

# Noisy Matrix Completion: Understanding Statistical Guarantees for Convex Relaxation via Nonconvex Optimization

Yuxin Chen\*    Yuejie Chi†    Jianqing Fan‡    Cong Ma‡    Yuling Yan‡

October 13, 2019

## Abstract

This paper studies noisy low-rank matrix completion: given partial and noisy entries of a large low-rank matrix, the goal is to estimate the underlying matrix faithfully and efficiently. Arguably one of the most popular paradigms to tackle this problem is convex relaxation, which achieves remarkable efficacy in practice. However, the theoretical support of this approach is still far from optimal in the noisy setting, falling short of explaining its empirical success.

We make progress towards demystifying the practical efficacy of convex relaxation vis-à-vis random noise. When the rank and the condition number of the unknown matrix are bounded by a constant, we demonstrate that the convex programming approach achieves near-optimal estimation errors — in terms of the Euclidean loss, the entrywise loss, and the spectral norm loss — for a wide range of noise levels. All of this is enabled by bridging convex relaxation with the nonconvex Burer–Monteiro approach, a seemingly distinct algorithmic paradigm that is provably robust against noise. More specifically, we show that an approximate critical point of the nonconvex formulation serves as an extremely tight approximation of the convex solution, thus allowing us to transfer the desired statistical guarantees of the nonconvex approach to its convex counterpart.

**Keywords:** matrix completion, minimaxity, stability, convex relaxation, nonconvex optimization, Burer–Monteiro approach.

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Convex relaxation: limitations of prior results	3
1.2	A detour: nonconvex optimization	3
1.3	Empirical evidence: convex and nonconvex solutions are often close	4
1.4	Models and main results	5
1.4.1	Models and assumptions	5
1.4.2	Theoretical guarantees: when both the rank and the condition number are constants	6
1.4.3	Theoretical guarantees: extensions to more general settings	8
<b>2</b>	<b>Strategy and novelty</b>	<b>10</b>
2.1	Exact duality	10
2.2	A candidate primal solution via nonconvex optimization	11
2.3	Approximate nonconvex optimizers	12
2.4	Construction of an approximate nonconvex optimizer	12
2.5	Properties of the nonconvex iterates	13
2.6	Proof of Theorem 2	14

---

Author names are sorted alphabetically.

\*Department of Electrical Engineering, Princeton University, Princeton, NJ 08544, USA; Email: [yuxin.chen@princeton.edu](mailto:yuxin.chen@princeton.edu).

†Department of Electrical and Computer Engineering, Carnegie Mellon University, Pittsburgh, PA 15213, USA; Email: [yuejiechi@cmu.edu](mailto:yuejiechi@cmu.edu).

‡Department of Operations Research and Financial Engineering, Princeton University, Princeton, NJ 08544, USA; Email: [{jqfan, congm, yulingy}@princeton.edu">jqfan, congm, yulingy}@princeton.edu](mailto).

<b>3</b>	<b>Prior art</b>	<b>15</b>
<b>4</b>	<b>Discussion</b>	<b>16</b>
<b>A</b>	<b>Preliminaries</b>	<b>23</b>
<b>B</b>	<b>Exact duality analysis</b>	<b>23</b>
<b>C</b>	<b>Connections between convex and nonconvex solutions</b>	<b>24</b>
C.1	Proof of Lemma 1	24
C.2	Proof of Lemma 2	26
C.2.1	Proof of Claim 2	28
C.2.2	Proof of Claim 3	31
C.3	Proof of Lemma 4	31
C.3.1	Proof of Lemma 7	32
C.3.2	Proof of Lemma 8	36
<b>D</b>	<b>Analysis of the nonconvex gradient descent algorithm</b>	<b>36</b>
D.1	Preliminaries and notations	39
D.2	Proof of Lemma 9	41
D.3	Proof of Lemma 10	42
D.4	Proof of Lemma 11	44
D.5	Proof of Lemma 12	48
D.6	Proof of Lemma 13	52
D.7	Proof of Lemma 14	55
D.8	Proof of Lemma 15	56
D.9	Proof of Lemma 16	58
D.10	Proof of Lemma 17	59
D.11	Proof of Lemma 18	60
D.12	Proof of the inequalities (31)	61
<b>E</b>	<b>Technical lemmas</b>	<b>62</b>

# 1 Introduction

Suppose we are interested in a large low-rank data matrix, but only get to observe a highly incomplete subset of its entries. Can we hope to estimate the underlying data matrix in a reliable manner? This problem, often dubbed as *low-rank matrix completion*, spans a diverse array of science and engineering applications (e.g. collaborative filtering [RS05], localization [SY07], system identification [LV09], magnetic resonance parameter mapping [ZPL15], joint alignment [CC18a]), and has inspired a flurry of research activities in the past decade. In the statistics literature, matrix completion also falls under the category of factor models with a large amount of missing data, which finds numerous statistical applications such as controlling false discovery rates for dependence data [Efr07, Efr10, FHG12, FKSZ19], factor-adjusted variable selection [KS11, FKW18], principal component regression [Jol82, BN06, PBHT08, FXY17], and large covariance matrix estimation [FLM13, FWZ19]. Recent years have witnessed the development of many tractable algorithms that come with statistical guarantees, with convex relaxation being one of the most popular paradigms [FHB04, CR09, CT10]. See [DR16, CC18b] for an overview of this topic.

This paper focuses on noisy low-rank matrix completion, assuming that the revealed entries are corrupted by a certain amount of noise. Setting the stage, consider the task of estimating a rank- $r$  data matrix  $\mathbf{M}^* = [M_{ij}^*]_{1 \leq i, j \leq n} \in \mathbb{R}^{n \times n}$ ,<sup>1</sup> and suppose that this needs to be performed on the basis of a subset of noisy entries

$$M_{ij} = M_{ij}^* + E_{ij}, \quad (i, j) \in \Omega, \quad (1)$$

---

<sup>1</sup>It is straightforward to rephrase our discussions to a general rectangular matrix of size  $n_1 \times n_2$ . The current paper sets  $n = n_1 = n_2$  throughout for simplicity of presentation.

where  $\Omega \subseteq \{1, \dots, n\} \times \{1, \dots, n\}$  denotes a set of indices, and  $E_{ij}$  stands for the additive noise at the location  $(i, j)$ . As we shall elaborate shortly, solving noisy matrix completion via convex relaxation, while practically exhibiting excellent stability (in terms of the estimation errors against noise), is far less understood theoretically compared to the noiseless setting.

## 1.1 Convex relaxation: limitations of prior results

Naturally, one would search for a low-rank solution that best fits the observed entries. One choice is the regularized least-squares formulation given by

$$\underset{\mathbf{Z} \in \mathbb{R}^{n \times n}}{\text{minimize}} \quad \frac{1}{2} \sum_{(i,j) \in \Omega} (Z_{ij} - M_{ij})^2 + \lambda \text{rank}(\mathbf{Z}), \quad (2)$$

where  $\lambda > 0$  is some regularization parameter. In words, this approach optimizes certain trade-off between the goodness of fit (through the squared loss expressed in the first term of (2)) and the low-rank structure (through the rank function in the second term of (2)). Due to computational intractability of rank minimization, we often resort to convex relaxation in order to obtain computationally feasible solutions. One notable example is the following convex program:

$$\underset{\mathbf{Z} \in \mathbb{R}^{n \times n}}{\text{minimize}} \quad g(\mathbf{Z}) \triangleq \frac{1}{2} \sum_{(i,j) \in \Omega} (Z_{ij} - M_{ij})^2 + \lambda \|\mathbf{Z}\|_*, \quad (3)$$

where  $\|\mathbf{Z}\|_*$  denotes the nuclear norm (i.e. the sum of singular values) of  $\mathbf{Z}$  — a convex surrogate for the rank function. A significant portion of existing theory supports the use of this paradigm in the noiseless setting: when  $E_{ij}$  vanishes for all  $(i, j) \in \Omega$ , the solution to (3) is known to be faithful (i.e. the estimation error becomes zero) even under near-minimal sample complexity [CR09, CP10, CT10, Gro11, Rec11, Che15].

By contrast, the performance of convex relaxation remains largely unclear when it comes to noisy settings (which are often more practically relevant). Candès and Plan [CP10] first studied the stability of an equivalent variant<sup>2</sup> of (3) against noise. The estimation error  $\|\mathbf{Z}_{\text{cvx}} - \mathbf{M}^*\|_F$  derived therein, of the solution  $\mathbf{Z}_{\text{cvx}}$  to (3), is significantly larger than the oracle lower bound. This does not explain well the effectiveness of (3) in practice. In fact, the numerical experiments reported in [CP10] already indicated that the performance of convex relaxation is far better than their theoretical bounds. This discrepancy between numerical performance and existing theoretical bounds gives rise to the following natural yet challenging questions: *Where does the convex program (3) stand in terms of its stability vis-à-vis additive noise? Can we establish statistical performance guarantees that match its practical effectiveness?*

We note in passing that several other convex relaxation formulations have been thoroughly analyzed for noisy matrix completion, most notably by Negahban and Wainwright [NW12] and by Koltchinskii et al. [KLT11]. These works have significantly advanced our understanding of the power of convex relaxation. However, the estimators studied therein, particularly the one in [KLT11], are quite different from the one (3) considered here; as a consequence, the analysis therein does not lead to improved statistical guarantees of (3). Moreover, the performance guarantees provided for these variants are also suboptimal when restricted to the class of “incoherent” or “de-localized” matrices, unless the magnitudes of the noise are fairly large. See Section 1.4 for more detailed discussions as well as numerical comparisons of these algorithms.

## 1.2 A detour: nonconvex optimization

While the focus of the current paper is convex relaxation, we take a moment to discuss a seemingly distinct algorithmic paradigm: nonconvex optimization, which turns out to be remarkably helpful in understanding convex relaxation. Inspired by the Burer–Monteiro approach [BM03], the nonconvex scheme starts by representing the rank- $r$  decision matrix (or parameters)  $\mathbf{Z}$  as  $\mathbf{Z} = \mathbf{X}\mathbf{Y}^\top$  via low-rank factors  $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times r}$ ,

<sup>2</sup>Technically, [CP10] deals with the constrained version of (3), which is equivalent to the Lagrangian form as in (3) with a proper choice of the regularization parameter.

and proceeds by solving the following nonconvex (regularized) least-squares problem [KMO10a]

$$\underset{\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times r}}{\text{minimize}} \quad \frac{1}{2} \sum_{(i,j) \in \Omega} [(\mathbf{X}\mathbf{Y}^\top)_{ij} - M_{ij}]^2 + \text{reg}(\mathbf{X}, \mathbf{Y}). \quad (4)$$

Here,  $\text{reg}(\cdot, \cdot)$  denotes a certain regularization term that promotes additional structural properties.

To see its intimate connection with the convex program (3), we make the following observation: if the solution to (3) has rank  $r$ , then it must coincide with the solution to

$$\underset{\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times r}}{\text{minimize}} \quad \frac{1}{2} \sum_{(i,j) \in \Omega} [(\mathbf{X}\mathbf{Y}^\top)_{ij} - M_{ij}]^2 + \underbrace{\frac{\lambda}{2} \|\mathbf{X}\|_{\text{F}}^2 + \frac{\lambda}{2} \|\mathbf{Y}\|_{\text{F}}^2}_{\text{reg}(\mathbf{X}, \mathbf{Y})}. \quad (5)$$

This can be easily verified by recognizing the elementary fact that

$$\|\mathbf{Z}\|_* = \inf_{\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times r}: \mathbf{X}\mathbf{Y}^\top = \mathbf{Z}} \left\{ \frac{1}{2} \|\mathbf{X}\|_{\text{F}}^2 + \frac{1}{2} \|\mathbf{Y}\|_{\text{F}}^2 \right\} \quad (6)$$

for any rank- $r$  matrix  $\mathbf{Z}$  [SS05, MHT10]. Note, however, that it is very challenging to predict when the key assumption in establishing this connection — namely, the rank- $r$  assumption of the solution to the convex program (3) — can possibly hold (and in particular, whether it can hold under minimal sample complexity requirement).

Despite the nonconvexity of (4), simple first-order optimization methods, in conjunction with proper initialization, are often effective in solving (4). Partial examples include gradient descent on manifold [KMO10a, KMO10b, WCCL16], gradient descent [SL16, MWCC17], and projected gradient descent [CW15, ZL16]. Apart from their practical efficiency, the nonconvex optimization approach is also appealing in theory. To begin with, algorithms tailored to (4) often enable exact recovery in the noiseless setting. Perhaps more importantly, for a wide range of noise settings, the nonconvex approach achieves appealing estimation accuracy [CW15, MWCC17], which could be significantly better than those bounds derived for convex relaxation discussed earlier. See [CLC19, CC18b] for a summary of recent results. Such intriguing statistical guarantees motivate us to take a closer inspection of the underlying connection between the two contrasting algorithmic frameworks.

### 1.3 Empirical evidence: convex and nonconvex solutions are often close

In order to obtain a better sense of the relationships between convex and nonconvex approaches, we begin by comparing the estimates returned by the two approaches via numerical experiments. Fix  $n = 1000$  and  $r = 5$ . We generate  $\mathbf{M}^* = \mathbf{X}^* \mathbf{Y}^{*\top}$ , where  $\mathbf{X}^*, \mathbf{Y}^* \in \mathbb{R}^{n \times r}$  are random orthonormal matrices. Each entry  $M_{ij}^*$  of  $\mathbf{M}^*$  is observed with probability  $p = 0.2$  independently, and then corrupted by an independent Gaussian noise  $E_{ij} \sim \mathcal{N}(0, \sigma^2)$ . Throughout the experiments, we set  $\lambda = 5\sigma\sqrt{np}$ . The convex program (3) is solved by the proximal gradient method [PB14], whereas we attempt solving the nonconvex formulation (5) by gradient descent with spectral initialization (see [CLC19] for details). Let  $\mathbf{Z}_{\text{cvx}}$  (resp.  $\mathbf{Z}_{\text{ncvx}} = \mathbf{X}_{\text{ncvx}} \mathbf{Y}_{\text{ncvx}}^\top$ ) be the solution returned by the convex program (3) (resp. the nonconvex program (5)). Figure 1 displays the relative estimation errors of both methods ( $\|\mathbf{Z}_{\text{cvx}} - \mathbf{M}^*\|_{\text{F}} / \|\mathbf{M}^*\|_{\text{F}}$  and  $\|\mathbf{Z}_{\text{ncvx}} - \mathbf{M}^*\|_{\text{F}} / \|\mathbf{M}^*\|_{\text{F}}$ ) as well as the relative distance  $\|\mathbf{Z}_{\text{cvx}} - \mathbf{Z}_{\text{ncvx}}\|_{\text{F}} / \|\mathbf{M}^*\|_{\text{F}}$  between the two estimates. The results are averaged over 20 independent trials.

Interestingly, the distance between the convex and the nonconvex solutions seems extremely small (e.g.  $\|\mathbf{Z}_{\text{cvx}} - \mathbf{Z}_{\text{ncvx}}\|_{\text{F}} / \|\mathbf{M}^*\|_{\text{F}}$  is typically below  $10^{-7}$ ); in comparison, the relative estimation errors of both  $\mathbf{Z}_{\text{cvx}}$  and  $\mathbf{Z}_{\text{ncvx}}$  are substantially larger. In other words, the estimate returned by the nonconvex approach serves as a remarkably accurate approximation of the convex solution. Given that the nonconvex approach is often guaranteed to achieve intriguing statistical guarantees vis-à-vis random noise [MWCC17], this suggests that the convex program is equally stable — a phenomenon that was not captured by prior theory [CP10]. *Can we leverage existing theory for the nonconvex scheme to improve the statistical analysis of the convex relaxation approach?*

Before continuing, we remark that the above numerical connection between convex relaxation (3) and nonconvex optimization (5) has already been observed multiple times in prior literature [Faz02, SS05, RFP10,

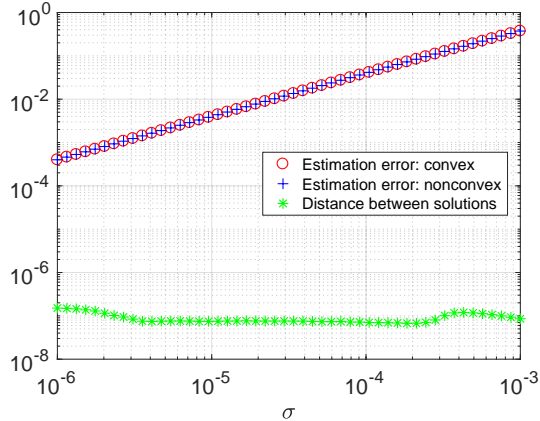


Figure 1: The relative estimation errors of both  $\mathbf{Z}_{\text{cvx}}$  (the estimate of the convex program (3)) and  $\mathbf{Z}_{\text{ncvx}}$  (the estimate returned by the nonconvex approach tailored to (5)) and the relative distance between them vs. the standard deviation  $\sigma$  of the noise. The results are reported for  $n = 1000$ ,  $r = 5$ ,  $p = 0.2$ ,  $\lambda = 5\sigma\sqrt{np}$  and are averaged over 20 independent trials.

[MHT10, KMO10b]. Nevertheless, all prior observations on this connection were either completely empirical, or provided in a way that does not lead to improved statistical error bounds of the convex paradigm (3). In fact, the difficulty in rigorously justifying the above numerical observations has been noted in the literature; see e.g. [KMO10b].<sup>3</sup>

## 1.4 Models and main results

The numerical experiments reported in Section 1.3 suggest an alternative route for analyzing convex relaxation for noisy matrix completion. If one can formally justify the proximity between the convex and the nonconvex solutions, then it is possible to propagate the appealing stability guarantees from the nonconvex scheme to the convex approach. As it turns out, this simple idea leads to significantly enhanced statistical guarantees for the convex program (3), which we formally present in this subsection.

### 1.4.1 Models and assumptions

Before proceeding, we introduce a few model assumptions that play a crucial role in our theory.

**Assumption 1.**

- (a) **(Random sampling)** Each index  $(i, j)$  belongs to the index set  $\Omega$  independently with probability  $p$ .
- (b) **(Random noise)** The noise matrix  $\mathbf{E} = [E_{ij}]_{1 \leq i, j \leq n}$  is composed of i.i.d. zero-mean sub-Gaussian random variables with sub-Gaussian norm at most  $\sigma > 0$ , i.e.  $\|E_{ij}\|_{\psi_2} \leq \sigma$  (see [Ver12, Definition 5.7]).

In addition, let  $\mathbf{M}^* = \mathbf{U}^* \mathbf{\Sigma}^* \mathbf{V}^{*\top}$  be the singular value decomposition (SVD) of  $\mathbf{M}^*$ , where  $\mathbf{U}^*, \mathbf{V}^* \in \mathbb{R}^{n \times r}$  consist of orthonormal columns and  $\mathbf{\Sigma}^* = \text{diag}(\sigma_1^*, \sigma_2^*, \dots, \sigma_r^*) \in \mathbb{R}^{r \times r}$  is a diagonal matrix obeying  $\sigma_{\max} \triangleq \sigma_1^* \geq \sigma_2^* \geq \dots \geq \sigma_r^* \triangleq \sigma_{\min}$ . Denote by  $\kappa \triangleq \sigma_{\max}/\sigma_{\min}$  the condition number of  $\mathbf{M}^*$ . We impose the following incoherence condition on  $\mathbf{M}^*$ , which is known to be crucial for reliable recovery of  $\mathbf{M}^*$  [CR09, Che15].

**Definition 1.** A rank- $r$  matrix  $\mathbf{M}^* \in \mathbb{R}^{n \times n}$  with SVD  $\mathbf{M}^* = \mathbf{U}^* \mathbf{\Sigma}^* \mathbf{V}^{*\top}$  is said to be  $\mu$ -incoherent if

$$\|\mathbf{U}^*\|_{2, \infty} \leq \sqrt{\frac{\mu}{n}} \|\mathbf{U}^*\|_{\text{F}} = \sqrt{\frac{\mu r}{n}} \quad \text{and} \quad \|\mathbf{V}^*\|_{2, \infty} \leq \sqrt{\frac{\mu}{n}} \|\mathbf{V}^*\|_{\text{F}} = \sqrt{\frac{\mu r}{n}}.$$

<sup>3</sup>The seminal work [KMO10b] by Keshavan, Montanari and Oh stated that “In view of the identity (6) it might be possible to use the results in this paper to prove stronger guarantees on the nuclear norm minimization approach. Unfortunately this implication is not immediate . . . Trying to establish such an implication, and clarifying the relation between the two approaches is nevertheless a promising research direction.”

Here,  $\|\mathbf{U}\|_{2,\infty}$  denotes the largest  $\ell_2$  norm of all rows of a matrix  $\mathbf{U}$ .

**Remark 1.** It is worth noting that several other conditions on the low-rank matrix have been proposed in the noisy setting. Examples include the spikiness condition [NW12] and the bounded  $\ell_\infty$  norm condition [KLT11]. However, these conditions alone are often unable to ensure identifiability of the true matrix even in the absence of noise.

#### 1.4.2 Theoretical guarantees: when both the rank and the condition number are constants

With these in place, we are positioned to present our improved statistical guarantees for convex relaxation. For convenience of presentation, we shall begin with a simple yet fundamentally important class of settings when the rank  $r$  and the condition number  $\kappa$  are both fixed constants. As it turns out, this class of problems arises in a variety of engineering applications. For example, in a fundamental problem in cryo-EM called angular synchronization [Sin11], one needs to deal with rank-2 or rank-3 matrices with  $\kappa = 1$ ; in a joint shape mapping problem that arises in computer graphics [HG13, CGH14], the matrix under consideration has low rank and a condition number equal to 1; and in structure from motion in computer vision [TK92], one often seeks to estimate a matrix with  $r \leq 3$  and a small condition number. Encouragingly, our theory delivers near-optimal statistical guarantees for such practically important scenarios.

**Theorem 1.** *Let  $\mathbf{M}^*$  be rank- $r$  and  $\mu$ -incoherent with a condition number  $\kappa$ , where the rank and the condition number satisfy  $r, \kappa = O(1)$ . Suppose that Assumption 1 holds and take  $\lambda = C_\lambda \sigma \sqrt{np}$  in (3) for some large enough constant  $C_\lambda > 0$ . Assume the sample size obeys  $n^2 p \geq C \mu^2 n \log^3 n$  for some sufficiently large constant  $C > 0$ , and the noise satisfies  $\sigma \lesssim \sqrt{\frac{np}{\mu^3 \log n}} \|\mathbf{M}^*\|_\infty$  for some sufficiently small constant  $c > 0$ . Then with probability exceeding  $1 - O(n^{-3})$ :*

1. Any minimizer  $\mathbf{Z}_{\text{cvx}}$  of (3) obeys

$$\|\mathbf{Z}_{\text{cvx}} - \mathbf{M}^*\|_{\text{F}} \lesssim \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \|\mathbf{M}^*\|_{\text{F}}; \quad \|\mathbf{Z}_{\text{cvx}} - \mathbf{M}^*\| \lesssim \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \|\mathbf{M}^*\|; \quad (7a)$$

$$\|\mathbf{Z}_{\text{cvx}} - \mathbf{M}^*\|_\infty \lesssim \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{\mu n \log n}{p}} \|\mathbf{M}^*\|_\infty. \quad (7b)$$

2. Letting  $\mathbf{Z}_{\text{cvx},r} \triangleq \arg \min_{\mathbf{Z}: \text{rank}(\mathbf{Z}) \leq r} \|\mathbf{Z} - \mathbf{Z}_{\text{cvx}}\|_{\text{F}}$  be the best rank- $r$  approximation of  $\mathbf{Z}_{\text{cvx}}$ , we have

$$\|\mathbf{Z}_{\text{cvx},r} - \mathbf{Z}_{\text{cvx}}\|_{\text{F}} \leq \frac{1}{n^3} \cdot \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \|\mathbf{M}^*\|, \quad (8)$$

and the error bounds in (7) continue to hold if  $\mathbf{Z}_{\text{cvx}}$  is replaced by  $\mathbf{Z}_{\text{cvx},r}$ .

**Remark 2.** Here and throughout,  $f(n) \lesssim g(n)$  or  $f(n) = O(g(n))$  means  $|f(n)|/|g(n)| \leq C$  for some constant  $C > 0$  when  $n$  is sufficiently large;  $f(n) \gtrsim g(n)$  means  $|f(n)|/|g(n)| \geq C$  for some constant  $C > 0$  when  $n$  is sufficiently large; and  $f(n) \asymp g(n)$  if and only if  $f(n) \lesssim g(n)$  and  $f(n) \gtrsim g(n)$ . In addition,  $\|\cdot\|_\infty$  denotes the entrywise  $\ell_\infty$  norm, whereas  $\|\cdot\|$  is the spectral norm.

**Remark 3.** The factor  $1/n^3$  in (8) can be replaced by  $1/n^c$  for an arbitrarily large fixed constant  $c > 0$  (e.g.  $c = 100$ ).

To explain the applicability of the above theorem, we first remark on the conditions required for this theorem to hold; for simplicity, we assume that  $\mu = O(1)$ .

- *Sample complexity.* To begin with, the sample size needs to exceed the order of  $n \text{poly} \log n$ , which is information-theoretically optimal up to some logarithmic term [CT10].
- *Noise size.* We then turn attention to the noise requirement, i.e.  $\sigma \lesssim \sqrt{\frac{np}{\log n}} \|\mathbf{M}^*\|_\infty$ . Note that under the sample size condition  $n^2 p \geq C n \log^3 n$ , the size of the noise in each entry is allowed to be substantially larger than the maximum entry in the matrix. In other words, the signal-to-noise ratio w.r.t. each observed

entry could be very small. According to prior literature (e.g. [KMO10b, Theorem 1.1] and [MWCC17, Theorem 2]), such noise conditions are typically required for spectral methods to perform noticeably better than random guessing.

Further, Theorem 1 has several important implications about the power of convex relaxation. The discussions below again concentrate on the case where  $\mu = O(1)$ .

- *Near-optimal stability guarantees.* Our results reveal that the Euclidean error of any convex optimizer  $\mathbf{Z}_{\text{cvx}}$  of (3) obeys

$$\|\mathbf{Z}_{\text{cvx}} - \mathbf{M}^*\|_{\text{F}} \lesssim \sigma \sqrt{n/p}, \quad (9)$$

implying that the performance of convex relaxation degrades gracefully as the signal-to-noise ratio decreases. This result matches the oracle lower bound derived in [CP10, Eq. (III.13)], which also improves upon their statistical guarantee. Specifically, Candès and Plan [CP10] provided a stability guarantee in the presence of arbitrary bounded noise. When applied to the random noise model assumed here, their results yield  $\|\mathbf{Z}_{\text{cvx}} - \mathbf{M}^*\|_{\text{F}} \lesssim \sigma n^{3/2}$ , which could be  $O(\sqrt{n^2 p})$  times more conservative than our bound (9).

- *Nearly low-rank structure of the convex solution.* In light of (8), the optimizer of the convex program (3) is almost, if not exactly, rank- $r$ . When the true rank  $r$  is known *a priori*, it is not uncommon for practitioners to return the rank- $r$  approximation of  $\mathbf{Z}_{\text{cvx}}$ . Our theorem formally justifies that there is no loss of statistical accuracy — measured in terms of either  $\|\cdot\|_{\text{F}}$  or  $\|\cdot\|_{\infty}$  — when performing the rank- $r$  projection operation.
- *Entrywise and spectral norm error control.* Moving beyond the Euclidean loss, our theory uncovers that the estimation errors of the convex optimizer are fairly spread out across all entries, thus implying near-optimal entrywise error control. This is a stronger form of error bounds, as an optimal Euclidean estimation accuracy alone does not preclude the possibility of the estimation errors being spiky and localized. Furthermore, the spectral norm error of the convex optimizer is also well-controlled. Figure 2 displays the relative estimation errors in both the  $\ell_{\infty}$  norm and the spectral norm, under the same setting as in Figure 1. As can be seen, both forms of estimation errors scale linearly with the noise level, corroborating our theory.
- *Implicit regularization.* As a byproduct of the entrywise error control, this result indicates that the additional constraint  $\|\mathbf{Z}\|_{\infty} \leq \alpha$  suggested by [NW12] is automatically satisfied and is hence unnecessary. In other words, the convex approach implicitly controls the spikiness of its entries, without resorting to explicit regularization. This is also confirmed by the numerical experiments reported in Figure 3, where we see that the estimation error of (3) and that of the constrained version considered in [NW12] are nearly identical.
- *Statistical guarantees for fast iterative optimization methods.* Various iterative algorithms have been developed to solve the nuclear norm regularized least-squares problem (3) up to an arbitrarily prescribed accuracy, examples including SVT (or proximal gradient methods) [CCS10], FPC [MGC11], SOFT-IMPUTE [MHT10], FISTA [BT09, TY10], to name just a few. Our theory immediately provides statistical guarantees for these algorithms. As we shall make precise in Section 2, any point  $\mathbf{Z}$  with  $g(\mathbf{Z}) \leq g(\mathbf{Z}_{\text{cvx}}) + \varepsilon$  (where  $g(\cdot)$  is defined in (3)) enjoys the same error bounds as in (7) (with  $\mathbf{Z}_{\text{cvx}}$  replaced by  $\mathbf{Z}$  in (7)), provided that  $\varepsilon > 0$  is sufficiently small. In other words, when these convex optimization algorithms converge w.r.t. the objective value, they are guaranteed to return a statistically reliable estimate.

To better understand our contributions, we take a moment to discuss two important but different convex programs studied in [NW12] and [KLT11]. To begin with, under a spikiness assumption on the low-rank matrix, Negahban and Wainwright [NW12] proposed to enforce an extra entrywise constraint  $\|\mathbf{Z}\|_{\infty} \leq \alpha$  when solving (3), in order to explicitly control the spikiness of the estimate. When applied to our model with  $r, \kappa, \mu \asymp 1$ , their results read (up to some logarithmic factor)

$$\|\hat{\mathbf{Z}} - \mathbf{M}^*\|_{\text{F}} \lesssim \max\{\sigma, \|\mathbf{M}^*\|_{\infty}\} \sqrt{n/p}, \quad (10)$$

where  $\hat{\mathbf{Z}}$  is the estimate returned by their modified convex algorithm. While this matches the optimal bound when  $\sigma \gtrsim \|\mathbf{M}^*\|_{\infty}$ , it becomes suboptimal when  $\sigma \ll \|\mathbf{M}^*\|_{\infty}$  (under our models). Moreover,

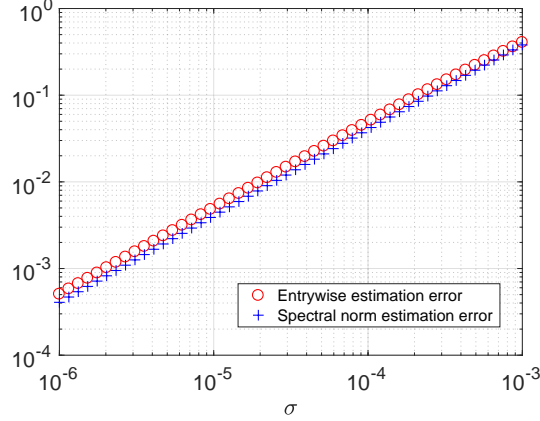


Figure 2: The relative estimation error of  $\mathbf{Z}_{\text{cvx}}$  measured by both  $\|\cdot\|_\infty$  (i.e.  $\|\mathbf{Z}_{\text{cvx}} - \mathbf{M}^*\|_\infty / \|\mathbf{M}^*\|_\infty$ ) and  $\|\cdot\|$  (i.e.  $\|\mathbf{Z}_{\text{cvx}} - \mathbf{M}^*\| / \|\mathbf{M}^*\|$ ) vs. the standard deviation  $\sigma$  of the noise. The results are reported for  $n = 1000$ ,  $r = 5$ ,  $p = 0.2$ ,  $\lambda = 5\sigma\sqrt{np}$  and are averaged over 20 independent trials.

as we have already discussed, the extra spikiness constraint becomes unnecessary in the regime considered herein. This also means that our result complements existing theory about the convex program in [NW12] by demonstrating its minimaxity for an additional range of noise. Another work by Koltchinskii et al. [KLT11] investigated a completely different convex algorithm, which is effectively a spectral method (namely, one round of soft singular value thresholding on a rescaled zero-padded data matrix). The algorithm is shown to be minimax optimal over the class of low-rank matrices with bounded  $\ell_\infty$  norm (note that this is very different from the set of incoherent matrices studied here). When specialized to our model, their error bound is the same as (10) (modulo some log factor), which also becomes suboptimal as  $\sigma$  decreases. As can be seen from the numerical experiments in Figure 3, the estimation error of this thresholding-based spectral algorithm does not decrease as the noise shrinks, and its performance seems uniformly outperformed by that of convex relaxation (3) and the constrained estimator in [NW12]. In fact, this is part of our motivation to pursue an improved theoretical understanding of the formulation (3).

Finally, we make note of a connection between our result and prior theory developed for the noiseless case. Specifically, when the noise vanishes (i.e.  $\sigma \rightarrow 0$ ), one can take a diminishing sequence of regularization parameters  $\{\lambda_k\}$  with  $\lambda_k \rightarrow 0$ , then the resulting estimation errors associated with this sequence should decrease to 0 as  $k \rightarrow \infty$  (which implies exact recovery in the limit of  $k$ ). This parallels the connection between Lasso in sparse linear regression and basis pursuit in compressed sensing.

### 1.4.3 Theoretical guarantees: extensions to more general settings

So far we have presented results when the true matrix has bounded rank and condition number, i.e.  $r, \kappa = O(1)$ . Our theory actually accommodates a significantly broader range of scenarios, where the rank and the condition number are both allowed to grow with the dimension  $n$ .

**Theorem 2.** *Let  $\mathbf{M}^*$  be rank- $r$  and  $\mu$ -incoherent with a condition number  $\kappa$ . Suppose Assumption 1 holds and take  $\lambda = C_\lambda \sigma \sqrt{np}$  in (3) for some large enough constant  $C_\lambda > 0$ . Assume the sample size obeys  $n^2 p \geq C \kappa^4 \mu^2 r^2 n \log^3 n$  for some sufficiently large constant  $C > 0$ , and the noise satisfies  $\sigma \sqrt{\frac{n}{p}} \leq c \frac{\sigma_{\min}}{\sqrt{\kappa^4 \mu r \log n}}$  for some sufficiently small constant  $c > 0$ . Then with probability exceeding  $1 - O(n^{-3})$ ,*

1. Any minimizer  $\mathbf{Z}_{\text{cvx}}$  of (3) obeys

$$\|\mathbf{Z}_{\text{cvx}} - \mathbf{M}^*\|_{\text{F}} \lesssim \kappa \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \|\mathbf{M}^*\|_{\text{F}}, \quad (11a)$$

$$\|\mathbf{Z}_{\text{cvx}} - \mathbf{M}^*\|_{\infty} \lesssim \sqrt{\kappa^3 \mu r} \cdot \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} \|\mathbf{M}^*\|_{\infty}, \quad (11b)$$



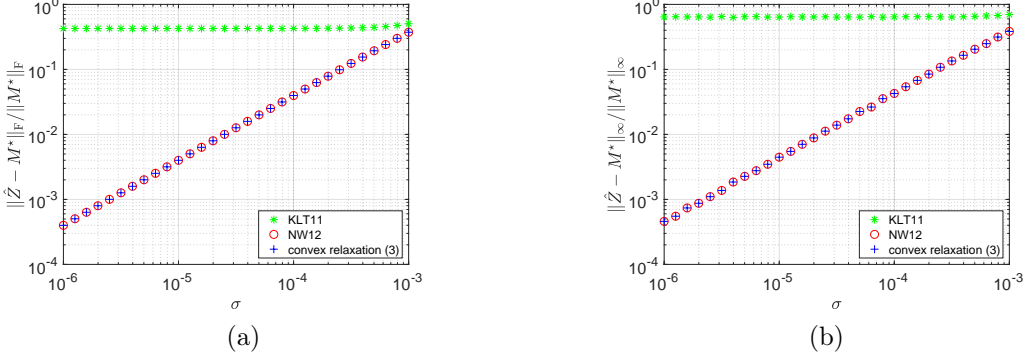


Figure 3: The relative estimation errors of  $\hat{\mathbf{Z}}$ , measured in terms of  $\ell_F$  and  $\ell_\infty$ , vs. the standard deviation  $\sigma$  of the noise. Here  $\hat{\mathbf{Z}}$  can be either the modified convex estimator in [KLT11], the constrained convex estimator in [NW12] or the vanilla convex estimator (3). The results are reported for  $n = 1000$ ,  $r = 5$ ,  $p = 0.2$ , and are averaged over 20 Monte-Carlo trials. For the modified convex estimator in [KLT11], we choose the regularization parameter  $\lambda$  therein to be  $1.5 \max\{\sigma, \|\mathbf{M}^*\|_\infty\} \sqrt{1/(n^3 p)}$ , as suggested by their theory. For the constrained one in [NW12], the regularization parameter  $\lambda$  is set to be  $5\sigma\sqrt{np}$  and the constraint  $\alpha$  is set to be  $\|\mathbf{M}^*\|_\infty$ . Both choices are recommended by [NW12]. As for (3), we set  $\lambda = 5\sigma\sqrt{np}$ .

$$\|\mathbf{Z}_{\text{cvx}} - \mathbf{M}^*\| \lesssim \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \|\mathbf{M}^*\|; \quad (11c)$$

2. Letting  $\mathbf{Z}_{\text{cvx},r} \triangleq \arg \min_{\mathbf{Z}: \text{rank}(\mathbf{Z}) \leq r} \|\mathbf{Z} - \mathbf{Z}_{\text{cvx}}\|_F$  be the best rank- $r$  approximation of  $\mathbf{Z}_{\text{cvx}}$ , we have

$$\|\mathbf{Z}_{\text{cvx},r} - \mathbf{Z}_{\text{cvx}}\|_F \leq \frac{1}{n^3} \cdot \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \|\mathbf{M}^*\|, \quad (12)$$

and the error bounds in (11) continue to hold if  $\mathbf{Z}_{\text{cvx}}$  is replaced by  $\mathbf{Z}_{\text{cvx},r}$ .

**Remark 4** (The noise condition). The incoherence condition (cf. Definition 1) guarantees that the largest entry  $\|\mathbf{M}^*\|_\infty$  of the matrix  $\mathbf{M}^*$  is no larger than  $\kappa \mu r \sigma_{\min} / n$ . As a result, the noise condition stated in Theorem 2 covers all scenarios obeying

$$\sigma \lesssim \sqrt{\frac{np}{\kappa^6 \mu^3 r^3 \log n}} \|\mathbf{M}^*\|_\infty.$$

Therefore, the typical size of the noise is allowed to be much larger than the size of the largest entry of  $\mathbf{M}^*$ , provided that  $p \gg \frac{\kappa^6 \mu^3 r^3 \log n}{n}$ . In particular, when  $r, \kappa = O(1)$ , this recovers the noise condition in Theorem 1.

Notably, the sample size condition for noisy matrix completion (i.e.  $n^2 p \geq C \kappa^4 \mu^2 r^2 n \log^3 n$ ) is more stringent than that in the noiseless setting (i.e.  $n^2 p \asymp nr \log^2 n$ ), and our statistical guarantees are likely suboptimal with respect to the dependency on  $r$  and  $\kappa$ . This sub-optimality is mainly due to the analysis of nonconvex optimization, a key ingredient of our analysis of convex relaxation. In fact, the state-of-the-art nonconvex analysis [KMO10b, CW15, MWCC17] requires the sample size to be much larger than the optimal one (e.g.  $n^2 p \gg npoly(r)poly(\kappa)$ ) even in the noiseless setting. It would certainly be interesting, and in fact important, to see whether it is possible to develop a theory with optimal dependency on  $r$  and  $\kappa$ . We leave this for future investigation.

Despite the above sub-optimality issue, implications similar to those of Theorem 1 hold for this general setting. To begin with, the nearly low-rank structure of the convex solution is preserved (cf. (12)). In addition, the estimation error of the convex estimate is spread out across entries (cf. (11b)), thus uncovering an implicit regularization phenomenon underlying convex relaxation (which implicitly regularizes the spikiness constraint on the solution). Last but not least, the upper bounds (11) and (12) continue to hold for approximate minimizers of the convex program (3), thus yielding statistical guarantees for numerous iterative algorithms aimed at minimizing (3).

## 2 Strategy and novelty

In this section, we introduce the strategy for proving our main theorem, i.e. Theorem 2. Theorem 1 follows immediately. Informally, the main technical difficulty stems from the lack of closed-form expressions for the primal solution to (3), which in turn makes it difficult to construct a dual certificate. This is in stark contrast to the noiseless setting, where one clearly anticipates the ground truth  $\mathbf{M}^*$  to be the primal solution; in fact, this is precisely why the analysis for the noisy case is significantly more challenging. Our strategy, as we shall detail below, mainly entails invoking an iterative nonconvex algorithm to “approximate” such a primal solution.

Before continuing, we introduce a few more notations. Let  $\mathcal{P}_\Omega(\cdot) : \mathbb{R}^{n \times n} \mapsto \mathbb{R}^{n \times n}$  represent the projection onto the subspace of matrices supported on  $\Omega$ , namely,

$$[\mathcal{P}_\Omega(\mathbf{Z})]_{ij} = \begin{cases} Z_{ij}, & \text{for } (i, j) \in \Omega \\ 0, & \text{otherwise} \end{cases} \quad (13)$$

for any matrix  $\mathbf{Z} \in \mathbb{R}^{n \times n}$ . For a rank- $r$  matrix  $\mathbf{M}$  with singular value decomposition  $\mathbf{U}\Sigma\mathbf{V}^\top$ , denote by  $T$  its tangent space, i.e.

$$T = \{\mathbf{U}\mathbf{A}^\top + \mathbf{B}\mathbf{V}^\top \mid \mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times r}\}. \quad (14)$$

Correspondingly, let  $\mathcal{P}_T(\cdot)$  be the orthogonal projection onto the subspace  $T$ , that is,

$$\mathcal{P}_T(\mathbf{Z}) = \mathbf{U}\mathbf{U}^\top \mathbf{Z} + \mathbf{Z}\mathbf{V}\mathbf{V}^\top - \mathbf{U}\mathbf{U}^\top \mathbf{Z}\mathbf{V}\mathbf{V}^\top \quad (15)$$

for any matrix  $\mathbf{Z} \in \mathbb{R}^{n \times n}$ . In addition, let  $T^\perp$  and  $\mathcal{P}_{T^\perp}(\cdot)$  denote the orthogonal complement of  $T$  and the projection onto  $T^\perp$ , respectively. With regards to the ground truth, we denote

$$\mathbf{X}^* = \mathbf{U}^*(\Sigma^*)^{1/2} \quad \text{and} \quad \mathbf{Y}^* = \mathbf{V}^*(\Sigma^*)^{1/2}. \quad (16)$$

The nonconvex problem (5) is equivalent to

$$\underset{\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times r}}{\text{minimize}} \quad f(\mathbf{X}, \mathbf{Y}) \triangleq \frac{1}{2p} \|\mathcal{P}_\Omega(\mathbf{X}\mathbf{Y}^\top - \mathbf{M})\|_F^2 + \frac{\lambda}{2p} \|\mathbf{X}\|_F^2 + \frac{\lambda}{2p} \|\mathbf{Y}\|_F^2, \quad (17)$$

where we have inserted an extra factor  $1/p$  (compared to (5)) to simplify the presentation of the analysis later on.

### 2.1 Exact duality

In order to analyze the convex program (3), it is natural to start with the first-order optimality condition. Specifically, suppose that  $\mathbf{Z} \in \mathbb{R}^{n \times n}$  is a (primal) solution to (3) with SVD  $\mathbf{Z} = \mathbf{U}\Sigma\mathbf{V}^\top$ .<sup>4</sup> As before, let  $T$  be the tangent space of  $\mathbf{Z}$ , and let  $T^\perp$  be the orthogonal complement of  $T$ . Then the first-order optimality condition for (3) reads: there exists a matrix  $\mathbf{W} \in T^\perp$  (called a dual certificate) such that

$$\frac{1}{\lambda} \mathcal{P}_\Omega(\mathbf{M} - \mathbf{Z}) = \mathbf{U}\mathbf{V}^\top + \mathbf{W}; \quad (18a)$$

$$\|\mathbf{W}\| \leq 1. \quad (18b)$$

This condition is not only necessary to certify the optimality of  $\mathbf{Z}$ , but also “almost sufficient” in guaranteeing the uniqueness of the solution  $\mathbf{Z}$ ; see Appendix B for in-depth discussions.

The challenge then boils down to identifying such a primal-dual pair  $(\mathbf{Z}, \mathbf{W})$  satisfying the optimality condition (18). For the noise-free case, the primal solution is clearly  $\mathbf{Z} = \mathbf{M}^*$  if exact recovery is to be expected; the dual certificate can then be either constructed exactly by the least-squares solution to a certain underdetermined linear system [CR09, CT10], or produced approximately via a clever golfing scheme pioneered by Gross [Gro11]. For the noisy case, however, it is often difficult to hypothesize on the primal solution  $\mathbf{Z}$ , as it depends on the random noise in a complicated way. In fact, the lack of a suitable guess of  $\mathbf{Z}$  (and hence  $\mathbf{W}$ ) was the major hurdle that prior works faced when carrying out the duality analysis.

<sup>4</sup>Here and below, we use  $\mathbf{Z}$  (rather than  $\mathbf{Z}_{\text{cvx}}$ ) for notational simplicity, whenever it is clear from the context.

## 2.2 A candidate primal solution via nonconvex optimization

Motivated by the numerical experiment in Section 1.3, we propose to examine whether the optimizer of the nonconvex problem (5) stays close to the solution to the convex program (3). Towards this, suppose that  $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times r}$  form a critical point of (5) with  $\text{rank}(\mathbf{X}) = \text{rank}(\mathbf{Y}) = r$ .<sup>5</sup> Then the first-order condition reads

$$\frac{1}{\lambda} \mathcal{P}_\Omega(\mathbf{M} - \mathbf{X}\mathbf{Y}^\top) \mathbf{Y} = \mathbf{X}; \quad (19a)$$

$$\frac{1}{\lambda} [\mathcal{P}_\Omega(\mathbf{M} - \mathbf{X}\mathbf{Y}^\top)]^\top \mathbf{X} = \mathbf{Y}. \quad (19b)$$

To develop some intuition about the connection between (18) and (19), let us take a look at the case with  $r = 1$ . Denote  $\mathbf{X} = \mathbf{x}$  and  $\mathbf{Y} = \mathbf{y}$  and assume that the two rank-1 factors are “balanced”, namely,  $\|\mathbf{x}\|_2 = \|\mathbf{y}\|_2 \neq 0$ . It then follows from (19) that  $\lambda^{-1} \mathcal{P}_\Omega(\mathbf{M} - \mathbf{x}\mathbf{y}^\top)$  has a singular value 1, whose corresponding left and right singular vectors are  $\mathbf{x}/\|\mathbf{x}\|_2$  and  $\mathbf{y}/\|\mathbf{y}\|_2$ , respectively. In other words, one can express

$$\frac{1}{\lambda} \mathcal{P}_\Omega(\mathbf{M} - \mathbf{x}\mathbf{y}^\top) = \frac{1}{\|\mathbf{x}\|_2 \|\mathbf{y}\|_2} \mathbf{x}\mathbf{y}^\top + \mathbf{W}, \quad (20)$$

where  $\mathbf{W}$  is orthogonal to the tangent space of  $\mathbf{x}\mathbf{y}^\top$ ; this is precisely the condition (18a). It remains to argue that (18b) is valid as well. Towards this end, the first-order condition (19) alone is insufficient, as there might be non-global critical points (e.g. saddle points) that are unable to approximate the convex solution well. Fortunately, as long as the candidate  $\mathbf{x}\mathbf{y}^\top$  is not far away from the ground truth  $\mathbf{M}^*$ , one can guarantee  $\|\mathbf{W}\| < 1$  as required in (18b).

The above informal argument about the link between the convex and the nonconvex problems can be rigorized. To begin with, we introduce the following conditions on the regularization parameter  $\lambda$ .

**Condition 1 (Regularization parameter).** The regularization parameter  $\lambda$  satisfies

- (a) **(Relative to noise)**  $\|\mathcal{P}_\Omega(\mathbf{E})\| < \lambda/8$ .
- (b) **(Relative to nonconvex solution)**  $\|\mathcal{P}_\Omega(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}^*) - p(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}^*)\| < \lambda/8$ .

**Remark 5.** Condition 1 requires that the regularization parameter  $\lambda$  should dominate a certain norm of the noise, as well as of the deviation of  $\mathbf{X}\mathbf{Y}^\top - \mathbf{M}^*$  from its mean  $p(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}^*)$ ; as will be seen shortly, the latter condition can be met when  $(\mathbf{X}, \mathbf{Y})$  is sufficiently close to  $(\mathbf{X}^*, \mathbf{Y}^*)$ .

With the above condition in place, the following result demonstrates that a critical point  $(\mathbf{X}, \mathbf{Y})$  of the nonconvex problem (5) readily translates to the unique minimizer of the convex program (3). This lemma is established in Appendix C.1.

**Lemma 1 (Exact nonconvex vs. convex optimizers).** *Suppose that  $(\mathbf{X}, \mathbf{Y})$  is a critical point of (5) satisfying  $\text{rank}(\mathbf{X}) = \text{rank}(\mathbf{Y}) = r$ , and the sampling operator  $\mathcal{P}_\Omega$  is injective when restricted to the elements of the tangent space  $T$  of  $\mathbf{X}\mathbf{Y}^\top$ , namely,*

$$\mathcal{P}_\Omega(\mathbf{H}) = \mathbf{0} \iff \mathbf{H} = \mathbf{0}, \quad \text{for all } \mathbf{H} \in T. \quad (21)$$

*Under Condition 1, the point  $\mathbf{Z} \triangleq \mathbf{X}\mathbf{Y}^\top$  is the unique minimizer of (3).*

In order to apply Lemma 1, one needs to locate a critical point of (5) that is sufficiently close to the truth, for which one natural candidate is the global optimizer of (5). The caveat, however, is the lack of theory characterizing directly the properties of the optimizer of (5). Instead, what is available in prior theory is the characterization of some iterative sequence (e.g. gradient descent iterates) aimed at solving (5). It is unclear from prior theory whether the iterative algorithm under study (e.g. gradient descent) converges to the global optimizer in the presence of noise. This leads to technical difficulty in justifying the proximity between the nonconvex optimizer and the convex solution via Lemma 1.

<sup>5</sup>Once again, we abuse the notation  $(\mathbf{X}, \mathbf{Y})$  (instead of using  $(\mathbf{X}_{\text{ncvx}}, \mathbf{Y}_{\text{ncvx}})$ ) for notational simplicity, whenever it is clear from the context.

### 2.3 Approximate nonconvex optimizers

Fortunately, perfect knowledge of the nonconvex optimizer is not pivotal. Instead, an approximate solution to the nonconvex problem (5) (or equivalently (17)) suffices to serve as a reasonably tight approximation of the convex solution. More precisely, we desire two factors  $(\mathbf{X}, \mathbf{Y})$  that result in nearly zero (rather than exactly zero) gradients:

$$\nabla_{\mathbf{X}} f(\mathbf{X}, \mathbf{Y}) \approx \mathbf{0} \quad \text{and} \quad \nabla_{\mathbf{Y}} f(\mathbf{X}, \mathbf{Y}) \approx \mathbf{0},$$

where  $f(\cdot, \cdot)$  is the nonconvex objective function as defined in (17). This relaxes the condition discussed in Lemma 1 (which only applies to critical points of (5) as opposed to approximate critical points). As it turns out, such points can be found via gradient descent tailored to (5). The sufficiency of the near-zero gradient condition is made possible by slightly strengthening the injectivity assumption (21), which is stated below.

**Condition 2 (Injectivity).** Let  $T$  be the tangent space of  $\mathbf{X}\mathbf{Y}^\top$ . There is a quantity  $c_{\text{inj}} > 0$  such that

$$p^{-1} \|\mathcal{P}_\Omega(\mathbf{H})\|_{\text{F}}^2 \geq c_{\text{inj}} \|\mathbf{H}\|_{\text{F}}^2, \quad \text{for all } \mathbf{H} \in T. \quad (22)$$

The following lemma states quantitatively how an approximate nonconvex optimizer serves as an excellent proxy of the convex solution, which we establish in Appendix C.2.

**Lemma 2 (Approximate nonconvex vs. convex optimizers).** *Suppose that  $(\mathbf{X}, \mathbf{Y})$  obeys*

$$\|\nabla f(\mathbf{X}, \mathbf{Y})\|_{\text{F}} \leq c \frac{\sqrt{c_{\text{inj}} p}}{\kappa} \cdot \frac{\lambda}{p} \sqrt{\sigma_{\min}} \quad (23)$$

for some sufficiently small constant  $c > 0$ . Further assume that any singular value of  $\mathbf{X}$  and  $\mathbf{Y}$  lies in  $[\sqrt{\sigma_{\min}}/2, \sqrt{2\sigma_{\max}}]$ . Then under Conditions 1 and 2, any minimizer  $\mathbf{Z}_{\text{cvx}}$  of (3) satisfies

$$\|\mathbf{X}\mathbf{Y}^\top - \mathbf{Z}_{\text{cvx}}\|_{\text{F}} \lesssim \frac{\kappa}{c_{\text{inj}}} \frac{1}{\sqrt{\sigma_{\min}}} \|\nabla f(\mathbf{X}, \mathbf{Y})\|_{\text{F}}. \quad (24)$$

**Remark 6.** In fact, this lemma continues to hold if  $\mathbf{Z}_{\text{cvx}}$  is replaced by any  $\mathbf{Z}$  obeying  $g(\mathbf{Z}) \leq g(\mathbf{X}\mathbf{Y}^\top)$ , where  $g(\cdot)$  is the objective function defined in (3) and  $\mathbf{X}$  and  $\mathbf{Y}$  are low-rank factors obeying conditions of Lemma 2. This is important in providing statistical guarantees for iterative methods like SVT [CCS10], FPC [MGC11], SOFT-IMPUTE [MHT10], FISTA [BT09], etc. To be more specific, suppose that  $(\mathbf{X}, \mathbf{Y})$  results in an approximate optimizer of (3), namely,  $g(\mathbf{X}\mathbf{Y}^\top) = g(\mathbf{Z}_{\text{cvx}}) + \varepsilon$  for some sufficiently small  $\varepsilon > 0$ . Then for any  $\mathbf{Z}$  obeying  $g(\mathbf{Z}) \leq g(\mathbf{X}\mathbf{Y}^\top) = g(\mathbf{Z}_{\text{cvx}}) + \varepsilon$ , one has

$$\|\mathbf{X}\mathbf{Y}^\top - \mathbf{Z}\|_{\text{F}} \lesssim \frac{\kappa}{c_{\text{inj}}} \frac{1}{\sqrt{\sigma_{\min}}} \|\nabla f(\mathbf{X}, \mathbf{Y})\|_{\text{F}}. \quad (25)$$

As a result, as long as the above-mentioned algorithms converge in terms of the objective value, they must return a solution obeying (25), which is exceedingly close to  $\mathbf{X}\mathbf{Y}^\top$  if  $\|\nabla f(\mathbf{X}, \mathbf{Y})\|_{\text{F}}$  is small.

It is clear from Lemma 2 that, as the size of the gradient  $\nabla f(\mathbf{X}, \mathbf{Y})$  gets smaller, the nonconvex estimate  $\mathbf{X}\mathbf{Y}^\top$  becomes an increasingly tighter approximation of any convex optimizer of (3), which is consistent with Lemma 1. In contrast to Lemma 1, due to the lack of strong convexity, a nonconvex estimate with a near-zero gradient does not imply the uniqueness of the optimizer of the convex program (3); rather, it indicates that any minimizer of (3) lies within a sufficiently small neighborhood surrounding  $\mathbf{X}\mathbf{Y}^\top$  (cf. (24)).

### 2.4 Construction of an approximate nonconvex optimizer

So far, Lemmas 1-2 are both deterministic results based on Condition 1. As we will soon see, under Assumption 1, we can derive simpler conditions that — with high probability — guarantee Condition 1. We start with Condition 1(a).

**Lemma 3.** *Suppose  $n^2 p \geq C n \log^2 n$  for some sufficiently large constant  $C > 0$ . Then with probability at least  $1 - O(n^{-10})$ , one has  $\|\mathcal{P}_\Omega(\mathbf{E})\| \lesssim \sigma \sqrt{np}$ . As a result, Condition 1 holds (i.e.  $\|\mathcal{P}_\Omega(\mathbf{E})\| < \lambda/8$ ) as long as  $\lambda = C_\lambda \sigma \sqrt{np}$  for some sufficiently large constant  $C_\lambda > 0$ .*

*Proof.* This follows from [CW15, Lemma 11] with a slight and straightforward modification to accommodate the asymmetric noise here. For brevity, we omit the proof.  $\square$

Turning attention to Condition 1(b) and Condition 2, we have the following lemma, the proof of which is deferred to Appendix C.3.

**Lemma 4.** *Under the assumptions of Theorem 2, with probability exceeding  $1 - O(n^{-10})$  we have*

$$\|\mathcal{P}_\Omega(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}^*) - p(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}^*)\| < \lambda/8 \quad (\text{Condition 1(b)})$$

$$\frac{1}{p} \|\mathcal{P}_\Omega(\mathbf{H})\|_{\text{F}}^2 \geq \frac{1}{32\kappa} \|\mathbf{H}\|_{\text{F}}^2, \quad \text{for all } \mathbf{H} \in T \quad (\text{Condition 2 with } c_{\text{inj}} = (32\kappa)^{-1})$$

hold simultaneously for all  $(\mathbf{X}, \mathbf{Y})$  obeying

$$\begin{aligned} & \max \left\{ \|\mathbf{X} - \mathbf{X}^*\|_{2,\infty}, \|\mathbf{Y} - \mathbf{Y}^*\|_{2,\infty} \right\} \\ & \leq C_\infty \kappa \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \max \left\{ \|\mathbf{X}^*\|_{2,\infty}, \|\mathbf{Y}^*\|_{2,\infty} \right\}. \end{aligned} \quad (26)$$

Here,  $T$  denotes the tangent space of  $\mathbf{X}\mathbf{Y}^\top$ , and  $C_\infty > 0$  is some absolute constant.

This lemma is a uniform result, namely, the bounds hold irrespective of the statistical dependency between  $(\mathbf{X}, \mathbf{Y})$  and  $\Omega$ . As a consequence, to demonstrate the proximity between the convex and the nonconvex solutions (cf. (24)), it remains to identify a point  $(\mathbf{X}, \mathbf{Y})$  with vanishingly small gradient (cf. (23)) that is sufficiently close to the truth (cf. (26)).

As we already alluded to previously, a simple gradient descent algorithm aimed at solving the nonconvex problem (5) might help us produce an approximate nonconvex optimizer. This procedure is summarized in Algorithm 1. Our hope is this: when initialized at the ground truth and run for sufficiently many iterations, the GD trajectory produced by Algorithm 1 will contain at least one approximate stationary point of (5) with the desired properties (23) and (26). We shall note that Algorithm 1 is *not practical* since it starts from the ground truth  $(\mathbf{X}^*, \mathbf{Y}^*)$ ; this is an auxiliary step mainly to simplify the theoretical analysis. While we can certainly make it practical by adopting spectral initialization as in [MWCC17, CLL19], it requires more lengthy proofs without further improving our statistical guarantees.

---

**Algorithm 1** Construction of an approximate primal solution.

---

**Initialization:**  $\mathbf{X}^0 = \mathbf{X}^*$ ;  $\mathbf{Y}^0 = \mathbf{Y}^*$ .

**Gradient updates:** for  $t = 0, 1, \dots, t_0 - 1$  do

$$\mathbf{X}^{t+1} = \mathbf{X}^t - \eta \nabla_{\mathbf{X}} f(\mathbf{X}^t, \mathbf{Y}^t) = \mathbf{X}^t - \frac{\eta}{p} \left( \mathcal{P}_\Omega(\mathbf{X}^t \mathbf{Y}^{t\top} - \mathbf{M}) \mathbf{Y}^t + \lambda \mathbf{X}^t \right); \quad (27a)$$

$$\mathbf{Y}^{t+1} = \mathbf{Y}^t - \eta \nabla_{\mathbf{Y}} f(\mathbf{X}^t, \mathbf{Y}^t) = \mathbf{Y}^t - \frac{\eta}{p} \left( [\mathcal{P}_\Omega(\mathbf{X}^t \mathbf{Y}^{t\top} - \mathbf{M})]^\top \mathbf{X}^t + \lambda \mathbf{Y}^t \right). \quad (27b)$$

Here,  $\eta > 0$  is the step size.

---

## 2.5 Properties of the nonconvex iterates

In this subsection, we will build upon the literature on nonconvex low-rank matrix completion to justify that the estimates returned by Algorithm 1 satisfy the requirement stated in (26). Our theory will be largely established upon the leave-one-out strategy introduced by Ma et al. [MWCC17], which is an effective analysis technique to control the  $\ell_{2,\infty}$  error of the estimates. This strategy has recently been extended by Chen et al. [CLL19] to the more general rectangular case with an improved sample complexity bound.

Before continuing, we introduce several useful notations. Notice that the matrix product of  $\mathbf{X}^*$  and  $\mathbf{Y}^{*\top}$  is invariant under global orthonormal transformation, namely, for any orthonormal matrix  $\mathbf{R} \in \mathbb{R}^{r \times r}$  one has

$\mathbf{X}^* \mathbf{R} (\mathbf{Y}^* \mathbf{R})^\top = \mathbf{X}^* \mathbf{Y}^{*\top}$ . Viewed in this light, we shall consider distance metrics modulo global rotation. In particular, the theory relies heavily on a specific global rotation matrix defined as follows

$$\mathbf{H}^t \triangleq \arg \min_{\mathbf{R} \in \mathcal{O}^{r \times r}} \left( \|\mathbf{X}^t \mathbf{R} - \mathbf{X}^*\|_{\text{F}}^2 + \|\mathbf{Y}^t \mathbf{R} - \mathbf{Y}^*\|_{\text{F}}^2 \right)^{1/2}, \quad (28)$$

where  $\mathcal{O}^{r \times r}$  is the set of  $r \times r$  orthonormal matrices.

We are now ready to present the performance guarantees for Algorithm 1.

**Lemma 5 (Quality of the nonconvex estimates).** *Instate the notation and hypotheses of Theorem 2. With probability at least  $1 - O(n^{-3})$ , the iterates  $\{(\mathbf{X}^t, \mathbf{Y}^t)\}_{0 \leq t \leq t_0}$  of Algorithm 1 satisfy*

$$\max \left\{ \|\mathbf{X}^t \mathbf{H}^t - \mathbf{X}^*\|_{\text{F}}, \|\mathbf{Y}^t \mathbf{H}^t - \mathbf{Y}^*\|_{\text{F}} \right\} \leq C_{\text{F}} \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{X}^*\|_{\text{F}}, \quad (29a)$$

$$\max \left\{ \|\mathbf{X}^t \mathbf{H}^t - \mathbf{X}^*\|, \|\mathbf{Y}^t \mathbf{H}^t - \mathbf{Y}^*\| \right\} \leq C_{\text{op}} \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{X}^*\|, \quad (29b)$$

$$\begin{aligned} & \max \left\{ \|\mathbf{X}^t \mathbf{H}^t - \mathbf{X}^*\|_{2, \infty}, \|\mathbf{Y}^t \mathbf{H}^t - \mathbf{Y}^*\|_{2, \infty} \right\} \\ & \leq C_{\infty} \kappa \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \max \left\{ \|\mathbf{X}^*\|_{2, \infty}, \|\mathbf{Y}^*\|_{2, \infty} \right\}, \end{aligned} \quad (29c)$$

$$\min_{0 \leq t < t_0} \|\nabla f(\mathbf{X}^t, \mathbf{Y}^t)\|_{\text{F}} \leq \frac{1}{n^5} \frac{\lambda}{p} \sqrt{\sigma_{\min}}, \quad (30)$$

where  $C_{\text{F}}, C_{\text{op}}, C_{\infty} > 0$  are some absolute constants, provided that  $\eta \asymp 1/(n\kappa^3\sigma_{\max})$  and that  $t_0 = n^{18}$ .

This lemma, which we establish in Appendix D, reveals that for a polynomially large number of iterations, all iterates of the gradient descent sequence — when initialized at the ground truth — remain fairly close to the true low-rank factors. This holds in terms of the estimation errors measured by the Frobenius norm, the spectral norm, and the  $\ell_{2, \infty}$  norm. In particular, the proximity in terms of the  $\ell_{2, \infty}$  norm error plays a pivotal role in implementing our analysis strategy (particularly Lemmas 2-4) described previously. In addition, this lemma (cf. (30)) guarantees the existence of a small-gradient point within this sequence  $\{(\mathbf{X}^t, \mathbf{Y}^t)\}_{0 \leq t \leq t_0}$ , a somewhat straightforward property of GD tailored to smooth problems [Nes12]. This in turn enables us to invoke Lemma 2.

As immediate consequences of Lemma 5, with high probability we have

$$\|\mathbf{X}^t \mathbf{Y}^{t\top} - \mathbf{M}^*\|_{\text{F}} \leq 3\kappa C_{\text{F}} \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{M}^*\|_{\text{F}} \quad (31a)$$

$$\|\mathbf{X}^t \mathbf{Y}^{t\top} - \mathbf{M}^*\|_{\infty} \leq 3C_{\infty} \sqrt{\kappa^3 \mu r} \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{M}^*\|_{\infty} \quad (31b)$$

$$\|\mathbf{X}^t \mathbf{Y}^{t\top} - \mathbf{M}^*\| \leq 3C_{\text{op}} \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{M}^*\| \quad (31c)$$

for all  $0 \leq t \leq t_0$ . The proof is deferred to Appendix D.12.

## 2.6 Proof of Theorem 2

Let  $t_* \triangleq \arg \min_{0 \leq t < t_0} \|\nabla f(\mathbf{X}^t, \mathbf{Y}^t)\|_{\text{F}}$ , and take  $(\mathbf{X}_{\text{ncvx}}, \mathbf{Y}_{\text{ncvx}}) = (\mathbf{X}^{t_*} \mathbf{H}^{t_*}, \mathbf{Y}^{t_*} \mathbf{H}^{t_*})$  (cf. (28)). It is straightforward to verify that  $(\mathbf{X}_{\text{ncvx}}, \mathbf{Y}_{\text{ncvx}})$  obeys (i) the small-gradient condition (23), and (ii) the proximity condition (26). We are now positioned to invoke Lemma 2: for any optimizer  $\mathbf{Z}_{\text{cvx}}$  of (3), one has

$$\begin{aligned} \|\mathbf{Z}_{\text{cvx}} - \mathbf{X}_{\text{ncvx}} \mathbf{Y}_{\text{ncvx}}^\top\|_{\text{F}} & \lesssim \frac{\kappa}{c_{\text{inj}}} \frac{1}{\sqrt{\sigma_{\min}}} \|\nabla f(\mathbf{X}_{\text{ncvx}}, \mathbf{Y}_{\text{ncvx}})\|_{\text{F}} \lesssim \frac{\kappa^2 \lambda}{n^5 p} \\ & = \frac{\kappa}{n^5} \frac{\lambda}{p \sigma_{\min}} (\kappa \sigma_{\min}) = \frac{\kappa}{n^5} \frac{\lambda}{p \sigma_{\min}} \|\mathbf{M}^*\| \end{aligned}$$

$$\lesssim \frac{1}{n^4} \frac{\lambda}{p \sigma_{\min}} \|\mathbf{M}^*\|. \quad (32)$$

The last line arises since  $n \gg \kappa$  — a consequence of the sample complexity condition  $np \gtrsim \kappa^4 \mu^2 r^2 \log^3 n$  (and hence  $n \geq np \gtrsim \kappa^4 \mu^2 r^2 \log^3 n \gg \kappa^4$ ). This taken collectively with the property (31) implies that

$$\begin{aligned} \|\mathbf{Z}_{\text{cvx}} - \mathbf{M}^*\|_{\text{F}} &\leq \|\mathbf{Z}_{\text{cvx}} - \mathbf{X}_{\text{ncvx}} \mathbf{Y}_{\text{ncvx}}^{\top}\|_{\text{F}} + \|\mathbf{X}_{\text{ncvx}} \mathbf{Y}_{\text{ncvx}}^{\top} - \mathbf{M}^*\|_{\text{F}} \\ &\lesssim \frac{1}{n^4} \frac{\lambda}{p \sigma_{\min}} \|\mathbf{M}^*\| + \kappa \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{M}^*\|_{\text{F}} \\ &\asymp \kappa \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{M}^*\|_{\text{F}}. \end{aligned}$$

In other words, since  $\mathbf{X}_{\text{ncvx}} \mathbf{Y}_{\text{ncvx}}^{\top}$  and  $\mathbf{Z}_{\text{ncvx}}$  are exceedingly close, the error  $\mathbf{Z}_{\text{cvx}} - \mathbf{M}^*$  is mainly accredited to  $\mathbf{X}_{\text{ncvx}} \mathbf{Y}_{\text{ncvx}}^{\top} - \mathbf{M}^*$ . Similar arguments lead to

$$\begin{aligned} \|\mathbf{Z}_{\text{cvx}} - \mathbf{M}^*\| &\lesssim \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{M}^*\|, \\ \|\mathbf{Z}_{\text{cvx}} - \mathbf{M}^*\|_{\infty} &\lesssim \sqrt{\kappa^3 \mu r} \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{M}^*\|_{\infty}. \end{aligned}$$

We are left with proving the properties of  $\mathbf{Z}_{\text{cvx},r}$ . Since  $\mathbf{Z}_{\text{cvx},r}$  is defined to be the best rank- $r$  approximation of  $\mathbf{Z}_{\text{cvx}}$ , one can invoke (32) to derive

$$\|\mathbf{Z}_{\text{cvx}} - \mathbf{Z}_{\text{cvx},r}\|_{\text{F}} \leq \|\mathbf{Z}_{\text{cvx}} - \mathbf{X}_{\text{ncvx}} \mathbf{Y}_{\text{ncvx}}^{\top}\|_{\text{F}} \lesssim \frac{1}{n^4} \frac{\lambda}{p \sigma_{\min}} \|\mathbf{M}^*\|,$$

from which (12) follows. Repeating the above calculations implies that (11) holds if  $\mathbf{Z}_{\text{cvx}}$  is replaced by  $\mathbf{Z}_{\text{cvx},r}$ , thus concluding the proof.

### 3 Prior art

Nuclear norm minimization, pioneered by the seminal works [RFP10, CR09, CT10, Faz02], has been a popular and principled approach to low-rank matrix recovery. In the noiseless setting, i.e.  $\mathbf{E} = \mathbf{0}$ , it amounts to solving the following constrained convex program

$$\text{minimize}_{\mathbf{Z} \in \mathbb{R}^{n \times n}} \|\mathbf{Z}\|_* \quad \text{subject to} \quad \mathcal{P}_{\Omega}(\mathbf{Z}) = \mathcal{P}_{\Omega}(\mathbf{M}^*), \quad (33)$$

which enjoys great theoretical success. Informally, this approach enables exact recovery of a rank- $r$  matrix  $\mathbf{M}^* \in \mathbb{R}^{n \times n}$  as soon as the sample size is about the order of  $nr$  — the intrinsic degrees of freedom of a rank- $r$  matrix [Gro11, Rec11, Che15]. In particular, Gross [Gro11] blazed a trail by developing an ingenious golfing scheme for dual construction — an analysis technique that has found applications far beyond matrix completion. When it comes to the noisy case, Candès and Plan [CP10] first studied the stability of convex programming when the noise is bounded and possibly adversarial, followed by [NW12] and [KLT11] using two modified convex programs. As we have already discussed, none of these papers provide optimal statistical guarantees under our model when  $r = O(1)$ . Other related papers such as [Klo14, CZ16] include similar estimation error bounds and suffer from similar sub-optimality issues.

Turning to nonconvex optimization, we note that this approach has recently received much attention for various low-rank matrix factorization problems, owing to its superior computational advantage compared to convex programming (e.g. [KMO10a, JNS13, CLS15, CC17, TBS<sup>+</sup>16, ZZLC17]). The convergence guarantees for matrix completion have been established for various algorithms such as gradient descent on manifold [KMO10a, KMO10b], alternating minimization [JNS13, Har14], gradient descent [SL16, MWCC17, WZG16, CLL19], and projected gradient descent [CW15], provided that a suitable initialization (like spectral initialization) is available [KMO10a, JNS13, SL16, MWCC17, CCF18]. Our work is mostly related to [MWCC17, CLL19], which studied (vanilla) gradient descent for nonconvex matrix completion. This algorithm was first analyzed by [MWCC17] via a leave-one-out argument — a technique that proves useful in

analyzing various statistical algorithms [EK15, SCC17, ZB18, CFMW19, AFWZ17, LMCC18, DC18, CCFM19]. In the absence of noise and omitting logarithmic factors, [MWCC17] showed that  $O(nr^3)$  samples are sufficient for vanilla GD to yield  $\varepsilon$  accuracy in  $O(\log \frac{1}{\varepsilon})$  iterations (without the need of extra regularization procedures); the sample complexity was further improved to  $O(nr^2)$  by [CLL19]. Apart from gradient descent, other nonconvex methods (e.g. [RS05, JMD10, WYZ12, JNS13, FRW11, Van13, LXY13, Har14, JKN16, RT<sup>+</sup>11, WCCL16, DC18, GAGG13, CX16, ZWL15]) and landscape / geometry properties have been investigated [GLM16, CL17, PKCS17, GJZ17, SXZ19]; these are, however, beyond the scope of the current paper.

Another line of works asserted that a large family of SDPs admits low-rank solutions [Bar95], which in turn motivates the Burer-Monteiro approach [BM03, BVB16]. When applied to matrix completion, however, the generic theoretical guarantees therein lead to conservative results. Take the noiseless case (33) for instance: these results revealed the existence of a solution of rank at most  $O(\sqrt{n^2p})$ , which however is often much larger the true rank (e.g. when  $r \asymp 1$  and  $p \asymp \text{poly} \log(n)/n$ , one has  $\sqrt{n^2p} \gg \sqrt{n} \gg r$ ). Moreover, this line of works does not imply that all solutions to the SDP of interest are (approximately) low-rank.

Finally, the connection between convex and nonconvex optimization has also been explored in line spectral estimation [LT18], although the context therein is drastically different from ours.

## 4 Discussion

This paper provides an improved statistical analysis for the natural convex program (3), without the need of enforcing additional spikiness constraint. Our theoretical analysis uncovers an intriguing connection between convex relaxation and nonconvex optimization, which we believe is applicable to many other problems beyond matrix completion. Having said that, our current theory leaves open a variety of important directions for future exploration. Here we sample a few interesting ones.

- *Improving dependency on  $r$  and  $\kappa$ .* While our theory is optimal when  $r$  and  $\kappa$  are both constants, it becomes increasingly looser as either  $r$  or  $\kappa$  grows. For instance, in the noiseless setting, it has been shown that the sample complexity for convex relaxation scales as  $O(nr)$  — linear in  $r$  and independent of  $\kappa$  — which is better than the current results. It is worth noting that existing theory for nonconvex matrix factorization typically falls short of providing optimal scaling in  $r$  and  $\kappa$  [KMO10a, SL16, CW15, MWCC17, CLL19]. Thus, tightening the dependency of sample complexity on  $r$  and  $\kappa$  might call for new analysis tools.
- *Approximate low-rank structure.* So far our theory is built upon the assumption that the ground-truth matrix  $M^*$  is exactly low-rank, which falls short of accommodating the more realistic scenario where  $M^*$  is only approximately low-rank. For the approximate low-rank case, it is not yet clear whether the nonconvex factorization approach can still serve as a tight proxy. In addition, the landscape of nonconvex optimization for the approximately low-rank case [CL17] might shed light on how to handle this case.
- *Extension to deterministic noise.* Our current theory — in particular, the leave-one-out analysis for the nonconvex approach — relies heavily on the randomness assumption (i.e. i.i.d. sub-Gaussian) of the noise. In order to justify the broad applicability of convex relaxation, it would be interesting to see whether one can generalize the theory to cover deterministic noise with bounded magnitudes.
- *Extension to structured matrix completion.* Many applications involve low-rank matrices that exhibit additional structures, enabling a further reduction of the sample complexity [FHB03, CC14, CWW19]. For instance, if a matrix is Hankel and low-rank, then the sample complexity can be  $O(n)$  times smaller than the generic low-rank case. The existing stability guarantee of Hankel matrix completion, however, is overly pessimistic compared to practical performance [CC14]. The analysis framework herein might be amenable to the study of Hankel matrix completion and help close the theory-practice gap.
- *Extension to robust PCA and blind deconvolution.* Moving beyond matrix completion, there are other problems that are concerned with recovering low-rank matrices. Notable examples include robust principal component analysis [CLMW11, CSPW11, CJSC13], blind deconvolution [ARR14, LS15] and blind demixing [LS17, JKS17]. The stability analyses of the convex relaxation approaches for these problems [ZLW<sup>+</sup>10, ARR14, LS17] often adopt a similar approach as [CP10], and consequently are sub-optimal. The insights from the present paper might promise tighter statistical guarantees for such problems.



Finally, we remark that the intimate link between convex and nonconvex optimization enables statistically optimal inference and uncertainty quantification for noisy matrix completion (e.g. construction of optimal confidence intervals for each missing entry). The interested readers are referred to our companion paper [CFMY19] for in-depth discussions.

## Acknowledgements

Y. Chen is supported in part by the AFOSR YIP award FA9550-19-1-0030, by the ARO grant W911NF-18-1-0303, by the ONR grant N00014-19-1-2120, by the NSF grants CCF-1907661 and IIS-1900140, and by the Princeton SEAS innovation award. Y. Chi is supported in part by ONR under the grants N00014-18-1-2142 and N00014-19-1-2404, by ARO under the grant W911NF-18-1-0303, and by NSF under the grants CAREER ECCS-1818571 and CCF-1806154. J. Fan is supported in part by NSF Grants DMS-1662139 and DMS-1712591, ONR grant N00014-19-1-2120, and NIH Grant R01-GM072611-12. This work was done in part while Y. Chen was visiting the Kavli Institute for Theoretical Physics (supported in part by NSF grant PHY-1748958). Y. Chen thanks Emmanuel Candès for motivating discussions about noisy matrix completion.

## References

- [AFWZ17] E. Abbe, J. Fan, K. Wang, and Y. Zhong. Entrywise eigenvector analysis of random matrices with low expected rank. *arXiv preprint arXiv:1709.09565*, 2017.
- [ARR14] A. Ahmed, B. Recht, and J. Romberg. Blind deconvolution using convex programming. *IEEE Transactions on Information Theory*, 60(3):1711–1732, 2014.
- [Bar95] A. I. Barvinok. Problems of distance geometry and convex properties of quadratic maps. *Discrete & Computational Geometry*, 13(2):189–202, 1995.
- [BM03] S. Burer and R. D. Monteiro. A nonlinear programming algorithm for solving semidefinite programs via low-rank factorization. *Mathematical Programming*, 95(2):329–357, 2003.
- [BN06] J. Bai and S. Ng. Confidence intervals for diffusion index forecasts and inference for factor-augmented regressions. *Econometrica*, 74(4):1133–1150, 2006.
- [BT09] A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM journal on imaging sciences*, 2(1):183–202, 2009.
- [BVB16] N. Boumal, V. Voroninski, and A. Bandeira. The non-convex burer-monteiro approach works on smooth semidefinite programs. In *NIPS*, pages 2757–2765, 2016.
- [CC14] Y. Chen and Y. Chi. Robust spectral compressed sensing via structured matrix completion. *IEEE Transactions on Information Theory*, 60(10):6576 – 6601, 2014.
- [CC17] Y. Chen and E. J. Candès. Solving random quadratic systems of equations is nearly as easy as solving linear systems. *Comm. Pure Appl. Math.*, 70(5):822–883, 2017.
- [CC18a] Y. Chen and E. Candès. The projected power method: An efficient algorithm for joint alignment from pairwise differences. *Comm. Pure and Appl. Math.*, 71(8):1648–1714, 2018.
- [CC18b] Y. Chen and Y. Chi. Harnessing structures in big data via guaranteed low-rank matrix estimation: Recent theory and fast algorithms via convex and nonconvex optimization. *IEEE Signal Processing Magazine*, 35(4):14–31, July 2018.
- [CCF18] Y. Chen, C. Cheng, and J. Fan. Asymmetry helps: Eigenvalue and eigenvector analyses of asymmetrically perturbed low-rank matrices. *arXiv preprint arXiv:1811.12804*, 2018.

- [CCFM19] Y. Chen, Y. Chi, J. Fan, and C. Ma. Gradient descent with random initialization: Fast global convergence for nonconvex phase retrieval. *Mathematical Programming*, 176(1-2):5–37, July 2019.
- [CCS10] J. F. Cai, E. J. Candès, and Z. Shen. A singular value thresholding algorithm for matrix completion. *SIAM Journal on Optimization*, 20(4):1956–1982, 2010.
- [CFMW19] Y. Chen, J. Fan, C. Ma, and K. Wang. Spectral method and regularized MLE are both optimal for top- $K$  ranking. *Annals of Statistics*, 47(4):2204–2235, August 2019.
- [CFMY19] Y. Chen, J. Fan, C. Ma, and Y. Yan. Inference and uncertainty quantification for noisy matrix completion. *arXiv preprint arXiv:1906.04159*, 2019.
- [CGH14] Y. Chen, L. J. Guibas, and Q. Huang. Near-optimal joint optimal matching via convex relaxation. *International Conference on Machine Learning (ICML)*, pages 100 – 108, June 2014.
- [Che15] Y. Chen. Incoherence-optimal matrix completion. *IEEE Transactions on Information Theory*, 61(5):2909–2923, 2015.
- [CJSC13] Y. Chen, A. Jalali, S. Sanghavi, and C. Caramanis. Low-rank matrix recovery from errors and erasures. *IEEE Transactions on Information Theory*, 59(7):4324–4337, 2013.
- [CL17] J. Chen and X. Li. Memory-efficient kernel PCA via partial matrix sampling and nonconvex optimization: a model-free analysis of local minima. *arXiv:1711.01742*, 2017.
- [CLC19] Y. Chi, Y. M. Lu, and Y. Chen. Nonconvex optimization meets low-rank matrix factorization: An overview. *IEEE Transactions on Signal Processing*, 67(20):5239 – 5269, October 2019.
- [CLL19] J. Chen, D. Liu, and X. Li. Nonconvex rectangular matrix completion via gradient descent without  $\ell_{2,\infty}$  regularization. *arXiv:1901.06116v1*, 2019.
- [CLMW11] E. Candès, X. Li, Y. Ma, and J. Wright. Robust principal component analysis? *Journal of ACM*, 58(3):11:1–11:37, Jun 2011.
- [CLS15] E. Candès, X. Li, and M. Soltanolkotabi. Phase retrieval via Wirtinger flow: Theory and algorithms. *IEEE Transactions on Information Theory*, 61(4):1985–2007, April 2015.
- [CP10] E. Candès and Y. Plan. Matrix completion with noise. *Proceedings of the IEEE*, 98(6):925–936, June 2010.
- [CR09] E. Candès and B. Recht. Exact matrix completion via convex optimization. *Foundations of Computational Mathematics*, 9(6):717–772, April 2009.
- [CSPW11] V. Chandrasekaran, S. Sanghavi, P. A. Parrilo, and A. S. Willsky. Rank-sparsity incoherence for matrix decomposition. *SIAM Journal on Optimization*, 21(2):572–596, 2011.
- [CT10] E. Candès and T. Tao. The power of convex relaxation: Near-optimal matrix completion. *IEEE Transactions on Information Theory*, 56(5):2053–2080, May 2010.
- [CW15] Y. Chen and M. J. Wainwright. Fast low-rank estimation by projected gradient descent: General statistical and algorithmic guarantees. *arXiv preprint arXiv:1509.03025*, 2015.
- [CWW19] J.-F. Cai, T. Wang, and K. Wei. Fast and provable algorithms for spectrally sparse signal reconstruction via low-rank hankel matrix completion. *Applied and Computational Harmonic Analysis*, 46(1):94–121, 2019.
- [CX16] Y. Cao and Y. Xie. Poisson matrix recovery and completion. *IEEE Transactions on Signal Processing*, 64(6):1609–1620, 2016.
- [CZ16] T. T. Cai and W.-X. Zhou. Matrix completion via max-norm constrained optimization. *Electronic Journal of Statistics*, 10(1):1493–1525, 2016.

- [DC18] L. Ding and Y. Chen. The leave-one-out approach for matrix completion: Primal and dual analysis. *arXiv preprint arXiv:1803.07554*, 2018.
- [DR16] M. A. Davenport and J. Romberg. An overview of low-rank matrix recovery from incomplete observations. *IEEE Journal of Selected Topics in Signal Processing*, 10(4):608–622, 2016.
- [Efr07] B. Efron. Correlation and large-scale simultaneous significance testing. *Journal of the American Statistical Association*, 102(477):93–103, 2007.
- [Efr10] B. Efron. Correlated z-values and the accuracy of large-scale statistical estimates. *Journal of the American Statistical Association*, 105(491):1042–1055, 2010.
- [EK15] N. El Karoui. On the impact of predictor geometry on the performance on high-dimensional ridge-regularized generalized robust regression estimators. *Probability Theory and Related Fields*, pages 1–81, 2015.
- [Faz02] M. Fazel. *Matrix rank minimization with applications*. PhD thesis, 2002.
- [FHB03] M. Fazel, H. Hindi, and S. P. Boyd. Log-det heuristic for matrix rank minimization with applications to Hankel and Euclidean distance matrices. *American Control Conference*, 2003.
- [FHB04] M. Fazel, H. Hindi, and S. Boyd. Rank minimization and applications in system theory. In *American Control Conference*, volume 4, pages 3273–3278, 2004.
- [FHG12] J. Fan, X. Han, and W. Gu. Estimating false discovery proportion under arbitrary covariance dependence. *Journal of the American Statistical Association*, 107(499):1019–1035, 2012.
- [FKSZ19] J. Fan, Y. Ke, Q. Sun, and W.-X. Zhou. Farmtest: Factor-adjusted robust multiple testing with approximate false discovery control. *Journal of American Statistical Association*, 2019+.
- [FKW18] J. Fan, Y. Ke, and K. Wang. Factor-adjusted regularized model selection. *arXiv preprint arXiv:1612.08490*, 2018.
- [FLM13] J. Fan, Y. Liao, and M. Mincheva. Large covariance estimation by thresholding principal orthogonal complements. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 75(4):603–680, 2013.
- [FRW11] M. Fornasier, H. Rauhut, and R. Ward. Low-rank matrix recovery via iteratively reweighted least squares minimization. *SIAM Journal on Optimization*, 21(4):1614–1640, 2011.
- [FWZ19] J. Fan, W. Wang, and Y. Zhong. Robust covariance estimation for approximate factor models. *Journal of econometrics*, 208(1):5–22, 2019.
- [FXY17] J. Fan, L. Xue, and J. Yao. Sufficient forecasting using factor models. *Journal of econometrics*, 201(2):292–306, 2017.
- [GAGG13] S. Gunasekar, A. Acharya, N. Gaur, and J. Ghosh. Noisy matrix completion using alternating minimization. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 194–209, 2013.
- [GJZ17] R. Ge, C. Jin, and Y. Zheng. No spurious local minima in nonconvex low rank problems: A unified geometric analysis. *arXiv preprint arXiv:1704.00708*, 2017.
- [GLM16] R. Ge, J. D. Lee, and T. Ma. Matrix completion has no spurious local minimum. In *Advances in Neural Information Processing Systems*, pages 2973–2981, 2016.
- [Gro11] D. Gross. Recovering low-rank matrices from few coefficients in any basis. *IEEE Transactions on Information Theory*, 57(3):1548–1566, March 2011.
- [Har14] M. Hardt. Understanding alternating minimization for matrix completion. In *Foundations of Computer Science (FOCS)*, pages 651–660, 2014.

- [HG13] Q.-X. Huang and L. Guibas. Consistent shape maps via semidefinite programming. In *Computer Graphics Forum*, volume 32, pages 177–186. Wiley Online Library, 2013.
- [JKN16] C. Jin, S. M. Kakade, and P. Netrapalli. Provable efficient online matrix completion via non-convex stochastic gradient descent. In *NIPS*, pages 4520–4528, 2016.
- [JKS17] P. Jung, F. Krahmer, and D. Stöger. Blind demixing and deconvolution at near-optimal rate. *IEEE Transactions on Information Theory*, 64(2):704–727, 2017.
- [JMD10] P. Jain, R. Meka, and I. S. Dhillon. Guaranteed rank minimization via singular value projection. In *Advances in Neural Information Processing Systems*, pages 937–945, 2010.
- [JNS13] P. Jain, P. Netrapalli, and S. Sanghavi. Low-rank matrix completion using alternating minimization. In *ACM symposium on Theory of computing*, pages 665–674, 2013.
- [Jol82] I. T. Jolliffe. A note on the use of principal components in regression. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 31(3):300–303, 1982.
- [Klo14] O. Klopp. Noisy low-rank matrix completion with general sampling distribution. *Bernoulli*, 20(1):282–303, 2014.
- [KLT11] V. Koltchinskii, K. Lounici, and A. B. Tsybakov. Nuclear-norm penalization and optimal rates for noisy low-rank matrix completion. *Ann. Statist.*, 39(5):2302–2329, 2011.
- [KMO10a] R. H. Keshavan, A. Montanari, and S. Oh. Matrix completion from a few entries. *IEEE Transactions on Information Theory*, 56(6):2980–2998, June 2010.
- [KMO10b] R. H. Keshavan, A. Montanari, and S. Oh. Matrix completion from noisy entries. *J. Mach. Learn. Res.*, 11:2057–2078, 2010.
- [KS11] A. Kneip and P. Sarda. Factor models and variable selection in high-dimensional regression analysis. *The Annals of Statistics*, 39(5):2410–2447, 2011.
- [Lan93] S. Lang. Real and functional analysis. *Springer-Verlag, New York*, 10:11–13, 1993.
- [LMCC18] Y. Li, C. Ma, Y. Chen, and Y. Chi. Nonconvex matrix factorization from rank-one measurements. *arXiv:1802.06286, accepted to AISTATS*, 2018.
- [LS15] S. Ling and T. Strohmer. Self-calibration and biconvex compressive sensing. *Inverse Problems*, 31(11):115002, 2015.
- [LS17] S. Ling and T. Strohmer. Blind deconvolution meets blind demixing: Algorithms and performance bounds. *IEEE Transactions on Information Theory*, 63(7):4497–4520, 2017.
- [LT18] Q. Li and G. Tang. Approximate support recovery of atomic line spectral estimation: A tale of resolution and precision. *Applied and Computational Harmonic Analysis*, 2018.
- [LV09] Z. Liu and L. Vandenberghe. Interior-point method for nuclear norm approximation with application to system identification. *SIAM Journal on Matrix Analysis and Applications*, 31(3):1235–1256, 2009.
- [LXY13] M.-J. Lai, Y. Xu, and W. Yin. Improved iteratively reweighted least squares for unconstrained smoothed  $\ell_q$  minimization. *SIAM Journal on Numerical Analysis*, 51(2):927–957, 2013.
- [MGC11] S. Ma, D. Goldfarb, and L. Chen. Fixed point and bregman iterative methods for matrix rank minimization. *Mathematical Programming*, 128(1-2):321–353, 2011.
- [MHT10] R. Mazumder, T. Hastie, and R. Tibshirani. Spectral regularization algorithms for learning large incomplete matrices. *Journal of machine learning research*, 11(Aug):2287–2322, 2010.

- [MWCC17] C. Ma, K. Wang, Y. Chi, and Y. Chen. Implicit regularization in nonconvex statistical estimation: Gradient descent converges linearly for phase retrieval, matrix completion and blind deconvolution. *arXiv preprint arXiv:1711.10467*, accepted to *Foundations of Computational Mathematics*, 2017.
- [Nes12] Y. Nesterov. How to make the gradients small. *Optima*, 88:10–11, 2012.
- [NW12] S. Negahban and M. Wainwright. Restricted strong convexity and weighted matrix completion: Optimal bounds with noise. *Journal of Machine Learning Research*, pages 1665–1697, May 2012.
- [PB14] N. Parikh and S. Boyd. Proximal algorithms. *Foundations and Trends® in Optimization*, 1(3):127–239, 2014.
- [PBHT08] D. Paul, E. Bair, T. Hastie, and R. Tibshirani. “Preconditioning” for feature selection and regression in high-dimensional problems. *The Annals of Statistics*, 36(4):1595–1618, 2008.
- [PKCS17] D. Park, A. Kyrillidis, C. Carmanis, and S. Sanghavi. Non-square matrix sensing without spurious local minima via the burer-monteiro approach. In *Artificial Intelligence and Statistics*, pages 65–74, 2017.
- [Rec11] B. Recht. A simpler approach to matrix completion. *Journal of Machine Learning Research*, 12(Dec):3413–3430, 2011.
- [RFP10] B. Recht, M. Fazel, and P. A. Parrilo. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM Review*, 52(3):471–501, 2010.
- [RS05] J. D. Rennie and N. Srebro. Fast maximum margin matrix factorization for collaborative prediction. In *International conference on Machine learning*, pages 713–719. ACM, 2005.
- [RT<sup>+</sup>11] A. Rohde, A. B. Tsybakov, et al. Estimation of high-dimensional low-rank matrices. *The Annals of Statistics*, 39(2):887–930, 2011.
- [SCC17] P. Sur, Y. Chen, and E. J. Candès. The likelihood ratio test in high-dimensional logistic regression is asymptotically a rescaled chi-square. *arXiv:1706.01191*, accepted to *Probability Theory and Related Fields*, 2017.
- [Sin11] A. Singer. Angular synchronization by eigenvectors and semidefinite programming. *Applied and computational harmonic analysis*, 30(1):20–36, 2011.
- [SL16] R. Sun and Z.-Q. Luo. Guaranteed matrix completion via non-convex factorization. *IEEE Transactions on Information Theory*, 62(11):6535–6579, 2016.
- [SS05] N. Srebro and A. Shraibman. Rank, trace-norm and max-norm. In *International Conference on Computational Learning Theory*, pages 545–560. Springer, 2005.
- [SXZ19] A. Shapiro, Y. Xie, and R. Zhang. Matrix completion with deterministic pattern: A geometric perspective. *IEEE Transactions on Signal Processing*, 67(4):1088–1103, 2019.
- [SY07] A. M.-C. So and Y. Ye. Theory of semidefinite programming for sensor network localization. *Mathematical Programming*, 109(2-3):367–384, 2007.
- [TBS<sup>+</sup>16] S. Tu, R. Boczar, M. Simchowitz, M. Soltanolkotabi, and B. Recht. Low-rank solutions of linear matrix equations via procrustes flow. In *International Conference on Machine Learning*, pages 964–973, 2016.
- [TK92] C. Tomasi and T. Kanade. Shape and motion from image streams under orthography: a factorization method. *International journal of computer vision*, 9(2):137–154, 1992.
- [Tro15] J. A. Tropp. An introduction to matrix concentration inequalities. *Found. Trends Mach. Learn.*, 8(1-2):1–230, May 2015.

- [TY10] K.-C. Toh and S. Yun. An accelerated proximal gradient algorithm for nuclear norm regularized linear least squares problems. *Pacific Journal of optimization*, 6(615-640):15, 2010.
- [Van13] B. Vandereycken. Low-rank matrix completion by riemannian optimization. *SIAM Journal on Optimization*, 23(2):1214–1236, 2013.
- [Ver12] R. Vershynin. Introduction to the non-asymptotic analysis of random matrices. *Compressed Sensing, Theory and Applications*, pages 210 – 268, 2012.
- [WCCL16] K. Wei, J.-F. Cai, T. Chan, and S. Leung. Guarantees of riemannian optimization for low rank matrix recovery. *SIAM Journal on Matrix Analysis and Applications*, 37(3):1198–1222, 2016.
- [WYZ12] Z. Wen, W. Yin, and Y. Zhang. Solving a low-rank factorization model for matrix completion by a nonlinear successive over-relaxation algorithm. *Mathematical Programming Computation*, 4(4):333–361, 2012.
- [WZG16] L. Wang, X. Zhang, and Q. Gu. A unified computational and statistical framework for nonconvex low-rank matrix estimation. *arXiv preprint arXiv:1610.05275*, 2016.
- [YPCC16] X. Yi, D. Park, Y. Chen, and C. Caramanis. Fast algorithms for robust PCA via gradient descent. In *NIPS*, pages 4152–4160, 2016.
- [ZB18] Y. Zhong and N. Boumal. Near-optimal bound for phase synchronization. *SIAM Journal on Optimization*, 2018.
- [ZL16] Q. Zheng and J. Lafferty. Convergence analysis for rectangular matrix completion using Burer-Monteiro factorization and gradient descent. *arXiv:1605.07051*, 2016.
- [ZLW<sup>+</sup>10] Z. Zhou, X. Li, J. Wright, E. Candès, and Y. Ma. Stable principal component pursuit. In *International Symposium on Information Theory*, pages 1518–1522, 2010.
- [ZPL15] T. Zhang, J. M. Pauly, and I. R. Levesque. Accelerating parameter mapping with a locally low rank constraint. *Magnetic resonance in medicine*, 73(2):655–661, 2015.
- [ZWL15] T. Zhao, Z. Wang, and H. Liu. A nonconvex optimization framework for low rank matrix estimation. In *NIPS*, pages 559–567, 2015.
- [ZZLC17] H. Zhang, Y. Zhou, Y. Liang, and Y. Chi. A nonconvex approach for phase retrieval: Reshaped wirtinger flow and incremental algorithms. *The Journal of Machine Learning Research*, 18(1):5164–5198, 2017.

## A Preliminaries

In this section, we gather a few notations and preliminary facts that are used throughout the proofs.

To begin with, in view of the incoherence assumption (cf. Definition 1), one has

$$\|\mathbf{X}^*\|_{2,\infty} \leq \sqrt{\mu r/n} \|\mathbf{X}^*\| \quad \text{and} \quad \|\mathbf{Y}^*\|_{2,\infty} \leq \sqrt{\mu r/n} \|\mathbf{Y}^*\|. \quad (34)$$

This follows from

$$\|\mathbf{X}^*\|_{2,\infty} = \|\mathbf{U}^* (\boldsymbol{\Sigma}^*)^{1/2}\|_{2,\infty} \leq \|\mathbf{U}^*\|_{2,\infty} \|(\boldsymbol{\Sigma}^*)^{1/2}\| \leq \sqrt{\mu r/n} \|\mathbf{X}^*\|.$$

The bound for  $\mathbf{Y}^*$  follows from the same argument. In addition, we write  $A \ll B$  (resp.  $A \gg B$ ) if there exists a sufficiently small (resp. large) constant  $c$  such that  $A \leq cB$  (resp.  $A \geq cB$ ).

Finally, for notational convenience, we shall often denote

$$\mathcal{P}_\Omega^{\text{debias}}(\mathbf{B}) \triangleq \mathcal{P}_\Omega(\mathbf{B}) - p\mathbf{B}, \quad \text{for all } \mathbf{B} \in \mathbb{R}^{n \times n}. \quad (35)$$

## B Exact duality analysis

We show in this section that why the first-order optimality condition is almost sufficient in guaranteeing the uniqueness of the optimizer. The argument is standard, see e.g. [CR09].

**Lemma 6.** *Let  $\mathbf{Z} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^\top$  be the SVD of  $\mathbf{Z} \in \mathbb{R}^{n \times n}$ . Denote by  $T$  be the tangent space of  $\mathbf{Z}$  and by  $T^\perp$  its orthogonal complement. Suppose that there exists  $\mathbf{W} \in T^\perp$  such that*

$$\frac{1}{\lambda} \mathcal{P}_\Omega(\mathbf{M} - \mathbf{Z}) = \mathbf{U}\mathbf{V}^\top + \mathbf{W}. \quad (36)$$

Then  $\mathbf{Z}$  is the unique minimizer of (3) if

1.  $\|\mathbf{W}\| < 1$ ;
2. The operator  $\mathcal{P}_\Omega(\cdot)$  restricted to elements in  $T$  is injective, i.e.  $\mathcal{P}_\Omega(\mathbf{H}) = \mathbf{0}$  implies  $\mathbf{H} = \mathbf{0}$  for any  $\mathbf{H} \in T$ .

*Proof of Lemma 6.* To begin with, the assumption of this lemma implies that

$$\mathbf{U}\mathbf{V}^\top + \mathbf{W} \in \partial\|\mathbf{Z}\|_*,$$

where  $\partial\|\mathbf{Z}\|_*$  denotes the subdifferential of  $\|\cdot\|_*$  at  $\mathbf{Z}$ . This combined with (36) reveals that

$$\frac{1}{\lambda} \mathcal{P}_\Omega(\mathbf{M} - \mathbf{Z}) \in \partial\|\mathbf{Z}\|_*, \quad (37)$$

thus indicating that  $\mathbf{Z}$  is a minimizer of the convex program (3).

Next, we justify the uniqueness of  $\mathbf{Z}$ . Before continuing, we record a fact regarding the minimizers of (3).

**Claim 1.** *Suppose that  $\mathbf{Z}_1$  and  $\mathbf{Z}_2$  are both minimizers of (3). Then one has  $\mathcal{P}_\Omega(\mathbf{Z}_1) = \mathcal{P}_\Omega(\mathbf{Z}_2)$ .*

With this claim at hand, every minimizer of (3) can be written as  $\mathbf{Z} + \mathbf{H}$  for some  $\mathbf{H}$  obeying  $\mathcal{P}_\Omega(\mathbf{H}) = \mathbf{0}$ . It then suffices to prove that for any  $\mathbf{H} \neq \mathbf{0}$ , one has  $g(\mathbf{Z} + \mathbf{H}) > g(\mathbf{Z})$ , where  $g(\cdot)$  is the objective function in (3). To this end, we note that

$$\begin{aligned} g(\mathbf{Z} + \mathbf{H}) &= \frac{1}{2} \|\mathcal{P}_\Omega(\mathbf{Z} + \mathbf{H} - \mathbf{M})\|_{\mathbb{F}}^2 + \lambda \|\mathbf{Z} + \mathbf{H}\|_* \\ &= \frac{1}{2} \|\mathcal{P}_\Omega(\mathbf{Z} - \mathbf{M})\|_{\mathbb{F}}^2 + \lambda \|\mathbf{Z} + \mathbf{H}\|_*, \end{aligned} \quad (38)$$

where the last relation follows from Claim 1 (i.e.  $\mathcal{P}_\Omega(\mathbf{H}) = \mathbf{0}$ ). Let  $\mathbf{S}$  be a subgradient of  $\|\cdot\|_*$  at point  $\mathbf{Z}$  obeying

$$\mathcal{P}_T(\mathbf{S}) = \mathbf{U}\mathbf{V}^\top, \quad \|\mathcal{P}_{T^\perp}(\mathbf{S})\| \leq 1 \quad \text{and} \quad \langle \mathcal{P}_{T^\perp}(\mathbf{S}), \mathcal{P}_{T^\perp}(\mathbf{H}) \rangle = \|\mathcal{P}_{T^\perp}(\mathbf{H})\|_*. \quad (39)$$

Using the convexity of  $\|\cdot\|_*$ , one can further lower bound (38) by

$$\begin{aligned}
g(\mathbf{Z} + \mathbf{H}) &\geq \frac{1}{2} \|\mathcal{P}_\Omega(\mathbf{Z} - \mathbf{M})\|_{\mathbb{F}}^2 + \lambda(\|\mathbf{Z}\|_* + \langle \mathbf{S}, \mathbf{H} \rangle) \\
&= g(\mathbf{Z}) + \lambda \langle \mathbf{S}, \mathbf{H} \rangle \\
&= g(\mathbf{Z}) + \lambda \langle \mathbf{UV}^\top + \mathbf{W}, \mathbf{H} \rangle + \lambda \langle \mathbf{S} - \mathbf{UV}^\top - \mathbf{W}, \mathbf{H} \rangle \\
&\stackrel{(i)}{=} g(\mathbf{Z}) + \lambda \langle \mathbf{S} - \mathbf{UV}^\top - \mathbf{W}, \mathbf{H} \rangle \\
&\stackrel{(ii)}{=} g(\mathbf{Z}) + \lambda \langle \mathcal{P}_{T^\perp}(\mathbf{S}) - \mathbf{W}, \mathbf{H} \rangle.
\end{aligned}$$

Here, (i) follows from our assumption that  $\mathbf{UV}^\top + \mathbf{W}$  is supported on  $\Omega$  (cf. (36)) and the fact that  $\mathcal{P}_\Omega(\mathbf{H}) = \mathbf{0}$ , and (ii) holds since  $\mathcal{P}_T(\mathbf{S}) = \mathbf{UV}^\top$  (cf. (39)). We can now expand the above expression as

$$\begin{aligned}
g(\mathbf{Z} + \mathbf{H}) &\geq g(\mathbf{Z}) + \lambda \langle \mathcal{P}_{T^\perp}(\mathbf{S}), \mathcal{P}_{T^\perp}(\mathbf{H}) \rangle - \lambda \langle \mathbf{W}, \mathcal{P}_{T^\perp}(\mathbf{H}) \rangle \\
&\geq g(\mathbf{Z}) + \lambda(1 - \|\mathbf{W}\|) \|\mathcal{P}_{T^\perp}(\mathbf{H})\|_*,
\end{aligned} \tag{40}$$

where the last inequality holds by using the last property of (39) and invoking the elementary inequality

$$\langle \mathbf{W}, \mathcal{P}_{T^\perp}(\mathbf{H}) \rangle \leq \|\mathbf{W}\| \|\mathcal{P}_{T^\perp}(\mathbf{H})\|_*.$$

Given that  $\mathbf{W}$  is assumed to obey  $\|\mathbf{W}\| < 1$ , one has  $g(\mathbf{Z} + \mathbf{H}) > g(\mathbf{Z})$  unless  $\mathcal{P}_{T^\perp}(\mathbf{H}) = \mathbf{0}$ . However, if  $\mathcal{P}_{T^\perp}(\mathbf{H}) = \mathbf{0}$  (and hence  $\mathbf{H} \in T$ ), then the injectivity assumption together with the fact that  $\mathcal{P}_\Omega(\mathbf{H}) = \mathbf{0}$  forces  $\mathbf{H} = \mathbf{0}$ . Consequently, any minimizer  $\mathbf{Z} + \mathbf{H}$  with  $\mathbf{H} \neq \mathbf{0}$  must satisfy  $g(\mathbf{Z} + \mathbf{H}) > g(\mathbf{Z})$ , which results in contradiction. This concludes the proof.  $\square$

*Proof of Claim 1.* Consider any minimizers  $\mathbf{Z}_1 \neq \mathbf{Z}_2$ , and suppose instead that  $\mathcal{P}_\Omega(\mathbf{Z}_1 - \mathbf{Z}_2) \neq \mathbf{0}$ . For any  $0 < \alpha < 1$ , define

$$\mathbf{Z}_\alpha \triangleq \alpha \mathbf{Z}_1 + (1 - \alpha) \mathbf{Z}_2.$$

Since  $\|\cdot\|_*$  is convex, we have

$$\begin{aligned}
g(\mathbf{Z}_\alpha) &= \frac{1}{2} \|\mathcal{P}_\Omega(\alpha \mathbf{Z}_1 + (1 - \alpha) \mathbf{Z}_2 - \mathbf{M})\|_{\mathbb{F}}^2 + \lambda \|\alpha \mathbf{Z}_1 + (1 - \alpha) \mathbf{Z}_2\|_* \\
&\leq \frac{1}{2} \|\mathcal{P}_\Omega(\alpha \mathbf{Z}_1 + (1 - \alpha) \mathbf{Z}_2 - \mathbf{M})\|_{\mathbb{F}}^2 + \alpha \lambda \|\mathbf{Z}_1\|_* + (1 - \alpha) \lambda \|\mathbf{Z}_2\|_*.
\end{aligned} \tag{41}$$

Furthermore, by the strong convexity of  $\|\cdot\|_{\mathbb{F}}^2$  we have

$$\begin{aligned}
g(\mathbf{Z}_\alpha) &< \frac{1}{2}(\alpha \|\mathcal{P}_\Omega(\mathbf{Z}_1 - \mathbf{M})\|_{\mathbb{F}}^2 + (1 - \alpha) \|\mathcal{P}_\Omega(\mathbf{Z}_2 - \mathbf{M})\|_{\mathbb{F}}^2) + \alpha \lambda \|\mathbf{Z}_1\|_* + (1 - \alpha) \lambda \|\mathbf{Z}_2\|_* \\
&= \alpha g(\mathbf{Z}_1) + (1 - \alpha) g(\mathbf{Z}_2) = g(\mathbf{Z}_1).
\end{aligned}$$

This contradicts the fact that  $\mathbf{Z}_1$  is a minimizer of (3), thus completing the proof.  $\square$

## C Connections between convex and nonconvex solutions

### C.1 Proof of Lemma 1

First of all, since  $(\mathbf{X}, \mathbf{Y})$  is a stationary point of (5), we have the first-order optimality conditions

$$\mathcal{P}_\Omega(\mathbf{M} - \mathbf{XY}^\top) \mathbf{Y} = \lambda \mathbf{X}; \tag{42a}$$

$$[\mathcal{P}_\Omega(\mathbf{M} - \mathbf{XY}^\top)]^\top \mathbf{X} = \lambda \mathbf{Y}. \tag{42b}$$

As an immediate consequence, one has

$$\mathbf{X}^\top \mathbf{X} = \lambda^{-1} \mathbf{X}^\top \mathcal{P}_\Omega(\mathbf{M} - \mathbf{XY}^\top) \mathbf{Y} = \mathbf{Y}^\top \mathbf{Y}. \tag{43}$$

In words, any stationary point  $(\mathbf{X}, \mathbf{Y})$  has ‘‘balanced’’ scale.



Let  $\mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^\top$  be the singular value decomposition of  $\mathbf{X}\mathbf{Y}^\top$  with  $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{n \times r}$  orthonormal and  $\boldsymbol{\Sigma} \in \mathbb{R}^{r \times r}$  diagonal. In view of the balanced scale of  $(\mathbf{X}, \mathbf{Y})$  (namely, (43)) and Lemma 20, we can write

$$\mathbf{X} = \mathbf{U}\boldsymbol{\Sigma}^{1/2}\mathbf{R} \quad \text{and} \quad \mathbf{Y} = \mathbf{V}\boldsymbol{\Sigma}^{1/2}\mathbf{R} \quad (44)$$

for some orthonormal matrix  $\mathbf{R} \in \mathbb{R}^{r \times r}$ . Substitution into (42) results in

$$\mathcal{P}_\Omega(\mathbf{M} - \mathbf{X}\mathbf{Y}^\top)\mathbf{V} = \lambda\mathbf{U}; \quad (45a)$$

$$[\mathcal{P}_\Omega(\mathbf{M} - \mathbf{X}\mathbf{Y}^\top)]^\top\mathbf{U} = \lambda\mathbf{V}, \quad (45b)$$

implying that the columns of  $\mathbf{U}$  (resp.  $\mathbf{V}$ ) are the left (resp. right) singular vectors of the matrix  $\mathcal{P}_\Omega(\mathbf{M} - \mathbf{X}\mathbf{Y}^\top)$ . We can therefore write

$$\frac{1}{\lambda}\mathcal{P}_\Omega(\mathbf{M} - \mathbf{X}\mathbf{Y}^\top) = \mathbf{U}\mathbf{V}^\top + \mathbf{W}, \quad (46)$$

where  $\mathbf{W} \in T^\perp$ ; recall that  $T$  is the tangent space of  $\mathbf{X}\mathbf{Y}^\top$  and also  $\mathbf{U}\mathbf{V}^\top$ . In view of Lemma 6, it suffices to show that  $\|\mathbf{W}\| < 1$ , which is the content of the rest of the proof.

One can rewrite  $\mathcal{P}_\Omega(\mathbf{M} - \mathbf{X}\mathbf{Y}^\top)$  as

$$\mathcal{P}_\Omega(\mathbf{M} - \mathbf{X}\mathbf{Y}^\top) = p(\mathbf{M}^* - \mathbf{X}\mathbf{Y}^\top) + \mathcal{P}_\Omega^{\text{debias}}(\mathbf{M}^* - \mathbf{X}\mathbf{Y}^\top) + \mathcal{P}_\Omega(\mathbf{E}).$$

Substitute this identity into (45) and rearrange terms to obtain

$$[p\mathbf{M}^* + \mathcal{P}_\Omega^{\text{debias}}(\mathbf{M}^* - \mathbf{X}\mathbf{Y}^\top) + \mathcal{P}_\Omega(\mathbf{E})]\mathbf{V} = \mathbf{U}(p\boldsymbol{\Sigma} + \lambda\mathbf{I}_r);$$

$$[p\mathbf{M}^* + \mathcal{P}_\Omega^{\text{debias}}(\mathbf{M}^* - \mathbf{X}\mathbf{Y}^\top) + \mathcal{P}_\Omega(\mathbf{E})]^\top\mathbf{U} = \mathbf{V}(p\boldsymbol{\Sigma} + \lambda\mathbf{I}_r).$$

These tell us that the columns of  $\mathbf{U}$  (resp.  $\mathbf{V}$ ) are the left (resp. right) singular vectors of the matrix

$$p\mathbf{M}^* + \mathcal{P}_\Omega^{\text{debias}}(\mathbf{M}^* - \mathbf{X}\mathbf{Y}^\top) + \mathcal{P}_\Omega(\mathbf{E}),$$

which is equivalent to saying that<sup>6</sup>

$$p\mathbf{M}^* + \mathcal{P}_\Omega^{\text{debias}}(\mathbf{M}^* - \mathbf{X}\mathbf{Y}^\top) + \mathcal{P}_\Omega(\mathbf{E}) = \mathbf{U}(p\boldsymbol{\Sigma} + \lambda\mathbf{I}_r)\mathbf{V}^\top + \lambda\mathbf{W}_2, \quad (47)$$

for some  $\mathbf{W}_2 \in T^\perp$ . One can then derive from (46) that

$$\begin{aligned} \mathbf{W} &\stackrel{(i)}{=} \frac{1}{\lambda}\mathcal{P}_{T^\perp}[\mathcal{P}_\Omega(\mathbf{M} - \mathbf{X}\mathbf{Y}^\top)] \\ &= \frac{1}{\lambda}\mathcal{P}_{T^\perp}[p\mathbf{M}^* - p\mathbf{X}\mathbf{Y}^\top + \mathcal{P}_\Omega^{\text{debias}}(\mathbf{M}^* - \mathbf{X}\mathbf{Y}^\top) + \mathcal{P}_\Omega(\mathbf{E})] \\ &\stackrel{(ii)}{=} \frac{1}{\lambda}\mathcal{P}_{T^\perp}[p\mathbf{M}^* + \mathcal{P}_\Omega^{\text{debias}}(\mathbf{M}^* - \mathbf{X}\mathbf{Y}^\top) + \mathcal{P}_\Omega(\mathbf{E})] \\ &\stackrel{(iii)}{=} \frac{1}{\lambda}\mathcal{P}_{T^\perp}[\mathbf{U}(p\boldsymbol{\Sigma} + \lambda\mathbf{I}_r)\mathbf{V}^\top + \lambda\mathbf{W}_2] \\ &\stackrel{(iv)}{=} \mathbf{W}_2, \end{aligned}$$

where (i), (ii) and (iv) arise from the facts that  $\mathbf{U}\mathbf{V}^\top \in T$ ,  $\mathbf{X}\mathbf{Y}^\top \in T$  and  $\mathbf{U}(p\boldsymbol{\Sigma} + \lambda\mathbf{I}_r)\mathbf{V}^\top \in T$ , respectively, and (iii) relies on the identity (47).

It then suffices to control  $\|\mathbf{W}_2\|$ . To this end, apply Weyl's inequality to (47) to obtain that: for  $r + 1 \leq i \leq n$ , the  $i$ th largest singular value of  $\mathbf{U}(p\boldsymbol{\Sigma} + \lambda\mathbf{I}_r)\mathbf{V}^\top + \lambda\mathbf{W}_2$  obeys

$$\begin{aligned} \sigma_i(\mathbf{U}(p\boldsymbol{\Sigma} + \lambda\mathbf{I}_r)\mathbf{V}^\top + \lambda\mathbf{W}_2) &\leq p\sigma_i(\mathbf{M}^*) + \|\mathcal{P}_\Omega^{\text{debias}}(\mathbf{M}^* - \mathbf{X}\mathbf{Y}^\top) + \mathcal{P}_\Omega(\mathbf{E})\| \\ &\leq \|\mathcal{P}_\Omega^{\text{debias}}(\mathbf{M}^* - \mathbf{X}\mathbf{Y}^\top)\| + \|\mathcal{P}_\Omega(\mathbf{E})\| \\ &< \lambda, \end{aligned}$$

where the second inequality comes from the fact that  $\mathbf{M}^*$  has rank  $r$  (so that  $\sigma_i(\mathbf{M}^*) = 0$  for  $r + 1 \leq i \leq n$ ) as well as the triangle inequality, and the last inequality follows from the assumptions of the lemma. Furthermore, it is seen that  $\mathbf{U}(p\boldsymbol{\Sigma} + \lambda\mathbf{I}_r)\mathbf{V}^\top$  has rank  $r$  and all of its singular values are at least  $\lambda$ . These facts taken collectively demonstrate that

$$\|\mathbf{W}\| = \|\mathbf{W}_2\| = \frac{1}{\lambda} \max_{r < i \leq n} \sigma_i(\mathbf{U}(p\boldsymbol{\Sigma} + \lambda\mathbf{I}_r)\mathbf{V}^\top + \lambda\mathbf{W}_2) < 1.$$

This together with Lemma 6 completes the proof.

<sup>6</sup>Here, the pre-factor  $\lambda$  is chosen to simplify the analysis later on.

## C.2 Proof of Lemma 2

We begin by collecting a few simple properties resulting from our assumptions. By definition, the gradient of  $f(\cdot, \cdot)$  in (17) is given by

$$\nabla f(\mathbf{X}, \mathbf{Y}) = \frac{1}{p} \begin{bmatrix} \mathcal{P}_\Omega(\mathbf{X}\mathbf{Y}^\top - \mathbf{M})\mathbf{Y} + \lambda\mathbf{X} \\ [\mathcal{P}_\Omega(\mathbf{X}\mathbf{Y}^\top - \mathbf{M})]^\top \mathbf{X} + \lambda\mathbf{Y} \end{bmatrix},$$

which together with the small-gradient assumption  $\|\nabla f(\mathbf{X}, \mathbf{Y})\|_F \leq c\lambda\sqrt{c_{\text{inj}}p\sigma_{\min}/\kappa^2}/p$  implies that

$$\|\mathcal{P}_\Omega(\mathbf{X}\mathbf{Y}^\top - \mathbf{M})\mathbf{Y} + \lambda\mathbf{X}\|_F \leq p\|\nabla f(\mathbf{X}, \mathbf{Y})\|_F \leq c\lambda\sqrt{c_{\text{inj}}p\sigma_{\min}/\kappa^2}; \quad (48a)$$

$$\|(\mathcal{P}_\Omega(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}))^\top \mathbf{X} + \lambda\mathbf{Y}\|_F \leq p\|\nabla f(\mathbf{X}, \mathbf{Y})\|_F \leq c\lambda\sqrt{c_{\text{inj}}p\sigma_{\min}/\kappa^2}. \quad (48b)$$

Throughout the proof, we let the SVD of  $\mathbf{X}\mathbf{Y}^\top$  be  $\mathbf{X}\mathbf{Y}^\top = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$ , and denote by  $T$  the tangent space of  $\mathbf{X}\mathbf{Y}^\top$  and by  $T^\perp$  its orthogonal complement. Additionally, our assumption regarding the singular values of  $\mathbf{X}$  and  $\mathbf{Y}$  implies that

$$\sigma_{\min}/2 \leq \sigma_{\min}(\mathbf{\Sigma}) \leq \sigma_{\max}(\mathbf{\Sigma}) \leq 2\sigma_{\max}. \quad (49)$$

This can be easily seen from the following two inequalities

$$\sigma_{\max}(\mathbf{\Sigma}) = \|\mathbf{X}\mathbf{Y}^\top\| \leq \|\mathbf{X}\| \|\mathbf{Y}\| \leq 2\sigma_{\max};$$

$$\sigma_{\min}(\mathbf{\Sigma}) = \sigma_{\min}(\mathbf{X}\mathbf{Y}^\top) \geq \sigma_{\min}(\mathbf{X}) \sigma_{\min}(\mathbf{Y}) \geq \sigma_{\min}/2.$$

Before proceeding, we record a claim that will prove useful in the subsequent analysis.

**Claim 2.** *Under the notations and assumptions of Lemma 2, one has*

$$\mathcal{P}_\Omega(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}) = -\lambda\mathbf{U}\mathbf{V}^\top + \mathbf{R}, \quad (50)$$

where  $\mathbf{R}$  is some residual matrix satisfying

$$\|\mathcal{P}_T(\mathbf{R})\|_F \leq 72\kappa \frac{p}{\sqrt{\sigma_{\min}}} \|\nabla f(\mathbf{X}, \mathbf{Y})\|_F \quad \text{and} \quad \|\mathcal{P}_{T^\perp}(\mathbf{R})\| < \lambda/2. \quad (51)$$

With Claim 2 in place, we are ready to prove Lemma 2. Let  $\mathbf{Z}_{\text{cvx}}$  be any minimizer of (3) and denote  $\mathbf{\Delta} \triangleq \mathbf{Z}_{\text{cvx}} - \mathbf{X}\mathbf{Y}^\top$ . The proof can be divided into the following steps.

- First, show that the difference  $\mathbf{\Delta}$  primarily lies in the tangent space of  $\mathbf{X}\mathbf{Y}^\top$ ; see (58).
- Next, utilize this property to connect  $\|\mathcal{P}_\Omega(\mathbf{\Delta})\|_F^2$  with the size of the gradient  $\nabla f(\mathbf{X}, \mathbf{Y})$ ; see (60).
- In the end, obtain a lower bound on  $\|\mathcal{P}_\Omega(\mathbf{\Delta})\|_F^2$  in terms of  $\|\mathbf{\Delta}\|_F$  using the injectivity property; see (61).

The desired upper bound on  $\|\mathbf{\Delta}\|_F$  advertised in the lemma then follows by combining these results. In what follows, we shall carry out these steps one by one.

1. The optimality of  $\mathbf{Z}_{\text{cvx}} = \mathbf{X}\mathbf{Y}^\top + \mathbf{\Delta}$  reveals that

$$\frac{1}{2} \|\mathcal{P}_\Omega(\mathbf{X}\mathbf{Y}^\top + \mathbf{\Delta} - \mathbf{M})\|_F^2 + \lambda \|\mathbf{X}\mathbf{Y}^\top + \mathbf{\Delta}\|_* \leq \frac{1}{2} \|\mathcal{P}_\Omega(\mathbf{X}\mathbf{Y}^\top - \mathbf{M})\|_F^2 + \lambda \|\mathbf{X}\mathbf{Y}^\top\|_*.$$

A little algebra allows us to rearrange terms as follows

$$\frac{1}{2} \|\mathcal{P}_\Omega(\mathbf{\Delta})\|_F^2 \leq -\langle \mathcal{P}_\Omega(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}), \mathbf{\Delta} \rangle + \lambda \|\mathbf{X}\mathbf{Y}^\top\|_* - \lambda \|\mathbf{X}\mathbf{Y}^\top + \mathbf{\Delta}\|_*. \quad (52)$$

In addition, it follows from the convexity of  $\|\cdot\|_*$  that

$$\|\mathbf{X}\mathbf{Y}^\top + \mathbf{\Delta}\|_* \geq \|\mathbf{X}\mathbf{Y}^\top\|_* + \langle \mathbf{U}\mathbf{V}^\top + \mathbf{W}, \mathbf{\Delta} \rangle \quad (53)$$

for any  $\mathbf{W} \in T^\perp$  obeying  $\|\mathbf{W}\| \leq 1$ , where  $\mathbf{UV}^\top + \mathbf{W}$  serves as a subgradient of  $\|\cdot\|_*$  at  $\mathbf{XY}^\top$ . In what follows, we shall pick  $\mathbf{W}$  such that  $\langle \mathbf{W}, \Delta \rangle = \|\mathcal{P}_{T^\perp}(\Delta)\|_*$ . Combining this with (52) and (53), we reach

$$\begin{aligned} \frac{1}{2} \|\mathcal{P}_\Omega(\Delta)\|_{\mathbb{F}}^2 &\leq -\langle \mathcal{P}_\Omega(\mathbf{XY}^\top - \mathbf{M}), \Delta \rangle - \lambda \langle \mathbf{UV}^\top, \Delta \rangle - \lambda \langle \mathbf{W}, \Delta \rangle \\ &= -\langle \mathcal{P}_\Omega(\mathbf{XY}^\top - \mathbf{M}), \Delta \rangle - \lambda \langle \mathbf{UV}^\top, \Delta \rangle - \lambda \|\mathcal{P}_{T^\perp}(\Delta)\|_*. \end{aligned} \quad (54)$$

This together with the decomposition (50) leads to

$$\begin{aligned} 0 &\leq \frac{1}{2} \|\mathcal{P}_\Omega(\Delta)\|_{\mathbb{F}}^2 \leq -\langle \mathbf{R}, \Delta \rangle - \lambda \|\mathcal{P}_{T^\perp}(\Delta)\|_* \\ &= -\langle \mathcal{P}_T(\mathbf{R}), \Delta \rangle - \langle \mathcal{P}_{T^\perp}(\mathbf{R}), \Delta \rangle - \lambda \|\mathcal{P}_{T^\perp}(\Delta)\|_*, \end{aligned} \quad (55)$$

and therefore

$$\langle \mathcal{P}_T(\mathbf{R}), \Delta \rangle + \langle \mathcal{P}_{T^\perp}(\mathbf{R}), \Delta \rangle + \lambda \|\mathcal{P}_{T^\perp}(\Delta)\|_* \leq 0. \quad (56)$$

In addition, elementary inequalities give

$$\begin{aligned} -\|\mathcal{P}_T(\mathbf{R})\|_{\mathbb{F}} \|\mathcal{P}_T(\Delta)\|_{\mathbb{F}} - \|\mathcal{P}_{T^\perp}(\mathbf{R})\| \|\mathcal{P}_{T^\perp}(\Delta)\|_* + \lambda \|\mathcal{P}_{T^\perp}(\Delta)\|_* \\ \leq \langle \mathcal{P}_T(\mathbf{R}), \Delta \rangle + \langle \mathcal{P}_{T^\perp}(\mathbf{R}), \Delta \rangle + \lambda \|\mathcal{P}_{T^\perp}(\Delta)\|_* \leq 0. \end{aligned}$$

From the condition (51) we have  $\|\mathcal{P}_{T^\perp}(\mathbf{R})\| \leq \lambda/2$ , and hence the above inequality gives

$$\|\mathcal{P}_T(\mathbf{R})\|_{\mathbb{F}} \|\mathcal{P}_T(\Delta)\|_{\mathbb{F}} \geq -\|\mathcal{P}_{T^\perp}(\mathbf{R})\| \|\mathcal{P}_{T^\perp}(\Delta)\|_* + \lambda \|\mathcal{P}_{T^\perp}(\Delta)\|_* \geq \frac{\lambda}{2} \|\mathcal{P}_{T^\perp}(\Delta)\|_*, \quad (57)$$

which together with the condition (51) on  $\|\mathcal{P}_T(\mathbf{R})\|_{\mathbb{F}}$  and the small gradient assumption (23) yields

$$\|\mathcal{P}_{T^\perp}(\Delta)\|_* \leq 144\kappa \frac{p}{\lambda\sqrt{\sigma_{\min}}} \|\nabla f(\mathbf{X}, \mathbf{Y})\|_{\mathbb{F}} \|\mathcal{P}_T(\Delta)\|_{\mathbb{F}} \leq 144c\sqrt{c_{\text{inj}}p} \|\mathcal{P}_T(\Delta)\|_{\mathbb{F}}. \quad (58)$$

This essentially means that  $\Delta$  lies primarily in the tangent space of  $\mathbf{XY}^\top$  for  $c$  sufficiently small. As an immediate consequence,

$$\|\mathcal{P}_{T^\perp}(\Delta)\|_{\mathbb{F}} \leq \|\mathcal{P}_{T^\perp}(\Delta)\|_* \leq 144c\sqrt{c_{\text{inj}}p} \|\mathcal{P}_T(\Delta)\|_{\mathbb{F}} \leq \|\mathcal{P}_T(\Delta)\|_{\mathbb{F}}, \quad (59)$$

as long as  $c$  is sufficiently small. Note that we also use the elementary fact that  $c_{\text{inj}} \leq 1/p$  (otherwise we will have the contradictory inequality  $p^{-1}\|\mathcal{P}_\Omega(\mathbf{H})\|_{\mathbb{F}}^2 \geq c_{\text{inj}}\|\mathbf{H}\|_{\mathbb{F}}^2 > p^{-1}\|\mathbf{H}\|_{\mathbb{F}}^2$ ).

2. Continue the upper bound in (55) to obtain

$$\begin{aligned} \frac{1}{2} \|\mathcal{P}_\Omega(\Delta)\|_{\mathbb{F}}^2 &\leq -\langle \mathcal{P}_T(\mathbf{R}), \Delta \rangle - \langle \mathcal{P}_{T^\perp}(\mathbf{R}), \Delta \rangle - \lambda \|\mathcal{P}_{T^\perp}(\Delta)\|_* \\ &\leq \|\mathcal{P}_T(\mathbf{R})\|_{\mathbb{F}} \|\mathcal{P}_T(\Delta)\|_{\mathbb{F}} - \frac{\lambda}{2} \|\mathcal{P}_{T^\perp}(\Delta)\|_*. \end{aligned}$$

Here, the last line uses the fact that  $-\langle \mathcal{P}_{T^\perp}(\mathbf{R}), \Delta \rangle \leq \|\mathcal{P}_{T^\perp}(\mathbf{R})\| \cdot \|\mathcal{P}_{T^\perp}(\Delta)\|_* \leq \frac{\lambda}{2} \|\mathcal{P}_{T^\perp}(\Delta)\|_*$ , which follows from (51). Therefore, using the condition (51) we reach

$$\frac{1}{2} \|\mathcal{P}_\Omega(\Delta)\|_{\mathbb{F}}^2 \leq \|\mathcal{P}_T(\mathbf{R})\|_{\mathbb{F}} \|\mathcal{P}_T(\Delta)\|_{\mathbb{F}} \leq 72\kappa \frac{p}{\sqrt{\sigma_{\min}}} \|\nabla f(\mathbf{X}, \mathbf{Y})\|_{\mathbb{F}} \|\Delta\|_{\mathbb{F}}. \quad (60)$$

3. We are left with lower bounding  $\|\mathcal{P}_\Omega(\Delta)\|_{\mathbb{F}}^2$ . Using the decomposition  $\Delta = \mathcal{P}_T(\Delta) + \mathcal{P}_{T^\perp}(\Delta)$ , we obtain

$$\begin{aligned} \frac{1}{\sqrt{p}} \|\mathcal{P}_\Omega(\Delta)\|_{\mathbb{F}} &= \frac{1}{\sqrt{p}} \|\mathcal{P}_\Omega \mathcal{P}_T(\Delta) + \mathcal{P}_\Omega \mathcal{P}_{T^\perp}(\Delta)\|_{\mathbb{F}} \geq \frac{1}{\sqrt{p}} \|\mathcal{P}_\Omega \mathcal{P}_T(\Delta)\|_{\mathbb{F}} - \frac{1}{\sqrt{p}} \|\mathcal{P}_\Omega \mathcal{P}_{T^\perp}(\Delta)\|_{\mathbb{F}} \\ &\geq \sqrt{c_{\text{inj}}} \|\mathcal{P}_T(\Delta)\|_{\mathbb{F}} - \frac{1}{\sqrt{p}} \|\mathcal{P}_{T^\perp}(\Delta)\|_{\mathbb{F}}, \end{aligned}$$

where the last inequality follows from the injectivity assumption (22). In addition, (58) implies

$$\frac{1}{\sqrt{p}} \|\mathcal{P}_{T^\perp}(\Delta)\|_{\mathbb{F}} \leq \frac{1}{\sqrt{p}} \|\mathcal{P}_{T^\perp}(\Delta)\|_* \leq \frac{1}{\sqrt{p}} 144c\sqrt{c_{\text{inj}}p} \|\mathcal{P}_T(\Delta)\|_{\mathbb{F}} \leq \frac{\sqrt{c_{\text{inj}}}}{2} \|\mathcal{P}_T(\Delta)\|_{\mathbb{F}}$$

as long as  $c$  is sufficiently small. As a result,

$$\frac{1}{\sqrt{p}} \|\mathcal{P}_\Omega(\Delta)\|_F \geq \frac{\sqrt{c_{\text{inj}}}}{2} \|\mathcal{P}_T(\Delta)\|_F.$$

In addition, by (59) we have

$$\|\Delta\|_F \leq \|\mathcal{P}_T(\Delta)\|_F + \|\mathcal{P}_{T^\perp}(\Delta)\|_F \leq 2 \|\mathcal{P}_T(\Delta)\|_F,$$

and therefore

$$\frac{1}{\sqrt{p}} \|\mathcal{P}_\Omega(\Delta)\|_F \geq \frac{\sqrt{c_{\text{inj}}}}{2} \|\mathcal{P}_T(\Delta)\|_F \geq \frac{\sqrt{c_{\text{inj}}}}{4} \|\Delta\|_F. \quad (61)$$

Taking (60) and (61) collectively yields

$$\frac{c_{\text{inj}}}{32} \|\Delta\|_F^2 \leq \frac{1}{2p} \|\mathcal{P}_\Omega(\Delta)\|_F^2 \leq 72\kappa \frac{1}{\sqrt{\sigma_{\min}}} \|\nabla f(\mathbf{X}, \mathbf{Y})\|_F \|\Delta\|_F,$$

thus indicating that

$$\|\Delta\|_F \lesssim \frac{\kappa}{c_{\text{inj}}} \frac{1}{\sqrt{\sigma_{\min}}} \|\nabla f(\mathbf{X}, \mathbf{Y})\|_F.$$

### C.2.1 Proof of Claim 2

Before proceeding to the proof of Claim 2, we state a useful fact; the proof is deferred to Appendix C.2.2.

**Claim 3.** *Instate the notations and assumptions in Lemma 2. Let  $\mathbf{U}\Sigma\mathbf{V}^\top$  be the SVD of  $\mathbf{X}\mathbf{Y}^\top$ . There exists an invertible matrix  $\mathbf{Q} \in \mathbb{R}^{r \times r}$  such that  $\mathbf{X} = \mathbf{U}\Sigma^{1/2}\mathbf{Q}$ ,  $\mathbf{Y} = \mathbf{V}\Sigma^{1/2}\mathbf{Q}^{-\top}$  and*

$$\|\Sigma_{\mathbf{Q}} - \Sigma_{\mathbf{Q}}^{-1}\|_F \leq 8\sqrt{\kappa} \frac{p}{\lambda\sqrt{\sigma_{\min}}} \|\nabla f(\mathbf{X}, \mathbf{Y})\|_F \leq 8c\sqrt{c_{\text{inj}}p/\kappa}, \quad (62)$$

where  $\mathbf{U}_{\mathbf{Q}}\Sigma_{\mathbf{Q}}\mathbf{V}_{\mathbf{Q}}^\top$  is the SVD of  $\mathbf{Q}$ .

In light of the assumptions (48), one has

$$\mathcal{P}_\Omega(\mathbf{X}\mathbf{Y}^\top - \mathbf{M})\mathbf{Y} = -\lambda\mathbf{X} + \mathbf{B}_1 \quad \text{and} \quad [\mathcal{P}_\Omega(\mathbf{X}\mathbf{Y}^\top - \mathbf{M})]^\top \mathbf{X} = -\lambda\mathbf{Y} + \mathbf{B}_2 \quad (63)$$

for some  $\mathbf{B}_1 \in \mathbb{R}^{n \times r}$  and  $\mathbf{B}_2 \in \mathbb{R}^{n \times r}$ , where  $\max\{\|\mathbf{B}_1\|_F, \|\mathbf{B}_2\|_F\} \leq p\|\nabla f(\mathbf{X}, \mathbf{Y})\|_F$ . Recall that

$$\mathcal{P}_\Omega(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}) = -\lambda\mathbf{U}\mathbf{V}^\top + \mathbf{R}. \quad (64)$$

In the sequel, we shall prove the upper bounds on both  $\|\mathcal{P}_T(\mathbf{R})\|_F$  and  $\|\mathcal{P}_{T^\perp}(\mathbf{R})\|_F$  separately.

1. From the definition of  $\mathcal{P}_T(\cdot)$  (see (15)), we have

$$\begin{aligned} \|\mathcal{P}_T(\mathbf{R})\|_F &= \|\mathbf{U}\mathbf{U}^\top \mathbf{R}(\mathbf{I} - \mathbf{V}\mathbf{V}^\top) + \mathbf{R}\mathbf{V}\mathbf{V}^\top\|_F \\ &\leq \|\mathbf{U}^\top \mathbf{R}(\mathbf{I} - \mathbf{V}\mathbf{V}^\top)\|_F + \|\mathbf{R}\mathbf{V}\|_F \\ &\leq \|\mathbf{U}^\top \mathbf{R}\|_F + \|\mathbf{R}\mathbf{V}\|_F. \end{aligned} \quad (65)$$

In addition, invoke Claim 3 to obtain

$$\mathbf{X} = \mathbf{U}\Sigma^{1/2}\mathbf{Q} \quad \text{and} \quad \mathbf{Y} = \mathbf{V}\Sigma^{1/2}\mathbf{Q}^{-\top} \quad (66)$$

for some invertible matrix  $\mathbf{Q} \in \mathbb{R}^{r \times r}$ , whose SVD  $\mathbf{U}_{\mathbf{Q}}\Sigma_{\mathbf{Q}}\mathbf{V}_{\mathbf{Q}}^\top$  obeys (62). Combine (63) and (64) to see

$$-\lambda\mathbf{U}\mathbf{V}^\top \mathbf{Y} + \mathbf{R}\mathbf{Y} = -\lambda\mathbf{X} + \mathbf{B}_1,$$

which together with (66) yields

$$\mathbf{R}\mathbf{V} = \lambda\mathbf{U}\Sigma^{1/2}(\mathbf{I}_r - \mathbf{Q}\mathbf{Q}^\top)\Sigma^{-1/2} + \mathbf{B}_1\mathbf{Q}^\top\Sigma^{-1/2}.$$

Apply the triangle inequality to get

$$\begin{aligned}\|\mathbf{R}\mathbf{V}\|_{\mathbb{F}} &\leq \|\lambda\mathbf{U}\boldsymbol{\Sigma}^{1/2}(\mathbf{I}_r - \mathbf{Q}\mathbf{Q}^\top)\boldsymbol{\Sigma}^{-1/2}\|_{\mathbb{F}} + \|\mathbf{B}_1\mathbf{Q}^\top\boldsymbol{\Sigma}^{-1/2}\|_{\mathbb{F}} \\ &\leq \lambda\|\boldsymbol{\Sigma}^{1/2}\|\|\boldsymbol{\Sigma}^{-1/2}\|\|\mathbf{Q}\mathbf{Q}^\top - \mathbf{I}_r\|_{\mathbb{F}} + \|\mathbf{Q}\| \|\boldsymbol{\Sigma}^{-1/2}\| \|\mathbf{B}_1\|_{\mathbb{F}}.\end{aligned}\quad (67)$$

In order to further upper bound (67), we first recognize that (49) implies

$$\|\boldsymbol{\Sigma}^{1/2}\| \leq \sqrt{2\sigma_{\max}}, \quad \text{and} \quad \|\boldsymbol{\Sigma}^{-1/2}\| = 1/\sqrt{\sigma_{\min}(\boldsymbol{\Sigma})} \leq \sqrt{2/\sigma_{\min}}.$$

Second, Claim 3 yields

$$\|\boldsymbol{\Sigma}_{\mathbf{Q}} - \boldsymbol{\Sigma}_{\mathbf{Q}}^{-1}\|_{\mathbb{F}} \leq 8\sqrt{\kappa} \frac{p}{\lambda\sqrt{\sigma_{\min}}} \|\nabla f(\mathbf{X}, \mathbf{Y})\|_{\mathbb{F}} \leq 8c\sqrt{c_{\text{inj}}p/\kappa} \ll 1,$$

with the proviso that  $c$  is sufficiently small. Here we have used the facts that  $c_{\text{inj}} \leq 1/p$  and that  $\kappa \geq 1$ . This in turn implies that  $\|\mathbf{Q}\| = \|\boldsymbol{\Sigma}_{\mathbf{Q}}\| \leq 2$ . Putting the above bounds together yields

$$\begin{aligned}\|\mathbf{R}\mathbf{V}\|_{\mathbb{F}} &\leq \lambda\sqrt{2\sigma_{\max}}\sqrt{\frac{2}{\sigma_{\min}}}\|\boldsymbol{\Sigma}_{\mathbf{Q}}^2 - \mathbf{I}_r\|_{\mathbb{F}} + 2\sqrt{\frac{2}{\sigma_{\min}}}p\|\nabla f(\mathbf{X}, \mathbf{Y})\|_{\mathbb{F}} \\ &\leq \lambda\sqrt{2\sigma_{\max}}\sqrt{\frac{2}{\sigma_{\min}}}\|\boldsymbol{\Sigma}_{\mathbf{Q}}\|\|\boldsymbol{\Sigma}_{\mathbf{Q}} - \boldsymbol{\Sigma}_{\mathbf{Q}}^{-1}\|_{\mathbb{F}} + 2\sqrt{\frac{2}{\sigma_{\min}}}p\|\nabla f(\mathbf{X}, \mathbf{Y})\|_{\mathbb{F}} \\ &\leq 2\lambda\sqrt{2\sigma_{\max}}\sqrt{\frac{2}{\sigma_{\min}}}8\sqrt{\kappa} \frac{p}{\lambda\sqrt{\sigma_{\min}}}\|\nabla f(\mathbf{X}, \mathbf{Y})\|_{\mathbb{F}} + 2\sqrt{\frac{2}{\sigma_{\min}}}p\|\nabla f(\mathbf{X}, \mathbf{Y})\|_{\mathbb{F}} \\ &\leq 36\kappa \frac{p}{\sqrt{\sigma_{\min}}}\|\nabla f(\mathbf{X}, \mathbf{Y})\|_{\mathbb{F}}.\end{aligned}$$

Similarly we can show that  $\|\mathbf{U}^\top \mathbf{R}\|_{\mathbb{F}} \leq 36\kappa p \|\nabla f(\mathbf{X}, \mathbf{Y})\|_{\mathbb{F}} / \sqrt{\sigma_{\min}}$ . These bounds together with (65) result in

$$\|\mathcal{P}_T(\mathbf{R})\|_{\mathbb{F}} \leq 72\kappa \frac{p}{\sqrt{\sigma_{\min}}}\|\nabla f(\mathbf{X}, \mathbf{Y})\|_{\mathbb{F}}. \quad (68)$$

2. We now move on to bounding  $\|\mathcal{P}_{T^\perp}(\mathbf{R})\|$ . In view of the definition of  $\mathcal{P}_\Omega^{\text{debias}}(\cdot)$  in (35), we can rearrange (63) to derive

$$\begin{aligned}[p\mathbf{M}^* + \mathcal{P}_\Omega(\mathbf{E}) - \mathcal{P}_\Omega^{\text{debias}}(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}^*)]\mathbf{Y} &= p\mathbf{X}\mathbf{Y}^\top\mathbf{Y} + \lambda\mathbf{X} - \mathbf{B}_1, \\ [p\mathbf{M}^* + \mathcal{P}_\Omega(\mathbf{E}) - \mathcal{P}_\Omega^{\text{debias}}(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}^*)]^\top\mathbf{X} &= p\mathbf{Y}\mathbf{X}^\top\mathbf{X} + \lambda\mathbf{Y} - \mathbf{B}_2.\end{aligned}$$

In view of the representation  $\mathbf{X} = \mathbf{U}\boldsymbol{\Sigma}^{1/2}\mathbf{Q}$  and  $\mathbf{Y} = \mathbf{V}\boldsymbol{\Sigma}^{1/2}\mathbf{Q}^{-\top}$ , the above identities are equivalent to

$$\begin{aligned}[p\mathbf{M}^* + \mathcal{P}_\Omega(\mathbf{E}) - \mathcal{P}_\Omega^{\text{debias}}(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}^*)]\mathbf{V} &= p\mathbf{U}\boldsymbol{\Sigma} + \lambda\mathbf{U}\boldsymbol{\Sigma}^{1/2}\mathbf{Q}\mathbf{Q}^\top\boldsymbol{\Sigma}^{-1/2} - \mathbf{B}_1\mathbf{Q}^\top\boldsymbol{\Sigma}^{-1/2}, \\ [p\mathbf{M}^* + \mathcal{P}_\Omega(\mathbf{E}) - \mathcal{P}_\Omega^{\text{debias}}(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}^*)]^\top\mathbf{U} &= p\mathbf{V}\boldsymbol{\Sigma} + \lambda\mathbf{V}\boldsymbol{\Sigma}^{1/2}\mathbf{Q}^{-\top}\mathbf{Q}^{-1}\boldsymbol{\Sigma}^{-1/2} - \mathbf{B}_2\mathbf{Q}^{-1}\boldsymbol{\Sigma}^{-1/2}.\end{aligned}$$

Letting

$$p\mathbf{M}^* + \mathcal{P}_\Omega(\mathbf{E}) - \mathcal{P}_\Omega^{\text{debias}}(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}^*) = p\mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^\top + \lambda\mathbf{U}\boldsymbol{\Sigma}^{1/2}\mathbf{Q}\mathbf{Q}^\top\boldsymbol{\Sigma}^{-1/2}\mathbf{V}^\top + \tilde{\mathbf{R}} \quad (69)$$

for some residual matrix  $\tilde{\mathbf{R}} \in \mathbb{R}^{n \times n}$ , we have

$$\begin{aligned}\mathcal{P}_{T^\perp}(\mathbf{R}) &\stackrel{(i)}{=} \mathcal{P}_{T^\perp}[\mathcal{P}_\Omega(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}^* - \mathbf{E})] \\ &\stackrel{(ii)}{=} \mathcal{P}_{T^\perp}[p(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}^*) + \mathcal{P}_\Omega^{\text{debias}}(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}^*) - \mathcal{P}_\Omega(\mathbf{E})] \\ &\stackrel{(iii)}{=} \mathcal{P}_{T^\perp}[p\mathbf{M}^* + \mathcal{P}_\Omega(\mathbf{E}) - \mathcal{P}_\Omega^{\text{debias}}(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}^*)] \\ &\stackrel{(iv)}{=} \mathcal{P}_{T^\perp}(\tilde{\mathbf{R}}),\end{aligned}\quad (70)$$

where (i) follows from the definition of  $\mathbf{R}$  and the fact that  $\mathbf{UV}^\top \in T$ , (ii) uses the definition of  $\mathcal{P}_\Omega^{\text{debias}}(\cdot)$ , (iii) relies on the fact that  $\mathbf{XY}^\top \in T$ , and (iv) applies (69) and the facts that  $\mathbf{U}\Sigma\mathbf{V}^\top \in T$  and that  $\mathbf{U}\Sigma^{1/2}\mathbf{Q}\mathbf{Q}^\top\Sigma^{-1/2}\mathbf{V}^\top \in T$ . Therefore, it suffices to bound  $\|\mathcal{P}_{T^\perp}(\tilde{\mathbf{R}})\|$ . Rewrite (69) as

$$p\mathbf{M}^* + \mathcal{P}_\Omega(\mathbf{E}) - \mathcal{P}_\Omega^{\text{debias}}(\mathbf{XY}^\top - \mathbf{M}^*) - \mathcal{P}_T(\tilde{\mathbf{R}}) = \mathbf{U}(p\Sigma + \lambda\Sigma^{1/2}\mathbf{Q}\mathbf{Q}^\top\Sigma^{-1/2})\mathbf{V}^\top + \mathcal{P}_{T^\perp}(\tilde{\mathbf{R}}). \quad (71)$$

Suppose for the moment that

$$\|\mathcal{P}_T(\tilde{\mathbf{R}})\| \leq \lambda/4. \quad (72)$$

This together with the assumptions that  $\|\mathcal{P}_\Omega^{\text{debias}}(\mathbf{XY}^\top - \mathbf{M}^*)\| < \lambda/8$  and  $\|\mathcal{P}_\Omega(\mathbf{E})\| < \lambda/8$  reveals that

$$\|\mathcal{P}_\Omega(\mathbf{E}) - \mathcal{P}_\Omega^{\text{debias}}(\mathbf{XY}^\top - \mathbf{M}^*) - \mathcal{P}_T(\tilde{\mathbf{R}})\| < \lambda/2. \quad (73)$$

By Weyl's inequality and the relations (71) and (73), one has

$$\begin{aligned} \sigma_i \left[ \mathbf{U}(p\Sigma + \lambda\Sigma^{1/2}\mathbf{Q}\mathbf{Q}^\top\Sigma^{-1/2})\mathbf{V}^\top + \mathcal{P}_{T^\perp}(\tilde{\mathbf{R}}) \right] &\leq \sigma_i(p\mathbf{M}^*) + \|\mathcal{P}_\Omega(\mathbf{E}) - \mathcal{P}_\Omega^{\text{debias}}(\mathbf{XY}^\top - \mathbf{M}^*) - \mathcal{P}_T(\tilde{\mathbf{R}})\| \\ &< p\sigma_i(\mathbf{M}^*) + \lambda/2 = \lambda/2 \end{aligned} \quad (74)$$

for any  $r+1 \leq i \leq n$ , where  $\sigma_i(\mathbf{A})$  denotes the  $i$ th largest singular value of a matrix  $\mathbf{A}$ . Here, we have used the fact that  $\mathbf{M}^*$  has rank  $r$  and hence  $\sigma_i(\mathbf{M}^*) = 0$  for any  $i > r$ . In addition, it is seen that

$$\begin{aligned} \|\Sigma^{1/2}\mathbf{Q}\mathbf{Q}^\top\Sigma^{-1/2} - \mathbf{I}_r\| &= \|\Sigma^{1/2}(\mathbf{Q}\mathbf{Q}^\top - \mathbf{I}_r)\Sigma^{-1/2}\| \\ &\leq \|\Sigma^{1/2}\| \|\Sigma^{-1/2}\| \|\mathbf{Q}\mathbf{Q}^\top - \mathbf{I}_r\|_{\text{F}}. \end{aligned}$$

Note that in (67), we have obtained

$$\|\Sigma^{1/2}\| \|\Sigma^{-1/2}\| \|\mathbf{Q}\mathbf{Q}^\top - \mathbf{I}_r\|_{\text{F}} \leq 2\sqrt{2\sigma_{\max}}\sqrt{2/\sigma_{\min}}8c\sqrt{c_{\text{inj}}p/\kappa} \leq 1/10$$

as long as  $c$  is sufficiently small, and hence  $\|\Sigma^{1/2}\mathbf{Q}\mathbf{Q}^\top\Sigma^{-1/2} - \mathbf{I}_r\| \leq 1/10$ . Therefore, for any  $1 \leq i \leq r$  we know that

$$\begin{aligned} \sigma_i \left[ \mathbf{U}(p\Sigma + \lambda\Sigma^{1/2}\mathbf{Q}\mathbf{Q}^\top\Sigma^{-1/2})\mathbf{V}^\top \right] &\geq \sigma_r \left[ \mathbf{U}(p\Sigma + \lambda\mathbf{I}_r + \lambda(\Sigma^{1/2}\mathbf{Q}\mathbf{Q}^\top\Sigma^{-1/2} - \mathbf{I}_r))\mathbf{V}^\top \right] \\ &\geq \sigma_r(p\Sigma + \lambda\mathbf{I}_r) - \lambda\|\Sigma^{1/2}\mathbf{Q}\mathbf{Q}^\top\Sigma^{-1/2} - \mathbf{I}_r\| \\ &\geq \lambda - \lambda\|\Sigma^{1/2}\mathbf{Q}\mathbf{Q}^\top\Sigma^{-1/2} - \mathbf{I}_r\| \\ &\geq \lambda - \lambda/10 > \lambda/2, \end{aligned}$$

where the second inequality results from Weyl's inequality. This combined with (70) and (74) yields

$$\|\mathcal{P}_{T^\perp}(\mathbf{R})\| = \|\mathcal{P}_{T^\perp}(\tilde{\mathbf{R}})\| < \lambda/2;$$

this happens because at least  $n-r$  singular values of  $\mathbf{U}(p\Sigma + \lambda\Sigma^{1/2}\mathbf{Q}\mathbf{Q}^\top\Sigma^{-1/2})\mathbf{V}^\top + \mathcal{P}_{T^\perp}(\tilde{\mathbf{R}})$  are no larger than  $\lambda/2$  and they cannot correspond to directions simultaneously in the column space spanned by  $\mathbf{U}$  and the row space spanned by  $\mathbf{V}^\top$ .

The proof is then complete by verifying (72). To this end, observe that

$$\tilde{\mathbf{R}}\mathbf{V} = -\mathbf{B}_1\mathbf{Q}^\top\Sigma^{-1/2}, \quad \tilde{\mathbf{R}}^\top\mathbf{U} = \lambda\mathbf{V}\Sigma^{1/2}\mathbf{Q}^{-\top}\mathbf{Q}^{-1}\Sigma^{-1/2} - \lambda\mathbf{V}\Sigma^{-1/2}\mathbf{Q}\mathbf{Q}^\top\Sigma^{1/2} - \mathbf{B}_2\mathbf{Q}^{-1}\Sigma^{-1/2}.$$

Then following similar technique used to bound  $\|\mathcal{P}_T(\mathbf{R})\|$ , we have

$$\|\mathcal{P}_T(\tilde{\mathbf{R}})\| \leq \|\mathcal{P}_T(\tilde{\mathbf{R}})\|_{\text{F}} \leq \|\mathbf{U}^\top\tilde{\mathbf{R}}\|_{\text{F}} + \|\tilde{\mathbf{R}}\mathbf{V}\|_{\text{F}} \lesssim c\sqrt{c_{\text{inj}}p}\lambda < \lambda/4 \quad (75)$$

as long as  $c$  is small enough.

### C.2.2 Proof of Claim 3

Let

$$\mathcal{P}_\Omega(\mathbf{X}\mathbf{Y}^\top - \mathbf{M})\mathbf{Y} + \lambda\mathbf{X} = \mathbf{B}_1 \quad \text{and} \quad [\mathcal{P}_\Omega(\mathbf{X}\mathbf{Y}^\top - \mathbf{M})]^\top \mathbf{X} + \lambda\mathbf{Y} = \mathbf{B}_2 \quad (76)$$

for some  $\mathbf{B}_1, \mathbf{B}_2 \in \mathbb{R}^{n \times r}$ . Clearly, it is seen from the assumption (48) that

$$\max\{\|\mathbf{B}_1\|_F, \|\mathbf{B}_2\|_F\} \leq p\|\nabla f(\mathbf{X}, \mathbf{Y})\|_F. \quad (77)$$

In addition, the identities (76) allow us to obtain

$$\begin{aligned} \|\mathbf{X}^\top \mathbf{X} - \mathbf{Y}^\top \mathbf{Y}\|_F &= \frac{1}{\lambda} \|\mathbf{X}^\top (\mathbf{B}_1 - \mathcal{P}_\Omega(\mathbf{X}\mathbf{Y}^\top - \mathbf{M})\mathbf{Y}) - (\mathbf{B}_2 - [\mathcal{P}_\Omega(\mathbf{X}\mathbf{Y}^\top - \mathbf{M})]^\top \mathbf{X})^\top \mathbf{Y}\|_F \\ &= \frac{1}{\lambda} \|\mathbf{X}^\top \mathbf{B}_1 - \mathbf{B}_2^\top \mathbf{Y}\|_F \\ &\leq \frac{1}{\lambda} \|\mathbf{X}\| \|\mathbf{B}_1\|_F + \frac{1}{\lambda} \|\mathbf{B}_2\|_F \|\mathbf{Y}\| \\ &\leq 2\frac{p}{\lambda} \sqrt{2\sigma_{\max}} \|\nabla f(\mathbf{X}, \mathbf{Y})\|_F. \end{aligned} \quad (78)$$

Here, the last line makes use of (77) and the assumption that  $\|\mathbf{X}\|, \|\mathbf{Y}\| \leq \sqrt{2\sigma_{\max}}$ . In view of Lemma 20, one can find an invertible  $\mathbf{Q}$  such that  $\mathbf{X} = \mathbf{U}\Sigma^{1/2}\mathbf{Q}$ ,  $\mathbf{Y} = \mathbf{V}\Sigma^{1/2}\mathbf{Q}^{-\top}$  and

$$\begin{aligned} \|\Sigma_{\mathbf{Q}} - \Sigma_{\mathbf{Q}}^{-1}\|_F &\leq \frac{1}{\sigma_{\min}(\Sigma)} \|\mathbf{X}^\top \mathbf{X} - \mathbf{Y}^\top \mathbf{Y}\|_F \\ &\stackrel{(i)}{\leq} \frac{2}{\sigma_{\min}} \cdot 2\frac{p}{\lambda} \sqrt{2\sigma_{\max}} \|\nabla f(\mathbf{X}, \mathbf{Y})\|_F \\ &\leq 8\sqrt{\kappa} \frac{p}{\lambda\sqrt{\sigma_{\min}}} \|\nabla f(\mathbf{X}, \mathbf{Y})\|_F \\ &\stackrel{(ii)}{\leq} 8c\sqrt{c_{\text{inj}}p/\kappa}, \end{aligned}$$

where  $\Sigma_{\mathbf{Q}}$  is a diagonal matrix consisting of all singular values of  $\mathbf{Q}$ . Here, (i) follows from (49) as well as the bound (78), and the last inequality (ii) uses the assumption (23). This completes the proof.

### C.3 Proof of Lemma 4

Lemma 4 consists of two parts, which we restate into the following two lemmas, namely Lemmas 7-8.

First of all, Lemma 7 demonstrates that as long as  $(\mathbf{X}, \mathbf{Y})$  is sufficiently close to  $(\mathbf{X}^*, \mathbf{Y}^*)$ , the operator  $\mathcal{P}_\Omega(\cdot)$  restricted to the tangent space  $T$  of  $\mathbf{X}\mathbf{Y}^\top$  is injective. The proof is deferred to Appendix C.3.1.

**Lemma 7.** *Suppose that the sample complexity obeys  $n^2p \geq C\mu rn \log n$  for some sufficiently large constant  $C > 0$ . Then with probability exceeding  $1 - O(n^{-10})$ ,*

$$\frac{1}{p} \|\mathcal{P}_\Omega(\mathbf{H})\|_F^2 \geq \frac{1}{32\kappa} \|\mathbf{H}\|_F^2, \quad \forall \mathbf{H} \in T$$

holds simultaneously for all  $(\mathbf{X}, \mathbf{Y})$  obeying

$$\max\{\|\mathbf{X} - \mathbf{X}^*\|_{2,\infty}, \|\mathbf{Y} - \mathbf{Y}^*\|_{2,\infty}\} \leq \frac{c}{\kappa\sqrt{n}} \|\mathbf{X}^*\|. \quad (79)$$

Here,  $c > 0$  is some sufficiently small constant, and  $T$  denotes the tangent space of  $\mathbf{X}\mathbf{Y}^\top$ .

**Remark 7.** In the prior literature, the injectivity of  $\mathcal{P}_\Omega(\cdot)$  has been mostly studied when restricted to a fixed tangent space independent of  $\Omega$  (see [CR09, Gro11]). In comparison, this lemma demonstrates that the injectivity property holds uniformly over a large set of tangent spaces. This allows one to handle tangent spaces that are statistically dependent on  $\Omega$ .

**Remark 8.** Note that the condition (79) on  $(\mathbf{X}, \mathbf{Y})$  is weaker than (26) under the assumptions of Lemma 4. To see this, if (26) holds, then one necessarily has

$$\begin{aligned}
\|\mathbf{X} - \mathbf{X}^*\|_{2,\infty} &\leq C_\infty \kappa \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \max \left\{ \|\mathbf{X}^*\|_{2,\infty}, \|\mathbf{Y}^*\|_{2,\infty} \right\} \\
&\stackrel{(i)}{\lesssim} C_\infty \kappa \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} \max \left\{ \|\mathbf{X}^*\|_{2,\infty}, \|\mathbf{Y}^*\|_{2,\infty} \right\} \\
&\stackrel{(ii)}{\leq} C_\infty \kappa \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} \sqrt{\frac{\mu r}{n}} \|\mathbf{X}^*\| \\
&\stackrel{(iii)}{\leq} \frac{c}{\kappa \sqrt{n}} \|\mathbf{X}^*\|.
\end{aligned}$$

Here, (i) follows from the choice  $\lambda \asymp \sigma \sqrt{np}$ ; (ii) relies on the incoherence assumption (34); and (iii) holds true under the noise condition  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll \frac{1}{\sqrt{\kappa^4 \mu r \log n}}$ . A similar bound holds for  $\|\mathbf{Y} - \mathbf{Y}^*\|_{2,\infty}$ .

The next lemma shows that for all  $(\mathbf{X}, \mathbf{Y})$  close to  $(\mathbf{X}^*, \mathbf{Y}^*)$ ,  $\mathcal{P}_\Omega(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}^*)$  is uniformly close to its expectation  $p(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}^*)$ . The proof can be found in Appendix C.3.2.

**Lemma 8.** Suppose that  $n^2 p \gg \kappa^4 \mu^2 r^2 n \log^2 n$  and  $\sigma \sqrt{n(\log n)/p} \ll \sigma_{\min}/\kappa$ . With probability exceeding  $1 - O(n^{-10})$ , one has

$$\|\mathcal{P}_\Omega(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}^*) - p(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}^*)\| < \lambda/8$$

simultaneously for any  $(\mathbf{X}, \mathbf{Y})$  obeying (26), provided that  $\lambda = C_\lambda \sigma \sqrt{np}$  for some constant  $C_\lambda > 0$ .

### C.3.1 Proof of Lemma 7

By definition, any  $\mathbf{H} \in T$  can be expressed as

$$\mathbf{H} = \mathbf{X}\mathbf{A}^\top + \mathbf{B}\mathbf{Y}^\top \tag{80}$$

for some  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times r}$ . Given that this is an underdetermined linear system of equations, there might be numerous  $(\mathbf{A}, \mathbf{B})$ 's compatible with (80). We take a specific choice as follows

$$\begin{aligned}
(\mathbf{A}, \mathbf{B}) &:= \arg \min_{(\tilde{\mathbf{A}}, \tilde{\mathbf{B}})} 0.5 \|\tilde{\mathbf{A}}\|_{\text{F}}^2 + 0.5 \|\tilde{\mathbf{B}}\|_{\text{F}}^2 \\
&\text{subject to } \mathbf{H} = \mathbf{X}\tilde{\mathbf{A}}^\top + \tilde{\mathbf{B}}\mathbf{Y}^\top.
\end{aligned} \tag{81}$$

which satisfies a property that plays an important role in the subsequent analysis:

$$\mathbf{X}^\top \mathbf{B} = \mathbf{A}^\top \mathbf{Y}. \tag{82}$$

To see this, consider the Lagrangian

$$\mathcal{L}(\tilde{\mathbf{A}}, \tilde{\mathbf{B}}, \mathbf{\Lambda}) := 0.5 \|\tilde{\mathbf{A}}\|_{\text{F}}^2 + 0.5 \|\tilde{\mathbf{B}}\|_{\text{F}}^2 + \langle \mathbf{\Lambda}, \mathbf{X}\tilde{\mathbf{A}}^\top + \tilde{\mathbf{B}}\mathbf{Y}^\top - \mathbf{H} \rangle.$$

Taking the derivatives w.r.t.  $\tilde{\mathbf{A}}$  and  $\tilde{\mathbf{B}}$  and setting them to zero yield

$$\mathbf{A} = -\mathbf{\Lambda}^\top \mathbf{X} \quad \text{and} \quad \mathbf{B} = -\mathbf{\Lambda} \mathbf{Y}$$

for some Lagrangian multiplier matrix  $\mathbf{\Lambda} \in \mathbb{R}^{n \times n}$ . The claim (82) then follows immediately.

The remaining proof consists of two steps.

- First, we would like to show that

$$\|\mathbf{H}\|_{\text{F}}^2 \leq 8\sigma_{\max} (\|\mathbf{A}\|_{\text{F}}^2 + \|\mathbf{B}\|_{\text{F}}^2). \tag{83}$$



- Second, we prove that

$$\frac{1}{2p} \|\mathcal{P}_\Omega(\mathbf{H})\|_F^2 = \frac{1}{2p} \|\mathcal{P}_\Omega(\mathbf{X}\mathbf{A}^\top + \mathbf{B}\mathbf{Y}^\top)\|_F^2 \geq \frac{\sigma_{\min}}{8} (\|\mathbf{A}\|_F^2 + \|\mathbf{B}\|_F^2). \quad (84)$$

Taking (83) and (84) together immediately yields the claimed bounds in the lemma. In what follows, we shall establish these two bounds separately.

1. Regarding the upper bound (83), it follows from elementary inequalities that

$$\begin{aligned} \|\mathbf{H}\|_F^2 &= \|\mathbf{X}\mathbf{A}^\top + \mathbf{B}\mathbf{Y}^\top\|_F^2 \leq 2(\|\mathbf{X}\mathbf{A}^\top\|_F^2 + \|\mathbf{B}\mathbf{Y}^\top\|_F^2) \\ &\leq 2(\|\mathbf{X}\|^2 \|\mathbf{A}\|_F^2 + \|\mathbf{Y}\|^2 \|\mathbf{B}\|_F^2) \\ &\leq 2 \max\{\|\mathbf{X}\|^2, \|\mathbf{Y}\|^2\} (\|\mathbf{A}\|_F^2 + \|\mathbf{B}\|_F^2). \end{aligned} \quad (85)$$

It then suffices to control  $\max\{\|\mathbf{X}\|, \|\mathbf{Y}\|\}$ . In view of the assumption (79), one has

$$\|\mathbf{X} - \mathbf{X}^*\| \leq \|\mathbf{X} - \mathbf{X}^*\|_F \leq \sqrt{n} \|\mathbf{X} - \mathbf{X}^*\|_{2,\infty} \leq \frac{c}{\kappa} \|\mathbf{X}^*\| \leq \|\mathbf{X}^*\|, \quad (86)$$

as long as  $c < 1$ . This together with the triangle inequality reveals that

$$\|\mathbf{X}\| \leq \|\mathbf{X}^*\| + \|\mathbf{X} - \mathbf{X}^*\| \leq 2\|\mathbf{X}^*\| \leq 2\sqrt{\sigma_{\max}}.$$

Similarly, one has  $\|\mathbf{Y}\| \leq 2\sqrt{\sigma_{\max}}$ . Substitution into (85) yields the desired upper bound (83).

2. We now move on to the lower bound (84). To this end, one first decomposes

$$\frac{1}{2p} \|\mathcal{P}_\Omega(\mathbf{X}\mathbf{A}^\top + \mathbf{B}\mathbf{Y}^\top)\|_F^2 = \underbrace{\frac{1}{2p} \|\mathcal{P}_\Omega(\mathbf{X}\mathbf{A}^\top + \mathbf{B}\mathbf{Y}^\top)\|_F^2}_{:=\alpha_1} - \frac{1}{2} \|\mathbf{X}\mathbf{A}^\top + \mathbf{B}\mathbf{Y}^\top\|_F^2 + \underbrace{\frac{1}{2} \|\mathbf{X}\mathbf{A}^\top + \mathbf{B}\mathbf{Y}^\top\|_F^2}_{:=\alpha_2}.$$

The basic idea is to demonstrate that (1)  $\alpha_2$  is bounded from below, and (2)  $\alpha_1$  is sufficiently small compared to  $\alpha_2$ .

- (a) We start by controlling  $\alpha_2$ , towards which we can expand

$$\alpha_2 = \frac{1}{2} \left( \|\mathbf{X}\mathbf{A}^\top\|_F^2 + \|\mathbf{B}\mathbf{Y}^\top\|_F^2 \right) + \text{Tr}(\mathbf{X}^\top \mathbf{B}\mathbf{Y}^\top \mathbf{A}).$$

The property  $\mathbf{X}^\top \mathbf{B} = \mathbf{A}^\top \mathbf{Y}$  (see (82)) implies that

$$\text{Tr}(\mathbf{X}^\top \mathbf{B}\mathbf{Y}^\top \mathbf{A}) = \|\mathbf{X}^\top \mathbf{B}\|_F^2 \geq 0 \quad \implies \quad \alpha_2 \geq \frac{1}{2} \left( \|\mathbf{X}\mathbf{A}^\top\|_F^2 + \|\mathbf{B}\mathbf{Y}^\top\|_F^2 \right).$$

Write  $\mathbf{\Delta}_X = \mathbf{X} - \mathbf{X}^*$  and  $\mathbf{\Delta}_Y = \mathbf{Y} - \mathbf{Y}^*$ . We have

$$\begin{aligned} \|\mathbf{X}\mathbf{A}^\top\|_F^2 &= \|(\mathbf{X}^* + \mathbf{\Delta}_X)\mathbf{A}^\top\|_F^2 = \|\mathbf{X}^*\mathbf{A}^\top\|_F^2 + \|\mathbf{\Delta}_X\mathbf{A}^\top\|_F^2 + 2\langle \mathbf{X}^*\mathbf{A}^\top, \mathbf{\Delta}_X\mathbf{A}^\top \rangle \\ &\geq \|\mathbf{X}^*\mathbf{A}^\top\|_F^2 - 2\|\mathbf{X}^*\mathbf{A}^\top\|_F \|\mathbf{\Delta}_X\mathbf{A}^\top\|_F \\ &\geq \|\mathbf{X}^*\mathbf{A}^\top\|_F^2 - 2\|\mathbf{X}^*\| \|\mathbf{\Delta}_X\| \|\mathbf{A}\|_F^2, \end{aligned}$$

where the second line arises from the Cauchy-Schwarz inequality. Recalling from (86) that  $\|\mathbf{\Delta}_X\| \leq c\|\mathbf{X}^*\|/\kappa$ , we arrive at

$$\|\mathbf{X}\mathbf{A}^\top\|_F^2 \geq \|\mathbf{X}^*\mathbf{A}^\top\|_F^2 - 2c\sigma_{\min} \|\mathbf{A}\|_F^2 \geq \|\mathbf{X}^*\mathbf{A}^\top\|_F^2 - \sigma_{\min} \|\mathbf{A}\|_F^2 / 100,$$

provided that  $c \leq 1/200$ . A similar bound holds for  $\|\mathbf{B}\mathbf{Y}^\top\|_F^2$ , thus leading to

$$\alpha_2 \geq \frac{1}{2} \left( \|\mathbf{X}^*\mathbf{A}^\top\|_F^2 + \|\mathbf{B}\mathbf{Y}^{\star\top}\|_F^2 \right) - \frac{1}{100} \sigma_{\min} \left( \|\mathbf{A}\|_F^2 + \|\mathbf{B}\|_F^2 \right).$$

(b) Next, we control  $\alpha_1$ . First, it is seen that

$$\begin{aligned}\mathbf{X}\mathbf{A}^\top + \mathbf{B}\mathbf{Y}^\top &= (\mathbf{X}^* + \Delta_{\mathbf{X}})\mathbf{A}^\top + \mathbf{B}(\mathbf{Y}^* + \Delta_{\mathbf{Y}})^\top \\ &= \mathbf{X}^*\mathbf{A}^\top + \mathbf{B}\mathbf{Y}^{*\top} + \Delta_{\mathbf{X}}\mathbf{A}^\top + \mathbf{B}\Delta_{\mathbf{Y}}^\top.\end{aligned}$$

As a result, we can expand  $\alpha_1$  as

$$\begin{aligned}\alpha_1 &= \frac{1}{2p} \|\mathcal{P}_\Omega(\mathbf{X}^*\mathbf{A}^\top + \mathbf{B}\mathbf{Y}^{*\top} + \Delta_{\mathbf{X}}\mathbf{A}^\top + \mathbf{B}\Delta_{\mathbf{Y}}^\top)\|_{\mathbb{F}}^2 - \frac{1}{2} \|\mathbf{X}^*\mathbf{A}^\top + \mathbf{B}\mathbf{Y}^{*\top} + \Delta_{\mathbf{X}}\mathbf{A}^\top + \mathbf{B}\Delta_{\mathbf{Y}}^\top\|_{\mathbb{F}}^2 \\ &= \frac{1}{2p} \underbrace{\|\mathcal{P}_\Omega(\mathbf{X}^*\mathbf{A}^\top + \mathbf{B}\mathbf{Y}^{*\top})\|_{\mathbb{F}}^2 - \frac{1}{2} \|\mathbf{X}^*\mathbf{A}^\top + \mathbf{B}\mathbf{Y}^{*\top}\|_{\mathbb{F}}^2}_{:=\gamma_1} \\ &\quad + \frac{1}{2p} \underbrace{\|\mathcal{P}_\Omega(\Delta_{\mathbf{X}}\mathbf{A}^\top)\|_{\mathbb{F}}^2 - \frac{1}{2} \|\Delta_{\mathbf{X}}\mathbf{A}^\top\|_{\mathbb{F}}^2}_{:=\gamma_2} + \frac{1}{2p} \underbrace{\|\mathcal{P}_\Omega(\mathbf{B}\Delta_{\mathbf{Y}}^\top)\|_{\mathbb{F}}^2 - \frac{1}{2} \|\mathbf{B}\Delta_{\mathbf{Y}}^\top\|_{\mathbb{F}}^2}_{:=\gamma_3} \\ &\quad + \frac{1}{p} \underbrace{\langle \mathcal{P}_\Omega(\Delta_{\mathbf{X}}\mathbf{A}^\top), \mathcal{P}_\Omega(\mathbf{B}\Delta_{\mathbf{Y}}^\top) \rangle - \langle \Delta_{\mathbf{X}}\mathbf{A}^\top, \mathbf{B}\Delta_{\mathbf{Y}}^\top \rangle}_{:=\gamma_4} \\ &\quad + \frac{1}{p} \underbrace{\langle \mathcal{P}_\Omega(\mathbf{X}^*\mathbf{A}^\top + \mathbf{B}\mathbf{Y}^{*\top}), \mathcal{P}_\Omega(\Delta_{\mathbf{X}}\mathbf{A}^\top + \mathbf{B}\Delta_{\mathbf{Y}}^\top) \rangle - \langle \mathbf{X}^*\mathbf{A}^\top + \mathbf{B}\mathbf{Y}^{*\top}, \Delta_{\mathbf{X}}\mathbf{A}^\top + \mathbf{B}\Delta_{\mathbf{Y}}^\top \rangle}_{:=\gamma_5}.\end{aligned}$$

i. Regarding  $\gamma_1$ , it follows from the bounds in [CR09, Section 4.2] that

$$|\gamma_1| \leq \frac{1}{64} \|\mathbf{X}^*\mathbf{A}^\top + \mathbf{B}\mathbf{Y}^{*\top}\|_{\mathbb{F}}^2 \leq \frac{1}{32} \left( \|\mathbf{X}^*\mathbf{A}^\top\|_{\mathbb{F}}^2 + \|\mathbf{B}\mathbf{Y}^{*\top}\|_{\mathbb{F}}^2 \right),$$

as long as  $np \gg \mu r \log n$ .

ii. Invoke Lemma 19 to show that

$$\begin{aligned}|\gamma_2| &\leq \frac{3n}{2} \|\Delta_{\mathbf{X}}\|_{2,\infty}^2 \|\mathbf{A}\|_{\mathbb{F}}^2 \leq \frac{3c^2}{2\kappa} \sigma_{\min} \|\mathbf{A}\|_{\mathbb{F}}^2 \leq \frac{1}{100} \sigma_{\min} \|\mathbf{A}\|_{\mathbb{F}}^2, \\ |\gamma_3| &\leq \frac{3n}{2} \|\Delta_{\mathbf{Y}}\|_{2,\infty}^2 \|\mathbf{B}\|_{\mathbb{F}}^2 \leq \frac{3c^2}{2\kappa} \sigma_{\min} \|\mathbf{B}\|_{\mathbb{F}}^2 \leq \frac{1}{100} \sigma_{\min} \|\mathbf{B}\|_{\mathbb{F}}^2,\end{aligned}$$

as long as  $n^2p \gg n \log n$  and  $c > 0$  is sufficiently small. Here we have utilized the assumption that  $\max\{\|\Delta_{\mathbf{X}}\|_{2,\infty}, \|\Delta_{\mathbf{Y}}\|_{2,\infty}\} \leq c\|\mathbf{X}^*\|/(\kappa\sqrt{n})$ .

iii. The term  $\gamma_4$  can be controlled via Lemma 21:

$$\begin{aligned}|\gamma_4| &\leq \left\| \frac{1}{p} \mathcal{P}_\Omega(\mathbf{1}\mathbf{1}^\top) - \mathbf{1}\mathbf{1}^\top \right\| \|\Delta_{\mathbf{X}}\|_{2,\infty} \|\mathbf{A}\|_{\mathbb{F}} \|\Delta_{\mathbf{Y}}\|_{2,\infty} \|\mathbf{B}\|_{\mathbb{F}} \\ &\lesssim \sqrt{\frac{n}{p}} \|\Delta_{\mathbf{X}}\|_{2,\infty} \|\mathbf{A}\|_{\mathbb{F}} \|\Delta_{\mathbf{Y}}\|_{2,\infty} \|\mathbf{B}\|_{\mathbb{F}},\end{aligned}$$

where the second line uses the bound  $\|p^{-1}\mathcal{P}_\Omega(\mathbf{1}\mathbf{1}^\top) - \mathbf{1}\mathbf{1}^\top\| \lesssim \sqrt{n/p}$  guaranteed by [KMO10a, Lemma 3.2]. Continue the upper bound to get

$$|\gamma_4| \stackrel{(i)}{\lesssim} n \frac{c^2}{\kappa^2 n} \sigma_{\max} \|\mathbf{A}\|_{\mathbb{F}} \|\mathbf{B}\|_{\mathbb{F}} \stackrel{(ii)}{\leq} \frac{c^2}{2\kappa} \sigma_{\min} \left( \|\mathbf{A}\|_{\mathbb{F}}^2 + \|\mathbf{B}\|_{\mathbb{F}}^2 \right) \stackrel{(iii)}{\leq} \frac{1}{100} \sigma_{\min} \left( \|\mathbf{A}\|_{\mathbb{F}}^2 + \|\mathbf{B}\|_{\mathbb{F}}^2 \right).$$

Here the first relation (i) arises from the assumption that  $np \gg 1$ . The second inequality (ii) applies the elementary inequality  $ab \leq (a^2 + b^2)/2$  and the last one (iii) holds with the proviso that  $c > 0$  is small enough.

- iv. The last term  $\gamma_5$  can be further decomposed into the sum of four terms. For brevity, we take one out as an example, namely the term

$$\frac{1}{p} \langle \mathcal{P}_\Omega(\mathbf{X}^* \mathbf{A}^\top), \mathcal{P}_\Omega(\Delta_{\mathbf{X}} \mathbf{A}^\top) \rangle - \langle \mathbf{X}^* \mathbf{A}^\top, \Delta_{\mathbf{X}} \mathbf{A}^\top \rangle.$$

Apply the triangle inequality to obtain

$$\begin{aligned} & \left| \frac{1}{p} \langle \mathcal{P}_\Omega(\mathbf{X}^* \mathbf{A}^\top), \mathcal{P}_\Omega(\Delta_{\mathbf{X}} \mathbf{A}^\top) \rangle - \langle \mathbf{X}^* \mathbf{A}^\top, \Delta_{\mathbf{X}} \mathbf{A}^\top \rangle \right| \\ & \leq \left| \frac{1}{p} \langle \mathcal{P}_\Omega(\mathbf{X}^* \mathbf{A}^\top), \mathcal{P}_\Omega(\Delta_{\mathbf{X}} \mathbf{A}^\top) \rangle \right| + |\langle \mathbf{X}^* \mathbf{A}^\top, \Delta_{\mathbf{X}} \mathbf{A}^\top \rangle| \\ & \leq \frac{1}{\sqrt{p}} \|\mathcal{P}_\Omega(\mathbf{X}^* \mathbf{A}^\top)\|_{\text{F}} \frac{1}{\sqrt{p}} \|\mathcal{P}_\Omega(\Delta_{\mathbf{X}} \mathbf{A}^\top)\|_{\text{F}} + \|\mathbf{X}^* \mathbf{A}^\top\|_{\text{F}} \|\Delta_{\mathbf{X}} \mathbf{A}^\top\|_{\text{F}}. \end{aligned}$$

In light of [CR09, Section 4.2] and [ZL16, Lemma 9], we have

$$\begin{aligned} \frac{1}{\sqrt{p}} \|\mathcal{P}_\Omega(\mathbf{X}^* \mathbf{A}^\top)\|_{\text{F}} & \leq 1.1 \|\mathbf{X}^* \mathbf{A}^\top\|_{\text{F}}; \\ \frac{1}{\sqrt{p}} \|\mathcal{P}_\Omega(\Delta_{\mathbf{X}} \mathbf{A}^\top)\|_{\text{F}} & \leq \sqrt{2n} \|\Delta_{\mathbf{X}}\|_{2,\infty} \|\mathbf{A}\|_{\text{F}}. \end{aligned}$$

Taking the above three bounds collectively yields

$$\begin{aligned} & \left| \frac{1}{p} \langle \mathcal{P}_\Omega(\mathbf{X}^* \mathbf{A}^\top), \mathcal{P}_\Omega(\Delta_{\mathbf{X}} \mathbf{A}^\top) \rangle - \langle \mathbf{X}^* \mathbf{A}^\top, \Delta_{\mathbf{X}} \mathbf{A}^\top \rangle \right| \\ & \leq 5 \|\mathbf{X}^* \mathbf{A}^\top\|_{\text{F}} \sqrt{n} \|\Delta_{\mathbf{X}}\|_{2,\infty} \|\mathbf{A}\|_{\text{F}} + \|\mathbf{X}^* \mathbf{A}^\top\|_{\text{F}} \|\Delta_{\mathbf{X}} \mathbf{A}^\top\|_{\text{F}} \\ & \leq 5\sqrt{n} \|\Delta_{\mathbf{X}}\|_{2,\infty} \|\mathbf{X}^*\| \|\mathbf{A}\|_{\text{F}}^2 + \sqrt{n} \|\Delta_{\mathbf{X}}\|_{2,\infty} \|\mathbf{X}^*\| \|\mathbf{A}\|_{\text{F}}^2 \\ & = 6\sqrt{n} \|\Delta_{\mathbf{X}}\|_{2,\infty} \|\mathbf{X}^*\| \|\mathbf{A}\|_{\text{F}}^2. \end{aligned}$$

Using the assumption that  $\|\Delta_{\mathbf{X}}\|_{2,\infty} \leq c \|\mathbf{X}^*\| / (\kappa \sqrt{n})$ , one has

$$\left| \frac{1}{p} \langle \mathcal{P}_\Omega(\mathbf{X}^* \mathbf{A}^\top), \mathcal{P}_\Omega(\Delta_{\mathbf{X}} \mathbf{A}^\top) \rangle - \langle \mathbf{X}^* \mathbf{A}^\top, \Delta_{\mathbf{X}} \mathbf{A}^\top \rangle \right| \lesssim \sqrt{n} \frac{c}{\kappa \sqrt{n}} \sigma_{\max} \|\mathbf{A}\|_{\text{F}}^2 \leq \frac{1}{100} \sigma_{\min} \|\mathbf{A}\|_{\text{F}}^2$$

for  $c > 0$  small enough. The same argument applies to the remaining three terms, resulting in

$$|\gamma_5| \leq \frac{1}{50} \sigma_{\min} \left( \|\mathbf{A}\|_{\text{F}}^2 + \|\mathbf{B}\|_{\text{F}}^2 \right).$$

- v. Combining the previous bounds on  $\gamma_1$  through  $\gamma_5$ , we arrive at

$$\begin{aligned} |\alpha_1| & \leq |\gamma_1| + |\gamma_2| + |\gamma_3| + |\gamma_4| + |\gamma_5| \\ & \leq \frac{1}{32} \left( \|\mathbf{X}^* \mathbf{A}^\top\|_{\text{F}}^2 + \|\mathbf{B} \mathbf{Y}^{*\top}\|_{\text{F}}^2 \right) + \frac{1}{25} \sigma_{\min} \left( \|\mathbf{A}\|_{\text{F}}^2 + \|\mathbf{B}\|_{\text{F}}^2 \right). \end{aligned}$$

- (c) Taking the preceding bounds on  $\alpha_1$  and  $\alpha_2$  collectively yields

$$\begin{aligned} & \frac{1}{2p} \|\mathcal{P}_\Omega(\mathbf{X} \mathbf{A}^\top + \mathbf{B} \mathbf{Y}^\top)\|_{\text{F}}^2 \geq \alpha_2 - |\alpha_1| \\ & \geq \frac{15}{32} \left( \|\mathbf{X}^* \mathbf{A}^\top\|_{\text{F}}^2 + \|\mathbf{B} \mathbf{Y}^{*\top}\|_{\text{F}}^2 \right) - \frac{1}{5} \sigma_{\min} \left( \|\mathbf{A}\|_{\text{F}}^2 + \|\mathbf{B}\|_{\text{F}}^2 \right) \\ & \geq \frac{15}{32} \sigma_{\min} \left( \|\mathbf{A}\|_{\text{F}}^2 + \|\mathbf{B}\|_{\text{F}}^2 \right) - \frac{1}{5} \sigma_{\min} \left( \|\mathbf{A}\|_{\text{F}}^2 + \|\mathbf{B}\|_{\text{F}}^2 \right) \\ & \geq \frac{1}{8} \sigma_{\min} \left( \|\mathbf{A}\|_{\text{F}}^2 + \|\mathbf{B}\|_{\text{F}}^2 \right). \end{aligned}$$

The proof is then complete.

### C.3.2 Proof of Lemma 8

To start with, we have

$$\mathbf{X}\mathbf{Y}^\top - \mathbf{M}^* = (\mathbf{X} - \mathbf{X}^*)\mathbf{Y}^\top + \mathbf{X}^*(\mathbf{Y} - \mathbf{Y}^*)^\top,$$

which together with the triangle inequality implies

$$\|\mathcal{P}_\Omega^{\text{debias}}(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}^*)\| \leq \|\mathcal{P}_\Omega^{\text{debias}}[(\mathbf{X} - \mathbf{X}^*)\mathbf{Y}^\top]\| + \|\mathcal{P}_\Omega^{\text{debias}}[\mathbf{X}^*(\mathbf{Y} - \mathbf{Y}^*)^\top]\|.$$

Apply [CL17, Lemma 4.5] to obtain

$$\begin{aligned} \|\mathcal{P}_\Omega^{\text{debias}}[(\mathbf{X} - \mathbf{X}^*)\mathbf{Y}^\top]\| &\leq \|\mathcal{P}_\Omega^{\text{debias}}(\mathbf{1}\mathbf{1}^\top)\| \|\mathbf{X} - \mathbf{X}^*\|_{2,\infty} \|\mathbf{Y}\|_{2,\infty} \\ &\lesssim \sqrt{np} \|\mathbf{X} - \mathbf{X}^*\|_{2,\infty} \|\mathbf{Y}\|_{2,\infty}, \end{aligned}$$

where the second line is due to  $\|\mathcal{P}_\Omega^{\text{debias}}(\mathbf{1}\mathbf{1}^\top)\| \lesssim \sqrt{np}$  (cf. [KMO10a, Lemma 3.2]). Similarly,

$$\|\mathcal{P}_\Omega^{\text{debias}}[\mathbf{X}^*(\mathbf{Y} - \mathbf{Y}^*)^\top]\| \lesssim \sqrt{np} \|\mathbf{Y} - \mathbf{Y}^*\|_{2,\infty} \|\mathbf{X}^*\|_{2,\infty}.$$

In addition, the assumption (26) yields

$$\begin{aligned} \|\mathbf{Y}\|_{2,\infty} &\leq \|\mathbf{Y} - \mathbf{Y}^*\|_{2,\infty} + \|\mathbf{Y}^*\|_{2,\infty} \\ &\leq C_\infty \kappa \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{Y}^*\|_{2,\infty} + \|\mathbf{Y}^*\|_{2,\infty} \\ &\leq 2 \|\mathbf{Y}^*\|_{2,\infty}, \end{aligned}$$

as long as  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} \ll 1/\kappa$  (recall that  $\lambda = C_\lambda \sigma \sqrt{np}$  for some constant  $C_\lambda > 0$ ). As a consequence, one obtains

$$\begin{aligned} \|\mathcal{P}_\Omega^{\text{debias}}(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}^*)\| &\lesssim \sqrt{np} \kappa \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{X}^*\|_{2,\infty} \|\mathbf{Y}^*\|_{2,\infty} \\ &\leq \sqrt{np} \kappa \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \frac{\mu r \sigma_{\max}}{n}, \end{aligned} \quad (87)$$

where the last inequality follows from the upper bound  $\max\{\|\mathbf{X}^*\|_{2,\infty}, \|\mathbf{Y}^*\|_{2,\infty}\} \leq \sqrt{\mu r \sigma_{\max}/n}$  (cf. (34)). Rearrange the right-hand side of (87) to reach

$$\|\mathcal{P}_\Omega^{\text{debias}}(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}^*)\| \lesssim \sigma \sqrt{np} \cdot \sqrt{\frac{\kappa^4 \mu^2 r^2 \log n}{np}} + \lambda \sqrt{\frac{\kappa^4 \mu^2 r^2}{np}} < \lambda/8,$$

where the last line holds because of the assumption  $n^2 p \gg \kappa^4 \mu^2 r^2 n \log n$  as well as the choice of  $\lambda$ .

## D Analysis of the nonconvex gradient descent algorithm

Lemma 5 shares similar spirit as [MWCC17, Theorem 2] and [CLL19, Lemma 3.5] with one difference: the nonconvex loss function (17) has an additional term  $\|\mathbf{X}\|_{\text{F}}^2 + \|\mathbf{Y}\|_{\text{F}}^2$  to balance the scale of  $\mathbf{X}$  and  $\mathbf{Y}$ . To simplify the presentation, we find it convenient to introduce a few notations. Denote

$$\mathbf{F}^t \triangleq \begin{bmatrix} \mathbf{X}^t \\ \mathbf{Y}^t \end{bmatrix} \in \mathbb{R}^{2n \times r} \quad \text{and} \quad \mathbf{F}^* \triangleq \begin{bmatrix} \mathbf{X}^* \\ \mathbf{Y}^* \end{bmatrix} \in \mathbb{R}^{2n \times r}. \quad (88)$$

It is easily seen from (28) that

$$\mathbf{H}^t = \arg \min_{\mathbf{R} \in \mathcal{O}^{r \times r}} \|\mathbf{F}^t \mathbf{R} - \mathbf{F}^*\|_{\text{F}}. \quad (89)$$

---

**Algorithm 2** Construction of the  $l$ th leave-one-out sequence.

---

**Initialization:**  $\mathbf{X}^{0,(l)} = \mathbf{X}^*$ ;  $\mathbf{Y}^{0,(l)} = \mathbf{Y}^*$ ; Set  $\mathbf{F}^{0,(l)} \triangleq \begin{bmatrix} \mathbf{X}^{0,(l)} \\ \mathbf{Y}^{0,(l)} \end{bmatrix}$ .

**Gradient updates:** for  $t = 0, 1, \dots, t_0 - 1$  do

$$\mathbf{F}^{t+1,(l)} \triangleq \begin{bmatrix} \mathbf{X}^{t+1,(l)} \\ \mathbf{Y}^{t+1,(l)} \end{bmatrix} = \begin{bmatrix} \mathbf{X}^{t,(l)} - \eta \nabla_{\mathbf{X}} f^{(l)}(\mathbf{X}^{t,(l)}, \mathbf{Y}^{t,(l)}) \\ \mathbf{Y}^{t,(l)} - \eta \nabla_{\mathbf{Y}} f^{(l)}(\mathbf{X}^{t,(l)}, \mathbf{Y}^{t,(l)}) \end{bmatrix},$$

where  $\eta > 0$  is the step size.

---

Similar to [MWCC17, CLL19], we resort to the leave-one-out sequences to control the  $\ell_2/\ell_\infty$  error. Specifically, for each  $1 \leq l \leq n$  (corresponding to row indices), we construct  $\{\mathbf{F}^{t,(l)}\}_{t \geq 0}$  to be the gradient descent iterates (see Algorithm 2) w.r.t. the following auxiliary loss function

$$f^{(l)}(\mathbf{X}, \mathbf{Y}) = \frac{1}{2p} \|\mathcal{P}_{\Omega_{-l,\cdot}}(\mathbf{X}\mathbf{Y}^\top - \mathbf{M})\|_{\text{F}}^2 + \frac{1}{2} \|\mathcal{P}_{l,\cdot}(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}^*)\|_{\text{F}}^2 + \frac{\lambda}{2p} \|\mathbf{X}\|_{\text{F}}^2 + \frac{\lambda}{2p} \|\mathbf{Y}\|_{\text{F}}^2. \quad (91)$$

Here  $\mathcal{P}_{\Omega_{-l,\cdot}}(\cdot)$  (resp.  $\mathcal{P}_{l,\cdot}(\cdot)$ ) denotes the orthogonal projection onto the space of matrices which are supported on the index set  $\Omega_{-l,\cdot} = \{(i, j) \in \Omega | i \neq l\}$  (resp.  $\{(i, j) | i = l\}$ ). Mathematically, we have for any matrix  $\mathbf{B} \in \mathbb{R}^{n \times n}$

$$[\mathcal{P}_{\Omega_{-l,\cdot}}(\mathbf{B})]_{ij} = \begin{cases} B_{ij}, & \text{if } (i, j) \in \Omega \text{ and } i \neq l, \\ 0, & \text{otherwise} \end{cases} \quad \text{and} \quad [\mathcal{P}_{l,\cdot}(\mathbf{B})]_{ij} = \begin{cases} B_{ij}, & \text{if } i = l, \\ 0, & \text{otherwise.} \end{cases} \quad (92)$$

Similarly, for each  $n+1 \leq l \leq 2n$  (with  $l-n$  corresponding to the column index), we define  $\{\mathbf{F}^{t,(l)}\}_{t \geq 0}$  to be the GD iterates (see Algorithm 2) operating on

$$f^{(l)}(\mathbf{X}, \mathbf{Y}) = \frac{1}{2p} \|\mathcal{P}_{\Omega_{\cdot, -(l-n)}}(\mathbf{X}\mathbf{Y}^\top - \mathbf{M})\|_{\text{F}}^2 + \frac{1}{2} \|\mathcal{P}_{\cdot, (l-n)}(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}^*)\|_{\text{F}}^2 + \frac{\lambda}{2p} \|\mathbf{X}\|_{\text{F}}^2 + \frac{\lambda}{2p} \|\mathbf{Y}\|_{\text{F}}^2,$$

where  $\mathcal{P}_{\Omega_{\cdot, -(l-n)}}(\cdot)$  and  $\mathcal{P}_{\cdot, (l-n)}(\cdot)$  are defined as

$$[\mathcal{P}_{\Omega_{\cdot, -(l-n)}}(\mathbf{B})]_{ij} = \begin{cases} B_{ij}, & \text{if } (i, j) \in \Omega \text{ and } j \neq l-n, \\ 0, & \text{otherwise} \end{cases} \quad \text{and} \quad [\mathcal{P}_{\cdot, (l-n)}(\mathbf{B})]_{ij} = \begin{cases} B_{ij}, & \text{if } j = l-n, \\ 0, & \text{otherwise,} \end{cases}$$

for any matrix  $\mathbf{B} \in \mathbb{R}^{n \times n}$ . The key ideas are: (1) the iterates are not perturbed by much when one drops a small number of samples (and hence  $\mathbf{F}^t$  and  $\mathbf{F}^{t,(l)}$  remain sufficiently close); (2) the auxiliary iterates  $\mathbf{F}^{t,(l)}$  are independent of the samples directly related to the  $l$ th row of  $\mathbf{M}$ , which in turn allows to exploit certain statistical independence to control the  $l$ th row of  $\mathbf{F}^{t,(l)}$  (and hence  $\mathbf{F}^t$ ). See [MWCC17, Section 5] for a detailed explanation. Last but not least, the step size is set to be  $\eta$ , and we take  $\mathbf{F}^{0,(l)} = \mathbf{F}^*$  for all  $1 \leq l \leq 2n$  (the same initialization as in Algorithm 1).

With the help of the leave-one-out sequences, we are ready to establish Lemma 5 in an inductive manner. Concretely we aim at proving that

$$\|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*\|_{\text{F}} \leq C_{\text{F}} \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{X}^*\|_{\text{F}}, \quad (93a)$$

$$\|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*\| \leq C_{\text{op}} \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{X}^*\|, \quad (93b)$$

$$\max_{1 \leq l \leq 2n} \|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)}\|_{\text{F}} \leq C_3 \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{F}^*\|_{2, \infty}, \quad (93c)$$

$$\max_{1 \leq l \leq 2n} \|(\mathbf{F}^{t,(l)} \mathbf{H}^{t,(l)} - \mathbf{F}^*)_{l,\cdot}\|_2 \leq C_4 \kappa \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{F}^*\|_{2, \infty}, \quad (93d)$$

$$\|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*\|_{2,\infty} \leq C_\infty \kappa \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{F}^*\|_{2,\infty}, \quad (93e)$$

$$\|\mathbf{X}^{t\top} \mathbf{X}^t - \mathbf{Y}^{t\top} \mathbf{Y}^t\|_{\mathbb{F}} \leq C_B \kappa \eta \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \sqrt{r} \sigma_{\max}^2 \quad (93f)$$

hold for all  $0 \leq t \leq t_0 = n^{18}$  and for some constants  $C_{\mathbb{F}}, C_{\text{op}}, C_3, C_4, C_\infty, C_B > 0$ , provided that  $\eta \asymp 1/(n\kappa^3\sigma_{\max})$ . In addition, we also intend to establish that

$$f(\mathbf{X}^t, \mathbf{Y}^t) \leq f(\mathbf{X}^{t-1}, \mathbf{Y}^{t-1}) - \frac{\eta}{2} \|\nabla f(\mathbf{X}^{t-1}, \mathbf{Y}^{t-1})\|_{\mathbb{F}}^2 \quad (94)$$

holds for all  $1 \leq t \leq t_0 = n^{18}$ . Here,  $\mathbf{H}^{t,(l)}$  and  $\mathbf{R}^{t,(l)}$  are rotation matrices defined as

$$\mathbf{H}^{t,(l)} \triangleq \arg \min_{\mathbf{R} \in \mathcal{O}^{r \times r}} \|\mathbf{F}^{t,(l)} \mathbf{R} - \mathbf{F}^*\|_{\mathbb{F}}; \quad (95a)$$

$$\mathbf{R}^{t,(l)} \triangleq \arg \min_{\mathbf{R} \in \mathcal{O}^{r \times r}} \|\mathbf{F}^{t,(l)} \mathbf{R} - \mathbf{F}^t \mathbf{H}^t\|_{\mathbb{F}}. \quad (95b)$$

Note that the induction hypotheses (93a), (93b) and (93e) readily imply the statements (29a), (29b) and (29c) in Lemma 5, respectively, whereas the last bound on the size of the gradient (30) follows from (94). We summarize the last connection in the following lemma, whose proof is in Appendix D.2.

**Lemma 9 (Small gradient (30)).** *Set  $\lambda = C_\lambda \sigma \sqrt{np}$  for some large constant  $C_\lambda > 0$ . Suppose that the sample size obeys  $n^2 p \gg \kappa \mu r n \log^2 n$  and that the noise satisfies  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll \frac{1}{\sqrt{\kappa^4 \mu r \log n}}$ . If the induction hypotheses (93) hold for all  $0 \leq t \leq t_0$  and that (94) holds for all  $1 \leq t \leq t_0$ , then*

$$\min_{0 \leq t < t_0} \|\nabla f(\mathbf{X}^t, \mathbf{Y}^t)\|_{\mathbb{F}} \leq \frac{1}{n^5} \frac{\lambda}{p} \sqrt{\sigma_{\min}},$$

as long as  $\eta \asymp 1/(n\kappa^3\sigma_{\max})$ .

The rest of this section is devoted to proving the hypotheses (93) and (94) via induction. We start with the base case, i.e.  $t = 0$ . All the induction hypotheses (93) are easily verified by noting that

$$\mathbf{F}^0 = \mathbf{F}^{0,(l)} = \mathbf{F}^*, \quad \text{for all } 1 \leq l \leq 2n.$$

We now proceed to the induction step, which are demonstrated via the following lemmas. All the proofs are in subsequent subsections.

**Lemma 10 (Frobenius norm error (93a)).** *Set  $\lambda = C_\lambda \sigma \sqrt{np}$  for some large constant  $C_\lambda > 0$ . Suppose that the sample size obeys  $n^2 p \gg \kappa \mu r n \log^2 n$  and the noise satisfies  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll \frac{1}{\sqrt{\kappa^4 \mu r \log n}}$ . If the iterates satisfy (93) at the  $t$ th iteration, then with probability at least  $1 - O(n^{-100})$ ,*

$$\|\mathbf{F}^{t+1} \mathbf{H}^{t+1} - \mathbf{F}^*\|_{\mathbb{F}} \leq C_{\mathbb{F}} \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{X}^*\|_{\mathbb{F}},$$

holds as long as  $0 < \eta \ll 1/(\kappa^{5/2}\sigma_{\max})$  and  $C_{\mathbb{F}} > 0$  is large enough.

**Lemma 11 (Spectral norm error (93b)).** *Set  $\lambda = C_\lambda \sigma \sqrt{np}$  for some large constant  $C_\lambda > 0$ . Suppose the sample size obeys  $n^2 p \gg \kappa^4 \mu^2 r^2 n \log^2 n$  and the noise satisfies  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll \frac{1}{\sqrt{\kappa^4 \log n}}$ . If the iterates satisfy (93) at the  $t$ th iteration, then with probability at least  $1 - O(n^{-100})$ ,*

$$\|\mathbf{F}^{t+1} \mathbf{H}^{t+1} - \mathbf{F}^*\| \leq C_{\text{op}} \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{X}^*\|$$

holds with the proviso that  $0 < \eta \ll 1/(\kappa^3\sigma_{\max}\sqrt{r})$  and that  $C_{\text{op}} \gg 1$ .

**Lemma 12 (Leave-one-out perturbation (93c)).** Set  $\lambda = C_\lambda \sigma \sqrt{np}$  for some large constant  $C_\lambda > 0$ . Suppose that the sample size satisfies  $n^2 p \gg \kappa^4 \mu^2 r^2 n \log^3 n$  and that the noise satisfies  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll \frac{1}{\sqrt{\kappa^4 \mu r \log n}}$ . If the iterates satisfy (93) at the  $t$ th iteration, then with probability at least  $1 - O(n^{-99})$ ,

$$\max_{1 \leq l \leq 2n} \|\mathbf{F}^{t+1} \mathbf{H}^{t+1} - \mathbf{F}^{t+1, (l)} \mathbf{R}^{t+1, (l)}\|_{\text{F}} \leq C_3 \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{F}^*\|_{2, \infty}$$

holds, provided that  $0 < \eta \ll 1/(\kappa^2 \sigma_{\max} n)$  and that  $C_3 > 0$  is some sufficiently large constant.

**Lemma 13 ( $\ell_2/\ell_\infty$  norm error of leave-one-out sequences (93d)).** Set  $\lambda = C_\lambda \sigma \sqrt{np}$  for some large constant  $C_\lambda > 0$ . Suppose that the sample size obeys  $n^2 p \gg \kappa^2 \mu^2 r^2 n \log^3 n$  and that the noise satisfies  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll \frac{1}{\sqrt{\kappa^2 \log n}}$ . If the iterates satisfy (93) at the  $t$ th iteration, then with probability at least  $1 - O(n^{-99})$ ,

$$\max_{1 \leq l \leq 2n} \|(\mathbf{F}^{t+1, (l)} \mathbf{H}^{t+1, (l)} - \mathbf{F}^*)_{l, \cdot}\|_2 \leq C_4 \kappa \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{F}^*\|_{2, \infty}$$

holds, provided that  $0 < \eta \ll 1/(\kappa^2 \sqrt{r} \sigma_{\max})$ ,  $C_{\text{op}} \gg 1$  and  $C_4 \gg C_{\text{op}}$ .

**Lemma 14 ( $\ell_2/\ell_\infty$  norm error (93e)).** Set  $\lambda = C_\lambda \sigma \sqrt{np}$  for some large constant  $C_\lambda > 0$ . Suppose that  $n \geq \mu r$  and that the noise satisfies  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll \frac{1}{\sqrt{\kappa^2 \log n}}$ . If the iterates satisfy (93) at the  $t$ th iteration, then with probability at least  $1 - O(n^{-99})$ ,

$$\|\mathbf{F}^{t+1} \mathbf{H}^{t+1} - \mathbf{F}^*\|_{2, \infty} \leq C_\infty \kappa \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{F}^*\|_{2, \infty},$$

holds provided that  $C_\infty \geq 5C_3 + C_4$ .

**Lemma 15 (Approximate balancedness (93f)).** Set  $\lambda = C_\lambda \sigma \sqrt{np}$  for some large constant  $C_\lambda > 0$ . Suppose that the sample size satisfies  $n^2 p \gg \kappa^2 \mu^2 r^2 n \log n$  and that the noise satisfies  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll \frac{1}{\sqrt{\kappa^2 \log n}}$ . If the iterates satisfy (93) at the  $t$ th iteration, then with probability at least  $1 - O(n^{-100})$ ,

$$\begin{aligned} \|\mathbf{X}^{t+1 \top} \mathbf{X}^{t+1} - \mathbf{Y}^{t+1 \top} \mathbf{Y}^{t+1}\|_{\text{F}} &\leq C_{\text{B}} \kappa \eta \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \sqrt{r} \sigma_{\max}^2, \\ \max_{1 \leq l \leq 2n} \|\mathbf{X}^{t+1, (l) \top} \mathbf{X}^{t+1, (l)} - \mathbf{Y}^{t+1, (l) \top} \mathbf{Y}^{t+1, (l)}\|_{\text{F}} &\leq C_{\text{B}} \kappa \eta \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \sqrt{r} \sigma_{\max}^2, \end{aligned}$$

holds for some sufficiently large constant  $C_{\text{B}} \gg C_{\text{op}}^2$ , provided that  $0 < \eta < 1/\sigma_{\min}$ .

**Lemma 16 (Decreasing of function values (94)).** Set  $\lambda = C_\lambda \sigma \sqrt{np}$  for some large constant  $C_\lambda > 0$ . Suppose that the noise satisfies  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll 1/\sqrt{r}$ . If the iterates satisfy (93) at the  $t$ th iteration, then with probability at least  $1 - O(n^{-99})$ ,

$$f(\mathbf{X}^{t+1}, \mathbf{Y}^{t+1}) \leq f(\mathbf{X}^t, \mathbf{Y}^t) - \frac{\eta}{2} \|\nabla f(\mathbf{X}^t, \mathbf{Y}^t)\|_{\text{F}}^2,$$

as long as  $\eta \ll 1/(\kappa n \sigma_{\max})$ .

## D.1 Preliminaries and notations

Before proceeding to the proofs, we collect a few useful facts and notations. To begin with, for any matrix  $\mathbf{A}$ , we denote by  $\mathbf{A}_l$ . (resp.  $\mathbf{A}_{\cdot, l}$ ) the  $l$ th row (reps. column) of  $\mathbf{A}$ .

Define an augmented loss function  $f_{\text{aug}}(\mathbf{X}, \mathbf{Y})$  to be

$$f_{\text{aug}}(\mathbf{X}, \mathbf{Y}) \triangleq \frac{1}{2p} \|\mathcal{P}_\Omega(\mathbf{X}\mathbf{Y}^\top - \mathbf{M})\|_F^2 + \frac{\lambda}{2p} \|\mathbf{X}\|_F^2 + \frac{\lambda}{2p} \|\mathbf{Y}\|_F^2 + \frac{1}{8} \|\mathbf{X}^\top \mathbf{X} - \mathbf{Y}^\top \mathbf{Y}\|_F^2. \quad (96)$$

As the name suggests, this new function augments the original loss function (cf. (17)) with an additional term  $\|\mathbf{X}^\top \mathbf{X} - \mathbf{Y}^\top \mathbf{Y}\|_F^2/8$ , which is commonly used in the literature of asymmetric low-rank matrix factorization to balance the scale of  $\mathbf{X}$  and  $\mathbf{Y}$  [TBS<sup>+</sup>16, YPCC16, CLL19]. We emphasize that, in contrast to aforementioned works, here our gradient descent algorithm (cf. Algorithm 1) operates on  $f(\cdot, \cdot)$  instead of  $f_{\text{aug}}(\cdot, \cdot)$ . The introduction of  $f_{\text{aug}}(\cdot, \cdot)$  is mainly to simplify the proof.

It is easily seen that the gradients of  $f_{\text{aug}}(\cdot, \cdot)$  are given by

$$\nabla_{\mathbf{X}} f_{\text{aug}}(\mathbf{X}, \mathbf{Y}) = \frac{1}{p} \mathcal{P}_\Omega(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}) \mathbf{Y} + \frac{\lambda}{p} \mathbf{X} + \frac{1}{2} \mathbf{X} (\mathbf{X}^\top \mathbf{X} - \mathbf{Y}^\top \mathbf{Y}); \quad (97a)$$

$$\nabla_{\mathbf{Y}} f_{\text{aug}}(\mathbf{X}, \mathbf{Y}) = \frac{1}{p} \mathcal{P}_\Omega(\mathbf{X}\mathbf{Y}^\top - \mathbf{M})^\top \mathbf{X} + \frac{\lambda}{p} \mathbf{Y} + \frac{1}{2} \mathbf{Y} (\mathbf{Y}^\top \mathbf{Y} - \mathbf{X}^\top \mathbf{X}). \quad (97b)$$

Correspondingly, define the difference between gradients of  $\nabla f(\mathbf{X}, \mathbf{Y})$  and  $\nabla f_{\text{aug}}(\mathbf{X}, \mathbf{Y})$  as follows

$$\nabla_{\mathbf{X}} f_{\text{diff}}(\mathbf{X}, \mathbf{Y}) = -\mathbf{X} (\mathbf{X}^\top \mathbf{X} - \mathbf{Y}^\top \mathbf{Y}) / 2; \quad (98a)$$

$$\nabla_{\mathbf{Y}} f_{\text{diff}}(\mathbf{X}, \mathbf{Y}) = -\mathbf{Y} (\mathbf{Y}^\top \mathbf{Y} - \mathbf{X}^\top \mathbf{X}) / 2, \quad (98b)$$

such that

$$\nabla_{\mathbf{X}} f(\mathbf{X}, \mathbf{Y}) = \nabla_{\mathbf{X}} f_{\text{aug}}(\mathbf{X}, \mathbf{Y}) + \nabla_{\mathbf{X}} f_{\text{diff}}(\mathbf{X}, \mathbf{Y}); \quad (99a)$$

$$\nabla_{\mathbf{Y}} f(\mathbf{X}, \mathbf{Y}) = \nabla_{\mathbf{Y}} f_{\text{aug}}(\mathbf{X}, \mathbf{Y}) + \nabla_{\mathbf{Y}} f_{\text{diff}}(\mathbf{X}, \mathbf{Y}). \quad (99b)$$

Regarding  $\mathbf{F}^*$ , simple algebra reveals that

$$\sigma_1(\mathbf{F}^*) = \|\mathbf{F}^*\| = \sqrt{2\sigma_{\max}}, \quad \sigma_r(\mathbf{F}^*) = \sqrt{2\sigma_{\min}}, \quad (100a)$$

$$\|\mathbf{F}^*\|_{2,\infty} = \max \{ \|\mathbf{X}^*\|_{2,\infty}, \|\mathbf{Y}^*\|_{2,\infty} \} \leq \sqrt{\mu r \sigma_{\max} / n}, \quad (100b)$$

where the last one follows from the incoherence assumption (34).

We start with a lemma that characterizes the local geometry of the nonconvex loss function, whose proof is given in Appendix D.10.

**Lemma 17.** *Set  $\lambda = C_\lambda \sigma \sqrt{np}$  for some constant  $C_\lambda > 0$ . Suppose that the sample size obeys  $n^2 p \geq C \kappa \mu r n \log^2 n$  for some sufficiently large constant  $C > 0$  and that the noise satisfies  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll 1$ . Recall the function  $f_{\text{aug}}(\cdot, \cdot)$  defined in (96). Then with probability at least  $1 - O(n^{-10})$ ,*

$$\begin{aligned} \text{vec}(\Delta)^\top \nabla^2 f_{\text{aug}}(\mathbf{X}, \mathbf{Y}) \text{vec}(\Delta) &\geq \frac{1}{10} \sigma_{\min} \|\Delta\|_F^2, \\ \max \{ \|\nabla^2 f_{\text{aug}}(\mathbf{X}, \mathbf{Y})\|, \|\nabla^2 f(\mathbf{X}, \mathbf{Y})\| \} &\leq 10\sigma_{\max} \end{aligned}$$

hold uniformly over all  $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times r}$  obeying

$$\left\| \begin{bmatrix} \mathbf{X} - \mathbf{X}^* \\ \mathbf{Y} - \mathbf{Y}^* \end{bmatrix} \right\|_{2,\infty} \leq \frac{1}{1000\kappa\sqrt{n}} \|\mathbf{X}^*\|$$

and all  $\Delta = \begin{bmatrix} \Delta_{\mathbf{X}} \\ \Delta_{\mathbf{Y}} \end{bmatrix} \in \mathbb{R}^{2n \times r}$  lying in the set

$$\left\{ \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{Y}_1 \end{bmatrix} \hat{H} - \begin{bmatrix} \mathbf{X}_2 \\ \mathbf{Y}_2 \end{bmatrix} \right\| \left\| \begin{bmatrix} \mathbf{X}_2 - \mathbf{X}^* \\ \mathbf{Y}_2 - \mathbf{Y}^* \end{bmatrix} \right\| \leq \frac{1}{500\kappa} \|\mathbf{X}^*\|, \hat{H} \triangleq \arg \min_{\mathbf{R} \in \mathcal{O}^{r \times r}} \left\| \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{Y}_1 \end{bmatrix} \mathbf{R} - \begin{bmatrix} \mathbf{X}_2 \\ \mathbf{Y}_2 \end{bmatrix} \right\|_F \right\}.$$

Last but not least, a few immediate consequences of (93) are gathered in the following lemma, whose proof is given in Appendix D.11.



**Lemma 18.** *We have the following four sets of consequences of the induction hypotheses (29).*

1. *Suppose that the sample size obeys  $n \gg \mu r \log n$ . If the  $t$ th iterates obey (93), then one has*

$$\left\| \mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)} - \mathbf{F}^* \right\|_{2,\infty} \leq (C_\infty \kappa + C_3) \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{F}^*\|_{2,\infty}, \quad (101a)$$

$$\left\| \mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)} - \mathbf{F}^* \right\| \leq 2C_{\text{op}} \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{X}^*\|. \quad (101b)$$

2. *Suppose that the noise satisfies  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll \frac{1}{\sqrt{\kappa^2 \log n}}$ . If the  $t$ th iterates obey (93), then one has*

$$\|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*\| \leq \|\mathbf{X}^*\|, \quad \|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*\|_{\text{F}} \leq \|\mathbf{X}^*\|_{\text{F}}, \quad \|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*\|_{2,\infty} \leq \|\mathbf{F}^*\|_{2,\infty}, \quad (102a)$$

$$\|\mathbf{F}^t\| \leq 2\|\mathbf{X}^*\|, \quad \|\mathbf{F}^t\|_{\text{F}} \leq 2\|\mathbf{X}^*\|_{\text{F}}, \quad \|\mathbf{F}^t\|_{2,\infty} \leq 2\|\mathbf{F}^*\|_{2,\infty}. \quad (102b)$$

3. *Suppose that  $n \gg \kappa^2 \mu r \log n$  and that  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll \frac{1}{\sqrt{\kappa^2 \log n}}$ . If the  $t$ th iterates obey (93), then we have*

$$\|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^{t,(l)} \mathbf{H}^{t,(l)}\|_{\text{F}} \leq 5\kappa \|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)}\|_{\text{F}}.$$

4. *Suppose that  $n \geq \kappa \mu$  and that  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll \frac{1}{\sqrt{\kappa^2 \log n}}$ . If the  $t$ th iterates obey (93), then (102) also holds for  $\mathbf{F}^{t,(l)} \mathbf{H}^{t,(l)}$ . In addition, one has*

$$\sigma_{\min}/2 \leq \sigma_{\min} \left( (\mathbf{Y}^{t,(l)} \mathbf{H}^{t,(l)})^\top \mathbf{Y}^{t,(l)} \mathbf{H}^{t,(l)} \right) \leq \sigma_{\max} \left( (\mathbf{Y}^{t,(l)} \mathbf{H}^{t,(l)})^\top \mathbf{Y}^{t,(l)} \mathbf{H}^{t,(l)} \right) \leq 2\sigma_{\max}.$$

## D.2 Proof of Lemma 9

Summing (94) from  $t = 1$  to  $t = t_0$  leads to a telescopic sum

$$f(\mathbf{X}^{t_0}, \mathbf{Y}^{t_0}) \leq f(\mathbf{X}^0, \mathbf{Y}^0) - \frac{\eta}{2} \sum_{t=0}^{t_0-1} \|\nabla f(\mathbf{X}^t, \mathbf{Y}^t)\|_{\text{F}}^2.$$

This further implies that

$$\min_{0 \leq t < t_0} \|\nabla f(\mathbf{X}^t, \mathbf{Y}^t)\|_{\text{F}} \leq \left\{ \frac{1}{t_0} \sum_{t=0}^{t_0-1} \|\nabla f(\mathbf{X}^t, \mathbf{Y}^t)\|_{\text{F}}^2 \right\}^{1/2} \leq \left\{ \frac{2}{\eta t_0} [f(\mathbf{X}^*, \mathbf{Y}^*) - f(\mathbf{X}^{t_0}, \mathbf{Y}^{t_0})] \right\}^{1/2}, \quad (103)$$

where we have used the assumption that  $(\mathbf{X}^0, \mathbf{Y}^0) = (\mathbf{X}^*, \mathbf{Y}^*)$ .

It remains to control  $f(\mathbf{X}^*, \mathbf{Y}^*) - f(\mathbf{X}^{t_0}, \mathbf{Y}^{t_0})$ . Towards this end, we can use the fact that  $f(\mathbf{X}, \mathbf{Y}) = f(\mathbf{X}\mathbf{R}, \mathbf{Y}\mathbf{R})$  for any  $\mathbf{R} \in \mathcal{O}^{r \times r}$  to obtain

$$f(\mathbf{F}^{t_0}) = f(\mathbf{F}^{t_0} \mathbf{H}^{t_0}) = f(\mathbf{F}^*) + \langle \nabla f(\mathbf{F}^*), \mathbf{F}^{t_0} \mathbf{H}^{t_0} - \mathbf{F}^* \rangle + \frac{1}{2} \text{vec}(\mathbf{F}^{t_0} \mathbf{H}^{t_0} - \mathbf{F}^*)^\top \nabla^2 f(\tilde{\mathbf{F}}) \text{vec}(\mathbf{F}^{t_0} \mathbf{H}^{t_0} - \mathbf{F}^*),$$

where  $\tilde{\mathbf{F}}$  lies in the line segment connecting  $\mathbf{F}^{t_0} \mathbf{H}^{t_0}$  and  $\mathbf{F}^*$ . Apply the triangle inequality to see

$$\begin{aligned} f(\mathbf{F}^*) - f(\mathbf{F}^{t_0}) &\leq \|\nabla f(\mathbf{F}^*)\|_{\text{F}} \|\mathbf{F}^{t_0} \mathbf{H}^{t_0} - \mathbf{F}^*\|_{\text{F}} - \frac{1}{2} \text{vec}(\mathbf{F}^{t_0} \mathbf{H}^{t_0} - \mathbf{F}^*)^\top \nabla^2 f(\tilde{\mathbf{F}}) \text{vec}(\mathbf{F}^{t_0} \mathbf{H}^{t_0} - \mathbf{F}^*) \\ &\leq \|\nabla f(\mathbf{F}^*)\|_{\text{F}} \|\mathbf{F}^{t_0} \mathbf{H}^{t_0} - \mathbf{F}^*\|_{\text{F}} + 5\sigma_{\max} \|\mathbf{F}^{t_0} \mathbf{H}^{t_0} - \mathbf{F}^*\|_{\text{F}}^2. \end{aligned}$$

Here the second line follows from the fact that  $\|\nabla^2 f(\tilde{\mathbf{F}})\| \leq 10\sigma_{\max}$ . To see this, use (93e) to obtain that

$$\|\tilde{\mathbf{F}} - \mathbf{F}^*\|_{2,\infty} \leq \|\mathbf{F}^{t_0} \mathbf{H}^{t_0} - \mathbf{F}^*\|_{2,\infty} \leq C_\infty \kappa \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{F}^*\|_{2,\infty}$$

$$\begin{aligned}
&\leq C_\infty \kappa \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \sqrt{\frac{\mu r}{n}} \sqrt{\sigma_{\max}} \\
&\leq \frac{1}{2000 \kappa \sqrt{n}} \sqrt{\sigma_{\max}}, \tag{104}
\end{aligned}$$

where the second line arises from the incoherence assumption (100b) and the last inequality holds as long as  $\lambda \asymp \sigma \sqrt{np}$  and  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll \frac{1}{\sqrt{\kappa^4 \mu r \log n}}$ . Apply Lemma 17 to conclude that  $\|\nabla^2 f(\tilde{\mathbf{F}})\| \leq 10\sigma_{\max}$ . Recognize that

$$\begin{aligned}
\|\nabla f(\mathbf{F}^*)\|_{\mathbb{F}} &\leq \|\nabla_{\mathbf{X}} f(\mathbf{F}^*)\|_{\mathbb{F}} + \|\nabla_{\mathbf{Y}} f(\mathbf{F}^*)\|_{\mathbb{F}} \\
&\leq \frac{1}{p} \|\mathcal{P}_\Omega(\mathbf{E}) \mathbf{Y}^*\|_{\mathbb{F}} + \frac{\lambda}{p} \|\mathbf{X}^*\|_{\mathbb{F}} + \frac{1}{p} \|\mathcal{P}_\Omega(\mathbf{E})^\top \mathbf{X}^*\|_{\mathbb{F}} + \frac{\lambda}{p} \|\mathbf{Y}^*\|_{\mathbb{F}} \\
&\leq \left( \frac{1}{p} \|\mathcal{P}_\Omega(\mathbf{E})\| + \frac{\lambda}{p} \right) (\|\mathbf{X}^*\|_{\mathbb{F}} + \|\mathbf{Y}^*\|_{\mathbb{F}}), \tag{105}
\end{aligned}$$

where we have used the fact that  $\nabla_{\mathbf{X}} f(\mathbf{F}^*) = \frac{1}{p} \mathcal{P}_\Omega(\mathbf{X}^* \mathbf{Y}^{*\top} - \mathbf{M}^* - \mathbf{E}) \mathbf{Y}^* + \frac{\lambda}{p} \mathbf{X}^* = -\frac{1}{p} \mathcal{P}_\Omega(\mathbf{E}) \mathbf{Y}^* + \frac{\lambda}{p} \mathbf{X}^*$  (similar expression holds true for  $\nabla_{\mathbf{Y}} f(\mathbf{F}^*)$ ). This together with Lemma 3 and the assumption that  $\lambda \asymp \sigma \sqrt{np}$  yields

$$\|\nabla f(\mathbf{F}^*)\|_{\mathbb{F}} \lesssim \left( \sigma \sqrt{\frac{n}{p}} + \frac{\lambda}{p} \right) \sqrt{r \sigma_{\max}} \asymp \frac{\lambda}{p} \sqrt{r \sigma_{\max}}. \tag{106}$$

The above bounds together with the induction hypothesis (93a) for  $t = t_0$  give

$$\begin{aligned}
f(\mathbf{F}^*) - f(\mathbf{F}^{t_0}) &\lesssim \frac{\lambda}{p} \sqrt{r \sigma_{\max}} \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{X}^*\|_{\mathbb{F}} + \sigma_{\max} \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right)^2 \|\mathbf{X}^*\|_{\mathbb{F}}^2 \\
&\lesssim r \kappa^2 \left( \frac{\lambda}{p} \right)^2,
\end{aligned}$$

where the last relation arises from  $\sigma \sqrt{np} \asymp \lambda$ . Substitution into (103) results in

$$\min_{0 \leq t < t_0} \|\nabla f(\mathbf{X}^t, \mathbf{Y}^t)\|_{\mathbb{F}} \lesssim \sqrt{\frac{1}{\eta t_0} r \kappa^2 \left( \frac{\lambda}{p} \right)^2} \leq \frac{1}{n^5} \frac{\lambda}{p} \sqrt{\sigma_{\min}},$$

provided that  $\eta \asymp 1/(n \kappa^3 \sigma_{\max})$ ,  $t_0 = n^{18}$  and that  $n \geq \kappa$ , which is a consequence of our sample complexity  $n \geq np \gg \kappa \mu r \log^2 n$ .

### D.3 Proof of Lemma 10

From the definitions of  $\mathbf{H}^{t+1}$  (cf. (89)),  $\nabla f_{\text{aug}}$  (cf. (97)) and  $\nabla f_{\text{diff}}$  (cf. (98)), we have

$$\begin{aligned}
\|\mathbf{F}^{t+1} \mathbf{H}^{t+1} - \mathbf{F}^*\|_{\mathbb{F}} &\leq \|\mathbf{F}^{t+1} \mathbf{H}^t - \mathbf{F}^*\|_{\mathbb{F}} = \|\mathbf{F}^t - \eta \nabla f(\mathbf{F}^t)\|_{\mathbb{F}} \|\mathbf{H}^t - \mathbf{F}^*\|_{\mathbb{F}} \\
&\stackrel{(i)}{=} \|\mathbf{F}^t \mathbf{H}^t - \eta \nabla f(\mathbf{F}^t \mathbf{H}^t) - \mathbf{F}^*\|_{\mathbb{F}} \\
&\stackrel{(ii)}{\leq} \underbrace{\|\mathbf{F}^t \mathbf{H}^t - \eta \nabla f_{\text{aug}}(\mathbf{F}^t \mathbf{H}^t) - [\mathbf{F}^* - \eta \nabla f_{\text{aug}}(\mathbf{F}^*)]\|_{\mathbb{F}}}_{:=\alpha_1} + \underbrace{\eta \|\nabla f_{\text{diff}}(\mathbf{F}^t \mathbf{H}^t)\|_{\mathbb{F}}}_{:=\alpha_2} + \underbrace{\eta \|\nabla f_{\text{aug}}(\mathbf{F}^*)\|_{\mathbb{F}}}_{:=\alpha_3}.
\end{aligned}$$

Here (i) uses the fact that  $\nabla f(\mathbf{F}\mathbf{R}) = \nabla f(\mathbf{F})\mathbf{R}$  for all  $\mathbf{R} \in \mathcal{O}^{r \times r}$ ; the last relation (ii) uses the decomposition (99) and the triangle inequality. In the following, we bound  $\alpha_1$ ,  $\alpha_2$  and  $\alpha_3$  in the reverse order.

1. First, regarding  $\alpha_3$ , since  $\mathbf{X}^{*\top} \mathbf{X}^* = \mathbf{Y}^{*\top} \mathbf{Y}^*$ , one has  $\eta \|\nabla f(\mathbf{F}^*)\|_{\mathbb{F}} = \eta \|\nabla f_{\text{aug}}(\mathbf{F}^*)\|_{\mathbb{F}}$ . Repeating our arguments for (105) and (106) gives

$$\alpha_3 = \eta \|\nabla f(\mathbf{F}^*)\|_{\mathbb{F}} \leq 4\eta \frac{\lambda}{p} \|\mathbf{X}^*\|_{\mathbb{F}}$$

as long as  $\lambda \asymp \sigma \sqrt{np}$ . Here the last inequality also relies on the fact that  $\|\mathbf{X}^*\|_{\mathbb{F}} = \|\mathbf{Y}^*\|_{\mathbb{F}}$ .

2. We now move on to  $\alpha_2$ , for which one has

$$\begin{aligned}\alpha_2 &\leq \frac{\eta}{2} (\|\mathbf{X}^t (\mathbf{X}^{t\top} \mathbf{X}^t - \mathbf{Y}^{t\top} \mathbf{Y}^t) \mathbf{H}^t\|_{\mathbb{F}} + \|\mathbf{Y}^t (\mathbf{Y}^{t\top} \mathbf{Y}^t - \mathbf{X}^{t\top} \mathbf{X}^t) \mathbf{H}^t\|_{\mathbb{F}}) \\ &\leq \frac{\eta}{2} (\|\mathbf{X}^t\| + \|\mathbf{Y}^t\|) \|\mathbf{X}^{t\top} \mathbf{X}^t - \mathbf{Y}^{t\top} \mathbf{Y}^t\|_{\mathbb{F}}.\end{aligned}$$

Utilize the fact that  $\max\{\|\mathbf{X}^t\|, \|\mathbf{Y}^t\|\} \leq \|\mathbf{F}^*\| \leq 2\|\mathbf{X}^*\|$  (see Lemma 18) and the induction hypothesis (93f) to obtain

$$\begin{aligned}\alpha_2 &\leq 2\eta\sqrt{\sigma_{\max}} \cdot C_{\text{B}}\kappa\eta \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}} \right) \sqrt{r}\sigma_{\max}^2 \\ &\leq (2C_{\text{B}}\kappa^{5/2}\eta\sigma_{\max})\sigma_{\min}\eta \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}} \right) \|\mathbf{X}^*\|_{\mathbb{F}} \\ &\leq \sigma_{\min}\eta \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}} \right) \|\mathbf{X}^*\|_{\mathbb{F}},\end{aligned}$$

where the second inequality uses  $\|\mathbf{X}^*\|_{\mathbb{F}} \geq \sqrt{r}\sigma_{\min}$  and the last one holds as long as  $2C_{\text{B}}\kappa^{5/2}\sigma_{\max}\eta \leq 1$ .

3. In the end, for  $\alpha_1$ , the fundamental theorem of calculus [Lan93, Chapter XIII, Theorem 4.2] reveals that

$$\begin{aligned}\text{vec} [\mathbf{F}^t \mathbf{H}^t - \eta \nabla f_{\text{aug}}(\mathbf{F}^t \mathbf{H}^t) - [\mathbf{F}^* - \eta \nabla f_{\text{aug}}(\mathbf{F}^*)]] \\ &= \text{vec} [\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*] - \eta \cdot \text{vec} [\nabla f_{\text{aug}}(\mathbf{F}^t \mathbf{H}^t) - \nabla f_{\text{aug}}(\mathbf{F}^*)] \\ &= \left( \mathbf{I}_{2nr} - \eta \underbrace{\int_0^1 \nabla^2 f_{\text{aug}}(\mathbf{F}(\tau)) \text{d}\tau}_{:=\mathbf{A}} \right) \text{vec}(\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*),\end{aligned}\tag{107}$$

where we denote  $\mathbf{F}(\tau) \triangleq \mathbf{F}^* + \tau(\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*)$  for all  $0 \leq \tau \leq 1$ . Taking the squared Euclidean norm of both sides of the equality (107) leads to

$$\begin{aligned}\alpha_1^2 &= \text{vec}(\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*)^\top (\mathbf{I}_{2nr} - \eta \mathbf{A})^2 \text{vec}(\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*) \\ &= \text{vec}(\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*)^\top (\mathbf{I}_{2nr} - 2\eta \mathbf{A} + \eta^2 \mathbf{A}^2) \text{vec}(\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*) \\ &\leq \|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*\|_{\mathbb{F}}^2 + \eta^2 \|\mathbf{A}\|^2 \|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*\|_{\mathbb{F}}^2 - 2\eta \text{vec}(\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*)^\top \mathbf{A} \text{vec}(\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*),\end{aligned}\tag{108}$$

where (108) results from the fact that

$$\text{vec}(\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*)^\top \mathbf{A} \text{vec}(\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*) \leq \|\mathbf{A}\|^2 \|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*\|_{\mathbb{F}}^2.$$

Applying the same argument as in (104), one gets for all  $0 \leq \tau \leq 1$ ,  $\|\mathbf{F}(\tau) - \mathbf{F}^*\|_{2,\infty} \leq \frac{1}{2000\kappa\sqrt{n}} \|\mathbf{X}^*\|$ . Invoke Lemma 17 with  $\mathbf{X} = \mathbf{X}^* + \tau(\mathbf{X}^t \mathbf{H}^t - \mathbf{X}^*)$ ,  $\mathbf{Y} = \mathbf{Y}^* + \tau(\mathbf{Y}^t \mathbf{H}^t - \mathbf{Y}^*)$ ,  $(\mathbf{X}_1, \mathbf{Y}_1) = (\mathbf{X}^t, \mathbf{Y}^t)$  and  $(\mathbf{X}_2, \mathbf{Y}_2) = (\mathbf{X}^*, \mathbf{Y}^*)$  to obtain  $\|\mathbf{A}\| \leq 10\sigma_{\max}$  and

$$\text{vec}(\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*)^\top \mathbf{A} \text{vec}(\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*) \geq \frac{1}{10} \sigma_{\min} \|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*\|_{\mathbb{F}}^2.$$

Putting these two bounds back to (108) yields

$$\alpha_1^2 \leq \left( 1 + 100\eta^2\sigma_{\max}^2 - \frac{1}{5}\eta\sigma_{\min} \right) \|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*\|_{\mathbb{F}}^2 \leq \left( 1 - \frac{\sigma_{\min}}{10}\eta \right) \|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*\|_{\mathbb{F}}^2.$$

Here the last relation holds as long as  $0 \leq \eta \leq 1/(1000\kappa\sigma_{\max})$ . As a result, we have

$$\alpha_1 \leq \left( 1 - \frac{\sigma_{\min}}{20}\eta \right) \|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*\|_{\mathbb{F}}.$$

Combine the above bounds on  $\alpha_1, \alpha_2$  and  $\alpha_3$  to conclude that

$$\begin{aligned} \|\mathbf{F}^{t+1}\mathbf{H}^{t+1} - \mathbf{F}^*\|_{\text{F}} &\leq \left(1 - \frac{\sigma_{\min}}{20}\eta\right) \|\mathbf{F}^t\mathbf{H}^t - \mathbf{F}^*\|_{\text{F}} + 4\eta\frac{\lambda}{p} \|\mathbf{X}^*\|_{\text{F}} + \eta\sigma_{\min} \left(\frac{\sigma}{\sigma_{\min}}\sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}}\right) \|\mathbf{X}^*\|_{\text{F}} \\ &\leq \left(1 - \frac{\sigma_{\min}}{20}\eta\right) C_{\text{F}} \left(\frac{\sigma}{\sigma_{\min}}\sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}}\right) \|\mathbf{X}^*\|_{\text{F}} + 4\eta\sigma_{\min}\frac{\lambda}{p\sigma_{\min}} \|\mathbf{X}^*\|_{\text{F}} + \eta\sigma_{\min} \left(\frac{\sigma}{\sigma_{\min}}\sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}}\right) \|\mathbf{X}^*\|_{\text{F}} \\ &\leq C_{\text{F}} \left(\frac{\sigma}{\sigma_{\min}}\sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}}\right) \|\mathbf{X}^*\|_{\text{F}}, \end{aligned}$$

provided that  $C_{\text{F}} > 0$  is large enough.

#### D.4 Proof of Lemma 11

To facilitate analysis, we define an auxiliary point  $\tilde{\mathbf{F}}^{t+1} \triangleq \begin{bmatrix} \tilde{\mathbf{X}}^{t+1} \\ \tilde{\mathbf{Y}}^{t+1} \end{bmatrix}$  as

$$\tilde{\mathbf{X}}^{t+1} = \mathbf{X}^t\mathbf{H}^t - \eta \left[ \frac{1}{p}\mathcal{P}_{\Omega}(\mathbf{X}^t\mathbf{Y}^{t\top} - \mathbf{M}^* - \mathbf{E})\mathbf{Y}^* + \frac{\lambda}{p}\mathbf{X}^* + \frac{1}{2}\mathbf{X}^*\mathbf{H}^{t\top}(\mathbf{X}^{t\top}\mathbf{X}^t - \mathbf{Y}^{t\top}\mathbf{Y}^t)\mathbf{H}^t \right]; \quad (109\text{a})$$

$$\tilde{\mathbf{Y}}^{t+1} = \mathbf{Y}^t\mathbf{H}^t - \eta \left[ \frac{1}{p}\mathcal{P}_{\Omega}(\mathbf{X}^t\mathbf{Y}^{t\top} - \mathbf{M}^* - \mathbf{E})^{\top}\mathbf{X}^* + \frac{\lambda}{p}\mathbf{Y}^* + \frac{1}{2}\mathbf{Y}^*\mathbf{H}^{t\top}(\mathbf{Y}^{t\top}\mathbf{Y}^t - \mathbf{X}^{t\top}\mathbf{X}^t)\mathbf{H}^t \right]. \quad (109\text{b})$$

Then the triangle inequality tells us that

$$\|\mathbf{F}^{t+1}\mathbf{H}^{t+1} - \mathbf{F}^*\| \leq \underbrace{\|\mathbf{F}^{t+1}\mathbf{H}^{t+1} - \tilde{\mathbf{F}}^{t+1}\|}_{:=\alpha_1} + \underbrace{\|\tilde{\mathbf{F}}^{t+1} - \mathbf{F}^*\|}_{:=\alpha_2}. \quad (110)$$

In what follows, we shall control  $\alpha_1$  and  $\alpha_2$  separately.

1. We start with  $\alpha_2$ . By the triangle inequality again we have

$$\begin{aligned} \alpha_2 &\leq \left\| \underbrace{\begin{bmatrix} \mathbf{X}^t\mathbf{H}^t - \eta \left[ (\mathbf{X}^t\mathbf{Y}^{t\top} - \mathbf{M}^*)\mathbf{Y}^* + \frac{1}{2}\mathbf{X}^*\mathbf{H}^{t\top}(\mathbf{X}^{t\top}\mathbf{X}^t - \mathbf{Y}^{t\top}\mathbf{Y}^t)\mathbf{H}^t \right] - \mathbf{X}^* \\ \mathbf{Y}^t\mathbf{H}^t - \eta \left[ (\mathbf{X}^t\mathbf{Y}^{t\top} - \mathbf{M}^*)^{\top}\mathbf{X}^* + \frac{1}{2}\mathbf{Y}^*\mathbf{H}^{t\top}(\mathbf{Y}^{t\top}\mathbf{Y}^t - \mathbf{X}^{t\top}\mathbf{X}^t)\mathbf{H}^t \right] - \mathbf{Y}^* \end{bmatrix}}_{:=\beta_1} \right\| \\ &\quad + \underbrace{\frac{\eta}{p} \left\| \begin{bmatrix} \mathcal{P}_{\Omega}(\mathbf{E})\mathbf{Y}^* \\ \mathcal{P}_{\Omega}(\mathbf{E})^{\top}\mathbf{X}^* \end{bmatrix} \right\|}_{:=\beta_2} + \eta\frac{\lambda}{p} \left\| \begin{bmatrix} \mathbf{X}^* \\ \mathbf{Y}^* \end{bmatrix} \right\| + \eta \left\| \underbrace{\begin{bmatrix} \frac{1}{p}\mathcal{P}_{\Omega}(\mathbf{X}^t\mathbf{Y}^{t\top} - \mathbf{M}^*)\mathbf{Y}^* - (\mathbf{X}^t\mathbf{Y}^{t\top} - \mathbf{M}^*)\mathbf{Y}^* \\ \frac{1}{p}[\mathcal{P}_{\Omega}(\mathbf{X}^t\mathbf{Y}^{t\top} - \mathbf{M}^*)]^{\top}\mathbf{X}^* - (\mathbf{X}^t\mathbf{Y}^{t\top} - \mathbf{M}^*)^{\top}\mathbf{X}^* \end{bmatrix}}_{:=\beta_3} \right\| \end{aligned}$$

Denote  $\Delta^t \triangleq \mathbf{F}^t\mathbf{H}^t - \mathbf{F}^* = \begin{bmatrix} \Delta_{\mathbf{X}}^t \\ \Delta_{\mathbf{Y}}^t \end{bmatrix}$ . The term  $\beta_1$  is the same as the term  $\alpha_2$  in [CLL19, Section 4.2]. Therefore we can adopt the bound therein to obtain

$$\beta_1 \leq (1 - \eta\sigma_{\min}) \|\Delta^t\| + 4\eta \|\Delta^t\|^2 \|\mathbf{X}^*\|.$$

Moving to  $\beta_2$ , one has

$$\beta_2 = \eta \left\| \begin{bmatrix} \frac{1}{p}\mathcal{P}_{\Omega}(\mathbf{E}) & \mathbf{0} \\ \mathbf{0} & \frac{1}{p}\mathcal{P}_{\Omega}(\mathbf{E})^{\top} \end{bmatrix} \begin{bmatrix} \mathbf{Y}^* \\ \mathbf{X}^* \end{bmatrix} \right\| + \eta\frac{\lambda}{p} \|\mathbf{F}^*\| \leq \frac{\eta}{p} \|\mathcal{P}_{\Omega}(\mathbf{E})\| \|\mathbf{F}^*\| + \eta\frac{\lambda}{p} \|\mathbf{F}^*\| \leq C\eta \left( \sigma\sqrt{\frac{n}{p}} + \frac{\lambda}{p} \right) \|\mathbf{X}^*\|$$

for some constant  $C > 0$ . Here the last inequality arises from Lemma 3 and the fact that  $\|\mathbf{F}^*\| = \sqrt{2}\|\mathbf{X}^*\|$  (cf. (100a)). We are now left with the term  $\beta_3$ , which is exactly the term  $\alpha_1$  in [CLL19, Section 4.2]. Reusing their results, we have

$$\begin{aligned} \beta_3 &\leq \frac{2\eta}{p} \|\mathbf{X}^*\| \|\mathcal{P}_{\Omega}(\mathbf{1}\mathbf{1}^{\top}) - p\mathbf{1}\mathbf{1}^{\top}\| (\|\Delta_{\mathbf{X}}^t\|_{2,\infty} \|\Delta_{\mathbf{Y}}^t\|_{2,\infty} + \|\Delta_{\mathbf{X}}^t\|_{2,\infty} \|\mathbf{Y}^*\|_{2,\infty} + \|\mathbf{X}^*\|_{2,\infty} \|\Delta_{\mathbf{Y}}^t\|_{2,\infty}) \\ &\lesssim \eta\sqrt{\frac{n}{p}} \|\Delta^t\|_{2,\infty} \|\mathbf{F}^*\|_{2,\infty} \|\mathbf{X}^*\|. \end{aligned}$$

The last line follows from the facts that  $\|\mathcal{P}_\Omega(\mathbf{1}\mathbf{1}^\top) - p\mathbf{1}\mathbf{1}^\top\| \lesssim \sqrt{np}$  (see [KMO10a, Lemma 3.2]) and that  $\max\{\|\Delta_{\mathbf{X}}^t\|_{2,\infty}, \|\Delta_{\mathbf{Y}}^t\|_{2,\infty}\} \leq \|\mathbf{F}^*\|_{2,\infty}$ , provided that  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll \frac{1}{\sqrt{\kappa^2 \log n}}$  (see Lemma 18). Combining the above three bounds gives

$$\begin{aligned} \alpha_2 &\leq (1 - \eta\sigma_{\min}) \|\Delta^t\| + 4\eta \|\Delta^t\|^2 \|\mathbf{X}^*\| + \tilde{C}\eta \left( \sigma\sqrt{\frac{n}{p}} + \frac{\lambda}{p} \right) \|\mathbf{X}^*\| + \tilde{C}\eta\sqrt{\frac{n}{p}} \|\Delta^t\|_{2,\infty} \|\mathbf{F}^*\|_{2,\infty} \|\mathbf{X}^*\| \\ &\leq \left(1 - \frac{\eta}{2}\sigma_{\min}\right) \|\Delta^t\| + \tilde{C}\eta \left( \sigma\sqrt{\frac{n}{p}} + \frac{\lambda}{p} \right) \|\mathbf{X}^*\| + \tilde{C}\eta\sqrt{\frac{n}{p}} \|\Delta^t\|_{2,\infty} \|\mathbf{F}^*\|_{2,\infty} \|\mathbf{X}^*\| \end{aligned} \quad (111)$$

for some sufficiently large constant  $\tilde{C} > 0$ . Here the second inequality arises from the condition

$$4 \|\Delta^t\| \|\mathbf{X}^*\| \leq \sigma_{\min}/2,$$

which would hold if  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll \frac{1}{\kappa}$ . An immediate consequence of (111) is that

$$\alpha_2 = \|\tilde{\mathbf{F}}^{t+1} - \mathbf{F}^*\| \leq (\sqrt{2}\kappa)^{-1} \|\mathbf{X}^*\|. \quad (112)$$

To see this, apply the induction hypotheses (93b) and (93e) to get

$$\begin{aligned} \alpha_2 &\leq \left(1 - \frac{\eta\sigma_{\min}}{2}\right) C_{\text{op}} \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}} \right) \|\mathbf{X}^*\| + \tilde{C}\eta \left( \sigma\sqrt{\frac{n}{p}} + \frac{\lambda}{p} \right) \|\mathbf{X}^*\| \\ &\quad + \tilde{C}\eta\sqrt{\frac{n}{p}} C_{\infty}\kappa \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p\sigma_{\min}} \right) \|\mathbf{F}^*\|_{2,\infty}^2 \|\mathbf{X}^*\| \\ &\stackrel{(i)}{\leq} C_{\text{op}} \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}} \right) \|\mathbf{X}^*\| \\ &\stackrel{(ii)}{\leq} (\sqrt{2}\kappa)^{-1} \|\mathbf{X}^*\|. \end{aligned} \quad (113)$$

Here (i) holds under the assumptions that  $C_{\text{op}} \gg \tilde{C} + \tilde{C}C_{\infty}\kappa^2 \sqrt{\frac{\mu^2 r^2 \log n}{np}}$  and that  $\|\mathbf{F}^*\|_{2,\infty} \leq \sqrt{\frac{\mu r}{n}} \|\mathbf{X}^*\|$  (cf. (100b)); (ii) arises since  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll 1/\kappa$  and  $\lambda \ll \sigma\sqrt{np}$ . Under the sample complexity  $n^2 p \gg \kappa^4 \mu^2 r^2 n \log n$ , the first condition can be simplified to  $C_{\text{op}} \gg 2\tilde{C} \gg 1$ .

2. Next we bound  $\alpha_1$ , towards which we first observe that

$$\alpha_1 = \|\mathbf{F}^{t+1} \mathbf{H}^{t+1} - \tilde{\mathbf{F}}^{t+1}\| = \|\mathbf{F}^{t+1} \mathbf{H}^t \mathbf{H}^{t\top} \mathbf{H}^{t+1} - \tilde{\mathbf{F}}^{t+1}\|.$$

It is straightforward to verify that  $\mathbf{H}^{t\top} \mathbf{H}^{t+1}$  is the best rotation matrix to align  $\mathbf{F}^{t+1} \mathbf{H}^t$  and  $\mathbf{F}^*$  (in the sense of (89)). Regarding  $\tilde{\mathbf{F}}^{t+1}$ , we obtain the following claim, which demonstrates that it is already aligned with  $\mathbf{F}^*$ , i.e.  $\mathbf{I}_r$  is the best rotation matrix to align  $\tilde{\mathbf{F}}^{t+1}$  and  $\mathbf{F}^*$ .

**Claim 4.** *Suppose (113) holds true, one has*

$$\mathbf{I}_r = \arg \min_{\mathbf{R} \in \mathcal{O}^{r \times r}} \|\tilde{\mathbf{F}}^{t+1} \mathbf{R} - \mathbf{F}^*\|_{\text{F}}.$$

Now we intend to apply Lemma 22 with

$$\mathbf{F}_0 = \mathbf{F}^*, \quad \mathbf{F}_1 = \tilde{\mathbf{F}}^{t+1}, \quad \mathbf{F}_2 = \mathbf{F}^{t+1} \mathbf{H}^t,$$

for which we need to check the two conditions therein. First, in view of (112), one has

$$\|\mathbf{F}_1 - \mathbf{F}_0\| \|\mathbf{F}_0\| = \|\tilde{\mathbf{F}}^{t+1} - \mathbf{F}^*\| \|\mathbf{F}^*\| \leq \frac{1}{\sqrt{2}\kappa} \|\mathbf{X}^*\| \|\mathbf{F}^*\| = \sigma_{\min} = \frac{1}{2} \sigma_r^2(\mathbf{F}_0).$$

Second, making use of the gradient update rules (27) and the decomposition (99), we obtain

$$\begin{aligned} \|\mathbf{F}^{t+1}\mathbf{H}^t - \tilde{\mathbf{F}}^{t+1}\| &= \left\| (\mathbf{F}^t - \eta \nabla f_{\text{aug}}(\mathbf{F}^t) - \eta \nabla f_{\text{diff}}(\mathbf{F}^t)) \mathbf{H}^t - \tilde{\mathbf{F}}^{t+1} \right\| \\ &\leq \underbrace{\left\| (\mathbf{F}^t - \eta \nabla f_{\text{aug}}(\mathbf{F}^t)) \mathbf{H}^t - \tilde{\mathbf{F}}^{t+1} \right\|}_{:=\theta_1} + \underbrace{\eta \left\| \nabla f_{\text{diff}}(\mathbf{F}^t) \right\|}_{:=\theta_2}. \end{aligned}$$

The term  $\theta_2$  has been controlled as  $\alpha_2$  in the proof of Lemma 10, where we obtained

$$\theta_2 \leq 2C_B \kappa \eta^2 \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \sqrt{r} \sigma_{\max}^2 \|\mathbf{X}^*\|.$$

We now move on to  $\theta_1$ , for which we have

$$\begin{aligned} \theta_1 &\leq \left\| \left( \mathbf{F}^t - \eta \left\{ \nabla f_{\text{aug}}(\mathbf{F}^t) - \eta \begin{bmatrix} \frac{1}{p} \mathcal{P}_\Omega(\mathbf{E}) \mathbf{Y}^t \\ \frac{1}{p} \mathcal{P}_\Omega(\mathbf{E})^\top \mathbf{X}^t \end{bmatrix} - \eta \frac{\lambda}{p} \begin{bmatrix} \mathbf{X}^t \\ \mathbf{Y}^t \end{bmatrix} \right\} \right) \mathbf{H}^t \right. \\ &\quad \left. - \tilde{\mathbf{F}}^{t+1} - \eta \begin{bmatrix} \frac{1}{p} \mathcal{P}_\Omega(\mathbf{E}) \mathbf{Y}^* \\ \frac{1}{p} \mathcal{P}_\Omega(\mathbf{E})^\top \mathbf{X}^* \end{bmatrix} - \eta \frac{\lambda}{p} \begin{bmatrix} \mathbf{X}^* \\ \mathbf{Y}^* \end{bmatrix} \right\| \\ &\quad \underbrace{:= \xi_1} \\ &\quad + \eta \underbrace{\left\| \begin{bmatrix} \frac{1}{p} \mathcal{P}_\Omega(\mathbf{E}) \mathbf{Y}^t \\ \frac{1}{p} \mathcal{P}_\Omega(\mathbf{E})^\top \mathbf{X}^t \end{bmatrix} \mathbf{H}^t + \frac{\lambda}{p} \begin{bmatrix} \mathbf{X}^t \\ \mathbf{Y}^t \end{bmatrix} \mathbf{H}^t - \begin{bmatrix} \frac{1}{p} \mathcal{P}_\Omega(\mathbf{E}) \mathbf{Y}^* \\ \frac{1}{p} \mathcal{P}_\Omega(\mathbf{E})^\top \mathbf{X}^* \end{bmatrix} - \frac{\lambda}{p} \begin{bmatrix} \mathbf{X}^* \\ \mathbf{Y}^* \end{bmatrix} \right\|}_{:=\xi_2}. \end{aligned}$$

Combining [CLL19, Equation (4.13)] and [KMO10a, Lemma 3.2] yields

$$\begin{aligned} \xi_1 &\lesssim \eta \sqrt{\frac{n}{p}} \left( \|\Delta_{\mathbf{X}}^t\|_{2,\infty} \|\mathbf{Y}^*\|_{2,\infty} + \|\Delta_{\mathbf{Y}}^t\|_{2,\infty} \|\mathbf{X}^*\|_{2,\infty} + \|\Delta_{\mathbf{X}}^t\|_{2,\infty} \|\Delta_{\mathbf{Y}}^t\|_{2,\infty} \right) \|\Delta^t\| \\ &\quad + \eta \left( \|\Delta_{\mathbf{X}}^t\| \|\mathbf{Y}^*\| + \|\Delta_{\mathbf{Y}}^t\| \|\mathbf{X}^*\| + \|\Delta_{\mathbf{X}}^t\| \|\Delta_{\mathbf{Y}}^t\| + 2 \|\mathbf{X}^*\| \|\Delta_{\mathbf{X}}^t\| + 2 \|\mathbf{Y}^*\| \|\Delta_{\mathbf{Y}}^t\| + \|\Delta_{\mathbf{X}}^t\|^2 + \|\Delta_{\mathbf{Y}}^t\|^2 \right) \|\Delta^t\| \\ &\lesssim \eta \sqrt{\frac{n}{p}} \|\Delta^t\|_{2,\infty} \|\mathbf{F}^*\|_{2,\infty} \|\Delta^t\| + \eta \|\Delta^t\|^2 \|\mathbf{X}^*\| \\ &\leq \frac{1}{15\kappa} \frac{\sigma_{\min}}{4} \eta \|\Delta^t\|. \end{aligned}$$

Here the penultimate inequality arises from the facts that  $\max\{\|\Delta_{\mathbf{X}}^t\|_{2,\infty}, \|\Delta_{\mathbf{Y}}^t\|_{2,\infty}\} \leq \|\Delta^t\|_{2,\infty} \leq \|\mathbf{F}^*\|_{2,\infty}$  and similarly  $\max\{\|\Delta_{\mathbf{X}}^t\|, \|\Delta_{\mathbf{Y}}^t\|\} \leq \|\Delta^t\| \leq \|\mathbf{X}^*\|$ ; see Lemma 18. In addition, the last line holds because of the induction hypotheses (93b) and (93e), provided that

$$C_\infty \kappa \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \sqrt{\frac{\mu^2 r^2 \log n}{np}} \ll \frac{1}{\kappa^2} \quad \text{and} \quad C_{\text{op}} \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll \frac{1}{\kappa^2}.$$

Again, the first condition would be guaranteed by the sample size condition  $n^2 p \gg \kappa^4 \mu^2 r^2 n \log n$  and the noise condition  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll 1/\kappa$ . Next, the term  $\xi_2$  can be easily controlled as follows

$$\begin{aligned} \xi_2 &\leq \eta \left\| \begin{bmatrix} \frac{1}{p} \mathcal{P}_\Omega(\mathbf{E}) (\mathbf{Y}^t \mathbf{H}^t - \mathbf{Y}^*) \\ \frac{1}{p} \mathcal{P}_\Omega(\mathbf{E})^\top (\mathbf{X}^t \mathbf{H}^t - \mathbf{X}^*) \end{bmatrix} \right\| + \eta \frac{\lambda}{p} \|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*\| \\ &\leq \tilde{C} \eta \left( \sigma \sqrt{\frac{n}{p}} + \frac{\lambda}{p} \right) \|\Delta^t\|, \end{aligned}$$

where the last line follows from the same argument for bounding  $\beta_2$  above. Taking the bounds on  $\theta_1$  and  $\theta_2$  collectively yields

$$\|\mathbf{F}^{t+1}\mathbf{H}^t - \tilde{\mathbf{F}}^{t+1}\| \leq \frac{1}{15\kappa} \frac{\sigma_{\min}}{4} \eta \|\Delta^t\| + \tilde{C} \eta \left( \sigma \sqrt{\frac{n}{p}} + \frac{\lambda}{p} \right) \|\Delta^t\| + 2C_B \kappa \eta^2 \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \sqrt{r} \sigma_{\max}^2 \|\mathbf{X}^*\|$$

$$\leq \frac{1}{5\kappa} \frac{\sigma_{\min}}{4} \eta \|\Delta^t\| + 2C_B \kappa \eta^2 \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \sqrt{r} \sigma_{\max}^2 \|\mathbf{X}^*\|. \quad (114)$$

The final inequality is true as long as  $\lambda \asymp \sigma \sqrt{np}$  and  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll \frac{1}{\kappa}$ . An immediate consequence of (114) is that

$$\|\tilde{\mathbf{F}}^{t+1} - \mathbf{F}^{t+1} \mathbf{H}^t\| \leq (2\sqrt{2}\kappa)^{-1} \|\mathbf{X}^*\|, \quad (115)$$

as long as  $\eta \ll 1/(C_B \kappa^2 \sigma_{\max} \sqrt{r})$ ,  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll \frac{1}{\kappa}$  and  $\lambda \asymp \sigma \sqrt{np}$ . As a result, one obtains

$$\|\mathbf{F}_1 - \mathbf{F}_2\| \|\mathbf{F}_0\| = \|\tilde{\mathbf{F}}^{t+1} - \mathbf{F}^{t+1} \mathbf{H}^t\| \|\mathbf{F}^*\| \leq (2\sqrt{2}\kappa)^{-1} \|\mathbf{X}^*\| \|\mathbf{F}^*\| = \sigma_{\min}/2 = \sigma_{\min}^2(\mathbf{F}_0)/4.$$

Armed with these two conditions, we can invoke Lemma 22 to obtain

$$\begin{aligned} \alpha_1 &= \|\tilde{\mathbf{F}}^{t+1} - \mathbf{F}^{t+1} \mathbf{H}^{t+1}\| \leq 5\kappa \|\tilde{\mathbf{F}}^{t+1} - \mathbf{F}^{t+1} \mathbf{H}^t\| \\ &\leq \frac{1}{4} \sigma_{\min} \eta \|\Delta^t\| + 10C_B \kappa^2 \eta^2 \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \sqrt{r} \sigma_{\max}^2 \|\mathbf{X}^*\| \\ &\leq \frac{1}{4} \sigma_{\min} \eta \|\Delta^t\| + \eta \left( \sigma \sqrt{\frac{n}{p}} + \frac{\lambda}{p} \right) \|\mathbf{X}^*\|, \end{aligned}$$

provided that  $\eta \ll 1/(C_B \kappa^3 \sigma_{\max} \sqrt{r})$ .

Combine the bounds on  $\alpha_1$  and  $\alpha_2$  to reach

$$\begin{aligned} &\|\mathbf{F}^{t+1} \mathbf{H}^{t+1} - \mathbf{F}^*\| \\ &\leq \left(1 - \frac{\eta}{2} \sigma_{\min}\right) \|\Delta^t\| + (\tilde{C} + 1) \eta \left( \sigma \sqrt{\frac{n}{p}} + \frac{\lambda}{p} \right) \|\mathbf{X}^*\| + \tilde{C} \eta \sqrt{\frac{n}{p}} \|\Delta^t\|_{2,\infty} \|\mathbf{F}^*\|_{2,\infty} \|\mathbf{X}^*\| + \frac{\sigma_{\min}}{4} \eta \|\Delta^t\| \\ &\leq \left(1 - \frac{\eta}{4} \sigma_{\min}\right) C_{\text{op}} \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{X}^*\| + (\tilde{C} + 1) \eta \left( \sigma \sqrt{\frac{n}{p}} + \frac{\lambda}{p} \right) \|\mathbf{X}^*\| \\ &\quad + \tilde{C} \eta \sqrt{\frac{n}{p}} C_{\infty} \kappa \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{F}^*\|_{2,\infty}^2 \|\mathbf{X}^*\| \\ &\leq C_{\text{op}} \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{X}^*\|, \end{aligned}$$

with the proviso that  $C_{\text{op}} \gg 1$  and  $n^2 p \gg \kappa^4 \mu^2 r^2 n \log n$ . Here the last line follows from the same argument as in bounding (113). This completes the proof.

*Proof of Claim 4.* In view of [MWCC17, Lemma 35], it suffices to show that  $\mathbf{F}^{*\top} \tilde{\mathbf{F}}^{t+1}$  is symmetric and positive semidefinite. Recognizing that  $\mathbf{F}^{*\top} \mathbf{F}^t \mathbf{H}^t = (\mathbf{X}^{*\top} \mathbf{X}^t + \mathbf{Y}^{*\top} \mathbf{Y}^t) \mathbf{H}^t$  is symmetric (see [MWCC17, Lemma 35]), it is straightforward to verify that  $\mathbf{F}^{*\top} \tilde{\mathbf{F}}^{t+1}$  is also symmetric (which we omit here for brevity). In addition, by (112) we have

$$\|\mathbf{F}^{*\top} \tilde{\mathbf{F}}^{t+1} - \mathbf{F}^{*\top} \mathbf{F}^*\| \leq \|\mathbf{F}^*\| \|\tilde{\mathbf{F}}^{t+1} - \mathbf{F}^*\| = \alpha_2 \|\mathbf{F}^*\| \leq \frac{1}{\sqrt{2}\kappa} \|\mathbf{X}^*\| \|\mathbf{F}^*\| = \sigma_{\min}.$$

Since  $\mathbf{F}^{*\top} \mathbf{F}^* = \mathbf{X}^{*\top} \mathbf{X}^* + \mathbf{Y}^{*\top} \mathbf{Y}^* = 2\Sigma^*$ , Weyl's inequality gives

$$\lambda_{\min}(\mathbf{F}^{*\top} \tilde{\mathbf{F}}^{t+1}) \geq 2\sigma_{\min} - \|\mathbf{F}^{*\top} \tilde{\mathbf{F}}^{t+1} - \mathbf{F}^{*\top} \mathbf{F}^*\| \geq \sigma_{\min} > 0,$$

where  $\lambda_{\min}(\mathbf{A})$  stands for the minimum eigenvalue of a matrix  $\mathbf{A}$ . To conclude,  $\mathbf{F}^{*\top} \tilde{\mathbf{F}}^{t+1}$  is both symmetric and positive semidefinite, thus establishing the claim.  $\square$

## D.5 Proof of Lemma 12

Without loss of generality, we consider the case when  $1 \leq l \leq n$ ; the case with  $n+1 \leq l \leq 2n$  can be derived in a similar way. From the definition of  $\mathbf{R}^{t+1,(l)}$  (cf. (95b)), we have

$$\|\mathbf{F}^{t+1} \mathbf{H}^{t+1} - \mathbf{F}^{t+1,(l)} \mathbf{R}^{t+1,(l)}\|_{\mathbb{F}} \leq \|\mathbf{F}^{t+1} \mathbf{H}^t - \mathbf{F}^{t+1,(l)} \mathbf{R}^{t,(l)}\|_{\mathbb{F}}.$$

The gradient update rules (27) and (90) give

$$\begin{aligned} & \mathbf{F}^{t+1} \mathbf{H}^t - \mathbf{F}^{t+1,(l)} \mathbf{R}^{t,(l)} \\ &= [\mathbf{F}^t - \eta \nabla f(\mathbf{F}^t)] \mathbf{H}^t - [\mathbf{F}^{t,(l)} - \eta \nabla f^{(l)}(\mathbf{F}^{t,(l)})] \mathbf{R}^{t,(l)} \\ &= \mathbf{F}^t \mathbf{H}^t - \eta \nabla f(\mathbf{F}^t \mathbf{H}^t) - [\mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)} - \eta \nabla f^{(l)}(\mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)})] \\ &= \underbrace{(\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)}) - \eta [\nabla f_{\text{aug}}(\mathbf{F}^t \mathbf{H}^t) - \nabla f_{\text{aug}}(\mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)})]}_{:=\mathbf{A}_1} - \eta \underbrace{[\nabla f_{\text{diff}}(\mathbf{F}^t \mathbf{H}^t) - \nabla f_{\text{diff}}(\mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)})]}_{:=\mathbf{A}_2} \\ & \quad + \eta \underbrace{[\nabla f^{(l)}(\mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)}) - \nabla f(\mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)})]}_{:=\mathbf{A}_3}, \end{aligned}$$

where we have used the facts that  $\nabla f(\mathbf{F})\mathbf{R} = \nabla f(\mathbf{F}\mathbf{R})$  and  $\nabla f^{(l)}(\mathbf{F})\mathbf{R} = \nabla f^{(l)}(\mathbf{F}\mathbf{R})$  for any orthonormal matrix  $\mathbf{R} \in \mathcal{O}^{r \times r}$ .

In what follows, we shall bound  $\mathbf{A}_1$ ,  $\mathbf{A}_2$  and  $\mathbf{A}_3$  sequentially.

1. The first term  $\mathbf{A}_1$  is similar to  $\alpha_1$  in the proof of Lemma 10. Going through the same derivations therein, we obtain

$$\|\mathbf{A}_1\|_{\mathbb{F}} \leq \left(1 - \frac{\sigma_{\min}}{20} \eta\right) \|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)}\|_{\mathbb{F}}, \quad (116)$$

provided that  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll \frac{1}{\kappa^4 \mu r \log n}$  and that  $0 \leq \eta \leq 1/(1000\kappa\sigma_{\max})$ .

2. Next, we turn attention to  $\mathbf{A}_2$ , which clearly obeys

$$\|\mathbf{A}_2\|_{\mathbb{F}} \leq \eta \|\nabla f_{\text{diff}}(\mathbf{F}^t \mathbf{H}^t)\|_{\mathbb{F}} + \eta \|\nabla f_{\text{diff}}(\mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)})\|_{\mathbb{F}}.$$

Recall from the term  $\alpha_2$  in the proof of Lemma 10 that

$$\eta \|\nabla f_{\text{diff}}(\mathbf{F}^t \mathbf{H}^t)\|_{\mathbb{F}} \leq 2C_{\text{B}}\kappa\eta^2 \left(\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}}\right) \sqrt{r}\sigma_{\max}^2 \|\mathbf{X}^*\|.$$

Applying Lemma 15 and going through the same derivation as in bounding  $\alpha_2$  in the proof of Lemma 10, one gets

$$\eta \|\nabla f_{\text{diff}}(\mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)})\|_{\mathbb{F}} \leq 2C_{\text{B}}\kappa\eta^2 \left(\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}}\right) \sqrt{r}\sigma_{\max}^2 \|\mathbf{X}^*\|.$$

Combine the above three inequalities to obtain

$$\begin{aligned} \|\mathbf{A}_2\|_{\mathbb{F}} &\leq 4C_{\text{B}}\kappa\eta^2 \left(\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}}\right) \sqrt{r}\sigma_{\max}^2 \|\mathbf{X}^*\| \\ &\leq 4\sqrt{n}C_{\text{B}}\kappa\eta^2 \left(\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}}\right) \sqrt{r}\sigma_{\max}^2 \|\mathbf{X}^*\|_{2,\infty} \\ &\leq \eta \left(\sigma \sqrt{\frac{n}{p}} + \frac{\lambda}{p}\right) \|\mathbf{F}^*\|_{2,\infty}. \end{aligned}$$

Here the second inequality arises from the elementary inequality  $\|\mathbf{X}^*\| \leq \sqrt{n}\|\mathbf{X}^*\|_{2,\infty}$ , whereas the last one holds true because of  $\|\mathbf{X}^*\|_{2,\infty} \leq \|\mathbf{F}^*\|_{2,\infty}$  and the condition that  $\eta \ll \frac{1}{n\kappa^2\sigma_{\max}}$ .



3. We are now left with  $\mathbf{A}_3$ . To this end, we first observe that

$$\mathbf{A}_3 = \eta \begin{bmatrix} \underbrace{\left[ \mathcal{P}_{l,\cdot} \left( \mathbf{X}^{t,(l)} \mathbf{Y}^{t,(l)\top} - \mathbf{M}^* \right) - p^{-1} \mathcal{P}_{\Omega_{l,\cdot}} \left( \mathbf{X}^{t,(l)} \mathbf{Y}^{t,(l)\top} - \mathbf{M}^* \right) \right] \mathbf{Y}^{t,(l)} \mathbf{R}^{t,(l)}}_{:=\mathbf{B}_1} + \underbrace{p^{-1} \mathcal{P}_{\Omega_{l,\cdot}} \left( \mathbf{E} \right) \mathbf{Y}^{t,(l)} \mathbf{R}^{t,(l)}}_{:=\mathbf{C}_1} \\ \underbrace{\left[ \mathcal{P}_{l,\cdot} \left( \mathbf{X}^{t,(l)} \mathbf{Y}^{t,(l)\top} - \mathbf{M}^* \right) - p^{-1} \mathcal{P}_{\Omega_{l,\cdot}} \left( \mathbf{X}^{t,(l)} \mathbf{Y}^{t,(l)\top} - \mathbf{M}^* \right) \right]^\top \mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)}}_{:=\mathbf{B}_2} + \underbrace{p^{-1} \mathcal{P}_{\Omega_{l,\cdot}} \left( \mathbf{E} \right)^\top \mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)}}_{:=\mathbf{C}_2} \end{bmatrix}.$$

The following claims allow one to bound  $\mathbf{B}_1, \mathbf{B}_2$  and  $\mathbf{C}_1, \mathbf{C}_2$ ; the proofs are deferred to the end of this subsection.

**Claim 5.** Suppose that  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll \frac{1}{\sqrt{\kappa^2 \log n}}$  and that  $np \gg \log^2 n$ . With probability at least  $1 - O(n^{-100})$ ,

$$\|\mathbf{B}_1\|_{\text{F}} \lesssim \sqrt{\frac{\mu^2 r^2 \log n}{np}} \|\mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)} - \mathbf{F}^*\|_{2,\infty} \sigma_{\max}. \quad (117)$$

**Claim 6.** Suppose that  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll \frac{1}{\sqrt{\kappa^2 \log n}}$  and that  $np \gg \log n$ . With probability at least  $1 - O(n^{-100})$ ,

$$\|\mathbf{B}_2\|_{\text{F}} \lesssim \sqrt{\frac{\mu^2 r^2 \log n}{np}} \|\mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)} - \mathbf{F}^*\|_{2,\infty} \sigma_{\max}. \quad (118)$$

**Claim 7.** Suppose that  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll \frac{1}{\sqrt{\kappa^2 \log n}}$  and that  $np \gg \log^3 n$ . With probability at least  $1 - O(n^{-100})$ ,

$$\max \{ \|\mathbf{C}_1\|_{\text{F}}, \|\mathbf{C}_2\|_{\text{F}} \} \lesssim \sigma \sqrt{\frac{n \log n}{p}} \|\mathbf{F}^*\|_{2,\infty}. \quad (119)$$

With these claims in place, one can readily obtain that

$$\begin{aligned} \|\mathbf{A}_3\|_{\text{F}} &\leq \eta (\|\mathbf{B}_1\|_{\text{F}} + \|\mathbf{B}_2\|_{\text{F}} + \|\mathbf{C}_1\|_{\text{F}} + \|\mathbf{C}_2\|_{\text{F}}) \\ &\lesssim \eta \sigma \sqrt{\frac{n \log n}{p}} \|\mathbf{F}^*\|_{2,\infty} + \eta \sqrt{\frac{\mu^2 r^2 \log n}{np}} \|\mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)} - \mathbf{F}^*\|_{2,\infty} \sigma_{\max} \\ &\leq \eta \sigma \sqrt{\frac{n \log n}{p}} \|\mathbf{F}^*\|_{2,\infty} + \eta \sqrt{\frac{\mu^2 r^2 \log n}{np}} (C_\infty \kappa + C_3) \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{F}^*\|_{2,\infty} \sigma_{\max}, \end{aligned}$$

where the last line follows from the induction hypotheses (93c) and (93e).

This together with the bounds on  $\mathbf{A}_1$  and  $\mathbf{A}_2$  gives: for some constant  $\tilde{C} > 0$ ,

$$\begin{aligned} \|\mathbf{F}^{t+1} \mathbf{H}^{t+1} - \mathbf{F}^{t+1,(l)} \mathbf{R}^{t+1,(l)}\|_{\text{F}} &\leq \|\mathbf{A}_1\|_{\text{F}} + \|\mathbf{A}_2\|_{\text{F}} + \|\mathbf{A}_3\|_{\text{F}} \\ &\leq \left( 1 - \frac{\sigma_{\min}}{20} \eta \right) \|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)}\|_{\text{F}} + \eta \left( \sigma \sqrt{\frac{n}{p}} + \frac{\lambda}{p} \right) \|\mathbf{F}^*\|_{2,\infty} \\ &\quad + \tilde{C} \eta \sigma \sqrt{\frac{n \log n}{p}} \|\mathbf{F}^*\|_{2,\infty} + \tilde{C} \eta \sqrt{\frac{\mu^2 r^2 \log n}{np}} (C_\infty \kappa + C_3) \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{F}^*\|_{2,\infty} \sigma_{\max} \\ &\stackrel{(i)}{\leq} \left( 1 - \frac{\sigma_{\min}}{20} \eta \right) C_3 \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{F}^*\|_{2,\infty} + \eta \left( \sigma \sqrt{\frac{n}{p}} + \frac{\lambda}{p} \right) \|\mathbf{F}^*\|_{2,\infty} \\ &\quad + \tilde{C} \eta \sigma \sqrt{\frac{n \log n}{p}} \|\mathbf{F}^*\|_{2,\infty} + \tilde{C} \eta \sqrt{\frac{\mu^2 r^2 \log n}{np}} (C_\infty \kappa + C_3) \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{F}^*\|_{2,\infty} \sigma_{\max} \end{aligned}$$

$$\stackrel{\text{(ii)}}{\leq} C_3 \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{F}^*\|_{2,\infty}$$

as claimed. Here, (i) invokes the induction hypothesis (93c), whereas (ii) holds as long as  $C_3$  is large enough and the sample size satisfies  $n^2 p \gg \kappa^4 \mu^2 r^2 n \log n$ .

*Proof of Claim 5.* For notational simplicity, we denote

$$\mathbf{C} \triangleq \mathbf{X}^{t,(l)} \mathbf{Y}^{t,(l)\top} - \mathbf{M}^* = \mathbf{X}^{t,(l)} \mathbf{Y}^{t,(l)\top} - \mathbf{X}^* \mathbf{Y}^{*\top}. \quad (120)$$

Since the Frobenius norm is unitarily invariant, we have

$$\|\mathbf{B}_1\|_{\text{F}} = \left\| \underbrace{[p^{-1} \mathcal{P}_{\Omega_{l,\cdot}}(\mathbf{C}) - \mathcal{P}_{l,\cdot}(\mathbf{C})]}_{:=\mathbf{W}} \mathbf{Y}^{t,(l)} \right\|_{\text{F}}.$$

All nonzero entries of the matrix  $\mathbf{W}$  reside in its  $l$ th row and therefore

$$p \|\mathbf{B}_1\|_{\text{F}} = \left\| \sum_{j=1}^n (\delta_{l,j} - p) C_{l,j} \mathbf{Y}_{j,\cdot}^{t,(l)} \right\|_2,$$

where  $\delta_{l,j} \triangleq \mathbf{1}_{\{(l,j) \in \Omega\}}$ . Notice that conditional on  $\mathbf{X}^{t,(l)}$  and  $\mathbf{Y}^{t,(l)}$ , the right-hand side is composed of a sum of independent random vectors, where the randomness comes from  $\{\delta_{l,j}\}_{1 \leq j \leq n}$ . It then follows that

$$\begin{aligned} L &\triangleq \max_{1 \leq j \leq n} \left\| (\delta_{l,j} - p) C_{l,j} \mathbf{Y}_{j,\cdot}^{t,(l)} \right\|_2 \leq \|\mathbf{C}\|_{\infty} \|\mathbf{Y}^{t,(l)}\|_{2,\infty} \stackrel{\text{(i)}}{\leq} 2 \|\mathbf{C}\|_{\infty} \|\mathbf{Y}^*\|_{2,\infty}, \\ V &\triangleq \left\| \sum_{j=1}^n \mathbb{E}[(\delta_{l,j} - p)^2] C_{l,j}^2 \mathbf{Y}_{j,\cdot}^{t,(l)} \mathbf{Y}_{j,\cdot}^{t,(l)\top} \right\| \leq p \|\mathbf{C}\|_{\infty}^2 \left\| \sum_{j=1}^n \mathbf{Y}_{j,\cdot}^{t,(l)} \mathbf{Y}_{j,\cdot}^{t,(l)\top} \right\| \\ &= p \|\mathbf{C}\|_{\infty}^2 \|\mathbf{Y}^{t,(l)}\|_{\text{F}}^2 \stackrel{\text{(ii)}}{\leq} 4p \|\mathbf{C}\|_{\infty}^2 \|\mathbf{Y}^*\|_{\text{F}}^2. \end{aligned}$$

Here, both (i) and (ii) arise from Lemma 18, as long as  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll \frac{1}{\sqrt{\kappa^2 \log n}}$ . The matrix Bernstein inequality [Tro15, Theorem 6.1.1] reveals that

$$\left\| \sum_{j=1}^n (\delta_{l,j} - p) C_{l,j} \mathbf{Y}_{j,\cdot}^{t,(l)} \right\|_2 \lesssim \sqrt{V \log n} + L \log n \lesssim \sqrt{p \|\mathbf{C}\|_{\infty}^2 \|\mathbf{Y}^*\|_{\text{F}}^2 \log n} + \|\mathbf{C}\|_{\infty} \|\mathbf{Y}^*\|_{2,\infty} \log n$$

with probability exceeding  $1 - O(n^{-100})$ . As a result, we arrive at

$$p \|\mathbf{B}_1\|_{\text{F}} \lesssim \sqrt{p \log n} \|\mathbf{C}\|_{\infty} \|\mathbf{Y}^*\|_{\text{F}} + \sqrt{np} \|\mathbf{C}\|_{\infty} \|\mathbf{Y}^*\|_{2,\infty} \quad (121)$$

as soon as  $np \gg \log^2 n$ .

To finish up, we make the observation that

$$\begin{aligned} \|\mathbf{C}\|_{\infty} &= \|\mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)} (\mathbf{Y}^{t,(l)} \mathbf{R}^{t,(l)})^{\top} - \mathbf{X}^* \mathbf{Y}^{*\top}\|_{\infty} \\ &\leq \left\| \left( \mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)} - \mathbf{X}^* \right) \left( \mathbf{Y}^{t,(l)} \mathbf{R}^{t,(l)} \right)^{\top} \right\|_{\infty} + \left\| \mathbf{X}^* \left( \mathbf{Y}^{t,(l)} \mathbf{R}^{t,(l)} - \mathbf{Y}^* \right)^{\top} \right\|_{\infty} \\ &\leq \left\| \mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)} - \mathbf{X}^* \right\|_{2,\infty} \left\| \mathbf{Y}^{t,(l)} \mathbf{R}^{t,(l)} \right\|_{2,\infty} + \|\mathbf{X}^*\|_{2,\infty} \left\| \mathbf{Y}^{t,(l)} \mathbf{R}^{t,(l)} - \mathbf{Y}^* \right\|_{2,\infty} \\ &\leq 3 \|\mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)} - \mathbf{F}^*\|_{2,\infty} \|\mathbf{F}^*\|_{2,\infty}, \end{aligned} \quad (122)$$

where the last line arises from Lemma 18. This combined with (121) gives

$$\|\mathbf{B}_1\|_{\text{F}} \lesssim \sqrt{\frac{\log n}{p}} \|\mathbf{C}\|_{\infty} \|\mathbf{Y}^*\|_{\text{F}} + \sqrt{\frac{n}{p}} \|\mathbf{C}\|_{\infty} \|\mathbf{Y}^*\|_{2,\infty}$$

$$\begin{aligned}
&\stackrel{(i)}{\lesssim} \sqrt{\frac{\log n}{p}} \|\mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)} - \mathbf{F}^*\|_{2,\infty} \|\mathbf{F}^*\|_{2,\infty} \|\mathbf{Y}^*\|_{\text{F}} + \sqrt{\frac{n}{p}} \|\mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)} - \mathbf{F}^*\|_{2,\infty} \|\mathbf{F}^*\|_{2,\infty}^2 \\
&\stackrel{(ii)}{\lesssim} \sqrt{\frac{\log n}{p}} \|\mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)} - \mathbf{F}^*\|_{2,\infty} \sqrt{\frac{\mu r^2}{n}} \sigma_{\max} + \sqrt{\frac{n}{p}} \|\mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)} - \mathbf{F}^*\|_{2,\infty} \frac{\mu r}{n} \sigma_{\max} \\
&\lesssim \sqrt{\frac{\mu^2 r^2 \log n}{np}} \|\mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)} - \mathbf{F}^*\|_{2,\infty} \sigma_{\max},
\end{aligned}$$

where (i) comes from (122), and (ii) makes use of the incoherence condition  $\|\mathbf{F}^*\|_{2,\infty} \leq \sqrt{\mu r \sigma_{\max}/n}$  and the fact that  $\|\mathbf{Y}^*\|_{\text{F}} \leq \sqrt{r \sigma_{\max}}$ .  $\square$

*Proof of Claim 6.* Instate the notation in proof of Claim 5. By the unitary invariance of Frobenius norm and the fact that all nonzero entries of the matrix  $\mathbf{W}$  reside in its  $l$ th row, we have

$$p \|\mathbf{B}_2\|_{\text{F}} = \left\| p \mathbf{W}^\top \mathbf{X}^{t,(l)} \right\|_{\text{F}} = \left\| \underbrace{\begin{bmatrix} (\delta_{l1} - p) C_{l1} \\ \vdots \\ (\delta_{ln} - p) C_{ln} \end{bmatrix}}_{:=\mathbf{b}} \mathbf{X}_{l,\cdot}^{t,(l)} \right\|_{\text{F}} = \|\mathbf{b}\|_2 \|\mathbf{X}_{l,\cdot}^{t,(l)}\|_2.$$

We can write  $\mathbf{b}$  as

$$\mathbf{b} = \sum_{j=1}^n \underbrace{e_j (\delta_{lj} - p) C_{lj}}_{:=\mathbf{u}_j} = \sum_{j=1}^n \mathbf{u}_j.$$

Note that for all  $j$ , one has

$$\begin{aligned}
L &\triangleq \max_{1 \leq j \leq n} \|\mathbf{u}_j\|_2 \leq \|\mathbf{C}\|_{\infty}, \\
V &\triangleq \left\| \sum_{j=1}^n \mathbb{E}[(\delta_{lj} - p)^2] C_{lj}^2 \mathbf{e}_j \mathbf{e}_j^\top \right\| \leq p \|\mathbf{C}\|_{\infty}^2 \left\| \sum_{j=1}^n \mathbf{e}_j \mathbf{e}_j^\top \right\| = np \|\mathbf{C}\|_{\infty}^2.
\end{aligned}$$

Then the matrix Bernstein inequality [Tro15, Theorem 6.1.1] reveals that

$$\begin{aligned}
\|\mathbf{b}\|_2 &\lesssim \sqrt{V \log n} + L \log n \lesssim \sqrt{np \log n} \|\mathbf{C}\|_{\infty} + \|\mathbf{C}\|_{\infty} \log n \\
&\lesssim \sqrt{np \log n} \|\mathbf{C}\|_{\infty} \\
&\lesssim \sqrt{np \log n} \|\mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)} - \mathbf{F}^*\|_{2,\infty} \|\mathbf{F}^*\|_{2,\infty}
\end{aligned}$$

with probability exceeding  $1 - O(n^{-100})$  as long as  $np \gg \log n$ . Here the last relation uses (122). Observe that  $\|\mathbf{X}^{t,(l)}\|_{2,\infty} \leq 2\|\mathbf{F}^*\|_{2,\infty}$  as long as  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll \frac{1}{\sqrt{\kappa^2 \log n}}$ ; see Lemma 18. Making use of the incoherence condition (100a) to get

$$\|\mathbf{B}_2\|_{\text{F}} \lesssim \sqrt{\frac{n \log n}{p}} \|\mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)} - \mathbf{F}^*\|_{2,\infty} \|\mathbf{F}^*\|_{2,\infty}^2 \lesssim \sqrt{\frac{\mu^2 r^2 \log n}{np}} \|\mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)} - \mathbf{F}^*\|_{2,\infty} \sigma_{\max}.$$

We can then conclude the proof.  $\square$

*Proof of Claim 7.* By the unitary invariance of the Frobenius norm, one has

$$\|\mathbf{C}_1\|_{\text{F}} = p^{-1} \|\mathcal{P}_{\Omega_{l,\cdot}}(\mathbf{E}) \mathbf{Y}^{t,(l)}\|_{\text{F}}.$$

Since the entries of  $\mathcal{P}_{\Omega_{l,\cdot}}(\mathbf{E})$  are all zero except those on the  $l$ th row, we have

$$p \|\mathbf{C}_1\|_{\text{F}} = \left\| \sum_{j=1}^n \underbrace{\delta_{lj} E_{lj} \mathbf{Y}_{j,\cdot}^{t,(l)}}_{:=\mathbf{u}_j} \right\|_2,$$

where we denote  $\delta_{lj} \triangleq \mathbf{1}_{(l,j) \in \Omega}$ . Since  $\mathbf{Y}^{t,(l)}$  is independent of  $\{\delta_{lj}\}_{1 \leq j \leq n}$  and  $\{E_{lj}\}_{1 \leq j \leq n}$ , the vectors  $\{\mathbf{u}_j\}_{1 \leq j \leq n}$  are independent conditioning on  $\mathbf{Y}^{t,(l)}$ . Therefore, from now on we shall condition on a fixed  $\mathbf{Y}^{t,(l)}$ . It is easy to verify that

$$\|\|\mathbf{u}_j\|_2\|_{\psi_1} \leq \|\mathbf{Y}^{t,(l)}\|_{2,\infty} \|\delta_{lj} E_{lj}\|_{\psi_1} \lesssim \sigma \|\mathbf{Y}^{t,(l)}\|_{2,\infty},$$

where  $\|\cdot\|_{\psi_1}$  denotes the sub-exponential norm [KLT11, Section 6]. Further, one can calculate

$$V := \left\| \mathbb{E} \left[ \sum_{j=1}^n (\delta_{lj} E_{lj})^2 \mathbf{Y}_{j,\cdot}^{t,(l)} \mathbf{Y}_{j,\cdot}^{t,(l)\top} \right] \right\| \lesssim p \sigma^2 \left\| \mathbb{E} \left[ \sum_{j=1}^n \mathbf{Y}_{j,\cdot}^{t,(l)} \mathbf{Y}_{j,\cdot}^{t,(l)\top} \right] \right\| = p \sigma^2 \|\mathbf{Y}^{t,(l)}\|_{\mathbb{F}}^2.$$

Invoke the matrix Bernstein inequality [KLT11, Proposition 2] to discover that with probability at least  $1 - O(n^{-100})$ ,

$$\begin{aligned} \left\| \sum_{j=1}^n \mathbf{u}_j \right\|_2 &\lesssim \sqrt{V \log n} + \left\| \|\mathbf{u}_j\|_2 \right\|_{\psi_1} \log^2 n \\ &\lesssim \sqrt{p \sigma^2 \|\mathbf{Y}^{t,(l)}\|_{\mathbb{F}}^2 \log n} + \sigma \|\mathbf{Y}^{t,(l)}\|_{2,\infty} \log^2 n \\ &\lesssim \sigma \sqrt{np \log n} \|\mathbf{Y}^{t,(l)}\|_{2,\infty} + \sigma \|\mathbf{Y}^{t,(l)}\|_{2,\infty} \log^2 n \\ &\lesssim \sigma \sqrt{np \log n} \|\mathbf{Y}^{t,(l)}\|_{2,\infty}, \end{aligned}$$

where the third inequality follows from  $\|\mathbf{Y}^{t,(l)}\|_{\mathbb{F}}^2 \leq n \|\mathbf{Y}^{t,(l)}\|_{2,\infty}^2$ , and the last inequality holds if  $np \gg \log^3 n$ . We then complete the proof by observing that  $\|\mathbf{Y}^{t,(l)}\|_{2,\infty} \leq 2 \|\mathbf{F}^*\|_{2,\infty}$  as long as  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll \frac{1}{\sqrt{\kappa^2 \log n}}$ ; see Lemma 18. The bound on  $\mathbf{C}_2$  follows from similar arguments to that used to bound  $\mathbf{B}_2$ .  $\square$

## D.6 Proof of Lemma 13

Without loss of generality, we assume  $1 \leq l \leq n$ . One can then decompose  $(\mathbf{F}^{t+1,(l)} \mathbf{H}^{t+1,(l)} - \mathbf{F}^*)_{l,\cdot}$  as

$$\begin{aligned} (\mathbf{F}^{t+1,(l)} \mathbf{H}^{t+1,(l)} - \mathbf{F}^*)_{l,\cdot} &= \mathbf{X}_{l,\cdot}^{t+1,(l)} \mathbf{H}^{t+1,(l)} - \mathbf{X}_{l,\cdot}^* \\ &= \left\{ \mathbf{X}_{l,\cdot}^{t,(l)} - \eta \left[ (\mathbf{X}^{t,(l)} \mathbf{Y}^{t,(l)\top} - \mathbf{M}^*)_{l,\cdot} \mathbf{Y}^{t,(l)} + \frac{\lambda}{p} \mathbf{X}_{l,\cdot}^{t,(l)} \right] \right\} \mathbf{H}^{t+1,(l)} - \mathbf{X}_{l,\cdot}^* \\ &= \mathbf{X}_{l,\cdot}^{t,(l)} \mathbf{H}^{t+1,(l)} - \mathbf{X}_{l,\cdot}^* - \eta \left[ (\mathbf{X}^{t,(l)} \mathbf{Y}^{t,(l)\top} - \mathbf{M}^*)_{l,\cdot} \mathbf{Y}^{t,(l)} + \frac{\lambda}{p} \mathbf{X}_{l,\cdot}^{t,(l)} \right] \mathbf{H}^{t+1,(l)} \\ &= \underbrace{\mathbf{X}_{l,\cdot}^{t,(l)} \mathbf{H}^{t,(l)} - \mathbf{X}_{l,\cdot}^* - \eta \left[ (\mathbf{X}^{t,(l)} \mathbf{Y}^{t,(l)\top} - \mathbf{M}^*)_{l,\cdot} \mathbf{Y}^{t,(l)} + \frac{\lambda}{p} \mathbf{X}_{l,\cdot}^{t,(l)} \right] \mathbf{H}^{t,(l)}}_{:= \mathbf{a}_1} \\ &\quad + \underbrace{\left\{ \mathbf{X}_{l,\cdot}^{t,(l)} \mathbf{H}^{t,(l)} - \eta \left[ (\mathbf{X}^{t,(l)} \mathbf{Y}^{t,(l)\top} - \mathbf{M}^*)_{l,\cdot} \mathbf{Y}^{t,(l)} + \frac{\lambda}{p} \mathbf{X}_{l,\cdot}^{t,(l)} \right] \mathbf{H}^{t,(l)} \right\} \left[ (\mathbf{H}^{t,(l)})^{-1} \mathbf{H}^{t+1,(l)} - \mathbf{I}_r \right]}_{:= \mathbf{a}_2}. \end{aligned}$$

Note that here  $\mathbf{a}_1$  and  $\mathbf{a}_2$  are  $r$ -dimensional row vectors. In the sequel, let us control  $\|\mathbf{a}_1\|_2$  and  $\|\mathbf{a}_2\|_2$  separately.

1. We begin with  $\mathbf{a}_1$ . For notational convenience, define  $\Delta^{t,(l)} \triangleq \begin{bmatrix} \Delta_{\mathbf{X}}^{t,(l)} \\ \Delta_{\mathbf{Y}}^{t,(l)} \end{bmatrix}$ , where  $\Delta_{\mathbf{X}}^{t,(l)} \triangleq \mathbf{X}^{t,(l)} \mathbf{H}^{t,(l)} - \mathbf{X}^*$  and  $\Delta_{\mathbf{Y}}^{t,(l)} \triangleq \mathbf{Y}^{t,(l)} \mathbf{H}^{t,(l)} - \mathbf{Y}^*$ . Then  $\mathbf{a}_1$  can be rewritten as

$$\begin{aligned} \mathbf{a}_1 &= (\Delta_{\mathbf{X}}^{t,(l)})_{l,\cdot} - \eta \left[ (\Delta_{\mathbf{X}}^{t,(l)})_{l,\cdot} (\mathbf{Y}^{t,(l)} \mathbf{H}^{t,(l)})^\top + \mathbf{X}_{l,\cdot}^* \Delta_{\mathbf{Y}}^{t,(l)\top} \right] \mathbf{Y}^{t,(l)} \mathbf{H}^{t,(l)} - \eta \frac{\lambda}{p} \mathbf{X}_{l,\cdot}^{t,(l)} \mathbf{H}^{t,(l)} \\ &= (\Delta_{\mathbf{X}}^{t,(l)})_{l,\cdot} \left[ \mathbf{I}_r - \eta (\mathbf{Y}^{t,(l)} \mathbf{H}^{t,(l)})^\top \mathbf{Y}^{t,(l)} \mathbf{H}^{t,(l)} \right] - \eta \mathbf{X}_{l,\cdot}^* \Delta_{\mathbf{Y}}^{t,(l)\top} \mathbf{Y}^{t,(l)} \mathbf{H}^{t,(l)} - \eta \frac{\lambda}{p} \mathbf{X}_{l,\cdot}^{t,(l)} \mathbf{H}^{t,(l)}, \end{aligned}$$

which together with the triangle inequality yields

$$\|\mathbf{a}_1\|_2 \leq \left\| \mathbf{I}_r - \eta (\mathbf{Y}^{t,(l)} \mathbf{H}^{t,(l)})^\top \mathbf{Y}^{t,(l)} \mathbf{H}^{t,(l)} \right\| \cdot \left\| (\Delta_{\mathbf{X}}^{t,(l)})_{l,\cdot} \right\|_2$$

$$+ \eta \|\mathbf{X}_{l,\cdot}^*\|_2 \cdot \|\Delta_{\mathbf{Y}}^{t,(l)}\| \cdot \|\mathbf{Y}^{t,(l)} \mathbf{H}^{t,(l)}\| + \eta \frac{\lambda}{p} \|\mathbf{X}^{t,(l)} \mathbf{H}^{t,(l)}\|_{2,\infty}.$$

In view of Lemma 18, we have

$$\begin{aligned} \sigma_{\min}/2 &\leq \sigma_{\min} \left[ (\mathbf{Y}^{t,(l)} \mathbf{H}^{t,(l)})^\top \mathbf{Y}^{t,(l)} \mathbf{H}^{t,(l)} \right] \leq \sigma_{\max} \left[ (\mathbf{Y}^{t,(l)} \mathbf{H}^{t,(l)})^\top \mathbf{Y}^{t,(l)} \mathbf{H}^{t,(l)} \right] \leq 2\sigma_{\max}, \\ \|\Delta_{\mathbf{Y}}^{t,(l)}\| &\leq \|\Delta^{t,(l)}\| \leq 2C_{\text{op}} \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}} \right) \|\mathbf{X}^*\|, \\ \|\mathbf{Y}^{t,(l)} \mathbf{H}^{t,(l)}\| &\leq \|\mathbf{F}^{t,(l)}\| \leq 2\|\mathbf{X}^*\| \quad \text{and} \\ \|\mathbf{X}^{t,(l)} \mathbf{H}^{t,(l)}\|_{2,\infty} &\leq \|\mathbf{F}^{t,(l)}\|_{2,\infty} \leq 2\|\mathbf{F}^*\|_{2,\infty}, \end{aligned} \quad (123)$$

provided that the sample size obeys  $n \gg \kappa\mu$  and that the noise satisfies  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll \frac{1}{\sqrt{\kappa^2 \log n}}$ . These allow us to further upper bound  $\|\mathbf{a}_1\|_2$  by

$$\|\mathbf{a}_1\|_2 \leq \left(1 - \frac{\eta\sigma_{\min}}{2}\right) \|(\Delta_{\mathbf{X}}^{t,(l)})_{l,\cdot}\|_2 + 4\eta C_{\text{op}} \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}} \right) \sigma_{\max} \|\mathbf{F}^*\|_{2,\infty} + 2\eta \frac{\lambda}{p} \|\mathbf{F}^*\|_{2,\infty},$$

as long as  $\eta \leq 1/(2\sigma_{\max})$ . As an immediate consequence,

$$\begin{aligned} \|\mathbf{a}_1\|_2 &\leq C_4 \kappa \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p\sigma_{\min}} \right) \|\mathbf{F}^*\|_{2,\infty} + 4\eta C_{\text{op}} \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}} \right) \sigma_{\max} \|\mathbf{F}^*\|_{2,\infty} + \frac{2\eta\lambda}{p} \|\mathbf{F}^*\|_{2,\infty} \\ &\leq \|\mathbf{F}^*\|_{2,\infty}, \end{aligned} \quad (124)$$

where the first inequality follows from the induction hypothesis (93d) and the last one holds as long as  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll \frac{1}{\sqrt{\kappa^2 \log n}}$  and  $\eta \ll 1/\sigma_{\max}$ .

2. Next, we turn attention to  $\|\mathbf{a}_2\|_2$ , which satisfies

$$\|\mathbf{a}_2\|_2 = \left\| (\mathbf{a}_1 + \mathbf{X}_{l,\cdot}^*) \left[ (\mathbf{H}^{t,(l)})^{-1} \mathbf{H}^{t+1,(l)} - \mathbf{I}_r \right] \right\|_2 \leq \left\| (\mathbf{H}^{t,(l)})^{-1} \mathbf{H}^{t+1,(l)} - \mathbf{I}_r \right\| \|\mathbf{a}_1 + \mathbf{X}_{l,\cdot}^*\|_2.$$

From (124), it is easily seen that

$$\|\mathbf{a}_1 + \mathbf{X}_{l,\cdot}^*\|_2 \leq \|\mathbf{a}_1\|_2 + \|\mathbf{F}^*\|_{2,\infty} \leq 2\|\mathbf{F}^*\|_{2,\infty}.$$

Regarding the term  $\|(\mathbf{H}^{t,(l)})^{-1} \mathbf{H}^{t+1,(l)} - \mathbf{I}_r\|$ , we find the following claim useful.

**Claim 8.** *With probability at least  $1 - O(n^{-100})$ , we have*

$$\left\| (\mathbf{H}^{t,(l)})^{-1} \mathbf{H}^{t+1,(l)} - \mathbf{I}_r \right\| \lesssim \eta \kappa C_{\text{op}}^2 \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}} \right)^2 \sigma_{\max} + \eta^2 C_{\text{B}} \kappa^2 \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}} \right) \sqrt{r} \sigma_{\max}^2,$$

provided that  $C_{\text{op}} \gg 1$ .

Finally, taking the bounds on  $\|\mathbf{a}_1\|_2$  and  $\|\mathbf{a}_2\|_2$  collectively yields that: for some absolute constant  $\tilde{C} > 0$ ,

$$\begin{aligned} \left\| (\mathbf{F}^{t+1,(l)} \mathbf{H}^{t+1,(l)} - \mathbf{F}^*)_{l,\cdot} \right\|_2 &\leq \|\mathbf{a}_1\|_2 + \|\mathbf{a}_2\|_2 \\ &\leq \left(1 - \frac{\eta}{2}\sigma_{\min}\right) C_4 \kappa \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p\sigma_{\min}} \right) \|\mathbf{F}^*\|_{2,\infty} + 4C_{\text{op}} \eta \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}} \right) \sigma_{\max} \|\mathbf{F}^*\|_{2,\infty} + \frac{2\eta\lambda}{p} \|\mathbf{F}^*\|_{2,\infty} \\ &\quad + \tilde{C} \eta \kappa C_{\text{op}}^2 \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}} \right)^2 \sigma_{\max} \|\mathbf{F}^*\|_{2,\infty} + \tilde{C} \eta^2 C_{\text{B}} \kappa^2 \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}} \right) \sqrt{r} \sigma_{\max}^2 \|\mathbf{F}^*\|_{2,\infty} \\ &\leq C_4 \kappa \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p\sigma_{\min}} \right) \|\mathbf{F}^*\|_{2,\infty}, \end{aligned}$$

provided that  $C_4 \gg C_{\text{op}}$ ,  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll 1/\kappa$  and  $\eta \ll 1/(\kappa^2 \sqrt{r} \sigma_{\max})$ . This finishes the proof of the lemma. It remains to establish Claim 8.

*Proof of Claim 8.* To facilitate analysis, we introduce an auxiliary point  $\tilde{\mathbf{F}}^{t+1} \triangleq \begin{bmatrix} \tilde{\mathbf{X}}^{t+1,(l)} \\ \tilde{\mathbf{Y}}^{t+1,(l)} \end{bmatrix}$  where

$$\begin{aligned} \tilde{\mathbf{X}}^{t+1,(l)} &= \mathbf{X}^{t,(l)} \mathbf{H}^{t,(l)} - \eta \left[ \frac{1}{p} \mathcal{P}_{\Omega_{-l,\cdot}} \left( \mathbf{X}^{t,(l)} \mathbf{Y}^{t,(l)\top} - \mathbf{M}^* - \mathbf{E} \right) + \mathcal{P}_{l,\cdot} \left( \mathbf{X}^{t,(l)} \mathbf{Y}^{t,(l)\top} - \mathbf{M}^* \right) \right] \mathbf{Y}^* \\ &\quad - \eta \frac{\lambda}{p} \mathbf{X}^* - \frac{\eta}{2} \mathbf{X}^* \mathbf{H}^{t,(l)\top} \left( \mathbf{X}^{t,(l)\top} \mathbf{X}^{t,(l)} - \mathbf{Y}^{t,(l)\top} \mathbf{Y}^{t,(l)} \right) \mathbf{H}^{t,(l)}, \\ \tilde{\mathbf{Y}}^{t+1,(l)} &= \mathbf{Y}^{t,(l)} \mathbf{H}^{t,(l)} - \eta \left[ \frac{1}{p} \mathcal{P}_{\Omega_{-l,\cdot}} \left( \mathbf{X}^{t,(l)} \mathbf{Y}^{t,(l)\top} - \mathbf{M}^* - \mathbf{E} \right) + \mathcal{P}_{l,\cdot} \left( \mathbf{X}^{t,(l)} \mathbf{Y}^{t,(l)\top} - \mathbf{M}^* \right) \right]^\top \mathbf{X}^* \\ &\quad - \eta \frac{\lambda}{p} \mathbf{Y}^* - \frac{\eta}{2} \mathbf{Y}^* \mathbf{H}^{t,(l)\top} \left( \mathbf{Y}^{t,(l)\top} \mathbf{Y}^{t,(l)} - \mathbf{X}^{t,(l)\top} \mathbf{X}^{t,(l)} \right) \mathbf{H}^{t,(l)}. \end{aligned}$$

We first claim that  $\mathbf{I}_r$  is the best rotation matrix to align  $\tilde{\mathbf{F}}^{t+1,(l)}$  and  $\mathbf{F}^*$ ; its proof is similar to that of Claim 4 and is hence omitted for brevity.

**Claim 9.** *One has*

$$\mathbf{I}_r = \arg \min_{\mathbf{R} \in \mathcal{O}^r} \left\| \tilde{\mathbf{F}}^{t+1,(l)} \mathbf{R} - \mathbf{F}^* \right\|_{\mathbb{F}} \quad \text{and} \quad \sigma_{\min} \left( \tilde{\mathbf{F}}^{t+1,(l)\top} \mathbf{F}^* \right) \geq \sigma_{\min} / 2.$$

With this claim at hand, we intend to invoke Lemma 23 with

$$\mathbf{S} = \tilde{\mathbf{F}}^{t+1,(l)\top} \mathbf{F}^*, \quad \mathbf{K} = \left( \mathbf{F}^{t+1,(l)} \mathbf{H}^{t,(l)} - \tilde{\mathbf{F}}^{t+1,(l)} \right)^\top \mathbf{F}^*$$

to get

$$\begin{aligned} \left\| \left( \mathbf{H}^{t,(l)} \right)^{-1} \mathbf{H}^{t+1,(l)} - \mathbf{I}_r \right\| &= \left\| \text{sgn}(\mathbf{S} + \mathbf{K}) - \text{sgn}(\mathbf{S}) \right\| \leq \frac{1}{\sigma_{\min}(\mathbf{S})} \|\mathbf{K}\| \\ &= \frac{1}{\sigma_{\min}(\tilde{\mathbf{F}}^{t+1,(l)\top} \mathbf{F}^*)} \left\| \left( \mathbf{F}^{t+1,(l)} \mathbf{H}^{t,(l)} - \tilde{\mathbf{F}}^{t+1,(l)} \right)^\top \mathbf{F}^* \right\| \\ &\leq \frac{2}{\sigma_{\min}} \left\| \mathbf{F}^{t+1,(l)} \mathbf{H}^{t,(l)} - \tilde{\mathbf{F}}^{t+1,(l)} \right\| \|\mathbf{F}^*\|, \end{aligned} \quad (125)$$

where the last line uses Claim 9. Here  $\text{sgn}(\mathbf{A}) = \mathbf{U}\mathbf{V}^\top$  for any matrix  $\mathbf{A}$  with SVD  $\mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$ . It then boils down to controlling  $\left\| \mathbf{F}^{t+1,(l)} \mathbf{H}^{t,(l)} - \tilde{\mathbf{F}}^{t+1,(l)} \right\|$ , for which we have

$$\mathbf{F}^{t+1,(l)} \mathbf{H}^{t,(l)} - \tilde{\mathbf{F}}^{t+1,(l)} = \eta \begin{bmatrix} \mathbf{B} & \mathbf{0} \\ \mathbf{0} & \mathbf{B}^\top \end{bmatrix} \begin{bmatrix} \Delta_{\mathbf{Y}}^{t,(l)} \\ \Delta_{\mathbf{X}}^{t,(l)} \end{bmatrix} + \frac{\eta}{2} \begin{bmatrix} \mathbf{X}^* \\ -\mathbf{Y}^* \end{bmatrix} \mathbf{H}^{t,(l)\top} \mathbf{C} \mathbf{H}^{t,(l)} - \eta \frac{\lambda}{p} \Delta^{t,(l)},$$

where we denote

$$\begin{aligned} \mathbf{B} &\triangleq -p^{-1} \mathcal{P}_{\Omega_{-l,\cdot}} \left( \mathbf{X}^{t,(l)} \mathbf{Y}^{t,(l)\top} - \mathbf{M}^* - \mathbf{E} \right) - \mathcal{P}_{l,\cdot} \left( \mathbf{X}^{t,(l)} \mathbf{Y}^{t,(l)\top} - \mathbf{M}^* \right); \\ \mathbf{C} &\triangleq \mathbf{X}^{t,(l)\top} \mathbf{X}^{t,(l)} - \mathbf{Y}^{t,(l)\top} \mathbf{Y}^{t,(l)}. \end{aligned}$$

This enables us to obtain

$$\left\| \mathbf{F}^{t+1,(l)} \mathbf{H}^{t,(l)} - \tilde{\mathbf{F}}^{t+1,(l)} \right\| \leq \eta \|\mathbf{B}\| \|\Delta^{t,(l)}\| + \frac{\eta}{2} \|\mathbf{F}^*\| \|\mathbf{C}\|_{\mathbb{F}} + \frac{\eta\lambda}{p} \|\Delta^{t,(l)}\|. \quad (126)$$

In view of Lemma 15, one has

$$\|\mathbf{C}\|_{\mathbb{F}} \leq C_{\mathbf{B}\kappa} \eta \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \sqrt{r} \sigma_{\max}^2. \quad (127)$$

We are left with bounding  $\|\mathbf{B}\|$ . Decompose  $\mathbf{B}$  into

$$\mathbf{B} = \underbrace{-\frac{1}{p} \mathcal{P}_{\Omega} \left( \mathbf{X}^{t,(l)} \mathbf{Y}^{t,(l)\top} - \mathbf{M}^* \right)}_{:=\mathbf{B}_1} + \underbrace{\frac{1}{p} \mathcal{P}_{\Omega_{l,\cdot}} \left( \mathbf{X}^{t,(l)} \mathbf{Y}^{t,(l)\top} - \mathbf{M}^* \right) - \mathcal{P}_{l,\cdot} \left( \mathbf{X}^{t,(l)} \mathbf{Y}^{t,(l)\top} - \mathbf{M}^* \right)}_{:=\mathbf{B}_2} + \underbrace{\frac{1}{p} \mathcal{P}_{\Omega_{-l,\cdot}} \left( \mathbf{E} \right)}_{:=\mathbf{B}_3}.$$

To control  $\mathbf{B}_1$ , following the same argument in Lemma 8, we see that

$$\begin{aligned}\|\mathbf{B}_1\| &\leq \|p^{-1}\mathcal{P}_\Omega(\mathbf{X}^{t,(l)}\mathbf{Y}^{t,(l)\top} - \mathbf{M}^*) - (\mathbf{X}^{t,(l)}\mathbf{Y}^{t,(l)\top} - \mathbf{M}^*)\| + \|\mathbf{X}^{t,(l)}\mathbf{Y}^{t,(l)\top} - \mathbf{M}^*\| \\ &\lesssim \sqrt{n/p}\|\mathbf{F}^{t,(l)}\mathbf{R}^{t,(l)} - \mathbf{F}^*\|_{2,\infty}\|\mathbf{F}^*\|_{2,\infty} + \|\mathbf{F}^{t,(l)}\mathbf{R}^{t,(l)} - \mathbf{F}^*\|\|\mathbf{F}^*\|.\end{aligned}$$

We now move on to  $\|\mathbf{B}_2\|$ , which is equal to  $\|\mathbf{b}\|_2/p$  defined in the proof of Claim 6, namely,

$$\|\mathbf{B}_2\| = \|\mathbf{b}\|_2/p \lesssim \sqrt{n \log n/p}\|\mathbf{F}^{t,(l)}\mathbf{R}^{t,(l)} - \mathbf{F}^*\|_{2,\infty}\|\mathbf{F}^*\|_{2,\infty}.$$

The last term  $\mathbf{B}_3$  can be easily bound via Lemma 3, that is,

$$\|\mathbf{B}_3\| \leq p^{-1}\|\mathcal{P}_\Omega(\mathbf{E})\| \lesssim \sigma\sqrt{n/p}.$$

Combining the above three bounds with Lemma 18, we arrive at

$$\begin{aligned}\|\mathbf{B}\| &\leq \|\mathbf{B}_1\| + \|\mathbf{B}_2\| + \|\mathbf{B}_3\| \\ &\lesssim \sqrt{\frac{n \log n}{p}}\|\mathbf{F}^{t,(l)}\mathbf{R}^{t,(l)} - \mathbf{F}^*\|_{2,\infty}\|\mathbf{F}^*\|_{2,\infty} + \|\mathbf{F}^{t,(l)}\mathbf{R}^{t,(l)} - \mathbf{F}^*\|\|\mathbf{F}^*\| + \sigma\sqrt{\frac{n}{p}} \\ &\lesssim \sqrt{\frac{n \log n}{p}}(C_\infty\kappa + C_3)\left(\frac{\sigma}{\sigma_{\min}}\sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p\sigma_{\min}}\right)\frac{\mu r}{n}\sigma_{\max} + 2C_{\text{op}}\left(\frac{\sigma}{\sigma_{\min}}\sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}}\right)\sigma_{\max} + \sigma\sqrt{\frac{n}{p}} \\ &\lesssim C_{\text{op}}\left(\frac{\sigma}{\sigma_{\min}}\sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}}\right)\sigma_{\max},\end{aligned}\tag{128}$$

provided that  $n^2p \gg \kappa^2\mu^2r^2n \log^2 n$  and  $C_{\text{op}} > 0$  is large enough. Taking (126), (127) and (128) collectively, we arrive at

$$\begin{aligned}\|\mathbf{F}^{t+1,(l)}\mathbf{H}^{t,(l)} - \tilde{\mathbf{F}}^{t+1,(l)}\| &\lesssim \eta C_{\text{op}}^2\left(\frac{\sigma}{\sigma_{\min}}\sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}}\right)^2\sigma_{\max}\|\mathbf{X}^*\| + \eta^2 C_{\text{B}}\kappa\left(\frac{\sigma}{\sigma_{\min}}\sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}}\right)\sqrt{r}\sigma_{\max}^2\|\mathbf{X}^*\| \\ &\quad + \eta\frac{\lambda}{p}C_{\text{op}}\left(\frac{\sigma}{\sigma_{\min}}\sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}}\right)\|\mathbf{X}^*\| \\ &\lesssim \eta C_{\text{op}}^2\left(\frac{\sigma}{\sigma_{\min}}\sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}}\right)^2\sigma_{\max}\|\mathbf{X}^*\| + \eta^2 C_{\text{B}}\kappa\left(\frac{\sigma}{\sigma_{\min}}\sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}}\right)\sqrt{r}\sigma_{\max}^2\|\mathbf{X}^*\|,\end{aligned}$$

provided that  $C_{\text{op}}$  is large enough. Here the last relation uses (123). Substitution into (125) yields

$$\|(\mathbf{H}^{t,(l)})^{-1}\mathbf{H}^{t+1,(l)} - \mathbf{I}_r\| \lesssim \eta\kappa C_{\text{op}}^2\left(\frac{\sigma}{\sigma_{\min}}\sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}}\right)^2\sigma_{\max} + \eta^2 C_{\text{B}}\kappa^2\left(\frac{\sigma}{\sigma_{\min}}\sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}}\right)\sqrt{r}\sigma_{\max}^2,$$

which concludes the proof.  $\square$

## D.7 Proof of Lemma 14

Fix any  $1 \leq l \leq 2n$ . Apply the triangle inequality to see that

$$\begin{aligned}\|(\mathbf{F}^{t+1}\mathbf{H}^{t+1} - \mathbf{F}^*)_{l,\cdot}\|_2 &\leq \|(\mathbf{F}^{t+1}\mathbf{H}^{t+1} - \mathbf{F}^{t+1,(l)}\mathbf{H}^{t+1,(l)})_{l,\cdot}\|_2 + \|(\mathbf{F}^{t+1,(l)}\mathbf{H}^{t+1,(l)} - \mathbf{F}^*)_{l,\cdot}\|_2 \\ &\leq \|\mathbf{F}^{t+1}\mathbf{H}^{t+1} - \mathbf{F}^{t+1,(l)}\mathbf{H}^{t+1,(l)}\|_{\text{F}} + C_4\kappa\left(\frac{\sigma}{\sigma_{\min}}\sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p\sigma_{\min}}\right)\|\mathbf{F}^*\|_{2,\infty},\end{aligned}\tag{129}$$

where the second line follows from Lemma 13. Apply Lemma 18 to the  $(t+1)$ th iterates to see that

$$\|\mathbf{F}^{t+1}\mathbf{H}^{t+1} - \mathbf{F}^{t+1,(l)}\mathbf{H}^{t+1,(l)}\|_{\text{F}} \leq 5\kappa\|\mathbf{F}^{t+1}\mathbf{H}^{t+1} - \mathbf{F}^{t+1,(l)}\mathbf{R}^{t+1,(l)}\|_{\text{F}}$$

$$\leq 5\kappa C_3 \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{F}^*\|_{2,\infty}. \quad (130)$$

Here the second line follows from Lemma 12. Combine (129) and (130) to reach

$$\begin{aligned} \|(\mathbf{F}^{t+1} \mathbf{H}^{t+1} - \mathbf{F}^*)_{l,\cdot}\|_2 &\leq 5\kappa C_3 \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{F}^*\|_{2,\infty} + C_4 \kappa \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{F}^*\|_{2,\infty} \\ &\leq C_\infty \kappa \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{F}^*\|_{2,\infty} \end{aligned}$$

as long as  $C_\infty \geq 5C_3 + C_4$ . The proof is then complete since this holds for all  $1 \leq l \leq 2n$ .

## D.8 Proof of Lemma 15

To simplify the notation hereafter, we denote

$$\mathbf{A}^t \triangleq \mathbf{X}^{t\top} \mathbf{X}^t - \mathbf{Y}^{t\top} \mathbf{Y}^t \quad \text{and} \quad \mathbf{A}^{t+1} \triangleq \mathbf{X}^{t+1\top} \mathbf{X}^{t+1} - \mathbf{Y}^{t+1\top} \mathbf{Y}^{t+1}.$$

In view of the gradient descent update rules (27), we have

$$\begin{aligned} \mathbf{X}^{t+1\top} \mathbf{X}^{t+1} &= \mathbf{X}^{t\top} \mathbf{X}^t - \eta [\mathbf{X}^{t\top} \nabla_{\mathbf{X}} f(\mathbf{X}^t, \mathbf{Y}^t) + \nabla_{\mathbf{X}} f(\mathbf{X}^t, \mathbf{Y}^t)^\top \mathbf{X}^t] + \eta^2 \nabla_{\mathbf{X}} f(\mathbf{X}^t, \mathbf{Y}^t)^\top \nabla_{\mathbf{X}} f(\mathbf{X}^t, \mathbf{Y}^t), \\ \mathbf{Y}^{t+1\top} \mathbf{Y}^{t+1} &= \mathbf{Y}^{t\top} \mathbf{Y}^t - \eta [\mathbf{Y}^{t\top} \nabla_{\mathbf{Y}} f(\mathbf{X}^t, \mathbf{Y}^t) + \nabla_{\mathbf{Y}} f(\mathbf{X}^t, \mathbf{Y}^t)^\top \mathbf{Y}^t] + \eta^2 \nabla_{\mathbf{Y}} f(\mathbf{X}^t, \mathbf{Y}^t)^\top \nabla_{\mathbf{Y}} f(\mathbf{X}^t, \mathbf{Y}^t). \end{aligned}$$

This gives rise to the following identity

$$\mathbf{A}^{t+1} = \mathbf{A}^t - \eta \mathbf{B}^t + \eta^2 \mathbf{C}^t, \quad (131)$$

where we denote

$$\begin{aligned} \mathbf{B}^t &\triangleq \mathbf{X}^{t\top} \nabla_{\mathbf{X}} f(\mathbf{X}^t, \mathbf{Y}^t) + \nabla_{\mathbf{X}} f(\mathbf{X}^t, \mathbf{Y}^t)^\top \mathbf{X}^t - \mathbf{Y}^{t\top} \nabla_{\mathbf{Y}} f(\mathbf{X}^t, \mathbf{Y}^t) - \nabla_{\mathbf{Y}} f(\mathbf{X}^t, \mathbf{Y}^t)^\top \mathbf{Y}^t, \\ \mathbf{C}^t &\triangleq \nabla_{\mathbf{X}} f(\mathbf{X}^t, \mathbf{Y}^t)^\top \nabla_{\mathbf{X}} f(\mathbf{X}^t, \mathbf{Y}^t) - \nabla_{\mathbf{Y}} f(\mathbf{X}^t, \mathbf{Y}^t)^\top \nabla_{\mathbf{Y}} f(\mathbf{X}^t, \mathbf{Y}^t). \end{aligned}$$

Denoting

$$\mathbf{D}^t \triangleq p^{-1} \mathcal{P}_\Omega(\mathbf{X}^t \mathbf{Y}^{t\top} - \mathbf{M}), \quad (132)$$

we have

$$\nabla_{\mathbf{X}} f(\mathbf{X}^t, \mathbf{Y}^t) = \mathbf{D}^t \mathbf{Y}^t + \frac{\lambda}{p} \mathbf{X}^t \quad \text{and} \quad \nabla_{\mathbf{Y}} f(\mathbf{X}^t, \mathbf{Y}^t) = \mathbf{D}^{t\top} \mathbf{X}^t + \frac{\lambda}{p} \mathbf{Y}^t.$$

With these in mind, a little calculation reveals that

$$\mathbf{B}^t = \mathbf{X}^{t\top} \mathbf{D}^t \mathbf{Y}^t + \mathbf{Y}^{t\top} \mathbf{D}^{t\top} \mathbf{X}^t + \frac{2\lambda}{p} \mathbf{X}^{t\top} \mathbf{X}^t - \mathbf{Y}^{t\top} \mathbf{D}^{t\top} \mathbf{X}^t - \mathbf{X}^{t\top} \mathbf{D}^t \mathbf{Y}^t - \frac{2\lambda}{p} \mathbf{Y}^{t\top} \mathbf{Y}^t = \frac{2\lambda}{p} \mathbf{A}^t$$

as well as

$$\begin{aligned} \mathbf{C}^t &= \left( \mathbf{D}^t \mathbf{Y}^t + \frac{\lambda}{p} \mathbf{X}^t \right)^\top \left( \mathbf{D}^t \mathbf{Y}^t + \frac{\lambda}{p} \mathbf{X}^t \right) - \left( \mathbf{D}^{t\top} \mathbf{X}^t + \frac{\lambda}{p} \mathbf{Y}^t \right)^\top \left( \mathbf{D}^{t\top} \mathbf{X}^t + \frac{\lambda}{p} \mathbf{Y}^t \right) \\ &= \mathbf{Y}^{t\top} \mathbf{D}^{t\top} \mathbf{D}^t \mathbf{Y}^t + \frac{\lambda}{p} \mathbf{Y}^{t\top} \mathbf{D}^{t\top} \mathbf{X}^t + \frac{\lambda}{p} \mathbf{X}^{t\top} \mathbf{D}^t \mathbf{Y}^t + \left( \frac{\lambda}{p} \right)^2 \mathbf{X}^{t\top} \mathbf{X}^t \\ &\quad - \mathbf{X}^{t\top} \mathbf{D}^t \mathbf{D}^{t\top} \mathbf{X}^t - \frac{\lambda}{p} \mathbf{X}^{t\top} \mathbf{D}^t \mathbf{Y}^t - \frac{\lambda}{p} \mathbf{Y}^{t\top} \mathbf{D}^{t\top} \mathbf{X}^t - \left( \frac{\lambda}{p} \right)^2 \mathbf{Y}^{t\top} \mathbf{Y}^t \\ &= (\mathbf{Y}^{t\top} \mathbf{D}^{t\top} \mathbf{D}^t \mathbf{Y}^t - \mathbf{X}^{t\top} \mathbf{D}^t \mathbf{D}^{t\top} \mathbf{X}^t) + \left( \frac{\lambda}{p} \right)^2 \mathbf{A}^t. \end{aligned}$$

Substituting the identities for  $\mathbf{B}^t$  and  $\mathbf{C}^t$  into (131) yields

$$\mathbf{A}^{t+1} = \mathbf{A}^t - 2\eta \frac{\lambda}{p} \mathbf{A}^t + \eta^2 (\mathbf{Y}^{t\top} \mathbf{D}^{t\top} \mathbf{D}^t \mathbf{Y}^t - \mathbf{X}^{t\top} \mathbf{D}^t \mathbf{D}^{t\top} \mathbf{X}^t) + \eta^2 \left( \frac{\lambda}{p} \right)^2 \mathbf{A}^t$$



$$= (1 - \lambda\eta/p)^2 \mathbf{A}^t + \eta^2 (\mathbf{Y}^{t\top} \mathbf{D}^{t\top} \mathbf{D}^t \mathbf{Y}^t - \mathbf{X}^{t\top} \mathbf{D}^t \mathbf{D}^{t\top} \mathbf{X}^t),$$

which together with the triangle inequality gives

$$\begin{aligned} \|\mathbf{A}^{t+1}\|_{\mathbb{F}} &\leq (1 - \lambda\eta/p)^2 \|\mathbf{A}^t\|_{\mathbb{F}} + \eta^2 \|\mathbf{Y}^{t\top} \mathbf{D}^{t\top} \mathbf{D}^t \mathbf{Y}^t - \mathbf{X}^{t\top} \mathbf{D}^t \mathbf{D}^{t\top} \mathbf{X}^t\|_{\mathbb{F}} \\ &\leq (1 - \lambda\eta/p) \|\mathbf{A}^t\|_{\mathbb{F}} + \eta^2 \|\mathbf{Y}^{t\top} \mathbf{D}^{t\top} \mathbf{D}^t \mathbf{Y}^t - \mathbf{X}^{t\top} \mathbf{D}^t \mathbf{D}^{t\top} \mathbf{X}^t\|_{\mathbb{F}}, \end{aligned}$$

as long as  $\lambda\eta/p < 1$  — a condition that is guaranteed by our assumptions on  $\lambda$  and  $\eta$ . It then boils down to controlling  $\|\mathbf{Y}^{t\top} \mathbf{D}^{t\top} \mathbf{D}^t \mathbf{Y}^t - \mathbf{X}^{t\top} \mathbf{D}^t \mathbf{D}^{t\top} \mathbf{X}^t\|_{\mathbb{F}}$ , which is supplied in the following claim.

**Claim 10.** *Suppose that the sample complexity satisfies  $n^2 p \gg \kappa^2 \mu^2 r^2 n \log n$  and that the noise satisfies  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll \frac{1}{\sqrt{\kappa^2 \log n}}$ , then one has*

$$\|\mathbf{Y}^{t\top} \mathbf{D}^{t\top} \mathbf{D}^t \mathbf{Y}^t - \mathbf{X}^{t\top} \mathbf{D}^t \mathbf{D}^{t\top} \mathbf{X}^t\|_{\mathbb{F}} \lesssim C_{\text{op}}^2 \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right)^2 \sqrt{r} \sigma_{\max}^3.$$

Taking the above bounds together, we arrive at for some constant  $\tilde{C} > 0$ ,

$$\begin{aligned} \|\mathbf{A}^{t+1}\|_{\mathbb{F}} &\leq \left(1 - \frac{\lambda}{p} \eta\right) \|\mathbf{A}^t\|_{\mathbb{F}} + \eta^2 \tilde{C} C_{\text{op}}^2 \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right)^2 \sqrt{r} \sigma_{\max}^3 \\ &\leq \left(1 - \frac{\lambda}{p} \eta\right) C_{\text{B}} \kappa \eta \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \sqrt{r} \sigma_{\max}^2 + \eta^2 \tilde{C} C_{\text{op}}^2 \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right)^2 \sqrt{r} \sigma_{\max}^3 \\ &\leq C_{\text{B}} \kappa \eta \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \sqrt{r} \sigma_{\max}^2, \end{aligned}$$

as long as  $\lambda \geq \sigma \sqrt{np}$  and  $C_{\text{B}} \gg C_{\text{op}}^2$ .

*Proof of Claim 10.* The triangle inequality yields

$$\begin{aligned} \|\mathbf{Y}^{t\top} \mathbf{D}^{t\top} \mathbf{D}^t \mathbf{Y}^t - \mathbf{X}^{t\top} \mathbf{D}^t \mathbf{D}^{t\top} \mathbf{X}^t\|_{\mathbb{F}} &\leq \|\mathbf{Y}^{t\top} \mathbf{D}^{t\top} \mathbf{D}^t \mathbf{Y}^t\|_{\mathbb{F}} + \|\mathbf{X}^{t\top} \mathbf{D}^t \mathbf{D}^{t\top} \mathbf{X}^t\|_{\mathbb{F}} \\ &\leq \|\mathbf{Y}^t\| \|\mathbf{D}^t\|^2 \|\mathbf{Y}^t\|_{\mathbb{F}} + \|\mathbf{X}^t\| \|\mathbf{D}^t\|^2 \|\mathbf{X}^t\|_{\mathbb{F}}. \end{aligned} \quad (133)$$

It is easy to see from Lemma 18 that

$$\|\mathbf{Y}^t\| \leq 2 \|\mathbf{Y}^*\|, \quad \|\mathbf{Y}^t\|_{\mathbb{F}} \leq 2 \|\mathbf{Y}^*\|_{\mathbb{F}}, \quad \|\mathbf{X}^t\| \leq 2 \|\mathbf{X}^*\| \quad \text{and} \quad \|\mathbf{X}^t\|_{\mathbb{F}} \leq 2 \|\mathbf{X}^*\|_{\mathbb{F}}$$

provided that  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll \frac{1}{\sqrt{\kappa^2 \log n}}$ . These allow us to further upper bound (133) as

$$\begin{aligned} \|\mathbf{Y}^{t\top} \mathbf{D}^{t\top} \mathbf{D}^t \mathbf{Y}^t - \mathbf{X}^{t\top} \mathbf{D}^t \mathbf{D}^{t\top} \mathbf{X}^t\|_{\mathbb{F}} &\leq 4 \|\mathbf{D}^t\|^2 \|\mathbf{Y}^*\| \|\mathbf{Y}^*\|_{\mathbb{F}} + 4 \|\mathbf{D}^t\|^2 \|\mathbf{X}^*\| \|\mathbf{X}^*\|_{\mathbb{F}} \\ &\leq 8 \|\mathbf{D}^t\|^2 \sqrt{r} \sigma_{\max}. \end{aligned} \quad (134)$$

It remains to bound  $\|\mathbf{D}^t\|$ . To this end, recall from (132) that

$$\|\mathbf{D}^t\| \leq p^{-1} \|\mathcal{P}_{\Omega}(\mathbf{E})\| + p^{-1} \|\mathcal{P}_{\Omega}^{\text{debias}}(\mathbf{X}^t \mathbf{Y}^{t\top} - \mathbf{M}^*)\| + \|\mathbf{X}^t \mathbf{Y}^{t\top} - \mathbf{M}^*\|.$$

In the sequel we shall bound these three terms sequentially. First, Lemma 3 tells us that  $\frac{1}{p} \|\mathcal{P}_{\Omega}(\mathbf{E})\| \lesssim \sigma \sqrt{\frac{n}{p}}$ . Next, repeating the arguments in the proof of Lemma 8 gives

$$\begin{aligned} \|\mathcal{P}_{\Omega}^{\text{debias}}(\mathbf{X}^t \mathbf{Y}^{t\top} - \mathbf{M}^*)\| &= \left\| \mathcal{P}_{\Omega}^{\text{debias}} \left[ \mathbf{X}^t \mathbf{H}^t (\mathbf{Y}^t \mathbf{H}^t)^{\top} - \mathbf{M}^* \right] \right\| \\ &\lesssim \sqrt{np} \left( \|\mathbf{X}^t \mathbf{H}^t - \mathbf{X}^*\|_{2,\infty} \|\mathbf{Y}^*\|_{2,\infty} + \|\mathbf{Y}^t \mathbf{H}^t - \mathbf{Y}^*\|_{2,\infty} \|\mathbf{X}^*\|_{2,\infty} \right), \end{aligned}$$

which together with the induction hypothesis (93e) yields

$$\begin{aligned} \frac{1}{p} \|\mathcal{P}_\Omega^{\text{debias}}(\mathbf{X}^t \mathbf{Y}^{t\top} - \mathbf{M}^*)\| &\lesssim \sqrt{\frac{n}{p}} C_\infty \kappa \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{F}^*\|_{2,\infty} \|\mathbf{F}^*\|_{2,\infty} \\ &\lesssim C_\infty \kappa \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \sqrt{\frac{\mu^2 r^2}{np}} \sigma_{\max}. \end{aligned}$$

Here the last relation uses the incoherence assumption  $\|\mathbf{F}^*\|_{2,\infty} \leq \sqrt{\mu r \sigma_{\max}/n}$  (cf. (100a)). Regarding  $\|\mathbf{X}^t \mathbf{Y}^{t\top} - \mathbf{M}^*\|$ , the triangle inequality reveals that

$$\begin{aligned} \|\mathbf{X}^t \mathbf{Y}^{t\top} - \mathbf{M}^*\| &= \|\mathbf{X}^t \mathbf{H}^t (\mathbf{Y}^t \mathbf{H}^t)^\top - \mathbf{M}^*\| \\ &\leq \|\mathbf{X}^t \mathbf{H}^t (\mathbf{Y}^t \mathbf{H}^t)^\top - \mathbf{X}^t \mathbf{H}^t \mathbf{Y}^{*\top}\| + \|\mathbf{X}^t \mathbf{H}^t \mathbf{Y}^{*\top} - \mathbf{X}^* \mathbf{Y}^{*\top}\| \\ &\leq \|\mathbf{X}^t \mathbf{H}^t\| \|\mathbf{Y}^t \mathbf{H}^t - \mathbf{Y}^*\| + \|\mathbf{X}^t \mathbf{H}^t - \mathbf{X}^*\| \|\mathbf{Y}^*\|. \end{aligned}$$

Combine the induction hypothesis (93b) and the fact that  $\|\mathbf{X}^t \mathbf{H}^t\| = \|\mathbf{X}^t\| \leq 2\|\mathbf{X}^*\|$  to reach

$$\|\mathbf{X}^t \mathbf{Y}^{t\top} - \mathbf{M}^*\| \lesssim C_{\text{op}} \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \sigma_{\max}.$$

Putting together the previous three bounds, we arrive at

$$\begin{aligned} \|\mathbf{D}^t\| &\lesssim \sigma \sqrt{\frac{n}{p}} + C_\infty \kappa \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \sqrt{\frac{\mu^2 r^2}{np}} \sigma_{\max} + C_{\text{op}} \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \sigma_{\max} \\ &\lesssim C_{\text{op}} \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \sigma_{\max} \end{aligned} \quad (135)$$

since  $np \gg \kappa^2 \mu^2 r^2 \log n$ . Putting (135) back to (134) leads to the claimed upper bound.

The upper bound on the leave-one-out sequences can be derived similarly. For brevity, we omit it.  $\square$

## D.9 Proof of Lemma 16

In light of the facts that  $f(\mathbf{F}\mathbf{R}) = f(\mathbf{F})$  and  $\nabla f(\mathbf{F}\mathbf{R}) = \nabla f(\mathbf{F})\mathbf{R}$  for any  $\mathbf{R} \in \mathcal{O}^{r \times r}$ , one has

$$\begin{aligned} f(\mathbf{F}^{t+1}) &= f(\mathbf{F}^{t+1} \mathbf{H}^t) = f([\mathbf{F}^t - \eta \nabla f(\mathbf{F}^t)] \mathbf{H}^t) \\ &= f(\mathbf{F}^t \mathbf{H}^t - \eta \nabla f(\mathbf{F}^t \mathbf{H}^t)) \\ &= f(\mathbf{F}^t \mathbf{H}^t) - \eta \langle \nabla f(\mathbf{F}^t \mathbf{H}^t), \nabla f(\mathbf{F}^t \mathbf{H}^t) \rangle + \frac{\eta^2}{2} \text{vec}(\nabla f(\mathbf{F}^t \mathbf{H}^t))^\top \nabla^2 f(\tilde{\mathbf{F}}) \text{vec}(\nabla f(\mathbf{F}^t \mathbf{H}^t)) \end{aligned}$$

for some  $\tilde{\mathbf{F}}$  which lies between  $\mathbf{F}^t \mathbf{H}^t$  and  $\mathbf{F}^t \mathbf{H}^t - \eta \nabla f(\mathbf{F}^t \mathbf{H}^t)$ . Suppose for the moment that

$$\|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*\|_{2,\infty} \leq \frac{1}{2000\kappa\sqrt{n}} \|\mathbf{X}^*\|, \quad (136a)$$

$$\|\mathbf{F}^t \mathbf{H}^t - \eta \nabla f(\mathbf{F}^t \mathbf{H}^t) - \mathbf{F}^*\|_{2,\infty} \leq \frac{1}{1000\kappa\sqrt{n}} \|\mathbf{X}^*\|. \quad (136b)$$

One can invoke Lemma 17 to obtain  $\|\nabla^2 f(\tilde{\mathbf{F}})\| \leq 10\sigma_{\max}$  and hence

$$\begin{aligned} f(\mathbf{F}^{t+1}) &\leq f(\mathbf{F}^t \mathbf{H}^t) - \eta \|\nabla f(\mathbf{F}^t \mathbf{H}^t)\|_{\mathbb{F}}^2 + 5\eta^2 \sigma_{\max} \|\nabla f(\mathbf{F}^t \mathbf{H}^t)\|_{\mathbb{F}}^2 \\ &= f(\mathbf{F}^t) - \eta \|\nabla f(\mathbf{F}^t)\|_{\mathbb{F}}^2 + 5\eta^2 \sigma_{\max} \|\nabla f(\mathbf{F}^t)\|_{\mathbb{F}}^2 \\ &\leq f(\mathbf{F}^t) - \frac{\eta}{2} \|\nabla f(\mathbf{F}^t)\|_{\mathbb{F}}^2. \end{aligned}$$

Here the equality uses again the facts that  $f(\mathbf{F}\mathbf{R}) = f(\mathbf{F})$  and  $\nabla f(\mathbf{F}\mathbf{R}) = \nabla f(\mathbf{F})\mathbf{R}$  for any  $\mathbf{R} \in \mathcal{O}^{r \times r}$  and the last inequality holds as long as  $\eta \leq \frac{1}{10\sigma_{\max}}$ . We are left with proving the aforementioned conditions (136). The first condition has been established in the proof of Lemma 9 and hence we concentrate on the second one, namely (136b). Apply the triangle inequality and the fundamental theorem of calculus [Lan93, Chapter XIII, Theorem 4.2] to obtain

$$\begin{aligned} \|\mathbf{F}^t \mathbf{H}^t - \eta \nabla f(\mathbf{F}^t \mathbf{H}^t) - \mathbf{F}^*\|_{2,\infty} &\leq \|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*\|_{2,\infty} + \eta \|\nabla f(\mathbf{F}^t \mathbf{H}^t) - \nabla f(\mathbf{F}^*)\|_{\mathbb{F}} + \eta \|\nabla f(\mathbf{F}^*)\|_{\mathbb{F}} \\ &\leq \|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*\|_{2,\infty} + \eta \left\| \int_0^1 \nabla^2 f(\mathbf{F}(\tau)) d\tau \text{vec}(\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*) \right\|_2 + \eta \|\nabla f(\mathbf{F}^*)\|_{\mathbb{F}}, \end{aligned}$$

where  $\mathbf{F}(\tau) \triangleq \mathbf{F}^* + \tau(\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*)$  for  $0 \leq \tau \leq 1$ . Following similar arguments in the proof of Lemma 10 and the proof of Lemma 9, one obtains

$$\begin{aligned} &\eta \left\| \int_0^1 \nabla^2 f(\mathbf{F}(\tau)) d\tau \text{vec}(\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*) \right\|_{\mathbb{F}} + \eta \|\nabla f(\mathbf{F}^*)\|_{\mathbb{F}} \\ &\leq \eta \cdot 10\sigma_{\max} \|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*\|_{\mathbb{F}} + \eta \frac{\lambda}{p} \sqrt{r\sigma_{\max}} \\ &\lesssim \eta \sigma_{\max} \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p\sigma_{\min}} \right) \sqrt{r} \|\mathbf{X}^*\| + \eta \sigma_{\min} \frac{\lambda}{p\sigma_{\min}} \sqrt{r} \|\mathbf{X}^*\| \leq \frac{1}{2000\kappa\sqrt{n}} \|\mathbf{X}^*\|. \end{aligned}$$

Here the middle inequality uses the induction hypothesis (93b) and the last relation holds true provided that  $\lambda \asymp \sigma\sqrt{np}$ ,  $\frac{\sigma}{\sigma_{\min}}\sqrt{\frac{n}{p}} \ll 1/\sqrt{r}$  and that  $\eta \ll 1/(\kappa n\sigma_{\max})$ . This proves the second condition and also the whole lemma.

## D.10 Proof of Lemma 17

We start by defining a new loss function

$$f_{\text{clean}}(\mathbf{X}, \mathbf{Y}) \triangleq \frac{1}{2p} \|\mathcal{P}_{\Omega}(\mathbf{X}\mathbf{Y}^{\top} - \mathbf{M}^*)\|_{\mathbb{F}}^2 + \frac{1}{8} \|\mathbf{X}^{\top} \mathbf{X} - \mathbf{Y}^{\top} \mathbf{Y}\|_{\mathbb{F}}^2;$$

compared with  $f_{\text{aug}}(\cdot, \cdot)$ , this new function  $f_{\text{clean}}(\cdot, \cdot)$  sets  $\lambda = 0$  and excludes the noise  $\mathbf{E}$  from consideration. It is straightforward to check that for any  $\Delta \in \mathbb{R}^{2n \times r}$ ,

$$\text{vec}(\Delta)^{\top} \nabla^2 f_{\text{aug}}(\mathbf{X}, \mathbf{Y}) \text{vec}(\Delta) = \text{vec}(\Delta)^{\top} \nabla^2 f_{\text{clean}}(\mathbf{X}, \mathbf{Y}) \text{vec}(\Delta) - \frac{2}{p} \langle \mathcal{P}_{\Omega}(\mathbf{E}), \Delta_{\mathbf{X}} \Delta_{\mathbf{Y}}^{\top} \rangle + \frac{\lambda}{p} \|\Delta\|_{\mathbb{F}}^2.$$

It has been proven in [CLL19, Lemma 3.2] that under the assumptions stated in Lemma 17, one has

$$\text{vec}(\Delta)^{\top} \nabla^2 f_{\text{clean}}(\mathbf{X}, \mathbf{Y}) \text{vec}(\Delta) \geq \frac{1}{5} \sigma_{\min} \|\Delta\|_{\mathbb{F}}^2 \quad \text{and} \quad \|\nabla^2 f_{\text{clean}}(\mathbf{X}, \mathbf{Y})\| \leq 5\sigma_{\max}. \quad (137)$$

It then boils down to controlling  $-\frac{2}{p} \langle \mathcal{P}_{\Omega}(\mathbf{E}), \Delta_{\mathbf{X}} \Delta_{\mathbf{Y}}^{\top} \rangle + \frac{\lambda}{p} \|\Delta\|_{\mathbb{F}}^2$ . To this end, one has

$$\left| \frac{1}{p} \langle \mathcal{P}_{\Omega}(\mathbf{E}), \Delta_{\mathbf{X}} \Delta_{\mathbf{Y}}^{\top} \rangle \right| \leq \left\| \frac{1}{p} \mathcal{P}_{\Omega}(\mathbf{E}) \right\| \|\Delta_{\mathbf{X}} \Delta_{\mathbf{Y}}^{\top}\|_* \lesssim \sigma \sqrt{\frac{n}{p}} \|\Delta\|_{\mathbb{F}}^2, \quad (138)$$

where the last relation holds due to Lemma 3 and the elementary fact about the nuclear norm (6), i.e.

$$2 \|\Delta_{\mathbf{X}} \Delta_{\mathbf{Y}}^{\top}\|_* \leq \|\Delta_{\mathbf{X}}\|_{\mathbb{F}}^2 + \|\Delta_{\mathbf{Y}}\|_{\mathbb{F}}^2 = \|\Delta\|_{\mathbb{F}}^2.$$

Regarding the term  $\lambda \|\Delta\|_{\mathbb{F}}^2/p$ , it is easy to see from the assumption  $\lambda \asymp \sigma\sqrt{np}$  that  $\frac{\lambda}{p} \|\Delta\|_{\mathbb{F}}^2 \asymp \sigma \sqrt{\frac{n}{p}} \|\Delta\|_{\mathbb{F}}^2$ . Combine the above two bounds and use the triangle inequality to reach

$$\left| -\frac{2}{p} \langle \mathcal{P}_{\Omega}(\mathbf{E}), \Delta_{\mathbf{X}} \Delta_{\mathbf{Y}}^{\top} \rangle + \frac{\lambda}{p} \|\Delta\|_{\mathbb{F}}^2 \right| \lesssim \sigma \sqrt{\frac{n}{p}} \|\Delta\|_{\mathbb{F}}^2 \leq \frac{1}{10} \sigma_{\min} \|\Delta\|_{\mathbb{F}}^2, \quad (139)$$

with the proviso that  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll 1$ . Taking (137) and (139) together immediately establishes the claims on  $\nabla^2 f_{\text{aug}}(\cdot, \cdot)$ .

Moving on to  $\nabla^2 f(\mathbf{X}, \mathbf{Y})$ , one has

$$\begin{aligned}
\text{vec}(\Delta)^\top \nabla^2 f(\mathbf{X}, \mathbf{Y}) \text{vec}(\Delta) &= \frac{1}{p} \|\mathcal{P}_\Omega(\mathbf{X}\Delta_Y^\top + \Delta_X \mathbf{Y}^\top)\|_F^2 + \frac{2}{p} \langle \mathcal{P}_\Omega(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}^* - \mathbf{E}), \Delta_X \Delta_Y^\top \rangle + \frac{\lambda}{p} \|\Delta\|_F^2 \\
&= \underbrace{\|\mathbf{X}\Delta_Y^\top + \Delta_X \mathbf{Y}^\top\|_F^2}_{:=\alpha_1} + 2 \underbrace{\langle \mathbf{X}\mathbf{Y}^\top - \mathbf{M}^*, \Delta_X \Delta_Y^\top \rangle}_{:=\alpha_2} + \underbrace{\left(-\frac{2}{p} \langle \mathcal{P}_\Omega(\mathbf{E}), \Delta_X \Delta_Y^\top \rangle + \frac{\lambda}{p} \|\Delta\|_F^2\right)}_{:=\alpha_3} \\
&\quad + \underbrace{\frac{1}{p} \|\mathcal{P}_\Omega(\mathbf{X}\Delta_Y^\top + \Delta_X \mathbf{Y}^\top)\|_F^2 + \frac{2}{p} \langle \mathcal{P}_\Omega(\mathbf{X}\mathbf{Y}^\top - \mathbf{M}^*), \Delta_X \Delta_Y^\top \rangle - \|\mathbf{X}\Delta_Y^\top + \Delta_X \mathbf{Y}^\top\|_F^2 - 2 \langle \mathbf{X}\mathbf{Y}^\top - \mathbf{M}^*, \Delta_X \Delta_Y^\top \rangle}_{:=\alpha_4}.
\end{aligned}$$

The term  $\alpha_4$  can be bounded by [CLL19, Equation A.4]

$$|\alpha_4| \leq \frac{1}{5} \sigma_{\min} (\|\Delta_X\|_F^2 + \|\Delta_Y\|_F^2) + \frac{1}{5} (\|\Delta_X \mathbf{Y}^{*\top}\|_F^2 + \|\mathbf{X}^* \Delta_Y^\top\|_F^2) \leq \frac{2}{5} \sigma_{\max} \|\Delta\|_F^2.$$

The term  $\alpha_3$  has been bounded in (139) where  $|\alpha_3| \leq \sigma_{\max} \|\Delta\|_F^2$  provided that  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll 1$ . The term  $\alpha_2$  can be written as

$$|\alpha_2| \leq 2 \|\mathbf{X}\mathbf{Y}^\top - \mathbf{M}^*\| \|\Delta_X \Delta_Y^\top\|_* \leq (\|\mathbf{X} - \mathbf{X}^*\| \|\mathbf{Y}\| + \|\mathbf{X}^*\| \|\mathbf{Y} - \mathbf{Y}^*\|) \|\Delta\|_F^2.$$

Since

$$\left\| \begin{bmatrix} \mathbf{X} - \mathbf{X}^* \\ \mathbf{Y} - \mathbf{Y}^* \end{bmatrix} \right\| \leq \left\| \begin{bmatrix} \mathbf{X} - \mathbf{X}^* \\ \mathbf{Y} - \mathbf{Y}^* \end{bmatrix} \right\|_F \leq \sqrt{2n} \left\| \begin{bmatrix} \mathbf{X} - \mathbf{X}^* \\ \mathbf{Y} - \mathbf{Y}^* \end{bmatrix} \right\|_{2,\infty} \leq \frac{1}{500\kappa} \|\mathbf{X}^*\|,$$

we immediately have

$$|\alpha_2| \leq \frac{3}{500\kappa} \sigma_{\max} \|\Delta\|_F^2 \leq \frac{1}{2} \sigma_{\max} \|\Delta\|_F^2.$$

The term  $\alpha_1$  can be bounded by

$$\alpha_1 \leq 2 \left( \|\mathbf{X}^* \Delta_Y^\top\|_F^2 + \|\Delta_X \mathbf{Y}^{*\top}\|_F^2 \right) \leq 2 \left( \|\mathbf{X}^*\|^2 \|\Delta_Y\|_F^2 + \|\mathbf{Y}^*\|^2 \|\Delta_X\|_F^2 \right) = 2 \sigma_{\max} \|\Delta\|_F^2.$$

Combining all these bounds yields

$$\text{vec}(\Delta)^\top \nabla^2 f(\mathbf{X}, \mathbf{Y}) \text{vec}(\Delta) \leq 10 \sigma_{\max} \|\Delta\|_F^2.$$

## D.11 Proof of Lemma 18

The first set of consequences (101) follows straightforwardly from the triangle inequality. For instance, combine the induction hypotheses (93c) and (93e) to obtain

$$\begin{aligned}
\|\mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)} - \mathbf{F}^*\|_{2,\infty} &\leq \|\mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)} - \mathbf{F}^t \mathbf{H}^t\|_{2,\infty} + \|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*\|_{2,\infty} \\
&\leq (C_\infty \kappa + C_3) \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{F}^*\|_{2,\infty}.
\end{aligned}$$

Similar bounds can be obtained for  $\|\mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)} - \mathbf{F}^*\|$  provided that  $n \gg \mu r \log n$ .

We continue to establish the second set of consequences namely (102). Since  $\|\cdot\|$  is unitarily invariant, one can apply the triangle inequality to get

$$\begin{aligned}
\|\mathbf{F}^t\| &= \|\mathbf{F}^t \mathbf{H}^t\| \leq \|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*\| + \|\mathbf{F}^*\| \\
&\stackrel{(i)}{\leq} C_{\text{op}} \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{X}^*\| + \sqrt{2} \|\mathbf{X}^*\| \stackrel{(ii)}{\leq} 2 \|\mathbf{X}^*\|.
\end{aligned}$$

Here (i) uses the induction hypothesis (93b) and the fact that  $\|\mathbf{F}^*\| = \sqrt{2} \|\mathbf{X}^*\|$ , and (ii) holds as long as  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll 1$ . Similarly one can obtain  $\|\mathbf{F}^t\|_F \leq 2 \|\mathbf{X}^*\|_F$  provided that  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll 1$  and  $\|\mathbf{F}^t\|_{2,\infty} \leq 2 \|\mathbf{F}^*\|_{2,\infty}$  as long as  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll 1/(\sqrt{\kappa^2 \log n})$ . Notice that along the way, we have also proven that

$$\|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*\| \leq \|\mathbf{X}^*\|, \quad \|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*\|_F \leq \|\mathbf{X}^*\|_F \quad \text{and} \quad \|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*\|_{2,\infty} \leq \|\mathbf{F}^*\|_{2,\infty}.$$

We now move on to  $\|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^{t,(l)} \mathbf{H}^{t,(l)}\|_{\mathbb{F}}$ , for which we intend to apply Lemma 22 to connect it with  $\|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)}\|_{\mathbb{F}}$ . First, in view of the induction hypothesis (93b), one has

$$\begin{aligned} \|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^*\| \|\mathbf{F}^*\| &\leq C_{\text{op}} \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{X}^*\| \|\mathbf{F}^*\| \\ &= \sqrt{2} C_{\text{op}} \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \sigma_{\max} \\ &\leq \sigma_r^2(\mathbf{F}^*)/2, \end{aligned}$$

where the equality arises since  $\|\mathbf{F}^*\| = \sqrt{2\sigma_{\max}}$  (see (100a)) and  $\|\mathbf{X}^*\| = \sqrt{\sigma_{\max}}$ , and the last line holds as long as  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll 1/\kappa$ . In addition, it follows from the induction hypothesis (93c) that

$$\begin{aligned} \|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)}\| \|\mathbf{F}^*\| &\leq \|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)}\|_{\mathbb{F}} \|\mathbf{F}^*\| \\ &\leq C_3 \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{F}^*\|_{2,\infty} \|\mathbf{F}^*\| \\ &\leq \sqrt{2} C_3 \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \sqrt{\frac{\mu r}{n}} \sigma_{\max} \\ &\leq \sigma_r^2(\mathbf{F}^*)/4, \end{aligned}$$

where the penultimate inequality arises from the facts that  $\|\mathbf{F}^*\|_{2,\infty} \leq \sqrt{\mu r \sigma_{\max}/n}$  and that  $\|\mathbf{F}^*\| = \sqrt{2\sigma_{\max}}$  (cf. (100)), and the last relation holds as long as  $(\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}}) \sqrt{\frac{\mu r}{n}} \ll 1/\kappa$ . Invoke Lemma 22 with

$$\mathbf{F}_0 = \mathbf{F}^*, \quad \mathbf{F}_1 = \mathbf{F}^t \mathbf{H}^t \quad \text{and} \quad \mathbf{F}_2 = \mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)}$$

to arrive at

$$\|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^{t,(l)} \mathbf{H}^{t,(l)}\|_{\mathbb{F}} \leq 5 \frac{\sigma_r^2(\mathbf{F}^*)}{\sigma_r^2(\mathbf{F}^*)} \|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)}\|_{\mathbb{F}} = 5\kappa \|\mathbf{F}^t \mathbf{H}^t - \mathbf{F}^{t,(l)} \mathbf{R}^{t,(l)}\|_{\mathbb{F}}.$$

The last set of consequences can be derived following similar arguments to that for establishing the first set. For brevity, we omit the proof.

## D.12 Proof of the inequalities (31)

We single out the proof of  $\|\mathbf{X}^t \mathbf{Y}^{t\top} - \mathbf{M}^*\|_{\infty}$ , whereas the proofs of  $\|\mathbf{X}^t \mathbf{Y}^{t\top} - \mathbf{M}^*\|_{\mathbb{F}}$  and  $\|\mathbf{X}^t \mathbf{Y}^{t\top} - \mathbf{M}^*\|$  follow from the same argument. Recognize the following decomposition

$$\mathbf{X}^t \mathbf{Y}^{t\top} - \mathbf{M}^* = (\mathbf{X}^t \mathbf{H}^t - \mathbf{X}^*) (\mathbf{Y}^t \mathbf{H}^t)^{\top} + \mathbf{X}^* (\mathbf{Y}^t \mathbf{H}^t - \mathbf{Y}^*)^{\top},$$

which together with the triangle inequality gives

$$\begin{aligned} \|\mathbf{X}^t \mathbf{Y}^{t\top} - \mathbf{M}^*\|_{\infty} &\leq \left\| (\mathbf{X}^t \mathbf{H}^t - \mathbf{X}^*) (\mathbf{Y}^t \mathbf{H}^t)^{\top} \right\|_{\infty} + \left\| \mathbf{X}^* (\mathbf{Y}^t \mathbf{H}^t - \mathbf{Y}^*)^{\top} \right\|_{\infty} \\ &\leq \|\mathbf{X}^t \mathbf{H}^t - \mathbf{X}^*\|_{2,\infty} \|\mathbf{Y}^t \mathbf{H}^t\|_{2,\infty} + \|\mathbf{X}^*\|_{2,\infty} \|\mathbf{Y}^t \mathbf{H}^t - \mathbf{Y}^*\|_{2,\infty}. \end{aligned}$$

In view of Lemma 18, one has  $\|\mathbf{Y}^t \mathbf{H}^t\|_{2,\infty} \leq 2\|\mathbf{F}^*\|_{2,\infty}$  as long as the noise obeys  $\frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \ll \frac{1}{\sqrt{\kappa^2 \log n}}$ . This further implies that

$$\begin{aligned} \|\mathbf{X}^t \mathbf{Y}^{t\top} - \mathbf{M}^*\|_{\infty} &\leq 3C_{\infty} \kappa \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{F}^*\|_{2,\infty} \|\mathbf{F}^*\|_{2,\infty} \\ &\leq 3C_{\infty} \sqrt{\kappa^3 \mu r} \left( \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} + \frac{\lambda}{p \sigma_{\min}} \right) \|\mathbf{M}^*\|_{\infty}, \end{aligned}$$

where the last relation is  $\|\mathbf{F}^*\|_{2,\infty} \|\mathbf{F}^*\|_{2,\infty} \leq \sqrt{\kappa\mu r} \|\mathbf{M}^*\|_\infty$ . To see this, one has for any  $1 \leq i \leq n$ ,

$$n \|\mathbf{M}^*\|_\infty^2 \geq \sum_{j=1}^n (M_{ij}^*)^2 = \mathbf{X}_{i,\cdot}^* \mathbf{Y}^{*\top} \mathbf{Y}^* \mathbf{X}_{i,\cdot}^{*\top} \geq \|\mathbf{X}_{i,\cdot}^*\|_2^2 \lambda_{\min}(\mathbf{Y}^{*\top} \mathbf{Y}^*) = \sigma_{\min} \|\mathbf{X}_{i,\cdot}^*\|_2^2.$$

Here  $\lambda_{\min}(\cdot)$  denotes the minimum eigenvalue. Since the inequality holds for all  $1 \leq i \leq n$ , we arrive at

$$\|\mathbf{X}^*\|_{2,\infty} \leq \sqrt{\frac{n}{\sigma_{\min}}} \|\mathbf{M}^*\|_\infty.$$

Similarly one can obtain  $\|\mathbf{Y}^*\|_{2,\infty} \leq \sqrt{n/\sigma_{\min}} \|\mathbf{M}^*\|_\infty$ , which further implies  $\|\mathbf{F}^*\|_{2,\infty} = \max\{\|\mathbf{X}^*\|_{2,\infty}, \|\mathbf{Y}^*\|_{2,\infty}\} \leq \sqrt{n/\sigma_{\min}} \|\mathbf{M}^*\|_\infty$ . As a result, we arrive at

$$\|\mathbf{F}^*\|_{2,\infty} \|\mathbf{F}^*\|_{2,\infty} \leq \sqrt{\frac{n}{\sigma_{\min}}} \|\mathbf{M}^*\|_\infty \cdot \sqrt{\frac{\mu r}{n}} \sqrt{\sigma_{\max}} \leq \sqrt{\kappa\mu r} \|\mathbf{M}^*\|_\infty.$$

Here we used the incoherence assumption (100a).

## E Technical lemmas

**Lemma 19.** *Suppose  $n^2 p \geq C n \log n$  for some sufficiently large constant  $C > 0$ . Then with probability exceeding  $1 - O(n^{-10})$ ,*

$$\left| p^{-1} \|\mathcal{P}_\Omega(\mathbf{A}\mathbf{B}^\top)\|_{\text{F}}^2 - \|\mathbf{A}\mathbf{B}^\top\|_{\text{F}}^2 \right| \leq 3n \min \left\{ \|\mathbf{A}\|_{2,\infty}^2 \|\mathbf{B}\|_{\text{F}}^2, \|\mathbf{B}\|_{2,\infty}^2 \|\mathbf{A}\|_{\text{F}}^2 \right\}$$

holds uniformly for all matrices  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times r}$ .

*Proof.* In view of [ZL16, Lemma 9], one has

$$p^{-1} \|\mathcal{P}_\Omega(\mathbf{A}\mathbf{B}^\top)\|_{\text{F}}^2 \leq 2n \min \left\{ \|\mathbf{A}\|_{2,\infty}^2 \|\mathbf{B}\|_{\text{F}}^2, \|\mathbf{B}\|_{2,\infty}^2 \|\mathbf{A}\|_{\text{F}}^2 \right\}$$

with high probability. In addition, simple algebra reveals that

$$\|\mathbf{A}\mathbf{B}^\top\|_{\text{F}}^2 \leq \|\mathbf{A}\|_{\text{F}}^2 \|\mathbf{B}\|_{\text{F}}^2 \leq n \|\mathbf{A}\|_{2,\infty}^2 \|\mathbf{B}\|_{\text{F}}^2$$

and, similarly,  $\|\mathbf{A}\mathbf{B}^\top\|_{\text{F}}^2 \leq n \|\mathbf{A}\|_{\text{F}}^2 \|\mathbf{B}\|_{2,\infty}^2$ . Combining the previous bounds with the triangle inequality establishes the claim.  $\square$

**Lemma 20.** *Let  $\mathbf{U}\Sigma\mathbf{V}^\top$  be the SVD of a rank- $r$  matrix  $\mathbf{X}\mathbf{Y}^\top$  with  $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times r}$ . Then there exists an invertible matrix  $\mathbf{Q} \in \mathbb{R}^{r \times r}$  such that  $\mathbf{X} = \mathbf{U}\Sigma^{1/2}\mathbf{Q}$  and  $\mathbf{Y} = \mathbf{V}\Sigma^{1/2}\mathbf{Q}^{-\top}$ . In addition, one has*

$$\|\Sigma_{\mathbf{Q}} - \Sigma_{\mathbf{Q}}^{-1}\|_{\text{F}} \leq \frac{1}{\sigma_{\min}(\Sigma)} \|\mathbf{X}^\top \mathbf{X} - \mathbf{Y}^\top \mathbf{Y}\|_{\text{F}}, \quad (140)$$

where  $\mathbf{U}_{\mathbf{Q}}\Sigma_{\mathbf{Q}}\mathbf{V}_{\mathbf{Q}}^\top$  is the SVD of  $\mathbf{Q}$ . In particular, if  $\mathbf{X}$  and  $\mathbf{Y}$  have balanced scale, i.e.  $\mathbf{X}^\top \mathbf{X} - \mathbf{Y}^\top \mathbf{Y} = \mathbf{0}$ , then  $\mathbf{Q}$  must be a rotation matrix.

*Proof.* The existence of  $\mathbf{Q}$  is trivial by setting

$$\mathbf{Q} = \Sigma^{-1/2} \mathbf{U}^\top \mathbf{X}.$$

To see this, one has

$$\mathbf{U}\Sigma^{1/2}\mathbf{Q} = \mathbf{U}\Sigma^{1/2}\Sigma^{-1/2}\mathbf{U}^\top \mathbf{X} = \mathbf{U}\mathbf{U}^\top \mathbf{X} = \mathbf{X},$$

where the last equality follows from the fact that the columns of  $\mathbf{U}$  are the left singular vectors of  $\mathbf{X}$ . The relation  $\mathbf{Y} = \mathbf{V}\Sigma^{1/2}\mathbf{Q}^{-\top}$  can also be verified by the identity

$$\mathbf{X}\mathbf{Y}^\top = \mathbf{U}\Sigma^{1/2}\mathbf{Q}\mathbf{Y}^\top = \mathbf{U}\Sigma\mathbf{V}^\top.$$

We now move on to proving the perturbation bound (140). In view of the SVD of  $\mathbf{Q}$ , i.e.  $\mathbf{Q} = \mathbf{U}_Q \boldsymbol{\Sigma}_Q \mathbf{V}_Q^\top$ , one can obtain

$$\begin{aligned} \mathbf{X}^\top \mathbf{X} - \mathbf{Y}^\top \mathbf{Y} &= \mathbf{Q}^\top \boldsymbol{\Sigma} \mathbf{Q} - \mathbf{Q}^{-1} \boldsymbol{\Sigma} \mathbf{Q}^{-\top} \\ &= \mathbf{V}_Q \boldsymbol{\Sigma}_Q \mathbf{U}_Q^\top \boldsymbol{\Sigma} \mathbf{U}_Q \boldsymbol{\Sigma}_Q \mathbf{V}_Q^\top - \mathbf{V}_Q \boldsymbol{\Sigma}_Q^{-1} \mathbf{U}_Q^\top \boldsymbol{\Sigma} \mathbf{U}_Q \boldsymbol{\Sigma}_Q^{-1} \mathbf{V}_Q^\top. \end{aligned}$$

Denote  $\mathbf{B} := \mathbf{U}_Q^\top \boldsymbol{\Sigma} \mathbf{U}_Q \succ 0$ . Then we have

$$\|\mathbf{X}^\top \mathbf{X} - \mathbf{Y}^\top \mathbf{Y}\|_{\mathbb{F}}^2 = \left\| \mathbf{V}_Q \boldsymbol{\Sigma}_Q \mathbf{B} \boldsymbol{\Sigma}_Q \mathbf{V}_Q^\top - \mathbf{V}_Q \boldsymbol{\Sigma}_Q^{-1} \mathbf{B} \boldsymbol{\Sigma}_Q^{-1} \mathbf{V}_Q^\top \right\|_{\mathbb{F}}^2 = \left\| \boldsymbol{\Sigma}_Q \mathbf{B} \boldsymbol{\Sigma}_Q - \boldsymbol{\Sigma}_Q^{-1} \mathbf{B} \boldsymbol{\Sigma}_Q^{-1} \right\|_{\mathbb{F}}^2.$$

Let  $\mathbf{C} = \boldsymbol{\Sigma}_Q \mathbf{B}^{1/2}$  and  $\mathbf{D} = \boldsymbol{\Sigma}_Q^{-1} \mathbf{B}^{1/2}$ , and denote  $\boldsymbol{\Delta} = \mathbf{C} - \mathbf{D}$ . One then has

$$\begin{aligned} \|\mathbf{X}^\top \mathbf{X} - \mathbf{Y}^\top \mathbf{Y}\|_{\mathbb{F}}^2 &= \|\mathbf{C} \mathbf{C}^\top - \mathbf{D} \mathbf{D}^\top\|_{\mathbb{F}}^2 = \|\mathbf{C} \boldsymbol{\Delta}^\top + \boldsymbol{\Delta} \mathbf{C}^\top - \boldsymbol{\Delta} \boldsymbol{\Delta}^\top\|_{\mathbb{F}}^2 \\ &= \text{Tr} (2\mathbf{C}^\top \mathbf{C} \boldsymbol{\Delta}^\top \boldsymbol{\Delta} + \boldsymbol{\Delta} \boldsymbol{\Delta}^\top \boldsymbol{\Delta} \boldsymbol{\Delta}^\top + 2\mathbf{C}^\top \boldsymbol{\Delta} \mathbf{C}^\top \boldsymbol{\Delta} - 4\mathbf{C}^\top \boldsymbol{\Delta} \boldsymbol{\Delta}^\top \boldsymbol{\Delta}) \\ &= \text{Tr} \left[ (\boldsymbol{\Delta}^\top \boldsymbol{\Delta} - \sqrt{2} \mathbf{C}^\top \boldsymbol{\Delta})^2 + (4 - 2\sqrt{2}) \mathbf{C}^\top (\mathbf{C} - \boldsymbol{\Delta}) \boldsymbol{\Delta}^\top \boldsymbol{\Delta} + (2\sqrt{2} - 1) \mathbf{C}^\top \mathbf{C} \boldsymbol{\Delta}^\top \boldsymbol{\Delta} \right]. \end{aligned}$$

Note that  $\mathbf{C}^\top \mathbf{D} = \mathbf{B}$  and that  $\mathbf{C}^\top \boldsymbol{\Delta} = \mathbf{C}^\top \mathbf{C} - \mathbf{C}^\top \mathbf{D} = \mathbf{C}^\top \mathbf{C} - \mathbf{B}$  is symmetric. One can continue the bound as

$$\begin{aligned} \|\mathbf{X}^\top \mathbf{X} - \mathbf{Y}^\top \mathbf{Y}\|_{\mathbb{F}}^2 &= \left\| \boldsymbol{\Delta}^\top \boldsymbol{\Delta} - \sqrt{2} \mathbf{C}^\top \boldsymbol{\Delta} \right\|_{\mathbb{F}}^2 + (4 - 2\sqrt{2}) \text{Tr} (\mathbf{B} \boldsymbol{\Delta} \boldsymbol{\Delta}^\top) + (2\sqrt{2} - 1) \|\mathbf{C} \boldsymbol{\Delta}^\top\|_{\mathbb{F}}^2 \\ &\geq \text{Tr} (\mathbf{B} \boldsymbol{\Delta} \boldsymbol{\Delta}^\top), \end{aligned}$$

where the inequality follows since  $4 - 2\sqrt{2} \geq 1$ . Write  $\mathbf{B} = \mathbf{B}^{1/2} \cdot \mathbf{B}^{1/2}$  to see

$$\begin{aligned} \|\mathbf{X}^\top \mathbf{X} - \mathbf{Y}^\top \mathbf{Y}\|_{\mathbb{F}}^2 &\geq \text{Tr} (\mathbf{B}^{1/2} \boldsymbol{\Delta} \boldsymbol{\Delta}^\top \mathbf{B}^{1/2}) = \|\mathbf{B}^{1/2} \boldsymbol{\Delta}\|_{\mathbb{F}}^2 \\ &= \|\mathbf{B}^{1/2} (\boldsymbol{\Sigma}_Q - \boldsymbol{\Sigma}_Q^{-1}) \mathbf{B}^{1/2}\|_{\mathbb{F}}^2 \\ &\geq \sigma_{\min}^2 (\mathbf{B}) \|\boldsymbol{\Sigma}_Q - \boldsymbol{\Sigma}_Q^{-1}\|_{\mathbb{F}}^2. \end{aligned}$$

Recognizing that  $\sigma_{\min}(\mathbf{B}) = \sigma_{\min}(\boldsymbol{\Sigma})$  finishes the proof of (140).

Combining  $\mathbf{X}^\top \mathbf{X} = \mathbf{Y}^\top \mathbf{Y}$  and (140) yields  $\|\boldsymbol{\Sigma}_Q - \boldsymbol{\Sigma}_Q^{-1}\|_{\mathbb{F}} = 0$ , which implies  $\boldsymbol{\Sigma}_Q = \mathbf{I}$ . Under this circumstance,  $\mathbf{Q} = \mathbf{U}_Q \boldsymbol{\Sigma}_Q \mathbf{V}_Q^\top = \mathbf{U}_Q \mathbf{V}_Q^\top$  is a rotation matrix. The proof is then complete.  $\square$

**Lemma 21.** For all  $\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D} \in \mathbb{R}^{n \times r}$ , one has

$$|\langle \mathcal{P}_\Omega (\mathbf{A} \mathbf{C}^\top), \mathcal{P}_\Omega (\mathbf{B} \mathbf{D}^\top) \rangle - p \langle \mathbf{A} \mathbf{C}^\top, \mathbf{B} \mathbf{D}^\top \rangle| \leq \|\mathcal{P}_\Omega (\mathbf{1} \mathbf{1}^\top) - p \mathbf{1} \mathbf{1}^\top\| \|\mathbf{A}\|_{2, \infty} \|\mathbf{B}\|_{\mathbb{F}} \|\mathbf{C}\|_{2, \infty} \|\mathbf{D}\|_{\mathbb{F}}.$$

*Proof.* This is a simple consequence of [CL17, Lemma 4.4], where they have shown

$$\begin{aligned} &|\langle \mathcal{P}_\Omega (\mathbf{A} \mathbf{C}^\top), \mathcal{P}_\Omega (\mathbf{B} \mathbf{D}^\top) \rangle - p \langle \mathbf{A} \mathbf{C}^\top, \mathbf{B} \mathbf{D}^\top \rangle| \\ &\leq \|\mathcal{P}_\Omega (\mathbf{1} \mathbf{1}^\top) - p \mathbf{1} \mathbf{1}^\top\| \sqrt{\sum_{k=1}^n \|\mathbf{A}_{k, \cdot}\|_2^2 \|\mathbf{B}_{k, \cdot}\|_2^2} \sqrt{\sum_{k=1}^n \|\mathbf{C}_{k, \cdot}\|_2^2 \|\mathbf{D}_{k, \cdot}\|_2^2}. \end{aligned}$$

Recognize that

$$\sum_{k=1}^n \|\mathbf{A}_{k, \cdot}\|_2^2 \|\mathbf{B}_{k, \cdot}\|_2^2 \leq \|\mathbf{A}\|_{2, \infty}^2 \sum_{k=1}^n \|\mathbf{B}_{k, \cdot}\|_2^2 = \|\mathbf{A}\|_{2, \infty}^2 \|\mathbf{B}\|_{\mathbb{F}}^2$$

and, similarly,  $\sum_k \|\mathbf{C}_{k, \cdot}\|_2^2 \|\mathbf{D}_{k, \cdot}\|_2^2 \leq \|\mathbf{C}\|_{2, \infty}^2 \|\mathbf{D}\|_{\mathbb{F}}^2$ . Putting these together concludes the proof.  $\square$

**Lemma 22.** Suppose  $\mathbf{F}_1, \mathbf{F}_2, \mathbf{F}_0 \in \mathbb{R}^{2n \times r}$  are three matrices such that

$$\|\mathbf{F}_1 - \mathbf{F}_0\| \|\mathbf{F}_0\| \leq \sigma_r^2 (\mathbf{F}_0) / 2 \quad \text{and} \quad \|\mathbf{F}_1 - \mathbf{F}_2\| \|\mathbf{F}_0\| \leq \sigma_r^2 (\mathbf{F}_0) / 4,$$

where  $\sigma_i(\mathbf{A})$  stands for the  $i$ th largest singular value of  $\mathbf{A}$ . Denote

$$\mathbf{R}_1 \triangleq \arg \min_{\mathbf{R} \in \mathcal{O}^{r \times r}} \|\mathbf{F}_1 \mathbf{R} - \mathbf{F}_0\|_{\text{F}} \quad \text{and} \quad \mathbf{R}_2 \triangleq \arg \min_{\mathbf{R} \in \mathcal{O}^{r \times r}} \|\mathbf{F}_2 \mathbf{R} - \mathbf{F}_0\|_{\text{F}}.$$

Then the following two inequalities hold true:

$$\|\mathbf{F}_1 \mathbf{R}_1 - \mathbf{F}_2 \mathbf{R}_2\| \leq 5 \frac{\sigma_1^2(\mathbf{F}_0)}{\sigma_r^2(\mathbf{F}_0)} \|\mathbf{F}_1 - \mathbf{F}_2\| \quad \text{and} \quad \|\mathbf{F}_1 \mathbf{R}_1 - \mathbf{F}_2 \mathbf{R}_2\|_{\text{F}} \leq 5 \frac{\sigma_1^2(\mathbf{F}_0)}{\sigma_r^2(\mathbf{F}_0)} \|\mathbf{F}_1 - \mathbf{F}_2\|_{\text{F}}.$$

*Proof.* This is the same as [MWCC17, Lemma 37]. □

**Lemma 23.** Let  $\mathbf{S} \in \mathbb{R}^{r \times r}$  be a nonsingular matrix. Then for any matrix  $\mathbf{K} \in \mathbb{R}^{r \times r}$  with  $\|\mathbf{K}\| \leq \sigma_{\min}(\mathbf{S})$ , one has

$$\|\text{sgn}(\mathbf{S} + \mathbf{K}) - \text{sgn}(\mathbf{S})\| \leq \frac{2}{\sigma_{r-1}(\mathbf{S}) + \sigma_r(\mathbf{S})} \|\mathbf{K}\|,$$

where  $\text{sgn}(\cdot)$  denotes the matrix sign function, i.e.  $\text{sgn}(\mathbf{A}) = \mathbf{U}\mathbf{V}^{\top}$  for a matrix  $\mathbf{A}$  with SVD  $\mathbf{U}\mathbf{\Sigma}\mathbf{V}^{\top}$ .

*Proof.* This is the same as [MWCC17, Lemma 36]. □