mathcal{R}_\Omega(\mathbf{X}^*)) \\ &\geq \sigma_r^* - 2\|\mathcal{R}_\Omega(\mathbf{X}^*) - \mathbf{X}^*\| \\ &\geq 0.5\sigma_r^*. \end{aligned} \quad (204)$$

Combining the above two inequalities, we conclude that, with probability at least $1 - \frac{1}{d^3}$, Condition 1 is satisfied with $c_0 = \frac{1}{2\kappa}$.

Appendix G. Concentration Inequalities for Matrix Completion

Recall that the sampling matrix $\Omega \in \mathbb{R}^{d \times d}$ is defined as

$$\Omega_{i,j} = \begin{cases} 1 & \text{if } (i, j) \in \Omega, \\ 0 & \text{otherwise.} \end{cases} \quad (205)$$

The following lemma characterizes the concentration behavior of Ω .

Lemma 25 (Adapted from (Vu, 2018, Lemma 8)) *Suppose the sampling rate satisfies $p \gtrsim \frac{\log(d)}{d}$. There is a universal constant $\Gamma > 0$ such that, with probability at least $1 - \frac{1}{d^3}$, we have*

$$\left\| \frac{\Omega + \Omega^\top}{2p} - \mathbf{J} \right\| \leq \Gamma \sqrt{\frac{d}{p}}. \quad (206)$$

Here \mathbf{J} is the all-one matrix.

The original result appeared in (Vu, 2018, Lemma 8) only holds for symmetric Bernoulli model, i.e., $\Omega = \Omega^\top$. However, we can easily extend it to the asymmetric case via the dilation trick (Tropp et al., 2015). Hence, we omit the details here. Next, we have the following extension to the leave-one-out sequences.

Corollary 26 *Suppose the sampling rate satisfies $p \gtrsim \frac{\log(d)}{d}$. Then, with probability at least $1 - \frac{1}{d^3}$, for all $1 \leq l \leq d$, we have*

$$\left\| \frac{\Omega^{(l)} + \Omega^{(l)\top}}{2p} - \mathbf{J} \right\| \leq \Gamma \sqrt{\frac{d}{p}}. \quad (207)$$

Proof of Corollary 26. Note that for any $1 \leq l \leq d$, the matrix $\frac{\Omega^{(l)} + \Omega^{(l)\top}}{2p} - \mathbf{J}$ can be derived from $\frac{\Omega + \Omega^\top}{2p} - \mathbf{J}$ by zeroing out the l -th row and column. Hence, the proof follows by invoking Lemma 34 in Lemma 25. \blacksquare

Lemma 27 ((Chen et al., 2020, Lemma A.1) and (Chen and Li, 2019, Lemma 8)) *For all A, B, C , and $D \in \mathbb{R}^{d \times r}$, we have*

$$\left| \left\langle (\mathcal{I} - \mathcal{R}_\Omega) (AC^\top), BD^\top \right\rangle \right| \leq \left\| \frac{\Omega + \Omega^\top}{2p} - \mathbf{J} \right\| \cdot \|A\|_{2,\infty} \|B\|_F \cdot \|C\|_F \|D\|_{2,\infty}. \quad (208)$$

Moreover, for all $1 \leq l \leq d$, we have

$$\left| \left\langle (\mathcal{I} - \mathcal{R}_{\Omega^{(l)}}) (AC^\top), BD^\top \right\rangle \right| \leq \left\| \frac{\Omega^{(l)} + \Omega^{(l)\top}}{2p} - \mathbf{J} \right\| \cdot \|A\|_{2,\infty} \|B\|_F \cdot \|C\|_F \|D\|_{2,\infty}. \quad (209)$$

As a special case, we have

Lemma 28 For any matrix $\mathbf{X} \in \mathbb{R}^{d \times d}$ with the form $\mathbf{X} = \mathbf{U}\mathbf{V}^\top$, we have

$$\|(\mathcal{I} - \mathcal{R}_{\Omega^{(l)}})(\mathbf{X})\| \leq \|(\mathcal{I} - \mathcal{R}_\Omega)(\mathbf{X})\| \leq \left\| \frac{\boldsymbol{\Omega} + \boldsymbol{\Omega}^\top}{2p} - \mathbf{J} \right\| \|\mathbf{U}\|_{2,\infty} \|\mathbf{V}\|_{2,\infty}, \quad \forall 1 \leq l \leq d. \quad (210)$$

Note that the above two results are deterministic. Combining them with Lemma 25 leads to the following high-probability results.

Corollary 29 ((Chen et al., 2020, Lemma 4.3 and Lemma A.1)) Suppose that the sampling rate satisfies $p \gtrsim \frac{\log(d)}{d}$. There exists a universal constant $C > 0$ such that, for any $\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D} \in \mathbb{R}^{d \times r}$, with probability at least $1 - \frac{1}{d^3}$, we have

$$\left| \left\langle (\mathcal{I} - \mathcal{R}_\Omega)(\mathbf{A}\mathbf{C}^\top), \mathbf{B}\mathbf{D}^\top \right\rangle \right| \leq \Gamma \sqrt{\frac{d}{p}} \cdot \|\mathbf{A}\|_{2,\infty} \|\mathbf{B}\|_{\text{F}} \cdot \|\mathbf{C}\|_{\text{F}} \|\mathbf{D}\|_{2,\infty}. \quad (211)$$

Moreover, with the same probability, for any $1 \leq l \leq d$, we have

$$\left| \left\langle (\mathcal{I} - \mathcal{R}_{\Omega^{(l)}})(\mathbf{A}\mathbf{C}^\top), \mathbf{B}\mathbf{D}^\top \right\rangle \right| \leq \Gamma \sqrt{\frac{d}{p}} \cdot \|\mathbf{A}\|_{2,\infty} \|\mathbf{B}\|_{\text{F}} \cdot \|\mathbf{C}\|_{\text{F}} \|\mathbf{D}\|_{2,\infty}. \quad (212)$$

Corollary 30 ((Chen and Li, 2019, Lemma 9)) Consider an arbitrary matrix $\mathbf{X} \in \mathbb{R}^{d \times d}$ decomposed as $\mathbf{X} = \mathbf{U}\mathbf{V}^\top$. There exists a universal constant C such that, with probability at least $1 - \frac{1}{d^3}$, we have

$$\left\| (\mathcal{I} - \mathcal{R}_{\Omega^{(l)}})(\mathbf{X}) \right\| \leq \|(\mathcal{I} - \mathcal{R}_\Omega)(\mathbf{X})\| \leq \Gamma \sqrt{\frac{d}{p}} \|\mathbf{U}\|_{2,\infty} \|\mathbf{V}\|_{2,\infty}. \quad (213)$$

Lastly, we provide a finer result for a matrix of the form $\mathbf{X} - \mathbf{Y}$.

Lemma 31 For two arbitrary symmetric matrices $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{d \times d}$, with probability at least $1 - \frac{1}{d^3}$, we have

$$\left\| (\mathcal{I} - \mathcal{R}_{\Omega^{(l)}})(\mathbf{X} - \mathbf{Y}) \right\| \leq \|(\mathcal{I} - \mathcal{R}_\Omega)(\mathbf{X} - \mathbf{Y})\| \leq \Gamma \sqrt{\frac{d}{p}} \|\mathbf{X} - \mathbf{Y}\| \left(\|\mathbf{U}_\mathbf{X}\|_{2,\infty}^2 + \|\mathbf{U}_\mathbf{Y}\|_{2,\infty}^2 \right). \quad (214)$$

Proof We denote $\mathbf{Z} = \mathbf{X} - \mathbf{Y}$ with an SVD $\mathbf{Z} = \mathbf{L}_\mathbf{Z} \boldsymbol{\Sigma}_\mathbf{Z} \mathbf{L}_\mathbf{Z}^\top$ (recall that \mathbf{Z} is symmetric). Hence, with probability at least $1 - \frac{1}{d^3}$, we have

$$\begin{aligned} \|(\mathcal{I} - \mathcal{R}_\Omega)(\mathbf{Z})\| &\stackrel{(a)}{\leq} \Gamma \sqrt{\frac{d}{p}} \|\mathbf{L}_\mathbf{Z} \boldsymbol{\Sigma}_\mathbf{Z}\|_{2,\infty} \|\mathbf{L}_\mathbf{Z}\|_{2,\infty} \\ &\stackrel{(b)}{\leq} \Gamma \sqrt{\frac{d}{p}} \|\boldsymbol{\Sigma}_\mathbf{Z}\| \|\mathbf{L}_\mathbf{Z}\|_{2,\infty}^2 \\ &= \Gamma \sqrt{\frac{d}{p}} \|\mathbf{Z}\| \|\mathbf{L}_\mathbf{Z}\|_{2,\infty}^2 \\ &\stackrel{(c)}{\leq} \Gamma \sqrt{\frac{d}{p}} \|\mathbf{Z}\| \left(\|\mathbf{L}_\mathbf{X}\|_{2,\infty}^2 + \|\mathbf{L}_\mathbf{Y}\|_{2,\infty}^2 \right). \end{aligned} \quad (215)$$

Here in (a), we apply Lemma 30 upon setting $U = L_Z \Sigma_Z$ and $V = L_Z$. In (b), we apply Lemma 33. Lastly, in (c), we apply Lemma 39. This completes the proof. \blacksquare

Appendix H. Auxiliary Lemmas

H.1. Concentration Inequalities

Lemma 32 (Bernstein's inequality) *Let X_1, \dots, X_n be independent zero-mean random variables. Suppose that $|X_i| \leq M$ almost surely, for all i and set $\nu^2 = \sum_{i=1}^n \text{Var}[X_i^2]$. Then, for all positive t ,*

$$\mathbb{P}\left(\left|\sum_{i=1}^n X_i\right| \geq t\right) \leq 2 \exp\left(-\frac{t^2}{2\nu^2 + \frac{2}{3}Mt}\right). \quad (216)$$

Or equivalently, with probability at least $1 - \delta$, one has

$$\left|\sum_{i=1}^n X_i\right| \leq 2\nu\sqrt{\log(2/\delta)} + \frac{4}{3}M \log(2/\delta). \quad (217)$$

H.2. Matrix Norm Inequalities

Lemma 33 ((Cape et al., 2019, Proposition 6.5)) *For $A \in \mathbb{R}^{d_1 \times d_2}$, and $B \in \mathbb{R}^{d_2 \times d_3}$, we have*

$$\|AB\|_{2,\infty} \leq \|A\|_{2,\infty} \|B\| \quad \text{and} \quad \|AB\|_{2,\infty} \leq \|A\|_\infty \|B\|_{2,\infty}. \quad (218)$$

Lemma 34 (Adapted from (Sun and Luo, 2016, Proposition A.3)) *For any matrix $A \in \mathbb{R}^{d_1 \times d_2}$, denote $A_{-i,\cdot}$ ($A_{\cdot,-i}$) as the matrix obtained by replacing the i -th row (column) of A by zeros, respectively. Then, we have*

$$\|A_{-i,\cdot}\| \leq \|A\| \quad \text{and} \quad \|A_{\cdot,-i}\| \leq \|A\|, \forall i \in [d_1], j \in [d_2]. \quad (219)$$

Lemma 35 ((Sun and Luo, 2016, Proposition A.4)) *For any two matrices $A \in \mathbb{R}^{d_1 \times d_2}$, $B \in \mathbb{R}^{d_2 \times d_3}$, we have*

$$\|AB\| \leq \|A\| \|B\| \quad \text{and} \quad \|AB\|_F \leq \|A\| \|B\|_F. \quad (220)$$

Furthermore, if $d_1 \geq d_2$, then

$$\sigma_{\min}(A) \|B\|_F \leq \|AB\|_F \quad \text{and} \quad \sigma_{\min}(A) \|B\| \leq \|AB\|. \quad (221)$$

Lemma 36 *For arbitrary matrices $U \in \mathbb{R}^{d_1 \times d_2}$, $\Sigma \in \mathbb{R}^{d_2 \times d_3}$ and $V \in \mathbb{R}^{d_3 \times d_4}$, we have*

$$\left\|U\Sigma V^\top\right\|_{\max} \leq \|\Sigma\| \|U\|_{2,\infty} \|V\|_{2,\infty}. \quad (222)$$

Proof By Cauchy-Schwartz inequality, we have

$$\begin{aligned} \left\|U\Sigma V^\top\right\|_{\max} &= \max_{i,j} \left| \sum_k [U\Sigma]_{i,k} V_{j,k} \right| \\ &\leq \max_{i,j} \|(U\Sigma)_{i,\cdot}\| \|V_{j,\cdot}\| \\ &= \|U\Sigma\|_{2,\infty} \|V\|_{2,\infty} \\ &\leq \|\Sigma\| \|U\|_{2,\infty} \|V\|_{2,\infty}. \end{aligned} \quad (223)$$

Here we apply Lemma 33 to derive the last inequality. This completes the proof. \blacksquare

Lemma 37 For any matrix $\mathbf{A} \in \mathbb{R}^{d_1 \times d_2}$, we have

$$\|\mathbf{A}\|_{2,\infty}^2 = \|\mathbf{A}\mathbf{A}^\top\|_{\max}. \quad (224)$$

Proof We write $\mathbf{a}_1, \dots, \mathbf{a}_{d_1}$ as the row vectors of \mathbf{A} . Then, we have

$$\|\mathbf{A}\|_{2,\infty}^2 = \max_{1 \leq i \leq d_1} \|\mathbf{a}_i\|^2. \quad (225)$$

On the other hand, we can write $\|\mathbf{A}\mathbf{A}^\top\|_{\max}$ as

$$\|\mathbf{A}\mathbf{A}^\top\|_{\max} = \max_{1 \leq i, j \leq d_1} |\langle \mathbf{a}_i, \mathbf{a}_j \rangle|. \quad (226)$$

First, we have $\|\mathbf{A}\mathbf{A}^\top\|_{\max} \leq \|\mathbf{A}\|_{2,\infty}^2$ since

$$\|\mathbf{A}\mathbf{A}^\top\|_{\max} = \max_{1 \leq i, j \leq d_1} |\langle \mathbf{a}_i, \mathbf{a}_j \rangle| \leq \max_{1 \leq i, j \leq d_1} \|\mathbf{a}_i\| \|\mathbf{a}_j\| = \max_{1 \leq i \leq d_1} \|\mathbf{a}_i\|^2 = \|\mathbf{A}\|_{2,\infty}^2. \quad (227)$$

Second, we have $\|\mathbf{A}\mathbf{A}^\top\|_{\max} \geq \|\mathbf{A}\|_{2,\infty}^2$ upon noting that

$$\|\mathbf{A}\mathbf{A}^\top\|_{\max} = \max_{1 \leq i, j \leq d_1} |\langle \mathbf{a}_i, \mathbf{a}_j \rangle| \geq \max_{i=j} |\langle \mathbf{a}_i, \mathbf{a}_i \rangle| = \max_{1 \leq i \leq d_1} \|\mathbf{a}_i\|^2 = \|\mathbf{A}\|_{2,\infty}^2. \quad (228)$$

Therefore, we derive that $\|\mathbf{A}\mathbf{A}^\top\|_{\max} = \|\mathbf{A}\|_{2,\infty}^2$, which completes the proof. \blacksquare

Lemma 38 For two PSD matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{d \times d}$ with $\mathbf{A} \preceq \mathbf{B}$, we have

$$\|\mathbf{A}\|_{\max} \leq \|\mathbf{B}\|_{\max}. \quad (229)$$

Proof The proof follows by the fact that for any PSD matrix \mathbf{A} , we have $\|\mathbf{A}\|_{\max} = \max_i \{A_{i,i}\}$. According to this fact, we immediately have

$$\|\mathbf{A}\|_{\max} = \max_i \{A_{i,i}\} \leq \max_i \{B_{i,i}\} = \|\mathbf{B}\|_{\max} \quad (230)$$

since $\mathbf{A} \preceq \mathbf{B}$. Now we turn to prove this fact. Note that we can write any PSD matrix \mathbf{A} as $\mathbf{A} = \mathbf{P}\mathbf{P}^\top$. Then, according to Lemma 37, we have

$$\|\mathbf{A}\|_{\max} = \|\mathbf{P}\|_{2,\infty}^2 = \max_i \|\mathbf{p}_i\|^2 = \max_i \{A_{i,i}\}. \quad (231)$$

Here we write $\{\mathbf{p}_i\}$ as the row vectors of \mathbf{P} . This completes the proof. \blacksquare

Lemma 39 For $\mathbf{Z} = \mathbf{X} - \mathbf{Y}$, we have

$$\|\mathbf{L}\mathbf{Z}\|_{2,\infty}^2 \leq \|\mathbf{L}\mathbf{X}\|_{2,\infty}^2 + \|\mathbf{L}\mathbf{Y}\|_{2,\infty}^2. \quad (232)$$

Proof We bound $\|\mathbf{L}_Z\|_{2,\infty}^2$ as follows

$$\begin{aligned}
 \|\mathbf{L}_Z\|_{2,\infty}^2 &\stackrel{(a)}{=} \left\| \mathbf{L}_Z \mathbf{L}_Z^\top \right\|_{\max} \\
 &\stackrel{(b)}{\leq} \left\| \mathbf{L}_X \mathbf{L}_X^\top + \mathbf{L}_Y \mathbf{L}_Y^\top \right\|_{\max} \\
 &\leq \left\| \mathbf{L}_X \mathbf{L}_X^\top \right\|_{\max} + \left\| \mathbf{L}_Y \mathbf{L}_Y^\top \right\|_{\max} \\
 &\stackrel{(c)}{=} \|\mathbf{L}_X\|_{2,\infty}^2 + \|\mathbf{L}_Y\|_{2,\infty}^2.
 \end{aligned} \tag{233}$$

Here (a) and (c) follow from Lemma 37. In (b), we apply Lemma 38 since $\mathbf{L}_Z \mathbf{L}_Z^\top \preceq \mathbf{L}_X \mathbf{L}_X^\top + \mathbf{L}_Y \mathbf{L}_Y^\top$. This is due to the fact that $\text{col}(\mathbf{Z}) \subseteq \text{col}(\mathbf{X}) \oplus \text{col}(\mathbf{Y})$ which leads to $\mathbf{L}_Z \mathbf{L}_Z^\top = \mathcal{P}_Z \preceq \mathcal{P}_X + \mathcal{P}_Y$. This completes the proof. \blacksquare

Lemma 40 ((Tu et al., 2016, Lemma 5.4)) *For any $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{d \times r}$ with $\sigma_r(\mathbf{X}) > 0$, we have*

$$\text{dist}^2(\mathbf{X}, \mathbf{Y}) \leq \frac{1}{2(\sqrt{2}-1)\sigma_r^2(\mathbf{X})} \left\| \mathbf{X} \mathbf{X}^\top - \mathbf{Y} \mathbf{Y}^\top \right\|_{\text{F}}^2. \tag{234}$$

Lemma 41 *Consider a matrix $\mathbf{U} \in \mathbb{R}^{d_1 \times d_2}$ and a diagonal matrix $\mathbf{\Sigma} \in \mathbb{R}^{d_2 \times d_2}$. We have*

$$\|\mathbf{U} \mathbf{\Sigma}\|_{2,\infty} \geq \sigma_r(\mathbf{\Sigma}) \|\mathbf{U}\|_{2,\infty}. \tag{235}$$

Proof We first write $\mathbf{\Sigma} = \text{Diag}\{\sigma_1, \dots, \sigma_{d_2}\}$. Next, note that

$$\|(U \mathbf{\Sigma})_{l,\cdot}\|^2 = \sum_{j=1}^{d_2} U_{l,j}^2 \sigma_j^2 \geq \sigma_r^2(\mathbf{\Sigma}) \sum_{j=1}^{d_2} U_{l,j}^2 = \sigma_r^2(\mathbf{\Sigma}) \|U_{l,\cdot}\|^2. \tag{236}$$

Hence, taking the maximum over index l on both sides, we immediately obtain

$$\|\mathbf{U} \mathbf{\Sigma}\|_{2,\infty} \geq \sigma_r(\mathbf{\Sigma}) \|\mathbf{U}\|_{2,\infty}. \tag{237}$$

\blacksquare

H.3. Other Useful Inequalities

Lemma 42 (Adapted from (Vershynin, 2018, Corollary 4.2.13)) *The covering number \mathcal{N}_ϵ of $\mathcal{B}_{\text{F}}^{r \times r}(R)$ satisfies the following inequality for any $0 < \epsilon \leq 1$:*

$$\mathcal{N}_\epsilon \leq \left(\frac{3R}{\epsilon} \right)^{r^2}. \tag{238}$$

Lemma 43 *For two orthogonal matrices $\mathbf{V}_1, \mathbf{V}_2 \in \mathcal{O}_{d \times r}$, we have*

$$\left\| \mathbf{V}_1 \mathbf{V}_1^\top - \mathbf{V}_2 \mathbf{V}_2^\top \right\| \leq 2 \|\mathbf{V}_1 - \mathbf{V}_2\| \quad \text{and} \quad \left\| \mathbf{V}_1 \mathbf{V}_1^\top - \mathbf{V}_2 \mathbf{V}_2^\top \right\|_{\text{F}} \leq 2 \|\mathbf{V}_1 - \mathbf{V}_2\|_{\text{F}}. \tag{239}$$

Proof Note that $\mathbf{V}_1 \mathbf{V}_1^\top - \mathbf{V}_2 \mathbf{V}_2^\top = \mathbf{V}_1 (\mathbf{V}_1 - \mathbf{V}_2)^\top + (\mathbf{V}_1 - \mathbf{V}_2) \mathbf{V}_2^\top$. Hence, we have

$$\left\| \mathbf{V}_1 \mathbf{V}_1^\top - \mathbf{V}_2 \mathbf{V}_2^\top \right\| \leq \|\mathbf{V}_1 - \mathbf{V}_2\| (\|\mathbf{V}_1\| + \|\mathbf{V}_2\|) \leq 2 \|\mathbf{V}_1 - \mathbf{V}_2\|. \quad (240)$$

Similarly, for the Frobenius norm, we also have

$$\left\| \mathbf{V}_1 \mathbf{V}_1^\top - \mathbf{V}_2 \mathbf{V}_2^\top \right\|_F \leq \|\mathbf{V}_1 - \mathbf{V}_2\|_F (\|\mathbf{V}_1\| + \|\mathbf{V}_2\|) \leq 2 \|\mathbf{V}_1 - \mathbf{V}_2\|_F. \quad (241)$$

This completes the proof. ■

Lemma 44 For arbitrary matrix $\mathbf{X} \in \mathbb{R}^{d_1 \times d_2}$ with $\text{rank}(\mathbf{X}) = r$ and $\mathbf{O} \in \mathcal{O}_{d_1 \times r}$, we have

$$\sigma_r(\mathbf{X}) \geq \sigma_r(\mathbf{O}^\top \mathbf{X}). \quad (242)$$

Moreover, if $\mathbf{O} = \mathbf{U}_\mathbf{X}$, we have $\sigma_i(\mathbf{X}) = \sigma_i(\mathbf{O}^\top \mathbf{X})$ for all $1 \leq i \leq r$.

Proof We first prove the special case. Suppose $\mathbf{O} = \mathbf{U}_\mathbf{X}$, then we have $\mathbf{O}^\top \mathbf{X} = \Sigma_\mathbf{X} \mathbf{V}_\mathbf{X}^\top$. Note that this is the SVD of $\mathbf{O}^\top \mathbf{X}$ with the singular value matrix $\Sigma_\mathbf{X}$. Hence, $\mathbf{O}^\top \mathbf{X}$ has the same singular values as \mathbf{X} . For the general case, it follows by

$$\begin{aligned} \sigma_r(\mathbf{O}^\top \mathbf{X}) &= \inf_{\mathbf{Y} \in \mathcal{M}_{\leq r-1}} \left\| \mathbf{O}^\top \mathbf{X} - \mathbf{Y} \right\|_F \\ &\leq \left\| \mathbf{O}^\top \mathbf{X} - \mathbf{O}^\top \mathbf{U}_\mathbf{X} \Sigma_{\mathbf{X}, -1} \mathbf{V}_\mathbf{X}^\top \right\|_F \\ &\leq \left\| \mathbf{O}^\top \sigma_r(\mathbf{X}) \mathbf{u}_{\mathbf{X}, -1} \mathbf{v}_{\mathbf{X}, -1}^\top \right\|_F \\ &\leq \sigma_r(\mathbf{X}). \end{aligned} \quad (243)$$

Here we define $\mathcal{M}_{\leq r-1}$ as the set of all matrices of rank at most $r - 1$. We also denote $\Sigma_{\mathbf{X}, -1} = \text{Diag}\{\sigma_1^*(\mathbf{X}), \dots, \sigma_{r-1}(\mathbf{X}), 0\}$. Finally, $\mathbf{u}_{\mathbf{X}, -1}, \mathbf{v}_{\mathbf{X}, -1}$ refer to the last columns of $\mathbf{U}_\mathbf{X}, \mathbf{V}_\mathbf{X}$, respectively. ■

Lemma 45 Consider two matrices $\mathbf{A} \in \mathbb{R}^{r \times r}, \mathbf{B} \in \mathbb{R}^{r \times d}$ where \mathbf{A} is invertible and $\mathbf{B} \neq \mathbf{0}$. We have

$$\sigma_r(\mathbf{AB}) \geq \sigma_r(\mathbf{A}) \sigma_r(\mathbf{B}). \quad (244)$$

Proof It directly follows by

$$\begin{aligned} \sigma_r(\mathbf{AB}) &= \min_{\|\mathbf{x}\|=1, \mathbf{x} \in \text{range}(\mathbf{B}^\top)} \|\mathbf{ABx}\| \\ &\geq \sigma_r(\mathbf{A}) \min_{\|\mathbf{x}\|=1, \mathbf{x} \in \text{range}(\mathbf{B}^\top)} \|\mathbf{Bx}\| \\ &= \sigma_r(\mathbf{A}) \sigma_r(\mathbf{B}). \end{aligned} \quad (245)$$

■

Lemma 46 For any $x \in [0, \frac{1}{2(r-1)})$ and $r > 1$, we have

$$(1+x)^r \leq 1+2rx. \quad (246)$$

Lemma 47 For the series $\{x_t\}_{t=0}^\infty$, the following two statements hold:

- Suppose that $x_{t+1} \leq Ax_t + B, \forall t \geq 0$ where $A > 0, A \neq 1$ and $x_0 + \frac{B}{A-1} \geq 0$. Then, we have

$$x_t \leq A^t \left(x_0 + \frac{B}{A-1} \right) - \frac{B}{A-1}. \quad (247)$$

- Suppose that $x_{t+1} \geq Ax_t - B, \forall t \geq 0$ where $A > 0$ and $x_0 - \frac{B}{A-1} \geq 0$. Then, we have

$$x_t \geq A^t \left(x_0 - \frac{B}{A-1} \right) + \frac{B}{A-1}. \quad (248)$$

Proof We first prove the first statement. Note that we can rewrite $x_{t+1} \leq Ax_t + B$ as $y_{t+1} \leq Ay_t$ where $y_t = x_t + \frac{B}{A-1}$. Note that $y_0 \geq 0, A > 0$ by our assumption. Hence, we derive $y_t \leq A^t y_0 = A^t \left(x_0 + \frac{B}{A-1} \right)$, which implies that $x_t = y_t - \frac{B}{A-1} \leq A^t \left(x_0 + \frac{B}{A-1} \right) - \frac{B}{A-1}$.

For the second statement, we first rewrite the condition as $y_{t+1} \geq Ay_t$ where $y_t = x_t - \frac{B}{A-1}$. Then, we have $y_t \geq A^t y_0$, which implies that $x_t = y_t + \frac{B}{A-1} \geq A^t y_0 + \frac{B}{A-1} = A^t \left(x_0 - \frac{B}{A-1} \right) + \frac{B}{A-1}$. ■