



# Implicit Regularization in Nonconvex Statistical Estimation: Gradient Descent Converges Linearly for Phase Retrieval, Matrix Completion, and Blind Deconvolution

Cong Ma<sup>1</sup> · Kaizheng Wang<sup>1</sup> · Yuejie Chi<sup>2</sup> · Yuxin Chen<sup>3</sup>

Received: 14 December 2017 / Revised: 8 May 2019 / Accepted: 18 June 2019 / Published online: 5 August 2019  
© The Author(s) 2019

## Abstract

Recent years have seen a flurry of activities in designing provably efficient nonconvex procedures for solving statistical estimation problems. Due to the highly nonconvex nature of the empirical loss, state-of-the-art procedures often require proper regularization (e.g., trimming, regularized cost, projection) in order to guarantee fast convergence. For vanilla procedures such as gradient descent, however, prior theory either recommends highly conservative learning rates to avoid overshooting, or completely lacks performance guarantees. This paper uncovers a striking phenomenon in nonconvex optimization: even in the absence of explicit regularization, gradient descent enforces proper regularization implicitly under various statistical models. In fact, gradient descent follows a trajectory staying within a basin that enjoys nice geometry, consisting of points incoherent with the sampling mechanism. This “implicit regularization” feature allows gradient descent to proceed in a far more aggressive fashion without overshooting, which in turn results in substantial computational savings. Focusing on three fundamental statistical estimation problems, i.e., phase retrieval, low-rank matrix completion, and blind deconvolution, we establish that gradient descent achieves near-optimal statistical and computational guarantees without explicit regularization. In particular, by marrying statistical modeling with generic optimization theory, we develop a general recipe for analyzing the trajectories of iterative algorithms via a leave-one-out perturbation argument. As a by-product, for noisy matrix completion, we demonstrate that gradient descent achieves near-optimal error control—measured entrywise and by the spectral norm—which might be of independent interest.

**Keywords** Nonconvex optimization · Gradient descent · Leave-one-out analysis · Phase retrieval · Matrix completion · Blind deconvolution

---

Communicated by Emmanuel J. Candès.

Extended author information available on the last page of the article

## Mathematics Subject Classification 90C26

### 1 Introduction

#### 1.1 Nonlinear Systems and Empirical Loss Minimization

A wide spectrum of science and engineering applications calls for solutions to a nonlinear system of equations. Imagine we have collected a set of data points  $\mathbf{y} = \{y_j\}_{1 \leq j \leq m}$ , generated by a nonlinear sensing system,

$$y_j \approx \mathcal{A}_j(\mathbf{x}^*), \quad 1 \leq j \leq m,$$

where  $\mathbf{x}^*$  is the unknown object of interest and the  $\mathcal{A}_j$ 's are certain nonlinear maps known *a priori*. Can we reconstruct the underlying object  $\mathbf{x}^*$  in a faithful yet efficient manner? Problems of this kind abound in information and statistical science, prominent examples including low-rank matrix recovery [19,64], robust principal component analysis [17,21], phase retrieval [20,59], neural networks [103,132], to name just a few.

In principle, it is possible to attempt reconstruction by searching for a solution that minimizes the empirical loss, namely

$$\text{minimize}_{\mathbf{x}} f(\mathbf{x}) = \sum_{j=1}^m |y_j - \mathcal{A}_j(\mathbf{x})|^2. \quad (1)$$

Unfortunately, this empirical loss minimization problem is, in many cases, nonconvex, making it NP-hard in general. This issue of nonconvexity comes up in, for example, several representative problems that epitomize the structures of nonlinear systems encountered in practice.<sup>1</sup>

- *Phase retrieval/solving quadratic systems of equations* Imagine we are asked to recover an unknown object  $\mathbf{x}^* \in \mathbb{R}^n$ , but are only given the square modulus of certain linear measurements about the object, with all sign/phase information of the measurements missing. This arises, for example, in X-ray crystallography [15], and in latent-variable models where the hidden variables are captured by the missing signs [33]. To fix ideas, assume we would like to solve for  $\mathbf{x}^* \in \mathbb{R}^n$  in the following quadratic system of  $m$  equations

$$y_j = (\mathbf{a}_j^\top \mathbf{x}^*)^2, \quad 1 \leq j \leq m,$$

<sup>1</sup> Here, we choose different pre-constants in front of the empirical loss in order to be consistent with the literature of the respective problems. In addition, we only introduce the problem in the noiseless case for simplicity of presentation.

where  $\{\mathbf{a}_j\}_{1 \leq j \leq m}$  are the known design vectors. One strategy is thus to solve the following problem

$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^n} \quad f(\mathbf{x}) = \frac{1}{4m} \sum_{j=1}^m \left[ y_j - (\mathbf{a}_j^\top \mathbf{x}) \right]^2. \tag{2}$$

- *Low-rank matrix completion* In many scenarios such as collaborative filtering, we wish to make predictions about all entries of an (approximately) low-rank matrix  $\mathbf{M}^* \in \mathbb{R}^{n \times n}$  (e.g., a matrix consisting of users’ ratings about many movies), yet only a highly incomplete subset of the entries are revealed to us [19]. For clarity of presentation, assume  $\mathbf{M}^*$  to be rank- $r$  ( $r \ll n$ ) and positive semidefinite (PSD), i.e.,  $\mathbf{M}^* = \mathbf{X}^* \mathbf{X}^{*\top}$  with  $\mathbf{X}^* \in \mathbb{R}^{n \times r}$ , and suppose we have only seen the entries

$$Y_{j,k} = M_{j,k}^* = (\mathbf{X}^* \mathbf{X}^{*\top})_{j,k}, \quad (j, k) \in \Omega$$

within some index subset  $\Omega$  of cardinality  $m$ . These entries can be viewed as nonlinear measurements about the low-rank factor  $\mathbf{X}^*$ . The task of completing the true matrix  $\mathbf{M}^*$  can then be cast as solving

$$\text{minimize}_{\mathbf{X} \in \mathbb{R}^{n \times r}} \quad f(\mathbf{X}) = \frac{n^2}{4m} \sum_{(j,k) \in \Omega} \left( Y_{j,k} - \mathbf{e}_j^\top \mathbf{X} \mathbf{X}^\top \mathbf{e}_k \right)^2, \tag{3}$$

where the  $\mathbf{e}_j$ ’s stand for the canonical basis vectors in  $\mathbb{R}^n$ .

- *Blind deconvolution/solving bilinear systems of equations* Imagine we are interested in estimating two signals of interest  $\mathbf{h}^*, \mathbf{x}^* \in \mathbb{C}^K$ , but only get to collect a few bilinear measurements about them. This problem arises from mathematical modeling of blind deconvolution [3,76], which frequently arises in astronomy, imaging, communications, etc. The goal is to recover two signals from their convolution. Put more formally, suppose we have acquired  $m$  bilinear measurements taking the following form

$$y_j = \mathbf{b}_j^H \mathbf{h}^* \mathbf{x}^{*H} \mathbf{a}_j, \quad 1 \leq j \leq m,$$

where  $\mathbf{a}_j, \mathbf{b}_j \in \mathbb{C}^K$  are distinct design vectors (e.g., Fourier and/or random design vectors) known *a priori* and  $\mathbf{b}_j^H$  denotes the conjugate transpose of  $\mathbf{b}_j$ . In order to reconstruct the underlying signals, one asks for solutions to the following problem

$$\text{minimize}_{\mathbf{h}, \mathbf{x} \in \mathbb{C}^K} \quad f(\mathbf{h}, \mathbf{x}) = \sum_{j=1}^m \left| y_j - \mathbf{b}_j^H \mathbf{h} \mathbf{x}^H \mathbf{a}_j \right|^2.$$

### 1.2 Nonconvex Optimization via Regularized Gradient Descent

First-order methods have been a popular heuristic in practice for solving nonconvex problems including (1). For instance, a widely adopted procedure is gradient descent, which follows the update rule

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \nabla f(\mathbf{x}^t), \quad t \geq 0, \quad (4)$$

where  $\eta_t$  is the learning rate (or step size) and  $\mathbf{x}^0$  is some proper initial guess. Given that it only performs a single gradient calculation  $\nabla f(\cdot)$  per iteration (which typically can be completed within near-linear time), this paradigm emerges as a candidate for solving large-scale problems. The concern is: whether  $\mathbf{x}^t$  converges to the global solution and, if so, how long it takes for convergence, especially since (1) is highly nonconvex.

Fortunately, despite the worst-case hardness, appealing convergence properties have been discovered in various statistical estimation problems, the blessing being that the statistical models help rule out ill-behaved instances. For the average case, the empirical loss often enjoys benign geometry, in a *local* region (or at least along certain directions) surrounding the global optimum. In light of this, an effective nonconvex iterative method typically consists of two stages:

1. a carefully designed initialization scheme (e.g., spectral method);
2. an iterative refinement procedure (e.g., gradient descent).

This strategy has recently spurred a great deal of interest, owing to its promise of achieving computational efficiency and statistical accuracy at once for a growing list of problems (e.g., [18,25,32,61,64,76,78,107]). However, rather than directly applying gradient descent (4), existing theory often suggests enforcing proper regularization. Such explicit regularization enables improved computational convergence by properly “stabilizing” the search directions. The following regularization schemes, among others, have been suggested to obtain or improve computational guarantees. We refer to these algorithms collectively as *Regularized Gradient Descent*.

- *Trimming/truncation*, which discards/truncates a subset of the gradient components when forming the descent direction. For instance, when solving quadratic systems of equations, one can modify the gradient descent update rule as

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \mathcal{T}(\nabla f(\mathbf{x}^t)), \quad (5)$$

where  $\mathcal{T}$  is an operator that effectively drops samples bearing too much influence on the search direction. This strategy [25,118,126] has been shown to enable exact recovery with linear-time computational complexity and optimal sample complexity.

- *Regularized loss*, which attempts to optimize a regularized empirical risk

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t (\nabla f(\mathbf{x}^t) + \nabla R(\mathbf{x}^t)), \quad (6)$$

where  $R(\mathbf{x})$  stands for an additional penalty term in the empirical loss. For example, in low-rank matrix completion  $R(\cdot)$  imposes penalty based on the  $\ell_2$  row norm [64,107] as well as the Frobenius norm [107] of the decision matrix, while in blind deconvolution, it penalizes the  $\ell_2$  norm as well as certain component-wise incoherence measure of the decision vectors [58,76,82].

**Table 1** Prior theory for gradient descent (with spectral initialization)

	Vanilla gradient descent			Regularized gradient descent		
	Sample complexity	Iteration complexity	Step size	Sample complexity	Iteration complexity	Type of regularization
Phase retrieval	$n \log n$	$n \log \frac{1}{\epsilon}$	$\frac{1}{n}$	$n$	$\log \frac{1}{\epsilon}$	Trimming [25,126]
Matrix completion	n/a	n/a	n/a	$nr^7$	$\frac{n}{r} \log \frac{1}{\epsilon}$	Regularized loss [107]
				$nr^2$	$r^2 \log \frac{1}{\epsilon}$	Projection [32,131]
Blind deconvolution	n/a	n/a	n/a	$K \text{poly} \log m$	$m \log \frac{1}{\epsilon}$	Regularized loss and projection [76]

- *Projection*, which projects the iterates onto certain sets based on prior knowledge, that is,

$$\mathbf{x}^{t+1} = \mathcal{P}(\mathbf{x}^t - \eta_t \nabla f(\mathbf{x}^t)), \tag{7}$$

where  $\mathcal{P}$  is a certain projection operator used to enforce, for example, incoherence properties. This strategy has been employed in both low-rank matrix completion [32,131] and blind deconvolution [76].

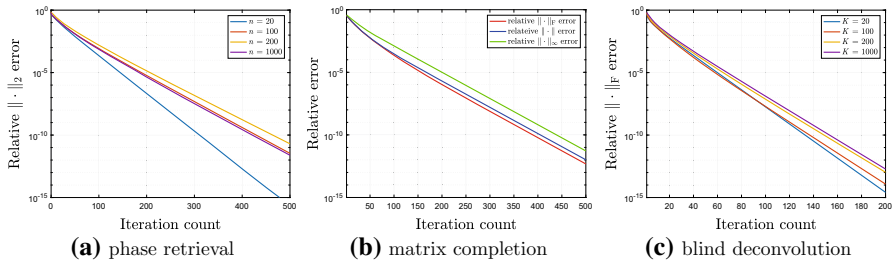
Equipped with such regularization procedures, existing works uncover appealing computational and statistical properties under various statistical models. Table 1 summarizes the performance guarantees derived in the prior literature; for simplicity, only orderwise results are provided.

**Remark 1** There is another role of regularization commonly studied in the literature, which exploits prior knowledge about the structure of the unknown object, such as sparsity to prevent over-fitting and improve statistical generalization ability. This is, however, not the focal point of this paper, since we are primarily pursuing solutions to (1) without imposing additional structures.

### 1.3 Regularization-Free Procedures?

The regularized gradient descent algorithms, while exhibiting appealing performance, usually introduce more algorithmic parameters that need to be carefully tuned based on the assumed statistical models. In contrast, vanilla gradient descent (cf. (4))—which is perhaps the very first method that comes into mind and requires minimal tuning parameters—is far less understood (cf. Table 1). Take matrix completion and blind deconvolution as examples: to the best of our knowledge, there is currently no theoretical guarantee derived for vanilla gradient descent.

The situation is better for phase retrieval: the local convergence of vanilla gradient descent, also known as Wirtinger flow (WF), has been investigated in [18,96]. Under i.i.d. Gaussian design and with near-optimal sample complexity, WF (combined with spectral initialization) provably achieves  $\epsilon$ -accuracy (in a relative sense)



**Fig. 1** **a** Relative  $\ell_2$  error of  $\mathbf{x}^t$  (modulo the global phase) versus iteration count for phase retrieval under i.i.d. Gaussian design, where  $m = 10n$  and  $\eta_t = 0.1$ . **b** Relative error of  $\mathbf{X}^t \mathbf{X}^{tT}$  (measured by  $\|\cdot\|_F, \|\cdot\|, \|\cdot\|_\infty$ ) versus iteration count for matrix completion, where  $n = 1000, r = 10, p = 0.1$ , and  $\eta_t = 0.2$ . **c** Relative error of  $\mathbf{h}^t \mathbf{x}^{tH}$  (measured by  $\|\cdot\|_F$ ) versus iteration count for blind deconvolution, where  $m = 10K$  and  $\eta_t = 0.5$

within  $O(n \log(1/\varepsilon))$  iterations. Nevertheless, the computational guarantee is significantly outperformed by the regularized version (called truncated Wirtinger flow [25]), which only requires  $O(\log(1/\varepsilon))$  iterations to converge with similar per-iteration cost. On closer inspection, the high computational cost of WF is largely due to the vanishingly small step size  $\eta_t = O(1/(n\|\mathbf{x}^*\|_2^2))$ —and hence slow movement—suggested by the theory [18]. While this is already the largest possible step size allowed in the theory published in [18], it is considerably more conservative than the choice  $\eta_t = O(1/\|\mathbf{x}^*\|_2^2)$  theoretically justified for the regularized version [25,126].

The lack of understanding and suboptimal results about vanilla gradient descent raise a very natural question: *Are regularization-free iterative algorithms inherently suboptimal when solving nonconvex statistical estimation problems of this kind?*

### 1.4 Numerical Surprise of Unregularized Gradient Descent

To answer the preceding question, it is perhaps best to first collect some numerical evidence. In what follows, we test the performance of vanilla gradient descent for phase retrieval, matrix completion, and blind deconvolution, using a *constant* step size. For all of these experiments, the initial guess is obtained by means of the standard spectral method. Our numerical findings are as follows:

- **Phase retrieval** For each  $n$ , set  $m = 10n$ , take  $\mathbf{x}^* \in \mathbb{R}^n$  to be a random vector with unit norm, and generate the design vectors  $\mathbf{a}_j \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n), 1 \leq j \leq m$ . Figure 1a illustrates the relative  $\ell_2$  error  $\min\{\|\mathbf{x}^t - \mathbf{x}^*\|_2, \|\mathbf{x}^t + \mathbf{x}^*\|_2\} / \|\mathbf{x}^*\|_2$  (modulo the unrecoverable global phase) versus the iteration count. The results are shown for  $n = 20, 100, 200, 1000$ , with the step size taken to be  $\eta_t = 0.1$  in all settings.
- **Matrix completion** Generate a random PSD matrix  $\mathbf{M}^* \in \mathbb{R}^{n \times n}$  with dimension  $n = 1000$ , rank  $r = 10$ , and all nonzero eigenvalues equal to one. Each entry of  $\mathbf{M}^*$  is observed independently with probability  $p = 0.1$ . Figure 1b plots the relative error  $\|\|\mathbf{X}^t \mathbf{X}^{tT} - \mathbf{M}^*\|\| / \|\|\mathbf{M}^*\|\|$  versus the iteration count, where  $\|\|\cdot\|\|$  can either be the Frobenius norm  $\|\cdot\|_F$ , the spectral norm  $\|\cdot\|$ , or the entrywise  $\ell_\infty$  norm  $\|\cdot\|_\infty$ . Here, we pick the step size as  $\eta_t = 0.2$ .

- *Blind deconvolution* For each  $K \in \{20, 100, 200, 1000\}$  and  $m = 10K$ , generate the design vectors  $\mathbf{a}_j \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \frac{1}{2}\mathbf{I}_K) + i\mathcal{N}(\mathbf{0}, \frac{1}{2}\mathbf{I}_K)$  for  $1 \leq j \leq m$  independently,<sup>2</sup> and the  $\mathbf{b}_j$ 's are drawn from a partial discrete Fourier transform (DFT) matrix (to be described in Sect. 3.3). The underlying signals  $\mathbf{h}^*, \mathbf{x}^* \in \mathbb{C}^K$  are produced as random vectors with unit norm. Figure 1c plots the relative error  $\|\mathbf{h}^t \mathbf{x}^{tH} - \mathbf{h}^* \mathbf{x}^{*H}\|_F / \|\mathbf{h}^* \mathbf{x}^{*H}\|_F$  versus the iteration count, with the step size taken to be  $\eta_t = 0.5$  in all settings.

In all of these numerical experiments, vanilla gradient descent enjoys remarkable linear convergence, always yielding an accuracy of  $10^{-5}$  (in a relative sense) within around 200 iterations. In particular, for the phase retrieval problem, the step size is taken to be  $\eta_t = 0.1$  although we vary the problem size from  $n = 20$  to  $n = 1000$ . The consequence is that the convergence rates experience little changes when the problem sizes vary. In comparison, the theory published in [18] seems overly pessimistic, as it suggests a diminishing step size inversely proportional to  $n$  and, as a result, an iteration complexity that worsens as the problem size grows.

In addition, it has been empirically observed in prior literature [25,76,127] that vanilla gradient descent performs comparably with the regularized counterpart for phase retrieval and blind deconvolution. To complete the picture, we further conduct experiments on matrix completion. In particular, we follow the experimental setup for matrix completion used above. We vary  $p$  from 0.01 to 0.1 with 51 logarithmically spaced points. For each  $p$ , we apply vanilla gradient descent, projected gradient descent [32] and gradient descent with additional regularization terms [107] with step size  $\eta = 0.2$  to 50 randomly generated instances. Successful recovery is declared if  $\|\mathbf{X}^t \mathbf{X}^{tT} - \mathbf{M}^*\|_F / \|\mathbf{M}^*\|_F \leq 10^{-5}$  in  $10^4$  iterations. Figure 2 reports the success rate versus the sampling rate. As can be seen, the phase transition of vanilla GD and that of GD with regularized cost are almost identical, whereas projected GD performs slightly better than the other two.

In short, the above empirical results are surprisingly positive yet puzzling. Why was the computational efficiency of vanilla gradient descent unexplained or substantially underestimated in prior theory?

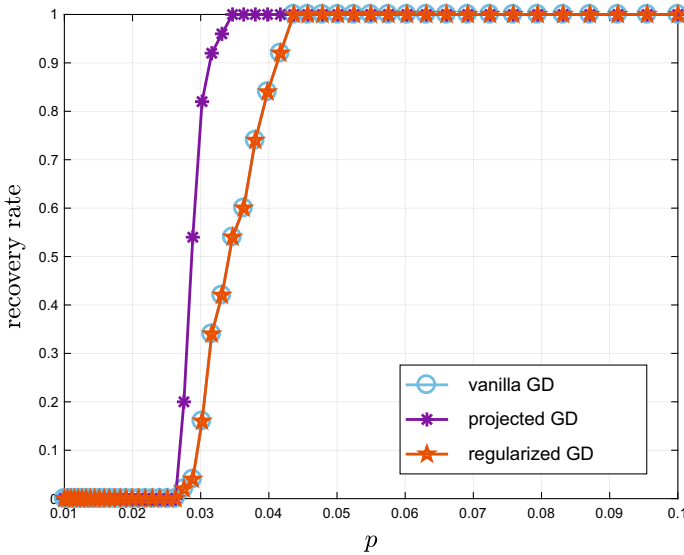
### 1.5 This Paper

The main contribution of this paper is toward demystifying the “unreasonable” effectiveness of regularization-free nonconvex iterative methods. As asserted in previous work, regularized gradient descent succeeds by properly enforcing/promoting certain incoherence conditions throughout the execution of the algorithm. In contrast, we discover that

*Vanilla gradient descent automatically forces the iterates to stay incoherent with the measurement mechanism, thus implicitly regularizing the search directions.*

This “implicit regularization” phenomenon is of fundamental importance, suggesting that vanilla gradient descent proceeds as if it were properly regularized. This

<sup>2</sup> Here and throughout,  $i$  represents the imaginary unit.



**Fig. 2** Success rate versus sampling rate  $p$  over 50 Monte Carlo trials for matrix completion with  $n = 1000$  and  $r = 10$

explains the remarkably favorable performance of unregularized gradient descent in practice. Focusing on the three representative problems mentioned in Sect. 1.1, our theory guarantees both statistical and computational efficiency of vanilla gradient descent under random designs and spectral initialization. With near-optimal sample complexity, to attain  $\epsilon$ -accuracy,

- *Phase retrieval (informal)* vanilla gradient descent converges in  $O(\log n \log \frac{1}{\epsilon})$  iterations;
- *Matrix completion (informal)* vanilla gradient descent converges in  $O(\log \frac{1}{\epsilon})$  iterations;
- *Blind deconvolution (informal)* vanilla gradient descent converges in  $O(\log \frac{1}{\epsilon})$  iterations.

In other words, gradient descent provably achieves (nearly) linear convergence in all of these examples. Throughout this paper, an algorithm is said to *converge (nearly) linearly* to  $\mathbf{x}^*$  in the noiseless case if the iterates  $\{\mathbf{x}^t\}$  obey

$$\text{dist}(\mathbf{x}^{t+1}, \mathbf{x}^*) \leq (1 - c) \text{dist}(\mathbf{x}^t, \mathbf{x}^*), \quad \forall t \geq 0$$

for some  $0 < c \leq 1$  that is (almost) independent of the problem size. Here,  $\text{dist}(\cdot, \cdot)$  can be any appropriate discrepancy measure.

As a by-product of our theory, gradient descent also provably controls the *entrywise* empirical risk uniformly across all iterations; for instance, this implies that vanilla gradient descent controls entrywise estimation error for the matrix completion task. Precise statements of these results are deferred to Sect. 3 and are briefly summarized in Table 2.



**Table 2** Prior theory versus our theory for vanilla gradient descent (with spectral initialization)

	Prior theory			Our theory		
	Sample complexity	Iteration complexity	Step size	Sample complexity	Iteration complexity	Step size
Phase retrieval	$n \log n$	$n \log (1/\varepsilon)$	$1/n$	$n \log n$	$\log n \log (1/\varepsilon)$	$1/\log n$
Matrix completion	n/a	n/a	n/a	$nr^3 \text{ poly } \log n$	$\log (1/\varepsilon)$	1
Blind deconvolution	n/a	n/a	n/a	$K \text{ poly } \log m$	$\log (1/\varepsilon)$	1

Notably, our study of implicit regularization suggests that the behavior of *nonconvex optimization* algorithms for statistical estimation needs to be examined in the context of *statistical models*, which induces an objective function as a finite sum. Our proof is accomplished via a leave-one-out perturbation argument, which is inherently tied to statistical models and leverages homogeneity across samples. Altogether, this allows us to localize benign landscapes for optimization and characterize finer dynamics not accounted for in generic gradient descent theory.

### 1.6 Notations

Before continuing, we introduce several notations used throughout the paper. First of all, boldfaced symbols are reserved for vectors and matrices. For any vector  $\mathbf{v}$ , we use  $\|\mathbf{v}\|_2$  to denote its Euclidean norm. For any matrix  $\mathbf{A}$ , we use  $\sigma_j(\mathbf{A})$  and  $\lambda_j(\mathbf{A})$  to denote its  $j$ th largest singular value and eigenvalue, respectively, and let  $\mathbf{A}_{j,\cdot}$  and  $\mathbf{A}_{\cdot,j}$  denote its  $j$ th row and  $j$ th column, respectively. In addition,  $\|\mathbf{A}\|$ ,  $\|\mathbf{A}\|_F$ ,  $\|\mathbf{A}\|_{2,\infty}$ , and  $\|\mathbf{A}\|_\infty$  stand for the spectral norm (i.e., the largest singular value), the Frobenius norm, the  $\ell_2/\ell_\infty$  norm (i.e., the largest  $\ell_2$  norm of the rows), and the entrywise  $\ell_\infty$  norm (the largest magnitude of all entries) of a matrix  $\mathbf{A}$ . Also,  $\mathbf{A}^\top$ ,  $\mathbf{A}^H$ , and  $\bar{\mathbf{A}}$  denote the transpose, the conjugate transpose, and the entrywise conjugate of  $\mathbf{A}$ , respectively.  $\mathbf{I}_n$  denotes the identity matrix with dimension  $n \times n$ . The notation  $\mathcal{O}^{n \times r}$  represents the set of all  $n \times r$  orthonormal matrices. The notation  $[n]$  refers to the set  $\{1, \dots, n\}$ . Also, we use  $\text{Re}(x)$  to denote the real part of a complex number  $x$ . Throughout the paper, we use the terms “samples” and “measurements” interchangeably.

Additionally, the standard notation  $f(n) = O(g(n))$  or  $f(n) \lesssim g(n)$  means that there exists a constant  $c > 0$  such that  $|f(n)| \leq c|g(n)|$ ,  $f(n) \gtrsim g(n)$  means that there exists a constant  $c > 0$  such that  $|f(n)| \geq c|g(n)|$ , and  $f(n) \asymp g(n)$  means that there exist constants  $c_1, c_2 > 0$  such that  $c_1|g(n)| \leq |f(n)| \leq c_2|g(n)|$ . Also,  $f(n) \gg g(n)$  means that there exists some large enough constant  $c > 0$  such that  $|f(n)| \geq c|g(n)|$ . Similarly,  $f(n) \ll g(n)$  means that there exists some sufficiently small constant  $c > 0$  such that  $|f(n)| \leq c|g(n)|$ .

## 2 Implicit Regularization: A Case Study

To reveal reasons behind the effectiveness of vanilla gradient descent, we first examine existing theory of gradient descent and identify the geometric properties that enable

linear convergence. We then develop an understanding as to why prior theory is conservative, and describe the phenomenon of implicit regularization that helps explain the effectiveness of vanilla gradient descent. To facilitate discussion, we will use the problem of solving random quadratic systems (phase retrieval) and Wirtinger flow as a case study, but our diagnosis applies more generally, as will be seen in later sections.

### 2.1 Gradient Descent Theory Revisited

In the convex optimization literature, there are two standard conditions about the objective function—strong convexity and smoothness—that allow for linear convergence of gradient descent.

**Definition 1** (*Strong convexity*) A twice continuously differentiable function  $f : \mathbb{R}^n \mapsto \mathbb{R}$  is said to be  $\alpha$ -strongly convex for  $\alpha > 0$  if

$$\nabla^2 f(\mathbf{x}) \succeq \alpha \mathbf{I}_n, \quad \forall \mathbf{x} \in \mathbb{R}^n.$$

**Definition 2** (*Smoothness*) A twice continuously differentiable function  $f : \mathbb{R}^n \mapsto \mathbb{R}$  is said to be  $\beta$ -smooth for  $\beta > 0$  if

$$\|\nabla^2 f(\mathbf{x})\| \leq \beta, \quad \forall \mathbf{x} \in \mathbb{R}^n.$$

It is well known that for an unconstrained optimization problem, if the objective function  $f$  is both  $\alpha$ -strongly convex and  $\beta$ -smooth, then vanilla gradient descent (4) enjoys  $\ell_2$  error contraction [9, Theorem 3.12], namely

$$\begin{aligned} \|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2 &\leq \left(1 - \frac{2}{\beta/\alpha + 1}\right) \|\mathbf{x}^t - \mathbf{x}^*\|_2, \quad \text{and} \quad \|\mathbf{x}^t - \mathbf{x}^*\|_2 \\ &\leq \left(1 - \frac{2}{\beta/\alpha + 1}\right)^t \|\mathbf{x}^0 - \mathbf{x}^*\|_2, \quad t \geq 0, \end{aligned} \tag{8}$$

as long as the step size is chosen as  $\eta_t = 2/(\alpha + \beta)$ . Here,  $\mathbf{x}^*$  denotes the global minimum. This immediately reveals the iteration complexity for gradient descent: the number of iterations taken to attain  $\epsilon$ -accuracy (in a relative sense) is bounded by

$$O\left(\frac{\beta}{\alpha} \log \frac{1}{\epsilon}\right).$$

In other words, the iteration complexity is dictated by and scales linearly with the condition number—the ratio  $\beta/\alpha$  of smoothness to strong convexity parameters.

Moving beyond convex optimization, one can easily extend the above theory to *nonconvex* problems with *local* strong convexity and smoothness. More precisely, suppose the objective function  $f$  satisfies

$$\nabla^2 f(\mathbf{x}) \succeq \alpha \mathbf{I} \quad \text{and} \quad \|\nabla^2 f(\mathbf{x})\| \leq \beta$$

over a local  $\ell_2$  ball surrounding the global minimum  $\mathbf{x}^*$ :

$$\mathcal{B}_\delta(\mathbf{x}) := \{\mathbf{x} \mid \|\mathbf{x} - \mathbf{x}^*\|_2 \leq \delta \|\mathbf{x}^*\|_2\}. \tag{9}$$

Then the contraction result (8) continues to hold, as long as the algorithm is seeded with an initial point that falls inside  $\mathcal{B}_\delta(\mathbf{x})$ .

### 2.2 Local Geometry for Solving Random Quadratic Systems

To invoke generic gradient descent theory, it is critical to characterize the local strong convexity and smoothness properties of the loss function. Take the problem of solving random quadratic systems (phase retrieval) as an example. Consider the i.i.d. Gaussian design in which  $\mathbf{a}_j \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ ,  $1 \leq j \leq m$ , and suppose without loss of generality that the underlying signal obeys  $\|\mathbf{x}^*\|_2 = 1$ . It is well known that  $\mathbf{x}^*$  is the unique minimizer—up to global phase—of (2) under this statistical model, provided that the ratio  $m/n$  of equations to unknowns is sufficiently large. The Hessian of the loss function  $f(\mathbf{x})$  is given by

$$\nabla^2 f(\mathbf{x}) = \frac{1}{m} \sum_{j=1}^m \left[ 3 \left( \mathbf{a}_j^\top \mathbf{x} \right)^2 - y_j \right] \mathbf{a}_j \mathbf{a}_j^\top. \tag{10}$$

- *Population-level analysis* Consider the case with an infinite number of equations or samples, i.e.,  $m \rightarrow \infty$ , where  $\nabla^2 f(\mathbf{x})$  converges to its expectation. Simple calculation yields that

$$\mathbb{E}[\nabla^2 f(\mathbf{x})] = 3 \left( \|\mathbf{x}\|_2^2 \mathbf{I}_n + 2\mathbf{x}\mathbf{x}^\top \right) - \left( \mathbf{I}_n + 2\mathbf{x}^* \mathbf{x}^{*\top} \right).$$

It is straightforward to verify that for any sufficiently small constant  $\delta > 0$ , one has the crude bound

$$\mathbf{I}_n \leq \mathbb{E}[\nabla^2 f(\mathbf{x})] \leq 10\mathbf{I}_n, \quad \forall \mathbf{x} \in \mathcal{B}_\delta(\mathbf{x}) : \|\mathbf{x} - \mathbf{x}^*\|_2 \leq \delta \|\mathbf{x}^*\|_2,$$

meaning that  $f$  is 1-strongly convex and 10-smooth within a local ball around  $\mathbf{x}^*$ . As a consequence, when we have infinite samples and an initial guess  $\mathbf{x}^0$  such that  $\|\mathbf{x}^0 - \mathbf{x}^*\|_2 \leq \delta \|\mathbf{x}^*\|_2$ , vanilla gradient descent with a constant step size converges to the global minimum within logarithmic iterations.

- *Finite-sample regime with  $m \asymp n \log n$*  Now that  $f$  exhibits favorable landscape in the population level, one thus hopes that the fluctuation can be well controlled so that the nice geometry carries over to the finite-sample regime. In the regime where  $m \asymp n \log n$  (which is the regime considered in [18]), the local strong convexity is still preserved, in the sense that

$$\nabla^2 f(\mathbf{x}) \geq (1/2) \cdot \mathbf{I}_n, \quad \forall \mathbf{x} : \|\mathbf{x} - \mathbf{x}^*\|_2 \leq \delta \|\mathbf{x}^*\|_2$$

occurs with high probability, provided that  $\delta > 0$  is sufficiently small (see [96,101] and Lemma 1). The smoothness parameter, however, is not well controlled. In fact, it can be as large as (up to logarithmic factors)<sup>3</sup>

$$\|\nabla^2 f(\mathbf{x})\| \lesssim n$$

even when we restrict attention to the local  $\ell_2$  ball (9) with  $\delta > 0$  being a fixed small constant. This means that the condition number  $\beta/\alpha$  (defined in Sect. 2.1) may scale as  $O(n)$ , leading to the step size recommendation

$$\eta_t \asymp 1/n,$$

and, as a consequence, a high iteration complexity  $O(n \log(1/\epsilon))$ . This underpins the analysis in [18].

In summary, the geometric properties of the loss function—even in the local  $\ell_2$  ball centering around the global minimum—are not as favorable as one anticipates, in particular in view of its population counterpart. A direct application of generic gradient descent theory leads to an overly conservative step size and a pessimistic convergence rate, unless the number of samples is enormously larger than the number of unknowns.

**Remark 2** Notably, due to Gaussian designs, the phase retrieval problem enjoys more favorable geometry compared to other nonconvex problems. In matrix completion and blind deconvolution, the Hessian matrices are rank-deficient even at the population level. In such cases, the above discussions need to be adjusted, e.g., strong convexity is only possible when we restrict attention to certain directions.

### 2.3 Which Region Enjoys Nicer Geometry?

Interestingly, our theory identifies a local region surrounding  $\mathbf{x}^*$  with a large diameter that enjoys much nicer geometry. This region does not mimic an  $\ell_2$  ball, but rather, the intersection of an  $\ell_2$  ball and a polytope. We term it the *region of incoherence and contraction* (RIC). For phase retrieval, the RIC includes all points  $\mathbf{x} \in \mathbb{R}^n$  obeying

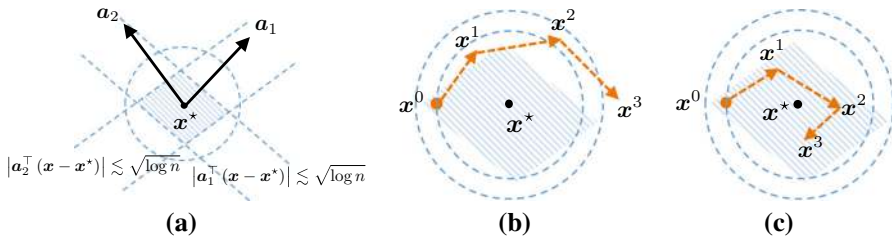
$$\|\mathbf{x} - \mathbf{x}^*\|_2 \leq \delta \|\mathbf{x}^*\|_2 \quad \text{and} \tag{11a}$$

$$\max_{1 \leq j \leq m} |\mathbf{a}_j^\top (\mathbf{x} - \mathbf{x}^*)| \lesssim \sqrt{\log n} \|\mathbf{x}^*\|_2, \tag{11b}$$

where  $\delta > 0$  is some small numerical constant. As will be formalized in Lemma 1, with high probability the Hessian matrix satisfies

$$(1/2) \cdot \mathbf{I}_n \leq \nabla^2 f(\mathbf{x}) \leq O(\log n) \cdot \mathbf{I}_n$$

<sup>3</sup> To demonstrate this, taking  $\mathbf{x} = \mathbf{x}^* + (\delta/\|\mathbf{a}_1\|_2) \cdot \mathbf{a}_1$  in (10), one can easily verify that, with high probability,  $\|\nabla^2 f(\mathbf{x})\| \geq |3(\mathbf{a}_1^\top \mathbf{x})^2 - y_1| \|\mathbf{a}_1 \mathbf{a}_1^\top\|/m - O(1) \gtrsim \delta^2 n^2/m \asymp \delta^2 n/\log n$ .



**Fig. 3** **a** The shaded region is an illustration of the incoherence region, which satisfies  $|a_j^\top(x - x^*)| \lesssim \sqrt{\log n}$  for all points  $x$  in the region. **b** When  $x^0$  resides in the desired region, we know that  $x^1$  remains within the  $\ell_2$  ball but might fall out of the incoherence region (the shaded region). Once  $x^1$  leaves the incoherence region, we lose control and may overshoot. **c** Our theory reveals that with high probability, all iterates will stay within the incoherence region, enabling fast convergence

simultaneously for  $x$  in the RIC. In words, the Hessian matrix is nearly well conditioned (with the condition number bounded by  $O(\log n)$ ), as long as (i) the iterate is not very far from the global minimizer (cf. (11a)) and (ii) the iterate remains incoherent<sup>4</sup> with respect to the sensing vectors (cf. (11b)). Another way to interpret the incoherence condition (11b) is that the empirical risk needs to be well controlled uniformly across all samples. See Fig. 3a for an illustration of the above region.

The following observation is thus immediate: one can safely adopt a far more aggressive step size (as large as  $\eta_t = O(1/\log n)$ ) to achieve acceleration, as long as the iterates stay within the RIC. This, however, fails to be guaranteed by generic gradient descent theory. To be more precise, if the current iterate  $x^t$  falls within the desired region, then in view of (8), we can ensure  $\ell_2$  error contraction after one iteration, namely

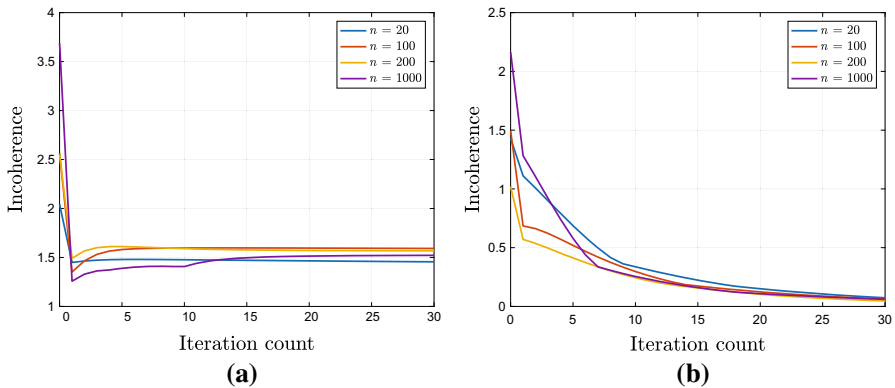
$$\|x^{t+1} - x^*\|_2 \leq \|x^t - x^*\|_2,$$

and hence  $x^{t+1}$  stays within the local  $\ell_2$  ball and hence satisfies (11a). However, it is not immediately obvious that  $x^{t+1}$  would still stay incoherent with the sensing vectors and satisfy (11b). If  $x^{t+1}$  leaves the RIC, it no longer enjoys the benign local geometry of the loss function, and the algorithm has to slow down in order to avoid overshooting. See Fig. 3b for a visual illustration. In fact, in almost all regularized gradient descent algorithms mentioned in Sect. 1.2, one of the main purposes of the proposed regularization procedures is to enforce such incoherence constraints.

### 2.4 Implicit Regularization

However, is regularization really necessary for the iterates to stay within the RIC? To answer this question, we plot in Fig. 4a (resp. Fig. 4b) the incoherence measure  $\frac{\max_j |a_j^\top x^t|}{\sqrt{\log n} \|x^t\|_2}$  (resp.  $\frac{\max_j |a_j^\top (x^t - x^*)|}{\sqrt{\log n} \|x^t - x^*\|_2}$ ) versus the iteration count in a typical Monte Carlo

<sup>4</sup> If  $x$  is aligned with (and hence very coherent with) one vector  $a_j$ , then with high probability one has  $|a_j^\top(x - x^*)| \gtrsim |a_j^\top x| \approx \sqrt{n} \|x\|_2$ , which is significantly larger than  $\sqrt{\log n} \|x\|_2$ .



**Fig. 4** The incoherence measure  $\frac{\max_{1 \leq j \leq m} |a_j^\top x^t|}{\sqrt{\log n} \|x^*\|_2}$  (in **a**) and  $\frac{\max_{1 \leq j \leq m} |a_j^\top (x^t - x^{t-1})|}{\sqrt{\log n} \|x^*\|_2}$  (in **b**) of the gradient iterates versus iteration count for the phase retrieval problem. The results are shown for  $n \in \{20, 100, 200, 1000\}$  and  $m = 10n$ , with the step size taken to be  $\eta_t = 0.1$ . The problem instances are generated in the same way as in Fig. 1a

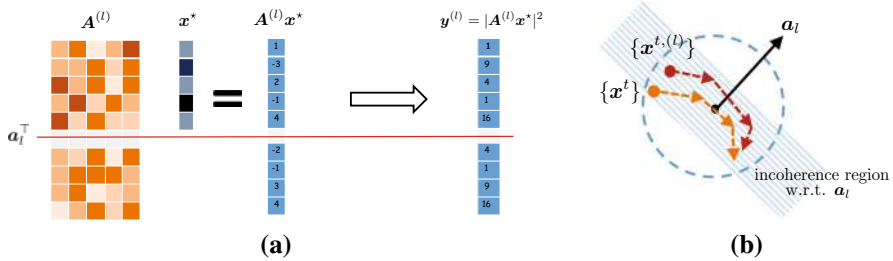
trial, generated in the same way as for Fig. 1a. Interestingly, the incoherence measure remains bounded by 2 for all iterations  $t > 1$ . This important observation suggests that one may adopt a substantially more aggressive step size throughout the whole algorithm.

The main objective of this paper is thus to provide a theoretical validation of the above empirical observation. As we will demonstrate shortly, with high probability all iterates along the execution of the algorithm (as well as the spectral initialization) are provably constrained within the RIC, implying fast convergence of vanilla gradient descent (cf. Fig. 3c). The fact that the iterates stay incoherent with the measurement mechanism automatically, without explicit enforcement, is termed “implicit regularization.”

### 2.5 A Glimpse of the Analysis: A Leave-One-Out Trick

In order to rigorously establish (11b) for all iterates, the current paper develops a powerful mechanism based on the leave-one-out perturbation argument, a trick rooted and widely used in probability and random matrix theory. Note that the iterate  $x^t$  is statistically dependent with the design vectors  $\{a_j\}$ . Under such circumstances, one often resorts to generic bounds like the Cauchy–Schwarz inequality, which would not yield a desirable estimate. To address this issue, we introduce a sequence of auxiliary iterates  $\{x^{t,(l)}\}$  for each  $1 \leq l \leq m$  (for analytical purposes only), obtained by running vanilla gradient descent using all but the  $l$ th sample. As one can expect, such auxiliary trajectories serve as extremely good surrogates of  $\{x^t\}$  in the sense that

$$x^t \approx x^{t,(l)}, \quad 1 \leq l \leq m, \quad t \geq 0, \tag{12}$$



**Fig. 5** Illustration of the leave-one-out sequence w.r.t.  $\mathbf{a}_l$ . **a** The sequence  $\{\mathbf{x}^{t,(l)}\}_{t \geq 0}$  is constructed without using the  $l$ th sample. **b** Since the auxiliary sequence  $\{\mathbf{x}^{t,(l)}\}$  is constructed without using  $\mathbf{a}_l$ , the leave-one-out iterates stay within the incoherence region w.r.t.  $\mathbf{a}_l$  with high probability. Meanwhile,  $\{\mathbf{x}^t\}$  and  $\{\mathbf{x}^{t,(l)}\}$  are expected to remain close as their construction differ only in a single sample

since their constructions only differ by a single sample. Most importantly, since  $\mathbf{x}^{t,(l)}$  is independent with the  $l$ th design vector, it is much easier to control its incoherence w.r.t.  $\mathbf{a}_l$  to the desired level:

$$|\mathbf{a}_l^\top (\mathbf{x}^{t,(l)} - \mathbf{x}^*)| \lesssim \sqrt{\log n} \|\mathbf{x}^*\|_2. \tag{13}$$

Combining (12) and (13) then leads to (11b). See Fig. 5 for a graphical illustration of this argument. Notably, this technique is very general and applicable to many other problems. We invite the readers to Sect. 5 for more details.

### 3 Main Results

This section formalizes the implicit regularization phenomenon underlying unregularized gradient descent and presents its consequences, namely near-optimal statistical and computational guarantees for phase retrieval, matrix completion, and blind deconvolution. Note that the discrepancy measure  $\text{dist}(\cdot, \cdot)$  may vary from problem to problem.

#### 3.1 Phase Retrieval

Suppose the  $m$  quadratic equations

$$y_j = (\mathbf{a}_j^\top \mathbf{x}^*)^2, \quad j = 1, 2, \dots, m \tag{14}$$

are collected using random design vectors, namely  $\mathbf{a}_j \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ , and the nonconvex problem to solve is

$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) := \frac{1}{4m} \sum_{j=1}^m \left[ (\mathbf{a}_j^\top \mathbf{x})^2 - y_j \right]^2. \tag{15}$$

The Wirtinger flow (WF) algorithm, first introduced in [18], is a combination of spectral initialization and vanilla gradient descent; see Algorithm 1.

---

**Algorithm 1** Wirtinger flow/gradient descent for phase retrieval

---

**Input:**  $\{a_j\}_{1 \leq j \leq m}$  and  $\{y_j\}_{1 \leq j \leq m}$ .

**Spectral initialization:** Let  $\lambda_1(Y)$  and  $\tilde{x}^0$  be the leading eigenvalue and eigenvector of

$$Y = \frac{1}{m} \sum_{j=1}^m y_j a_j a_j^\top, \tag{16}$$

respectively, and set  $x^0 = \sqrt{\lambda_1(Y)/3} \tilde{x}^0$ .

**Gradient updates:** for  $t = 0, 1, 2, \dots, T - 1$  do

$$x^{t+1} = x^t - \eta_t \nabla f(x^t). \tag{17}$$


---

Recognizing that the global phase/sign is unrecoverable from quadratic measurements, we introduce the  $\ell_2$  distance modulo the global phase as follows

$$\text{dist}(x, x^*) := \min \{ \|x - x^*\|_2, \|x + x^*\|_2 \}. \tag{18}$$

Our finding is summarized in the following theorem.

**Theorem 1** *Let  $x^* \in \mathbb{R}^n$  be a fixed vector. Suppose  $a_j \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I_n)$  for each  $1 \leq j \leq m$  and  $m \geq c_0 n \log n$  for some sufficiently large constant  $c_0 > 0$ . Assume the step size obeys  $\eta_t \equiv \eta = c_1 / (\log n \cdot \|x_0\|_2^2)$  for any sufficiently small constant  $c_1 > 0$ . Then there exist some absolute constants  $0 < \varepsilon < 1$  and  $c_2 > 0$  such that with probability at least  $1 - O(mn^{-5})$ , Algorithm 1 satisfies that for all  $t \geq 0$ ,*

$$\text{dist}(x^t, x^*) \leq \varepsilon (1 - \eta \|x^*\|_2^2 / 2)^t \|x^*\|_2, \tag{19a}$$

$$\max_{1 \leq j \leq m} |a_j^\top (x^t - x^*)| \leq c_2 \sqrt{\log n} \|x^*\|_2. \tag{19b}$$

Theorem 1 reveals a few intriguing properties of Algorithm 1.

- *Implicit regularization* Theorem 1 asserts that the incoherence properties are satisfied throughout the execution of the algorithm (see (19b)), which formally justifies the implicit regularization feature we hypothesized.
- *Near-constant step size* Consider the case where  $\|x^*\|_2 = 1$ . Theorem 1 establishes near-linear convergence of WF with a substantially more aggressive step size  $\eta \asymp 1/\log n$ . Compared with the choice  $\eta \lesssim 1/n$  admissible in [18, Theorem 3.3], Theorem 1 allows WF/GD to attain  $\varepsilon$ -accuracy within  $O(\log n \log(1/\varepsilon))$  iterations. The resulting computational complexity of the algorithm is

$$O\left(mn \log n \log \frac{1}{\varepsilon}\right),$$



which significantly improves upon the result  $O(mn^2 \log(1/\epsilon))$  derived in [18]. As a side note, if the sample size further increases to  $m \asymp n \log^2 n$ , then a constant step size  $\eta \asymp 1$  is also feasible, resulting in an iteration complexity  $\log(1/\epsilon)$ . This follows since with high probability, the entire trajectory resides within a more refined incoherence region  $\max_j |\mathbf{a}_j^\top (\mathbf{x}^t - \mathbf{x}^*)| \lesssim \|\mathbf{x}^*\|_2$ . We omit the details here.

- *Incoherence of spectral initialization* We have also demonstrated in Theorem 1 that the initial guess  $\mathbf{x}^0$  falls within the RIC and is hence nearly orthogonal to all design vectors. This provides a finer characterization of spectral initialization, in comparison with prior theory that focuses primarily on the  $\ell_2$  accuracy [18,90]. We expect our leave-one-out analysis to accommodate other variants of spectral initialization studied in the literature [12,25,83,88,118].

**Remark 3** As it turns out, a carefully designed initialization is not pivotal in enabling fast convergence. In fact, randomly initialized gradient descent provably attains  $\epsilon$ -accuracy in  $O(\log n + \log \frac{1}{\epsilon})$  iterations; see [27] for details.

### 3.2 Low-Rank Matrix Completion

Let  $\mathbf{M}^* \in \mathbb{R}^{n \times n}$  be a positive semidefinite matrix<sup>5</sup> with rank  $r$ , and suppose its eigendecomposition is

$$\mathbf{M}^* = \mathbf{U}^* \mathbf{\Sigma}^* \mathbf{U}^{*\top}, \tag{20}$$

where  $\mathbf{U}^* \in \mathbb{R}^{n \times r}$  consists of orthonormal columns and  $\mathbf{\Sigma}^*$  is an  $r \times r$  diagonal matrix with eigenvalues in a descending order, i.e.,  $\sigma_{\max} = \sigma_1 \geq \dots \geq \sigma_r = \sigma_{\min} > 0$ . Throughout this paper, we assume the condition number  $\kappa := \sigma_{\max}/\sigma_{\min}$  is bounded by a fixed constant, independent of the problem size (i.e.,  $n$  and  $r$ ). Denoting  $\mathbf{X}^* = \mathbf{U}^* (\mathbf{\Sigma}^*)^{1/2}$  allows us to factorize  $\mathbf{M}^*$  as

$$\mathbf{M}^* = \mathbf{X}^* \mathbf{X}^{*\top}. \tag{21}$$

Consider a random sampling model such that each entry of  $\mathbf{M}^*$  is observed independently with probability  $0 < p \leq 1$ , i.e., for  $1 \leq j \leq k \leq n$ ,

$$Y_{j,k} = \begin{cases} \mathbf{M}_{j,k}^* + E_{j,k}, & \text{with probability } p, \\ 0, & \text{else,} \end{cases} \tag{22}$$

where the entries of  $\mathbf{E} = [E_{j,k}]_{1 \leq j \leq k \leq n}$  are independent sub-Gaussian noise with sub-Gaussian norm  $\sigma$  (see [116, Definition 5.7]). We denote by  $\Omega$  the set of locations being sampled, and  $\mathcal{P}_\Omega(\mathbf{Y})$  represents the projection of  $\mathbf{Y}$  onto the set of matrices supported in  $\Omega$ . We note here that the sampling rate  $p$ , if not known, can be faithfully estimated by the sample proportion  $|\Omega|/n^2$ .

<sup>5</sup> Here, we assume  $\mathbf{M}^*$  to be positive semidefinite to simplify the presentation, but note that our analysis easily extends to asymmetric low-rank matrices.

To fix ideas, we consider the following nonconvex optimization problem

$$\text{minimize}_{\mathbf{X} \in \mathbb{R}^{n \times r}} \quad f(\mathbf{X}) := \frac{1}{4p} \sum_{(j,k) \in \Omega} (\mathbf{e}_j^\top \mathbf{X} \mathbf{X}^\top \mathbf{e}_k - Y_{j,k})^2. \quad (23)$$

The vanilla gradient descent algorithm (with spectral initialization) is summarized in Algorithm 2.

---

**Algorithm 2** Vanilla gradient descent for matrix completion (with spectral initialization)

---

**Input:**  $Y = [Y_{j,k}]_{1 \leq j,k \leq n}$ ,  $r, p$ .

**Spectral initialization:** Let  $U^0 \Sigma^0 U^{0\top}$  be the rank- $r$  eigendecomposition of

$$M^0 := \frac{1}{p} \mathcal{P}_\Omega(Y) = \frac{1}{p} \mathcal{P}_\Omega(M^* + E),$$

and set  $X^0 = U^0 (\Sigma^0)^{1/2}$ .

**Gradient updates:** for  $t = 0, 1, 2, \dots, T - 1$  do

$$X^{t+1} = X^t - \eta_t \nabla f(X^t). \quad (24)$$


---

Before proceeding to the main theorem, we first introduce a standard incoherence parameter required for matrix completion [19].

**Definition 3** (*Incoherence for matrix completion*) A rank- $r$  matrix  $M^*$  with eigendecomposition  $M^* = U^* \Sigma^* U^{*\top}$  is said to be  $\mu$ -incoherent if

$$\|U^*\|_{2,\infty} \leq \sqrt{\frac{\mu}{n}} \|U^*\|_F = \sqrt{\frac{\mu r}{n}}. \quad (25)$$

In addition, recognizing that  $X^*$  is identifiable only up to orthogonal transformation, we define the optimal transform from the  $t$ th iterate  $X^t$  to  $X^*$  as

$$\widehat{H}^t := \operatorname{argmin}_{R \in \mathcal{O}^{r \times r}} \|X^t R - X^*\|_F, \quad (26)$$

where  $\mathcal{O}^{r \times r}$  is the set of  $r \times r$  orthonormal matrices. With these definitions in place, we have the following theorem.

**Theorem 2** *Let  $M^*$  be a rank- $r$ ,  $\mu$ -incoherent PSD matrix, and its condition number  $\kappa$  is a fixed constant. Suppose the sample size satisfies  $n^2 p \geq C \mu^3 r^3 n \log^3 n$  for some sufficiently large constant  $C > 0$ , and the noise satisfies*

$$\sigma \sqrt{\frac{n}{p}} \ll \frac{\sigma_{\min}}{\sqrt{\kappa^3 \mu r \log^3 n}}. \quad (27)$$

With probability at least  $1 - O(n^{-3})$ , the iterates of Algorithm 2 satisfy

$$\|X^t \widehat{H}^t - X^*\|_F \leq \left( C_4 \rho^t \mu r \frac{1}{\sqrt{np}} + C_1 \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \right) \|X^*\|_F, \tag{28a}$$

$$\|X^t \widehat{H}^t - X^*\|_{2,\infty} \leq \left( C_5 \rho^t \mu r \sqrt{\frac{\log n}{np}} + C_8 \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} \right) \|X^*\|_{2,\infty}, \tag{28b}$$

$$\|X^t \widehat{H}^t - X^*\| \leq \left( C_9 \rho^t \mu r \frac{1}{\sqrt{np}} + C_{10} \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \right) \|X^*\| \tag{28c}$$

for all  $0 \leq t \leq T = O(n^5)$ , where  $C_1, C_4, C_5, C_8, C_9$ , and  $C_{10}$  are some absolute positive constants and  $1 - (\sigma_{\min}/5) \cdot \eta \leq \rho < 1$ , provided that  $0 < \eta_t \equiv \eta \leq 2 / (25\kappa\sigma_{\max})$ .

Theorem 2 provides the first theoretical guarantee of unregularized gradient descent for matrix completion, demonstrating near-optimal statistical accuracy and computational complexity.

- Implicit regularization** In Theorem 2, we bound the  $\ell_2/\ell_\infty$  error of the iterates in a uniform manner via (28b). Note that  $\|X - X^*\|_{2,\infty} = \max_j \|e_j^\top (X - X^*)\|_2$ , which implies the iterates remain incoherent with the sensing vectors throughout and have small incoherence parameters (cf. (25)). In comparison, prior works either include a penalty term on  $\{\|e_j^\top X\|_2\}_{1 \leq j \leq n}$  [64,107] and/or  $\|X\|_F$  [107] to encourage an incoherent and/or low-norm solution, or add an extra projection operation to enforce incoherence [32,131]. Our results demonstrate that such explicit regularization is unnecessary.
- Constant step size** Without loss of generality, we may assume that  $\sigma_{\max} = \|M^*\| = O(1)$ , which can be done by choosing proper scaling of  $M^*$ . Hence, we have a constant step size  $\eta_t \asymp 1$ . Actually, it is more convenient to consider the scale-invariant parameter  $\rho$ : Theorem 2 guarantees linear convergence of the vanilla gradient descent at a constant rate  $\rho$ . Remarkably, the convergence occurs with respect to three different unitarily invariant norms: the Frobenius norm  $\|\cdot\|_F$ , the  $\ell_2/\ell_\infty$  norm  $\|\cdot\|_{2,\infty}$ , and the spectral norm  $\|\cdot\|$ . As far as we know, the latter two are established for the first time. Note that our result even improves upon that for regularized gradient descent; see Table 1.
- Near-optimal sample complexity** When the rank  $r = O(1)$ , vanilla gradient descent succeeds under a near-optimal sample complexity  $n^2 p \gtrsim npoly \log n$ , which is statistically optimal up to some logarithmic factor.
- Near-minimal Euclidean error** In view of (28a), as  $t$  increases, the Euclidean error of vanilla GD converges to

$$\|X^t \widehat{H}^t - X^*\|_F \lesssim \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \|X^*\|_F, \tag{29}$$

which coincides with the theoretical guarantee in [32, Corollary 1] and matches the minimax lower bound established in [67,89].

- *Near-optimal entrywise error* The  $\ell_2/\ell_\infty$  error bound (28b) immediately yields entrywise control of the empirical risk. Specifically, as soon as  $t$  is sufficiently large (so that the first term in (28b) is negligible), we have

$$\begin{aligned} \|X^t X^{t\top} - M^*\|_\infty &\leq \|X^t \widehat{H}^t (X^t \widehat{H}^t - X^*)^\top\|_\infty + \|(X^t \widehat{H}^t - X^*) X^{*\top}\|_\infty \\ &\leq \|X^t \widehat{H}^t\|_{2,\infty} \|X^t \widehat{H}^t - X^*\|_{2,\infty} + \|X^t \widehat{H}^t - X^*\|_{2,\infty} \|X^*\|_{2,\infty} \\ &\lesssim \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} \|M^*\|_\infty, \end{aligned}$$

where the last line follows from (28b) as well as the facts that  $\|X^t \widehat{H}^t - X^*\|_{2,\infty} \leq \|X^*\|_{2,\infty}$  and  $\|M^*\|_\infty = \|X^*\|_{2,\infty}^2$ . Compared with the Euclidean loss (29), this implies that when  $r = O(1)$ , the entrywise error of  $X^t X^{t\top}$  is uniformly spread out across all entries. As far as we know, this is the first result that reveals near-optimal entrywise error control for noisy matrix completion using nonconvex optimization, without resorting to sample splitting.

**Remark 4** Theorem 2 remains valid if the total number  $T$  of iterations obeys  $T = n^{O(1)}$ . In the noiseless case where  $\sigma = 0$ , the theory allows arbitrarily large  $T$ .

Finally, we report the empirical statistical accuracy of vanilla gradient descent in the presence of noise. Figure 6 displays the squared relative error of vanilla gradient descent as a function of the signal-to-noise ratio (SNR), where the SNR is defined to be

$$\text{SNR} := \frac{\sum_{(j,k) \in \Omega} (M_{j,k}^*)^2}{\sum_{(j,k) \in \Omega} \text{Var}(E_{j,k})} \approx \frac{\|M^*\|_F^2}{n^2 \sigma^2}, \tag{30}$$

and the relative error is measured in terms of the square of the metrics as in (28) as well as the squared entrywise prediction error. Both the relative error and the SNR are shown on a dB scale (i.e.,  $10 \log_{10}(\text{SNR})$  and  $10 \log_{10}(\text{squared relative error})$  are plotted). The results are averaged over 20 independent trials. As one can see from the plot, the squared relative error scales inversely proportional to the SNR, which is consistent with our theory.<sup>6</sup>

### 3.3 Blind Deconvolution

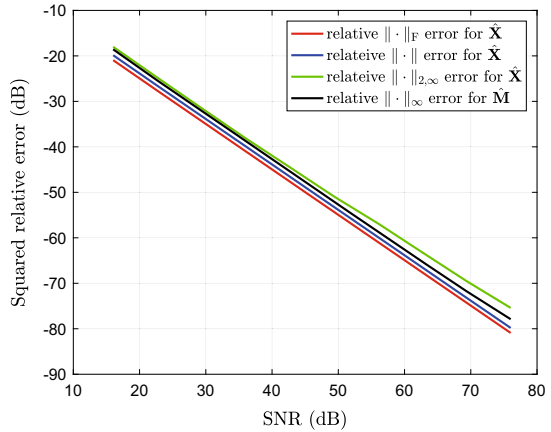
Suppose we have collected  $m$  bilinear measurements

$$y_j = \mathbf{b}_j^H \mathbf{h}^* \mathbf{x}^{*H} \mathbf{a}_j, \quad 1 \leq j \leq m, \tag{31}$$

where  $\mathbf{a}_j$  follows a complex Gaussian distribution, i.e.,  $\mathbf{a}_j \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \frac{1}{2} \mathbf{I}_K)$  +  $i\mathcal{N}(\mathbf{0}, \frac{1}{2} \mathbf{I}_K)$  for  $1 \leq j \leq m$ , and  $\mathbf{B} := [\mathbf{b}_1, \dots, \mathbf{b}_m]^H \in \mathbb{C}^{m \times K}$  is formed by the first

<sup>6</sup> Note that when  $M^*$  is well conditioned and when  $r = O(1)$ , one can easily check that  $\text{SNR} \approx (\|M^*\|_F^2) / (n^2 \sigma^2) \asymp \sigma_{\min}^2 / (n^2 \sigma^2)$ , and our theory says that the squared relative error bound is proportional to  $\sigma^2 / \sigma_{\min}^2$ .

**Fig. 6** Squared relative error of the estimate  $\widehat{\mathbf{X}}$  (measured by  $\|\cdot\|_F, \|\cdot\|, \|\cdot\|_{2,\infty}$  modulo global transformation) and  $\widehat{\mathbf{M}} = \widehat{\mathbf{X}}\widehat{\mathbf{X}}^\top$  (measured by  $\|\cdot\|_\infty$ ) versus SNR for noisy matrix completion, where  $n = 500$ ,  $r = 10$ ,  $p = 0.1$ , and  $\eta_t = 0.2$ . Here  $\widehat{\mathbf{X}}$  denotes the estimate returned by Algorithm 2 after convergence. The results are averaged over 20 independent Monte Carlo trials



$K$  columns of a unitary discrete Fourier transform (DFT) matrix  $\mathbf{F} \in \mathbb{C}^{m \times m}$  obeying  $\mathbf{F}\mathbf{F}^H = \mathbf{I}_m$  (see Appendix D.3.2 for a brief introduction to DFT matrices). This setup models blind deconvolution, where the two signals under convolution belong to known low-dimensional subspaces of dimension  $K$  [3].<sup>7</sup> In particular, the partial DFT matrix  $\mathbf{B}$  plays an important role in image blind deblurring. In this subsection, we consider solving the following nonconvex optimization problem

$$\text{minimize}_{\mathbf{h}, \mathbf{x} \in \mathbb{C}^K} f(\mathbf{h}, \mathbf{x}) = \sum_{j=1}^m \left| \mathbf{b}_j^H \mathbf{h} \mathbf{x}^H \mathbf{a}_j - y_j \right|^2. \tag{32}$$

The (Wirtinger) gradient descent algorithm (with spectral initialization) is summarized in Algorithm 3; here,  $\nabla_{\mathbf{h}} f(\mathbf{h}, \mathbf{x})$  and  $\nabla_{\mathbf{x}} f(\mathbf{h}, \mathbf{x})$  stand for the Wirtinger gradient and are given in (77) and (78), respectively; see [18, Section 6] for a brief introduction to Wirtinger calculus.

It is self-evident that  $\mathbf{h}^*$  and  $\mathbf{x}^*$  are only identifiable up to global scaling, that is, for any nonzero  $\alpha \in \mathbb{C}$ ,

$$\mathbf{h}^* \mathbf{x}^{*H} = \frac{1}{\alpha} \mathbf{h}^* (\alpha \mathbf{x}^*)^H.$$

In light of this, we will measure the discrepancy between

$$\mathbf{z} := \begin{bmatrix} \mathbf{h} \\ \mathbf{x} \end{bmatrix} \in \mathbb{C}^{2K} \quad \text{and} \quad \mathbf{z}^* := \begin{bmatrix} \mathbf{h}^* \\ \mathbf{x}^* \end{bmatrix} \in \mathbb{C}^{2K} \tag{33}$$

via the following function

$$\text{dist}(\mathbf{z}, \mathbf{z}^*) := \min_{\alpha \in \mathbb{C}} \sqrt{\left\| \frac{1}{\alpha} \mathbf{h} - \mathbf{h}^* \right\|_2^2 + \|\alpha \mathbf{x} - \mathbf{x}^*\|_2^2}. \tag{34}$$

<sup>7</sup> For simplicity, we have set the dimensions of the two subspaces equal, and it is straightforward to extend our results to the case of unequal subspace dimensions.

**Algorithm 3** Vanilla gradient descent for blind deconvolution (with spectral initialization)

**Input:**  $\{a_j\}_{1 \leq j \leq m}$ ,  $\{b_j\}_{1 \leq j \leq m}$  and  $\{y_j\}_{1 \leq j \leq m}$ .

**Spectral initialization:** Let  $\sigma_1(\mathbf{M})$ ,  $\check{\mathbf{h}}^0$  and  $\check{\mathbf{x}}^0$  be the leading singular value, left and right singular vectors of

$$\mathbf{M} := \sum_{j=1}^m y_j \mathbf{b}_j \mathbf{a}_j^H,$$

respectively. Set  $\mathbf{h}^0 = \sqrt{\sigma_1(\mathbf{M})} \check{\mathbf{h}}^0$  and  $\mathbf{x}^0 = \sqrt{\sigma_1(\mathbf{M})} \check{\mathbf{x}}^0$ .

**Gradient updates:** for  $t = 0, 1, 2, \dots, T - 1$  do

$$\begin{bmatrix} \mathbf{h}^{t+1} \\ \mathbf{x}^{t+1} \end{bmatrix} = \begin{bmatrix} \mathbf{h}^t \\ \mathbf{x}^t \end{bmatrix} - \eta \begin{bmatrix} \frac{1}{\|\mathbf{x}^t\|_2^2} \nabla_{\mathbf{h}} f(\mathbf{h}^t, \mathbf{x}^t) \\ \frac{1}{\|\mathbf{h}^t\|_2^2} \nabla_{\mathbf{x}} f(\mathbf{h}^t, \mathbf{x}^t) \end{bmatrix}. \tag{35}$$

Before proceeding, we need to introduce the incoherence parameter [3,76], which is crucial for blind deconvolution, whose role is similar to the incoherence parameter (cf. Definition 3) in matrix completion.

**Definition 4** (*Incoherence for blind deconvolution*) Let the incoherence parameter  $\mu$  of  $\mathbf{h}^*$  be the smallest number such that

$$\max_{1 \leq j \leq m} \left| \mathbf{b}_j^H \mathbf{h}^* \right| \leq \frac{\mu}{\sqrt{m}} \|\mathbf{h}^*\|_2. \tag{36}$$

The incoherence parameter describes the spectral flatness of the signal  $\mathbf{h}^*$ . With this definition in place, we have the following theorem, where for identifiability we assume that  $\|\mathbf{h}^*\|_2 = \|\mathbf{x}^*\|_2$ .

**Theorem 3** Suppose the number of measurements obeys  $m \geq C\mu^2 K \log^9 m$  for some sufficiently large constant  $C > 0$ , and suppose the step size  $\eta > 0$  is taken to be some sufficiently small constant. Then there exist constants  $c_1, c_2, C_1, C_3, C_4 > 0$  such that with probability exceeding  $1 - c_1 m^{-5} - c_1 m e^{-c_2 K}$ , the iterates in Algorithm 3 satisfy

$$\text{dist}(\mathbf{z}^t, \mathbf{z}^*) \leq C_1 \left(1 - \frac{\eta}{16}\right)^t \frac{1}{\log^2 m} \|\mathbf{z}^*\|_2, \tag{37a}$$

$$\max_{1 \leq l \leq m} \left| \mathbf{a}_l^H (\alpha^t \mathbf{x}^t - \mathbf{x}^*) \right| \leq C_3 \frac{1}{\log^{1.5} m} \|\mathbf{x}^*\|_2, \tag{37b}$$

$$\max_{1 \leq l \leq m} \left| \mathbf{b}_l^H \frac{1}{\alpha^t} \mathbf{h}^t \right| \leq C_4 \frac{\mu}{\sqrt{m}} \log^2 m \|\mathbf{h}^*\|_2 \tag{37c}$$

for all  $t \geq 0$ . Here, we denote  $\alpha^t$  as the alignment parameter,

$$\alpha^t := \arg \min_{\alpha \in \mathbb{C}} \left\| \frac{1}{\alpha} \mathbf{h}^t - \mathbf{h}^* \right\|_2^2 + \|\alpha \mathbf{x}^t - \mathbf{x}^*\|_2^2. \tag{38}$$

Theorem 3 provides the first theoretical guarantee of unregularized gradient descent for blind deconvolution at a near-optimal statistical and computational complexity. A few remarks are in order.

- *Implicit regularization* Theorem 3 reveals that the unregularized gradient descent iterates remain incoherent with the sampling mechanism (see (37b) and (37c)). Recall that prior works operate upon a regularized cost function with an additional penalty term that regularizes the global scaling  $\{\|\mathbf{h}\|_2, \|\mathbf{x}\|_2\}$  and the incoherence  $\{\|\mathbf{b}_j^H \mathbf{h}\|\}_{1 \leq j \leq m}$  [58,76,82]. In comparison, our theorem implies that it is unnecessary to regularize either the incoherence or the scaling ambiguity, which is somewhat surprising. This justifies the use of regularization-free (Wirtinger) gradient descent for blind deconvolution.
- *Constant step size* Compared to the step size  $\eta_t \lesssim 1/m$  suggested in [76] for regularized gradient descent, our theory admits a substantially more aggressive step size (i.e.,  $\eta_t \asymp 1$ ) even without regularization. Similar to phase retrieval, the computational efficiency is boosted by a factor of  $m$ , attaining  $\epsilon$ -accuracy within  $O(\log(1/\epsilon))$  iterations (vs.  $O(m \log(1/\epsilon))$  iterations in prior theory).
- *Near-optimal sample complexity* It is demonstrated that vanilla gradient descent succeeds at a near-optimal sample complexity up to logarithmic factors, although our requirement is slightly worse than [76] which uses explicit regularization. Notably, even under the sample complexity herein, the iteration complexity given in [76] is still  $O(m/\text{poly}(\log(m)))$ .
- *Incoherence of spectral initialization* As in phase retrieval, Theorem 3 demonstrates that the estimates returned by the spectral method are incoherent with respect to both  $\{\mathbf{a}_j\}$  and  $\{\mathbf{b}_j\}$ . In contrast, [76] recommends a projection operation (via a linear program) to enforce incoherence of the initial estimates, which is dispensable according to our theory.
- *Contraction in  $\|\cdot\|_F$*  It is easy to check that the Frobenius norm error satisfies  $\|\mathbf{h}^t \mathbf{x}^{tH} - \mathbf{h}^* \mathbf{x}^{*H}\|_F \lesssim \text{dist}(\mathbf{z}^t, \mathbf{z}^*)$ , and therefore, Theorem 3 corroborates the empirical results shown in Fig. 1c.

## 4 Related Work

Solving nonlinear systems of equations has received much attention in the past decade. Rather than directly attacking the nonconvex formulation, convex relaxation lifts the object of interest into a higher-dimensional space and then attempts recovery via semidefinite programming (e.g., [3,19,20,94]). This has enjoyed great success in both theory and practice. Despite appealing statistical guarantees, semidefinite programming is in general prohibitively expensive when processing large-scale datasets.

Nonconvex approaches, on the other end, have been under extensive study in the last few years, due to their computational advantages. There is a growing list of statistical estimation problems for which nonconvex approaches are guaranteed to find global optimal solutions, including but not limited to phase retrieval [18,25,90], low-rank matrix sensing and completion [7,32,48,115,130], blind deconvolution and self-calibration [72,76,78,82], dictionary learning [106], tensor decomposition [49],

joint alignment [24], learning shallow neural networks [103,132], robust subspace learning [34,74,86,91]. In several problems [40,48,49,75,77,86,87,105,106], it is further suggested that the optimization landscape is benign under sufficiently large sample complexity, in the sense that all local minima are globally optimal, and hence non-convex iterative algorithms become promising in solving such problems. See [37] for a recent overview. Below we review the three problems studied in this paper in more detail. Some state-of-the-art results are summarized in Table 1.

- *Phase retrieval* Candès et al. proposed *PhaseLift* [20] to solve the quadratic systems of equations based on convex programming. Specifically, it lifts the decision variable  $\mathbf{x}^*$  into a rank-one matrix  $\mathbf{X}^* = \mathbf{x}^* \mathbf{x}^{*\top}$  and translates the quadratic constraints of  $\mathbf{x}^*$  in (14) into linear constraints of  $\mathbf{X}^*$ . By dropping the rank constraint, the problem becomes convex [11,16,20,29,113]. Another convex program PhaseMax [5,41,50,53] operates in the natural parameter space via linear programming, provided that an anchor vector is available. On the other hand, alternating minimization [90] with sample splitting has been shown to enjoy much better computational guarantee. In contrast, Wirtinger flow [18] provides the first global convergence result for nonconvex methods without sample splitting, whose statistical and computational guarantees are later improved by [25] via an adaptive truncation strategy. Several other variants of WF are also proposed [12,68,102], among which an amplitude-based loss function has been investigated [117–119,127]. In particular, [127] demonstrates that the amplitude-based loss function has a better curvature, and vanilla gradient descent can indeed converge with a constant step size at the orderwise optimal sample complexity. A small sample of other nonconvex phase retrieval methods include [6,10,22,36,43,47,92,98,100,109,122], which are beyond the scope of this paper.
- *Matrix completion* Nuclear norm minimization was studied in [19] as a convex relaxation paradigm to solve the matrix completion problem. Under certain incoherence conditions imposed upon the ground truth matrix, exact recovery is guaranteed under near-optimal sample complexity [14,23,38,51,93]. Concurrently, several works [54,55,60,61,63–65,71,110,123,129,129] tackled the matrix completion problem via nonconvex approaches. In particular, the seminal work by Keshavan et al. [64,65] pioneered the two-stage approach that is widely adopted by later works. Sun and Luo [107] demonstrated the convergence of gradient descent type methods for noiseless matrix completion with a regularized nonconvex loss function. Instead of penalizing the loss function, [32,131] employed projection to enforce the incoherence condition throughout the execution of the algorithm. To the best of our knowledge, no rigorous guarantees have been established for matrix completion without explicit regularization. A notable exception is [63], which uses unregularized stochastic gradient descent for matrix completion in the online setting. However, the analysis is performed with fresh samples in each iteration. Our work closes the gap and makes the first contribution toward understanding implicit regularization in gradient descent without sample splitting. In addition, entrywise eigenvector perturbation has been studied by [1,26,60] in order to analyze the spectral algorithms for matrix completion, which helps us establish theoretical guarantees for the spectral initialization step. Finally, it has recently been shown



that the analysis of nonconvex gradient descent in turn yields near-optimal statistical guarantees for convex relaxation in the context of noisy matrix completion; see [28,31].

- *Blind deconvolution* In [3], Ahmed et al. first proposed to invoke similar lifting ideas for blind deconvolution, which translates the bilinear measurements (31) into a system of linear measurements of a rank-one matrix  $\mathbf{X}^* = \mathbf{h}^* \mathbf{x}^{*H}$ . Near-optimal performance guarantees have been established for convex relaxation [3]. Under the same model, Li et al. [76] proposed a regularized gradient descent algorithm that directly optimizes the nonconvex loss function (32) with a few regularization terms that account for scaling ambiguity and incoherence. In [58], a Riemannian steepest descent method is developed that removes the regularization for scaling ambiguity, although they still need to regularize for incoherence. In [2], a linear program is proposed but requires exact knowledge of the signs of the signals. Blind deconvolution has also been studied for other models—interested readers are referred to [35,72,73,81,82,120,128].

On the other hand, our analysis framework is based on a leave-one-out perturbation argument. This technique has been widely used to analyze high-dimensional problems with random designs, including but not limited to robust M-estimation [44,45], statistical inference for sparse regression [62], likelihood ratio test in logistic regression [108], phase synchronization [1,133], ranking from pairwise comparisons [30], community recovery [1], and covariance sketching [79]. In particular, this technique results in tight performance guarantees for the generalized power method [133], the spectral method [1,30], and convex programming approaches [30,44,108,133]; however, it has not been applied to analyze nonconvex optimization algorithms.

Finally, we note that the notion of implicit regularization—broadly defined—arises in settings far beyond the models and algorithms considered herein. For instance, it has been conjectured that in matrix factorization, over-parameterized stochastic gradient descent effectively enforces certain norm constraints, allowing it to converge to a minimal-norm solution as long as it starts from the origin [52]. The stochastic gradient methods have also been shown to implicitly enforce Tikhonov regularization in several statistical learning settings [80]. More broadly, this phenomenon seems crucial in enabling efficient training of deep neural networks [104,125].

## 5 A General Recipe for Trajectory Analysis

In this section, we sketch a general recipe for establishing performance guarantees of gradient descent, which conveys the key idea for proving the main results of this paper. The main challenge is to demonstrate that appropriate incoherence conditions are preserved throughout the trajectory of the algorithm. This requires exploiting statistical independence of the samples in a careful manner, in conjunction with generic optimization theory. Central to our approach is a leave-one-out perturbation argument, which allows to decouple the statistical dependency while controlling the component-wise incoherence measures.

### General Recipe (a leave-one-out analysis)

- Step 1:** characterize restricted strong convexity and smoothness of  $f$ , and identify the region of incoherence and contraction (RIC).
- Step 2:** introduce leave-one-out sequences  $\{X^{t,(l)}\}$  and  $\{H^{t,(l)}\}$  for each  $l$ , where  $\{X^{t,(l)}\}$  (resp.  $\{H^{t,(l)}\}$ ) is independent of any sample involving  $\phi_l$  (resp.  $\psi_l$ );
- Step 3:** establish the incoherence condition for  $\{X^t\}$  and  $\{H^t\}$  via induction. Suppose the iterates satisfy the claimed conditions in the  $t$ th iteration:
- show, via restricted strong convexity, that the true iterates  $(X^{t+1}, H^{t+1})$  and the leave-one-out version  $(X^{t+1,(l)}, H^{t+1,(l)})$  are exceedingly close;
  - use statistical independence to show that  $X^{t+1,(l)} - X^*$  (resp.  $H^{t+1,(l)} - H^*$ ) is incoherent w.r.t.  $\phi_l$  (resp.  $\psi_l$ ); namely,  $\|\phi_l^H(X^{t+1,(l)} - X^*)\|_2$  and  $\|\psi_l^H(H^{t+1,(l)} - H^*)\|_2$  are both well controlled;
  - combine the bounds to establish the desired incoherence condition concerning  $\max_l \|\phi_l^H(X^{t+1} - X^*)\|_2$  and  $\max_l \|\psi_l^H(H^{t+1} - H^*)\|_2$ .

## 5.1 General Model

Consider the following problem where the samples are collected in a bilinear/quadratic form as

$$y_j = \psi_j^H H^* X^{*H} \phi_j, \quad 1 \leq j \leq m, \quad (39)$$

where the objects of interest  $H^*$ ,  $X^* \in \mathbb{C}^{n \times r}$  or  $\mathbb{R}^{n \times r}$  might be vectors or tall matrices taking either real or complex values. The design vectors  $\{\psi_j\}$  and  $\{\phi_j\}$  are in either  $\mathbb{C}^n$  or  $\mathbb{R}^n$ , and can be either random or deterministic. This model is quite general and entails all three examples in this paper as special cases:

- *Phase retrieval:*  $H^* = X^* = x^* \in \mathbb{R}^n$ , and  $\psi_j = \phi_j = a_j$ ;
- *Matrix completion:*  $H^* = X^* \in \mathbb{R}^{n \times r}$  and  $\psi_j, \phi_j \in \{e_1, \dots, e_n\}$ ;
- *Blind deconvolution:*  $H^* = h^* \in \mathbb{C}^K$ ,  $X^* = x^* \in \mathbb{C}^K$ ,  $\phi_j = a_j$ , and  $\psi_j = b_j$ .

For this setting, the empirical loss function is given by

$$f(Z) := f(H, X) = \frac{1}{m} \sum_{j=1}^m \left| \psi_j^H H X^H \phi_j - y_j \right|^2,$$

where we denote  $Z = (H, X)$ . To minimize  $f(Z)$ , we proceed with vanilla gradient descent

$$Z^{t+1} = Z^t - \eta \nabla f(Z^t), \quad \forall t \geq 0$$

following a standard spectral initialization, where  $\eta$  is the step size. As a remark, for complex-valued problems, the gradient (resp. Hessian) should be understood as the Wirtinger gradient (resp. Hessian).

It is clear from (39) that  $\mathbf{Z}^* = (\mathbf{H}^*, \mathbf{X}^*)$  can only be recovered up to certain global ambiguity. For clarity of presentation, we assume in this section that such ambiguity has already been taken care of via proper global transformation.

### 5.2 Outline of the Recipe

We are now positioned to outline the general recipe, which entails the following steps.

- *Step 1: characterizing local geometry in the RIC* Our first step is to characterize a region  $\mathcal{R}$ —which we term as the *region of incoherence and contraction* (RIC)—such that the Hessian matrix  $\nabla^2 f(\mathbf{Z})$  obeys strong convexity and smoothness,

$$\mathbf{0} \prec \alpha \mathbf{I} \preceq \nabla^2 f(\mathbf{Z}) \preceq \beta \mathbf{I}, \quad \forall \mathbf{Z} \in \mathcal{R}, \tag{40}$$

or at least along certain directions (i.e., restricted strong convexity and smoothness), where  $\beta/\alpha$  scales slowly (or even remains bounded) with the problem size. As revealed by optimization theory, this geometric property (40) immediately implies linear convergence with the contraction rate  $1 - O(\alpha/\beta)$  for a properly chosen step size  $\eta$ , as long as all iterates stay within the RIC.

A natural question then arises: What does the RIC  $\mathcal{R}$  look like? As it turns out, the RIC typically contains all points such that the  $\ell_2$  error  $\|\mathbf{Z} - \mathbf{Z}^*\|_F$  is not too large and

$$\begin{aligned} \text{(incoherence)} \quad & \max_j \|\phi_j^H(\mathbf{X} - \mathbf{X}^*)\|_2 \text{ and } \max_j \|\psi_j^H(\mathbf{H} - \mathbf{H}^*)\|_2 \\ & \text{are well controlled.} \end{aligned} \tag{41}$$

In the three examples, the above incoherence condition translates to:

- *Phase retrieval*:  $\max_j |\mathbf{a}_j^T(\mathbf{x} - \mathbf{x}^*)|$  is well controlled;
  - *Matrix completion*:  $\|\mathbf{X} - \mathbf{X}^*\|_{2,\infty}$  is well controlled;
  - *Blind deconvolution*:  $\max_j |\mathbf{a}_j^T(\mathbf{x} - \mathbf{x}^*)|$  and  $\max_j |\mathbf{b}_j^T(\mathbf{h} - \mathbf{h}^*)|$  are well controlled.
- *Step 2: introducing the leave-one-out sequences* To justify that no iterates leave the RIC, we rely on the construction of auxiliary sequences. Specifically, for each  $l$ , produce an auxiliary sequence  $\{\mathbf{Z}^{t,(l)} = (\mathbf{X}^{t,(l)}, \mathbf{H}^{t,(l)})\}$  such that  $\mathbf{X}^{t,(l)}$  (resp.  $\mathbf{H}^{t,(l)}$ ) is independent of any sample involving  $\phi_l$  (resp.  $\psi_l$ ). As an example, suppose that the  $\phi_l$ 's and the  $\psi_l$ 's are independently and randomly generated. Then for each  $l$ , one can consider a leave-one-out loss function

$$f^{(l)}(\mathbf{Z}) := \frac{1}{m} \sum_{j:j \neq l} \left| \psi_j^H \mathbf{H} \mathbf{X}^H \phi_j - y_j \right|^2$$

that discards the  $l$ th sample. One further generates  $\{\mathbf{Z}^{t,(l)}\}$  by running vanilla gradient descent w.r.t. this auxiliary loss function, with a spectral initialization that similarly discards the  $l$ th sample. Note that this procedure is only introduced to facilitate analysis and is never implemented in practice.

- *Step 3: establishing the incoherence condition* We are now ready to establish the incoherence condition with the assistance of the auxiliary sequences. Usually, the proof proceeds by induction, where our goal is to show that the next iterate remains within the RIC, given that the current one does.

- *Step 3(a): proximity between the original and the leave-one-out iterates* As one can anticipate,  $\{\mathbf{Z}^t\}$  and  $\{\mathbf{Z}^{t,(l)}\}$  remain “glued” to each other along the whole trajectory, since their constructions differ by only a single sample. In fact, as long as the initial estimates stay sufficiently close, their gaps will never explode. To intuitively see why, use the fact  $\nabla f(\mathbf{Z}^t) \approx \nabla f^{(l)}(\mathbf{Z}^t)$  to discover that

$$\begin{aligned} \mathbf{Z}^{t+1} - \mathbf{Z}^{t+1,(l)} &= \mathbf{Z}^t - \eta \nabla f(\mathbf{Z}^t) - (\mathbf{Z}^{t,(l)} - \eta \nabla f^{(l)}(\mathbf{Z}^{t,(l)})) \\ &\approx \mathbf{Z}^t - \mathbf{Z}^{t,(l)} - \eta \nabla^2 f(\mathbf{Z}^t)(\mathbf{Z}^t - \mathbf{Z}^{t,(l)}), \end{aligned}$$

which together with the strong convexity condition implies  $\ell_2$  contraction

$$\|\mathbf{Z}^{t+1} - \mathbf{Z}^{t+1,(l)}\|_F \approx \|(I - \eta \nabla^2 f(\mathbf{Z}^t))(\mathbf{Z}^t - \mathbf{Z}^{t,(l)})\|_F \leq \|\mathbf{Z}^t - \mathbf{Z}^{t,(l)}\|_2.$$

Indeed, (restricted) strong convexity is crucial in controlling the size of leave-one-out perturbations.

- *Step 3(b): incoherence condition of the leave-one-out iterates* The fact that  $\mathbf{Z}^{t+1}$  and  $\mathbf{Z}^{t+1,(l)}$  are exceedingly close motivates us to control the incoherence of  $\mathbf{Z}^{t+1,(l)} - \mathbf{X}^*$  instead, for  $1 \leq l \leq m$ . By construction,  $\mathbf{X}^{t+1,(l)}$  (resp.  $\mathbf{H}^{t+1,(l)}$ ) is statistically *independent* of any sample involving the design vector  $\phi_l$  (resp.  $\psi_l$ ), a fact that typically leads to a more friendly analysis for controlling  $\|\phi_l^H(\mathbf{X}^{t+1,(l)} - \mathbf{X}^*)\|_2$  and  $\|\psi_l^H(\mathbf{H}^{t+1,(l)} - \mathbf{H}^*)\|_2$ .
- *Step 3(c): combining the bounds* With these results in place, apply the triangle inequality to obtain

$$\|\phi_l^H(\mathbf{X}^{t+1} - \mathbf{X}^*)\|_2 \leq \|\phi_l\|_2 \|\mathbf{X}^{t+1} - \mathbf{X}^{t+1,(l)}\|_F + \|\phi_l^H(\mathbf{X}^{t+1,(l)} - \mathbf{X}^*)\|_2,$$

where the first term is controlled in Step 3(a) and the second term is controlled in Step 3(b). The term  $\|\psi_l^H(\mathbf{H}^{t+1} - \mathbf{H}^*)\|_2$  can be bounded similarly. By choosing the bounds properly, this establishes the incoherence condition for all  $1 \leq l \leq m$  as desired.

## 6 Analysis for Phase Retrieval

In this section, we instantiate the general recipe presented in Sect. 5 to phase retrieval and prove Theorem 1. Similar to the section 7.1 in [18], we are going to use  $\eta_t =$

$c_1/(\log n \cdot \|\mathbf{x}^*\|_2^2)$  instead of  $c_1/(\log n \cdot \|\mathbf{x}_0\|_2^2)$  as the step size for analysis. This is because with high probability,  $\|\mathbf{x}_0\|_2$  and  $\|\mathbf{x}^*\|_2$  are rather close in the relative sense. Without loss of generality, we assume throughout this section that  $\|\mathbf{x}^*\|_2 = 1$  and

$$\text{dist}(\mathbf{x}^0, \mathbf{x}^*) = \|\mathbf{x}^0 - \mathbf{x}^*\|_2 \leq \|\mathbf{x}^0 + \mathbf{x}^*\|_2. \tag{42}$$

In addition, the gradient and the Hessian of  $f(\cdot)$  for this problem (see (15)) are given, respectively, by

$$\nabla f(\mathbf{x}) = \frac{1}{m} \sum_{j=1}^m \left[ \left( \mathbf{a}_j^\top \mathbf{x} \right)^2 - y_j \right] \left( \mathbf{a}_j^\top \mathbf{x} \right) \mathbf{a}_j, \tag{43}$$

$$\nabla^2 f(\mathbf{x}) = \frac{1}{m} \sum_{j=1}^m \left[ 3 \left( \mathbf{a}_j^\top \mathbf{x} \right)^2 - y_j \right] \mathbf{a}_j \mathbf{a}_j^\top, \tag{44}$$

which are useful throughout the proof.

### 6.1 Step 1: Characterizing Local Geometry in the RIC

#### 6.1.1 Local Geometry

We start by characterizing the region that enjoys both strong convexity and the desired level of smoothness. This is supplied in the following lemma, which plays a crucial role in the subsequent analysis.

**Lemma 1** (Restricted strong convexity and smoothness for phase retrieval) *Fix any sufficiently small constant  $C_1 > 0$  and any sufficiently large constant  $C_2 > 0$ , and suppose the sample complexity obeys  $m \geq c_0 n \log n$  for some sufficiently large constant  $c_0 > 0$ . With probability at least  $1 - O(mn^{-10})$ ,*

$$\nabla^2 f(\mathbf{x}) \succeq (1/2) \cdot \mathbf{I}_n$$

holds simultaneously for all  $\mathbf{x} \in \mathbb{R}^n$  satisfying  $\|\mathbf{x} - \mathbf{x}^*\|_2 \leq 2C_1$ , and

$$\nabla^2 f(\mathbf{x}) \preceq (5C_2(10 + C_2) \log n) \cdot \mathbf{I}_n$$

holds simultaneously for all  $\mathbf{x} \in \mathbb{R}^n$  obeying

$$\|\mathbf{x} - \mathbf{x}^*\|_2 \leq 2C_1, \tag{45a}$$

$$\max_{1 \leq j \leq m} \left| \mathbf{a}_j^\top (\mathbf{x} - \mathbf{x}^*) \right| \leq C_2 \sqrt{\log n}. \tag{45b}$$

**Proof** See Appendix A.1. □

In words, Lemma 1 reveals that the Hessian matrix is positive definite and (almost) well conditioned, if one restricts attention to the set of points that are (i) not far away from the truth (cf. (45a)) and (ii) incoherent with respect to the measurement vectors  $\{\mathbf{a}_j\}_{1 \leq j \leq m}$  (cf. (45b)).

### 6.1.2 Error Contraction

As we point out before, the nice local geometry enables  $\ell_2$  contraction, which we formalize below.

**Lemma 2** *There exists an event that does not depend on  $t$  and has probability  $1 - O(mn^{-10})$ , such that when it happens and  $\mathbf{x}^t$  obeys conditions (45), one has*

$$\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2 \leq (1 - \eta/2) \|\mathbf{x}^t - \mathbf{x}^*\|_2 \tag{46}$$

provided that the step size satisfies  $0 < \eta \leq 1/[5C_2(10 + C_2)\log n]$ .

**Proof** This proof applies the standard argument when establishing the  $\ell_2$  error contraction of gradient descent for strongly convex and smooth functions. See Appendix A.2. □

With the help of Lemma 2, we can turn the proof of Theorem 1 into ensuring that the trajectory  $\{\mathbf{x}^t\}_{0 \leq t \leq n}$  lies in the RIC specified by (47).<sup>8</sup> This is formally stated in the next lemma.

**Lemma 3** *Suppose for all  $0 \leq t \leq T_0 := n$ , the trajectory  $\{\mathbf{x}^t\}$  falls within the region of incoherence and contraction (termed the RIC), namely*

$$\|\mathbf{x}^t - \mathbf{x}^*\|_2 \leq C_1, \tag{47a}$$

$$\max_{1 \leq l \leq m} |\mathbf{a}_l^\top (\mathbf{x}^t - \mathbf{x}^*)| \leq C_2 \sqrt{\log n}, \tag{47b}$$

then the claims in Theorem 1 hold true. Here and throughout this section,  $C_1, C_2 > 0$  are two absolute constants as specified in Lemma 1.

**Proof** See Appendix A.3. □

## 6.2 Step 2: Introducing the Leave-One-Out Sequences

In comparison with the  $\ell_2$  error bound (47a) that captures the overall loss, the incoherence hypothesis (47b)—which concerns sample-wise control of the empirical risk—is more complicated to establish. This is partly due to the statistical dependence between  $\mathbf{x}^t$  and the sampling vectors  $\{\mathbf{a}_l\}$ . As described in the general recipe, the key idea is

---

<sup>8</sup> Here, we deliberately change  $2C_1$  in (45a) to  $C_1$  in the definition of the RIC (47a) to ensure the correctness of the analysis.

the introduction of a *leave-one-out* version of the WF iterates, which removes a single measurement from consideration.

To be precise, for each  $1 \leq l \leq m$ , we define the leave-one-out empirical loss function as

$$f^{(l)}(\mathbf{x}) := \frac{1}{4m} \sum_{j:j \neq l} \left[ \left( \mathbf{a}_j^\top \mathbf{x} \right)^2 - y_j \right]^2, \tag{48}$$

and the auxiliary trajectory  $\{\mathbf{x}^{t,(l)}\}_{t \geq 0}$  is constructed by running WF w.r.t.  $f^{(l)}(\mathbf{x})$ . In addition, the spectral initialization  $\mathbf{x}^{0,(l)}$  is computed based on the rescaled leading eigenvector of the leave-one-out data matrix

$$\mathbf{Y}^{(l)} := \frac{1}{m} \sum_{j:j \neq l} y_j \mathbf{a}_j \mathbf{a}_j^\top. \tag{49}$$

Clearly, the entire sequence  $\{\mathbf{x}^{t,(l)}\}_{t \geq 0}$  is independent of the  $l$ th sampling vector  $\mathbf{a}_l$ . This auxiliary procedure is formally described in Algorithm 4.

**Algorithm 4** The  $l$ th leave-one-out sequence for phase retrieval

**Input:**  $\{\mathbf{a}_j\}_{1 \leq j \leq m, j \neq l}$  and  $\{y_j\}_{1 \leq j \leq m, j \neq l}$ .

**Spectral initialization:** let  $\lambda_1(\mathbf{Y}^{(l)})$  and  $\tilde{\mathbf{x}}^{0,(l)}$  be the leading eigenvalue and eigenvector of

$$\mathbf{Y}^{(l)} = \frac{1}{m} \sum_{j:j \neq l} y_j \mathbf{a}_j \mathbf{a}_j^\top,$$

respectively, and set

$$\mathbf{x}^{0,(l)} = \begin{cases} \sqrt{\lambda_1(\mathbf{Y}^{(l)})/3} \tilde{\mathbf{x}}^{0,(l)}, & \text{if } \|\tilde{\mathbf{x}}^{0,(l)} - \mathbf{x}^*\|_2 \leq \|\tilde{\mathbf{x}}^{0,(l)} + \mathbf{x}^*\|_2, \\ -\sqrt{\lambda_1(\mathbf{Y}^{(l)})/3} \tilde{\mathbf{x}}^{0,(l)}, & \text{else.} \end{cases}$$

**Gradient updates:** for  $t = 0, 1, 2, \dots, T - 1$  do

$$\mathbf{x}^{t+1,(l)} = \mathbf{x}^{t,(l)} - \eta_t \nabla f^{(l)}(\mathbf{x}^{t,(l)}). \tag{50}$$

**6.3 Step 3: Establishing the Incoherence Condition by Induction**

As revealed by Lemma 3, it suffices to prove that the iterates  $\{\mathbf{x}^t\}_{0 \leq t \leq T_0}$  satisfies (47) with high probability. Our proof will be inductive in nature. For the sake of clarity, we list all the induction hypotheses:

$$\|\mathbf{x}^t - \mathbf{x}^*\|_2 \leq C_1, \tag{51a}$$

$$\max_{1 \leq l \leq m} \|\mathbf{x}^t - \mathbf{x}^{t,(l)}\|_2 \leq C_3 \sqrt{\frac{\log n}{n}} \tag{51b}$$

$$\max_{1 \leq j \leq m} \left| \mathbf{a}_j^\top (\mathbf{x}^t - \mathbf{x}^*) \right| \leq C_2 \sqrt{\log n}. \tag{51c}$$

Here  $C_3 > 0$  is some universal constant. For any  $t \geq 0$ , define  $\mathcal{E}_t$  to be the event where the conditions in (51) hold for the  $t$ th iteration. According to Lemma 2, there exists some event  $\mathcal{E}$  with probability  $1 - O(mn^{-10})$  such that on  $\mathcal{E}_t \cap \mathcal{E}$  one has

$$\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2 \leq C_1. \tag{52}$$

This subsection is devoted to establishing (51b) and (51c) for the  $(t + 1)$ th iteration, assuming that (51) holds true up to the  $t$ th iteration. We defer the justification of the base case (i.e., initialization at  $t = 0$ ) to Sect. 6.4.

- *Step 3(a): proximity between the original and the leave-one-out iterates* The leave-one-out sequence  $\{\mathbf{x}^{t,(l)}\}$  behaves similarly to the true WF iterates  $\{\mathbf{x}^t\}$  while maintaining statistical independence with  $\mathbf{a}_l$ , a key fact that allows us to control the incoherence of  $l$ th leave-one-out sequence w.r.t.  $\mathbf{a}_l$ . We will formally quantify the gap between  $\mathbf{x}^{t+1}$  and  $\mathbf{x}^{t+1,(l)}$  in the following lemma, which establishes the induction in (51b).

**Lemma 4** *Suppose that the sample size obeys  $m \geq Cn \log n$  for some sufficiently large constant  $C > 0$  and that the step size obeys  $0 < \eta < 1/[5C_2(10 + C_2) \log n]$ . Then on some event  $\mathcal{E}_{t+1,1} \subseteq \mathcal{E}_t$  obeying  $\mathbb{P}(\mathcal{E}_t \cap \mathcal{E}_{t+1,1}^c) = O(mn^{-10})$ , one has*

$$\max_{1 \leq l \leq m} \|\mathbf{x}^{t+1} - \mathbf{x}^{t+1,(l)}\|_2 \leq C_3 \sqrt{\frac{\log n}{n}}. \tag{53}$$

**Proof** The proof relies heavily on the restricted strong convexity (see Lemma 1) and is deferred to Appendix A.4. □

- *Step 3(b): incoherence of the leave-one-out iterates* By construction,  $\mathbf{x}^{t+1,(l)}$  is statistically independent of the sampling vector  $\mathbf{a}_l$ . One can thus invoke the standard Gaussian concentration results and the union bound to derive that on an event  $\mathcal{E}_{t+1,2} \subseteq \mathcal{E}_t$  obeying  $\mathbb{P}(\mathcal{E}_t \cap \mathcal{E}_{t+1,2}^c) = O(mn^{-10})$ ,

$$\begin{aligned} \max_{1 \leq l \leq m} \left| \mathbf{a}_l^\top (\mathbf{x}^{t+1,(l)} - \mathbf{x}^*) \right| &\leq 5\sqrt{\log n} \|\mathbf{x}^{t+1,(l)} - \mathbf{x}^*\|_2 \\ &\stackrel{(i)}{\leq} 5\sqrt{\log n} \left( \|\mathbf{x}^{t+1,(l)} - \mathbf{x}^{t+1}\|_2 + \|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2 \right) \\ &\stackrel{(ii)}{\leq} 5\sqrt{\log n} \left( C_3 \sqrt{\frac{\log n}{n}} + C_1 \right) \\ &\leq C_4 \sqrt{\log n} \end{aligned} \tag{54}$$



holds for some constant  $C_4 \geq 6C_1 > 0$  and  $n$  sufficiently large. Here, (i) comes from the triangle inequality and (ii) arises from the proximity bound (53) and the condition (52).

- *Step 3(c): combining the bounds* We are now prepared to establish (51c) for the  $(t + 1)$ th iteration. Specifically,

$$\begin{aligned} \max_{1 \leq l \leq m} \left| \mathbf{a}_l^\top (\mathbf{x}^{t+1} - \mathbf{x}^*) \right| &\leq \max_{1 \leq l \leq m} \left| \mathbf{a}_l^\top (\mathbf{x}^{t+1} - \mathbf{x}^{t+1,(l)}) \right| + \max_{1 \leq l \leq m} \left| \mathbf{a}_l^\top (\mathbf{x}^{t+1,(l)} - \mathbf{x}^*) \right| \\ &\stackrel{(i)}{\leq} \max_{1 \leq l \leq m} \|\mathbf{a}_l\|_2 \|\mathbf{x}^{t+1} - \mathbf{x}^{t+1,(l)}\|_2 + C_4 \sqrt{\log n} \\ &\stackrel{(ii)}{\leq} \sqrt{6n} \cdot C_3 \sqrt{\frac{\log n}{n}} + C_4 \sqrt{\log n} \leq C_2 \sqrt{\log n}, \end{aligned} \tag{55}$$

where (i) follows from the Cauchy–Schwarz inequality and (54), the inequality (ii) is a consequence of (53) and (98), and the last inequality holds as long as  $C_2/(C_3 + C_4)$  is sufficiently large. From the deduction above we easily get  $\mathbb{P}(\mathcal{E}_t \cap \mathcal{E}_{t+1}^c) = O(mn^{-10})$ .

Using mathematical induction and the union bound, we establish (51) for all  $t \leq T_0 = n$  with high probability. This in turn concludes the proof of Theorem 1, as long as the hypotheses are valid for the base case.

### 6.4 The Base Case: Spectral Initialization

In the end, we return to verify the induction hypotheses for the base case ( $t = 0$ ), i.e., the spectral initialization obeys (51). The following lemma justifies (51a) by choosing  $\delta$  sufficiently small.

**Lemma 5** *Fix any small constant  $\delta > 0$ , and suppose  $m > c_0 n \log n$  for some large constant  $c_0 > 0$ . Consider the two vectors  $\mathbf{x}^0$  and  $\tilde{\mathbf{x}}^0$  as defined in Algorithm 1, and suppose without loss of generality that (42) holds. Then with probability exceeding  $1 - O(n^{-10})$ , one has*

$$\|\mathbf{Y} - \mathbb{E}[\mathbf{Y}]\| \leq \delta, \tag{56}$$

$$\|\mathbf{x}^0 - \mathbf{x}^*\|_2 \leq 2\delta \quad \text{and} \quad \|\tilde{\mathbf{x}}^0 - \mathbf{x}^*\|_2 \leq \sqrt{2}\delta. \tag{57}$$

**Proof** This result follows directly from the Davis–Kahan  $\sin\Theta$  theorem. See Appendix A.5. □

We then move on to justifying (51b), the proximity between the original and leave-one-out iterates for  $t = 0$ .

**Lemma 6** *Suppose  $m > c_0 n \log n$  for some large constant  $c_0 > 0$ . Then with probability at least  $1 - O(mn^{-10})$ , one has*

$$\max_{1 \leq l \leq m} \|\mathbf{x}^0 - \mathbf{x}^{0,(l)}\|_2 \leq C_3 \sqrt{\frac{\log n}{n}}. \tag{58}$$

**Proof** This is also a consequence of the Davis–Kahan  $\sin\Theta$  theorem. See Appendix A.6. □

The final claim (51c) can be proved using the same argument as in deriving (55) and hence is omitted.

## 7 Analysis for Matrix Completion

In this section, we instantiate the general recipe presented in Sect. 5 to matrix completion and prove Theorem 2. Before continuing, we first gather a few useful facts regarding the loss function in (23). The gradient of it is given by

$$\nabla f(\mathbf{X}) = \frac{1}{p} \mathcal{P}_\Omega \left[ \mathbf{X}\mathbf{X}^\top - (\mathbf{M}^\star + \mathbf{E}) \right] \mathbf{X}. \tag{59}$$

We define the expected gradient (with respect to the sampling set  $\Omega$ ) to be

$$\nabla F(\mathbf{X}) = \left[ \mathbf{X}\mathbf{X}^\top - (\mathbf{M}^\star + \mathbf{E}) \right] \mathbf{X}$$

and also the (expected) gradient without noise to be

$$\nabla f_{\text{clean}}(\mathbf{X}) = \frac{1}{p} \mathcal{P}_\Omega \left( \mathbf{X}\mathbf{X}^\top - \mathbf{M}^\star \right) \mathbf{X} \quad \text{and} \quad \nabla F_{\text{clean}}(\mathbf{X}) = \left( \mathbf{X}\mathbf{X}^\top - \mathbf{M}^\star \right) \mathbf{X}. \tag{60}$$

In addition, we need the Hessian  $\nabla^2 f_{\text{clean}}(\mathbf{X})$ , which is represented by an  $nr \times nr$  matrix. Simple calculations reveal that for any  $\mathbf{V} \in \mathbb{R}^{nr \times r}$ ,

$$\begin{aligned} \text{vec}(\mathbf{V})^\top \nabla^2 f_{\text{clean}}(\mathbf{X}) \text{vec}(\mathbf{V}) &= \frac{1}{2p} \left\| \mathcal{P}_\Omega \left( \mathbf{V}\mathbf{X}^\top + \mathbf{X}\mathbf{V}^\top \right) \right\|_F^2 \\ &\quad + \frac{1}{p} \left\langle \mathcal{P}_\Omega \left( \mathbf{X}\mathbf{X}^\top - \mathbf{M}^\star \right), \mathbf{V}\mathbf{V}^\top \right\rangle, \end{aligned} \tag{61}$$

where  $\text{vec}(\mathbf{V}) \in \mathbb{R}^{nr}$  denotes the vectorization of  $\mathbf{V}$ .

### 7.1 Step 1: Characterizing Local Geometry in the RIC

#### 7.1.1 Local Geometry

The first step is to characterize the region where the empirical loss function enjoys restricted strong convexity and smoothness in an appropriate sense. This is formally stated in the following lemma.

**Lemma 7** (Restricted strong convexity and smoothness for matrix completion) *Suppose that the sample size obeys  $n^2 p \geq C\kappa^2 \mu r n \log n$  for some sufficiently large*

constant  $C > 0$ . Then with probability at least  $1 - O(n^{-10})$ , the Hessian  $\nabla^2 f_{\text{clean}}(\mathbf{X})$  as defined in (61) obeys

$$\text{vec}(\mathbf{V})^\top \nabla^2 f_{\text{clean}}(\mathbf{X}) \text{vec}(\mathbf{V}) \geq \frac{\sigma_{\min}}{2} \|\mathbf{V}\|_F^2 \quad \text{and} \quad \left\| \nabla^2 f_{\text{clean}}(\mathbf{X}) \right\| \leq \frac{5}{2} \sigma_{\max} \tag{62}$$

for all  $\mathbf{X}$  and  $\mathbf{V} = \mathbf{Y}\mathbf{H}_Y - \mathbf{Z}$ , with  $\mathbf{H}_Y := \arg \min_{\mathbf{R} \in \mathcal{O}^{r \times r}} \|\mathbf{Y}\mathbf{R} - \mathbf{Z}\|_F$ , satisfying:

$$\|\mathbf{X} - \mathbf{X}^*\|_{2,\infty} \leq \epsilon \|\mathbf{X}^*\|_{2,\infty}, \tag{63a}$$

$$\|\mathbf{Z} - \mathbf{X}^*\| \leq \delta \|\mathbf{X}^*\|, \tag{63b}$$

where  $\epsilon \ll 1/\sqrt{\kappa^3 \mu r \log^2 n}$  and  $\delta \ll 1/\kappa$ .

**Proof** See Appendix B.1. □

Lemma 7 reveals that the Hessian matrix is well conditioned in a neighborhood close to  $\mathbf{X}^*$  that remains incoherent measured in the  $\ell_2/\ell_\infty$  norm (cf. (63a)), and along directions that point toward points which are not far away from the truth in the spectral norm (cf. (63b)).

**Remark 5** The second condition (63b) is characterized using the spectral norm  $\|\cdot\|$ , while in previous works this is typically presented in the Frobenius norm  $\|\cdot\|_F$ . It is also worth noting that the Hessian matrix—even in the infinite-sample and noiseless case—is rank-deficient and cannot be positive definite. As a result, we resort to the form of strong convexity by restricting attention to certain directions (see the conditions on  $\mathbf{V}$ ).

### 7.1.2 Error Contraction

Our goal is to demonstrate the error bounds (28) measured in three different norms. Notably, as long as the iterates satisfy (28) at the  $t$ th iteration, then  $\|\mathbf{X}^t \widehat{\mathbf{H}}^t - \mathbf{X}^*\|_{2,\infty}$  is sufficiently small. Under our sample complexity assumption,  $\mathbf{X}^t \widehat{\mathbf{H}}^t$  satisfies the  $\ell_2/\ell_\infty$  condition (63a) required in Lemma 7. Consequently, we can invoke Lemma 7 to arrive at the following error contraction result.

**Lemma 8** (Contraction w.r.t. the Frobenius norm) *Suppose that  $n^2 p \geq C\kappa^3 \mu^3 r^3 n \log^3 n$  for some sufficiently large constant  $C > 0$ , and the noise satisfies (27). There exists an event that does not depend on  $t$  and has probability  $1 - O(n^{-10})$ , such that when it happens and (28a), (28b) hold for the  $t$ th iteration, one has*

$$\|\mathbf{X}^{t+1} \widehat{\mathbf{H}}^{t+1} - \mathbf{X}^*\|_F \leq C_4 \rho^{t+1} \mu r \frac{1}{\sqrt{np}} \|\mathbf{X}^*\|_F + C_1 \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \|\mathbf{X}^*\|_F$$

provided that  $0 < \eta \leq 2/(25\kappa\sigma_{\max})$ ,  $1 - (\sigma_{\min}/4) \cdot \eta \leq \rho < 1$ , and  $C_1$  is sufficiently large.

**Proof** The proof is built upon Lemma 7. See Appendix B.2. □

Further, if the current iterate satisfies all three conditions in (28), then we can derive a stronger sense of error contraction, namely contraction in terms of the spectral norm.

**Lemma 9** (Contraction w.r.t. the spectral norm) *Suppose  $n^2 p \geq C\kappa^3 \mu^3 r^3 n \log^3 n$  for some sufficiently large constant  $C > 0$ , and the noise satisfies (27). There exists an event that does not depend on  $t$  and has probability  $1 - O(n^{-10})$ , such that when it happens and (28) holds for the  $t$ th iteration, one has*

$$\|X^{t+1} \widehat{H}^{t+1} - X^*\| \leq C_9 \rho^{t+1} \mu r \frac{1}{\sqrt{np}} \|X^*\| + C_{10} \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \|X^*\| \tag{64}$$

provided that  $0 < \eta \leq 1/(2\sigma_{\max})$  and  $1 - (\sigma_{\min}/3) \cdot \eta \leq \rho < 1$ .

**Proof** The key observation is this: the iterate that proceeds according to the population-level gradient reduces the error w.r.t.  $\|\cdot\|$ , namely

$$\|X^t \widehat{H}^t - \eta \nabla F_{\text{clean}}(X^t \widehat{H}^t) - X^*\| < \|X^t \widehat{H}^t - X^*\|,$$

as long as  $X^t \widehat{H}^t$  is sufficiently close to the truth. Notably, the orthonormal matrix  $\widehat{H}^t$  is still chosen to be the one that minimizes the  $\|\cdot\|_F$  distance (as opposed to  $\|\cdot\|$ ), which yields a symmetry property  $X^{*\top} X^t \widehat{H}^t = (X^t \widehat{H}^t)^\top X^*$ , crucial for our analysis. See Appendix B.3 for details. □

### 7.2 Step 2: Introducing the Leave-One-Out Sequences

In order to establish the incoherence properties (28b) for the entire trajectory, which is difficult to deal with directly due to the complicated statistical dependence, we introduce a collection of *leave-one-out* versions of  $\{X^t\}_{t \geq 0}$ , denoted by  $\{X^{t,(l)}\}_{t \geq 0}$  for each  $1 \leq l \leq n$ . Specifically,  $\{X^{t,(l)}\}_{t \geq 0}$  is the iterates of gradient descent operating on the auxiliary loss function

$$f^{(l)}(X) := \frac{1}{4p} \left\| \mathcal{P}_{\Omega^{-l}} [XX^\top - (M^* + E)] \right\|_F^2 + \frac{1}{4} \left\| \mathcal{P}_l (XX^\top - M^*) \right\|_F^2. \tag{65}$$

Here,  $\mathcal{P}_{\Omega_l}$  (resp.  $\mathcal{P}_{\Omega^{-l}}$  and  $\mathcal{P}_l$ ) represents the orthogonal projection onto the subspace of matrices which vanish outside of the index set  $\Omega_l := \{(i, j) \in \Omega \mid i = l \text{ or } j = l\}$  (resp.  $\Omega^{-l} := \{(i, j) \in \Omega \mid i \neq l, j \neq l\}$  and  $\{(i, j) \mid i = l \text{ or } j = l\}$ ); that is, for any matrix  $M$ ,

$$\begin{aligned} [\mathcal{P}_{\Omega_l}(M)]_{i,j} &= \begin{cases} M_{i,j}, & \text{if } (i = l \text{ or } j = l) \text{ and } (i, j) \in \Omega, \\ 0, & \text{else,} \end{cases} & (66) \\ [\mathcal{P}_{\Omega^{-l}}(M)]_{i,j} &= \begin{cases} M_{i,j}, & \text{if } i \neq l \text{ and } j \neq l \text{ and } (i, j) \in \Omega \\ 0, & \text{else} \end{cases} \quad \text{and} \end{aligned}$$

$$[\mathcal{P}_l(\mathbf{M})]_{i,j} = \begin{cases} 0, & \text{if } i \neq l \text{ and } j \neq l, \\ M_{i,j}, & \text{if } i = l \text{ or } j = l. \end{cases} \tag{67}$$

The gradient of the leave-one-out loss function (65) is given by

$$\nabla f^{(l)}(\mathbf{X}) = \frac{1}{p} \mathcal{P}_{\Omega^{-l}} \left[ \mathbf{X} \mathbf{X}^\top - (\mathbf{M}^* + \mathbf{E}) \right] \mathbf{X} + \mathcal{P}_l \left( \mathbf{X} \mathbf{X}^\top - \mathbf{M}^* \right) \mathbf{X}. \tag{68}$$

The full algorithm to obtain the leave-one-out sequence  $\{\mathbf{X}^{t,(l)}\}_{t \geq 0}$  (including spectral initialization) is summarized in Algorithm 5.

---

**Algorithm 5** The  $l$ th leave-one-out sequence for matrix completion

---

**Input:**  $\mathbf{Y} = [Y_{i,j}]_{1 \leq i,j \leq n}$ ,  $\mathbf{M}_{:,l}^*$ ,  $\mathbf{M}_{l,:}^*$ ,  $r$ ,  $p$ .

**Spectral initialization:** Let  $U^{0,(l)} \Sigma^{(l)} U^{0,(l)\top}$  be the top- $r$  eigendecomposition of

$$\mathbf{M}^{(l)} := \frac{1}{p} \mathcal{P}_{\Omega^{-l}}(\mathbf{Y}) + \mathcal{P}_l(\mathbf{M}^*) = \frac{1}{p} \mathcal{P}_{\Omega^{-l}}(\mathbf{M}^* + \mathbf{E}) + \mathcal{P}_l(\mathbf{M}^*)$$

with  $\mathcal{P}_{\Omega^{-l}}$  and  $\mathcal{P}_l$  defined in (67), and set  $\mathbf{X}^{0,(l)} = U^{0,(l)}(\Sigma^{(l)})^{1/2}$ .

**Gradient updates:** for  $t = 0, 1, 2, \dots, T - 1$  do

$$\mathbf{X}^{t+1,(l)} = \mathbf{X}^{t,(l)} - \eta_t \nabla f^{(l)}(\mathbf{X}^{t,(l)}). \tag{69}$$


---

**Remark 6** Rather than simply dropping all samples in the  $l$ th row/column, we replace the  $l$ th row/column with their respective population means. In other words, the leave-one-out gradient forms an unbiased surrogate for the true gradient, which is particularly important in ensuring high estimation accuracy.

### 7.3 Step 3: Establishing the Incoherence Condition by Induction

We will continue the proof of Theorem 2 in an inductive manner. As seen in Sect. 7.1.2, the induction hypotheses (28a) and (28c) hold for the  $(t + 1)$ th iteration as long as (28) holds at the  $t$ th iteration. Therefore, we are left with proving the incoherence hypothesis (28b) for all  $0 \leq t \leq T = O(n^5)$ . For clarity of analysis, it is crucial to maintain a list of induction hypotheses, which includes a few more hypotheses that complement (28), and is given below.

$$\|\mathbf{X}^t \widehat{\mathbf{H}}^t - \mathbf{X}^*\|_F \leq \left( C_4 \rho^t \mu r \frac{1}{\sqrt{np}} + C_1 \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \right) \|\mathbf{X}^*\|_F, \tag{70a}$$

$$\|\mathbf{X}^t \widehat{\mathbf{H}}^t - \mathbf{X}^*\|_{2,\infty} \leq \left( C_5 \rho^t \mu r \sqrt{\frac{\log n}{np}} + C_8 \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} \right) \|\mathbf{X}^*\|_{2,\infty}, \tag{70b}$$

$$\|X^t \widehat{H}^t - X^*\| \leq \left( C_9 \rho^t \mu r \frac{1}{\sqrt{np}} + C_{10} \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \right) \|X^*\|, \tag{70c}$$

$$\max_{1 \leq l \leq n} \|X^t \widehat{H}^t - X^{t,(l)} R^{t,(l)}\|_F \leq \left( C_3 \rho^t \mu r \sqrt{\frac{\log n}{np}} + C_7 \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} \right) \|X^*\|_{2,\infty}, \tag{70d}$$

$$\max_{1 \leq l \leq n} \|(X^{t,(l)} \widehat{H}^{t,(l)} - X^*)_{l,\cdot}\|_2 \leq \left( C_2 \rho^t \mu r \frac{1}{\sqrt{np}} + C_6 \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} \right) \|X^*\|_{2,\infty} \tag{70e}$$

hold for some absolute constants  $0 < \rho < 1$  and  $C_1, \dots, C_{10} > 0$ . Here,  $\widehat{H}^{t,(l)}$  and  $R^{t,(l)}$  are orthonormal matrices defined by

$$\widehat{H}^{t,(l)} := \arg \min_{R \in O^{r \times r}} \|X^{t,(l)} R - X^*\|_F, \tag{71}$$

$$R^{t,(l)} := \arg \min_{R \in O^{r \times r}} \|X^{t,(l)} R - X^t \widehat{H}^t\|_F. \tag{72}$$

Clearly, the first three hypotheses (70a)–(70c) constitute the conclusion of Theorem 2, i.e., (28). The last two hypotheses (70d) and (70e) are auxiliary properties connecting the true iterates and the auxiliary leave-one-out sequences. Moreover, we summarize below several immediate consequences of (70), which will be useful throughout.

**Lemma 10** *Suppose  $n^2 p \geq C \kappa^3 \mu^2 r^2 n \log n$  for some sufficiently large constant  $C > 0$ , and the noise satisfies (27). Under hypotheses (70), one has*

$$\|X^t \widehat{H}^t - X^{t,(l)} \widehat{H}^{t,(l)}\|_F \leq 5\kappa \|X^t \widehat{H}^t - X^{t,(l)} R^{t,(l)}\|_F, \tag{73a}$$

$$\|X^{t,(l)} \widehat{H}^{t,(l)} - X^*\|_F \leq \|X^{t,(l)} R^{t,(l)} - X^*\|_F \leq \left\{ 2C_4 \rho^t \mu r \frac{1}{\sqrt{np}} + 2C_1 \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \right\} \|X^*\|_F, \tag{73b}$$

$$\|X^{t,(l)} R^{t,(l)} - X^*\|_{2,\infty} \leq \left\{ (C_3 + C_5) \rho^t \mu r \sqrt{\frac{\log n}{np}} + (C_8 + C_7) \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} \right\} \|X^*\|_{2,\infty}, \tag{73c}$$

$$\|X^{t,(l)} \widehat{H}^{t,(l)} - X^*\| \leq \left\{ 2C_9 \rho^t \mu r \frac{1}{\sqrt{np}} + 2C_{10} \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \right\} \|X^*\|. \tag{73d}$$

In particular, (73a) follows from hypotheses (70c) and (70d).

**Proof** See Appendix B.4. □

In the sequel, we follow the general recipe outlined in Sect. 5 to establish the induction hypotheses. We only need to establish (70b), (70d), and (70e) for the  $(t + 1)$ th iteration, since (70a) and (70c) are established in Sect. 7.1.2. Specifically, we resort to the leave-one-out iterates by showing that: first, the true and the auxiliary iterates remain exceedingly close throughout; second, the  $l$ th leave-one-out sequence stays incoherent with  $e_l$  due to statistical independence.

- *Step 3(a): proximity between the original and the leave-one-out iterates* We demonstrate that  $X^{t+1}$  is well approximated by  $X^{t+1,(l)}$ , up to proper orthonormal

transforms. This is precisely the induction hypothesis (70d) for the  $(t + 1)$ th iteration.

**Lemma 11** *Suppose the sample complexity satisfies  $n^2 p \geq C\kappa^4 \mu^3 r^3 n \log^3 n$  for some sufficiently large constant  $C > 0$ , and the noise satisfies (27). Let  $\mathcal{E}_t$  be the event where the hypotheses in (70) hold for the  $t$ th iteration. Then on some event  $\mathcal{E}_{t+1,1} \subseteq \mathcal{E}_t$  obeying  $\mathbb{P}(\mathcal{E}_t \cap \mathcal{E}_{t+1,1}^c) = O(n^{-10})$ , we have*

$$\begin{aligned} \left\| \mathbf{X}^{t+1} \widehat{\mathbf{H}}^{t+1} - \mathbf{X}^{t+1,(l)} \mathbf{R}^{t+1,(l)} \right\|_{\text{F}} &\leq C_3 \rho^{t+1} \mu r \sqrt{\frac{\log n}{np}} \|\mathbf{X}^*\|_{2,\infty} \\ &\quad + C_7 \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} \|\mathbf{X}^*\|_{2,\infty} \end{aligned} \tag{74}$$

provided that  $0 < \eta \leq 2/(25\kappa\sigma_{\max})$ ,  $1 - (\sigma_{\min}/5) \cdot \eta \leq \rho < 1$ , and  $C_7 > 0$  is sufficiently large.

**Proof** The fact that this difference is well controlled relies heavily on the benign geometric property of the Hessian revealed by Lemma 7. Two important remarks are in order: (1) both points  $\mathbf{X}^t \widehat{\mathbf{H}}^t$  and  $\mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)}$  satisfy (63a); (2) the difference  $\mathbf{X}^t \widehat{\mathbf{H}}^t - \mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)}$  forms a valid direction for restricted strong convexity. These two properties together allow us to invoke Lemma 7. See Appendix B.5.  $\square$

- *Step 3(b): incoherence of the leave-one-out iterates* Given that  $\mathbf{X}^{t+1,(l)}$  is sufficiently close to  $\mathbf{X}^{t+1}$ , we turn our attention to establishing the incoherence of this surrogate  $\mathbf{X}^{t+1,(l)}$  w.r.t.  $\mathbf{e}_l$ . This amounts to proving the induction hypothesis (70e) for the  $(t + 1)$ th iteration.

**Lemma 12** *Suppose the sample complexity meets  $n^2 p \geq C\kappa^3 \mu^3 r^3 n \log^3 n$  for some sufficiently large constant  $C > 0$ , and the noise satisfies (27). Let  $\mathcal{E}_t$  be the event where the hypotheses in (70) hold for the  $t$ th iteration. Then on some event  $\mathcal{E}_{t+1,2} \subseteq \mathcal{E}_t$  obeying  $\mathbb{P}(\mathcal{E}_t \cap \mathcal{E}_{t+1,2}^c) = O(n^{-10})$ , we have*

$$\left\| (\mathbf{X}^{t+1,(l)} \widehat{\mathbf{H}}^{t+1,(l)} - \mathbf{X}^*)_{l,\cdot} \right\|_2 \leq C_2 \rho^{t+1} \mu r \frac{1}{\sqrt{np}} \|\mathbf{X}^*\|_{2,\infty} + C_6 \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} \|\mathbf{X}^*\|_{2,\infty} \tag{75}$$

so long as  $0 < \eta \leq 1/\sigma_{\max}$ ,  $1 - (\sigma_{\min}/3) \cdot \eta \leq \rho < 1$ ,  $C_2 \gg \kappa C_9$ , and  $C_6 \gg \kappa C_{10}/\sqrt{\log n}$ .

**Proof** The key observation is that  $\mathbf{X}^{t+1,(l)}$  is statistically independent from any sample in the  $l$ th row/column of the matrix. Since there are an order of  $np$  samples in each row/column, we obtain enough information that helps establish the desired incoherence property. See Appendix B.6.  $\square$

- *Step 3(c): combining the bounds* Inequalities (70d) and (70e) taken collectively allow us to establish the induction hypothesis (70b). Specifically, for every  $1 \leq l \leq n$ , write

$$\begin{aligned} (X^{t+1} \widehat{H}^{t+1} - X^*)_{l,\cdot} &= (X^{t+1} \widehat{H}^{t+1} - X^{t+1,(l)} \widehat{H}^{t+1,(l)})_{l,\cdot} \\ &\quad + (X^{t+1,(l)} \widehat{H}^{t+1,(l)} - X^*)_{l,\cdot}, \end{aligned}$$

and the triangle inequality gives

$$\begin{aligned} \|(X^{t+1} \widehat{H}^{t+1} - X^*)_{l,\cdot}\|_2 &\leq \|X^{t+1} \widehat{H}^{t+1} - X^{t+1,(l)} \widehat{H}^{t+1,(l)}\|_F \\ &\quad + \|(X^{t+1,(l)} \widehat{H}^{t+1,(l)} - X^*)_{l,\cdot}\|_2. \end{aligned} \tag{76}$$

The second term has already been bounded by (75). Since we have established the induction hypotheses (70c) and (70d) for the  $(t + 1)$ th iteration, the first term can be bounded by (73a) for the  $(t + 1)$ th iteration, i.e.,

$$\|X^{t+1} \widehat{H}^{t+1} - X^{t+1,(l)} \widehat{H}^{t+1,(l)}\|_F \leq 5\kappa \|X^{t+1} \widehat{H}^{t+1} - X^{t+1,(l)} R^{t+1,(l)}\|_F.$$

Plugging the above inequality, (74), and (75) into (76), we have

$$\begin{aligned} &\|X^{t+1} \widehat{H}^{t+1} - X^*\|_{2,\infty} \\ &\leq 5\kappa \left( C_3 \rho^{t+1} \mu r \sqrt{\frac{\log n}{np}} \|X^*\|_{2,\infty} + \frac{C_7}{\sigma_{\min}} \sigma \sqrt{\frac{n \log n}{p}} \|X^*\|_{2,\infty} \right) \\ &\quad + C_2 \rho^{t+1} \mu r \frac{1}{\sqrt{np}} \|X^*\|_{2,\infty} + \frac{C_6}{\sigma_{\min}} \sigma \sqrt{\frac{n \log n}{p}} \|X^*\|_{2,\infty} \\ &\leq C_5 \rho^{t+1} \mu r \sqrt{\frac{\log n}{np}} \|X^*\|_{2,\infty} + \frac{C_8}{\sigma_{\min}} \sigma \sqrt{\frac{n \log n}{p}} \|X^*\|_{2,\infty} \end{aligned}$$

as long as  $C_5/(\kappa C_3 + C_2)$  and  $C_8/(\kappa C_7 + C_6)$  are sufficiently large. This establishes the induction hypothesis (70b). From the deduction above, we see  $\mathcal{E}_t \cap \mathcal{E}_{t+1}^c = O(n^{-10})$  and thus finish the proof.

### 7.4 The Base Case: Spectral Initialization

Finally, we return to check the base case; namely, we aim to show that the spectral initialization satisfies the induction hypotheses (70a)–(70e) for  $t = 0$ . This is accomplished via the following lemma.

**Lemma 13** *Suppose the sample size obeys  $n^2 p \geq C \mu^2 r^2 n \log n$  for some sufficiently large constant  $C > 0$ , the noise satisfies (27), and  $\kappa = \sigma_{\max}/\sigma_{\min} \asymp 1$ . Then with probability at least  $1 - O(n^{-10})$ , the claims in (70a)–(70e) hold simultaneously for  $t = 0$ .*

**Proof** This follows by invoking the Davis–Kahan  $\sin \Theta$  theorem [39] as well as the entrywise eigenvector perturbation analysis in [1]. We defer the proof to Appendix B.7. □



### 8 Analysis for Blind Deconvolution

In this section, we instantiate the general recipe presented in Sect. 5 to blind deconvolution and prove Theorem 3. Without loss of generality, we assume throughout that  $\|\mathbf{h}^*\|_2 = \|\mathbf{x}^*\|_2 = 1$ .

Before presenting the analysis, we first gather some simple facts about the empirical loss function in (32). Recall the definition of  $\mathbf{z}$  in (33), and for notational simplicity, we write  $f(\mathbf{z}) = f(\mathbf{h}, \mathbf{x})$ . Since  $\mathbf{z}$  is complex-valued, we need to resort to Wirtinger calculus; see [18, Section 6] for a brief introduction. The Wirtinger gradient of (32) with respect to  $\mathbf{h}$  and  $\mathbf{x}$  are given, respectively, by

$$\nabla_{\mathbf{h}} f(\mathbf{z}) = \nabla_{\mathbf{h}} f(\mathbf{h}, \mathbf{x}) = \sum_{j=1}^m \left( \mathbf{b}_j^H \mathbf{h} \mathbf{x}^H \mathbf{a}_j - y_j \right) \mathbf{b}_j \mathbf{a}_j^H \mathbf{x}; \tag{77}$$

$$\nabla_{\mathbf{x}} f(\mathbf{z}) = \nabla_{\mathbf{x}} f(\mathbf{h}, \mathbf{x}) = \sum_{j=1}^m \overline{\left( \mathbf{b}_j^H \mathbf{h} \mathbf{x}^H \mathbf{a}_j - y_j \right)} \mathbf{a}_j \mathbf{b}_j^H \mathbf{h}. \tag{78}$$

It is worth noting that the formal Wirtinger gradient contains  $\nabla_{\bar{\mathbf{h}}} f(\mathbf{h}, \mathbf{x})$  and  $\nabla_{\bar{\mathbf{x}}} f(\mathbf{h}, \mathbf{x})$  as well. Nevertheless, since  $f(\mathbf{h}, \mathbf{x})$  is a real-valued function, the following identities always hold

$$\nabla_{\mathbf{h}} f(\mathbf{h}, \mathbf{x}) = \overline{\nabla_{\bar{\mathbf{h}}} f(\mathbf{h}, \mathbf{x})} \quad \text{and} \quad \nabla_{\mathbf{x}} f(\mathbf{h}, \mathbf{x}) = \overline{\nabla_{\bar{\mathbf{x}}} f(\mathbf{h}, \mathbf{x})}.$$

In light of these observations, one often omits the gradient with respect to the conjugates; correspondingly, the gradient update rule (35) can be written as

$$\mathbf{h}^{t+1} = \mathbf{h}^t - \frac{\eta}{\|\mathbf{x}^t\|_2^2} \sum_{j=1}^m \left( \mathbf{b}_j^H \mathbf{h}^t \mathbf{x}^{tH} \mathbf{a}_j - y_j \right) \mathbf{b}_j \mathbf{a}_j^H \mathbf{x}^t, \tag{79a}$$

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \frac{\eta}{\|\mathbf{h}^t\|_2^2} \sum_{j=1}^m \overline{\left( \mathbf{b}_j^H \mathbf{h}^t \mathbf{x}^{tH} \mathbf{a}_j - y_j \right)} \mathbf{a}_j \mathbf{b}_j^H \mathbf{h}^t. \tag{79b}$$

We can also compute the Wirtinger Hessian of  $f(\mathbf{z})$  as follows,

$$\nabla^2 f(\mathbf{z}) = \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^H & \mathbf{A} \end{bmatrix}, \tag{80}$$

where

$$\mathbf{A} = \begin{bmatrix} \sum_{j=1}^m \left| \mathbf{a}_j^H \mathbf{x} \right|^2 \mathbf{b}_j \mathbf{b}_j^H & \sum_{j=1}^m \left( \mathbf{b}_j^H \mathbf{h} \mathbf{x}^H \mathbf{a}_j - y_j \right) \mathbf{b}_j \mathbf{a}_j^H \\ \sum_{j=1}^m \left[ \left( \mathbf{b}_j^H \mathbf{h} \mathbf{x}^H \mathbf{a}_j - y_j \right) \mathbf{b}_j \mathbf{a}_j^H \right]^H & \sum_{j=1}^m \left| \mathbf{b}_j^H \mathbf{h} \right|^2 \mathbf{a}_j \mathbf{a}_j^H \end{bmatrix} \in \mathbb{C}^{2K \times 2K};$$

$$\mathbf{B} = \begin{bmatrix} \mathbf{0} & \sum_{j=1}^m \mathbf{b}_j \mathbf{b}_j^H \mathbf{h} \left( \mathbf{a}_j \mathbf{a}_j^H \mathbf{x} \right)^T \\ \sum_{j=1}^m \mathbf{a}_j \mathbf{a}_j^H \mathbf{x} \left( \mathbf{b}_j \mathbf{b}_j^H \mathbf{h} \right)^T & \mathbf{0} \end{bmatrix} \in \mathbb{C}^{2K \times 2K}.$$

Last but not least, we say  $(\mathbf{h}_1, \mathbf{x}_1)$  is aligned with  $(\mathbf{h}_2, \mathbf{x}_2)$ , if the following holds,

$$\|\mathbf{h}_1 - \mathbf{h}_2\|_2^2 + \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 = \min_{\alpha \in \mathbb{C}} \left\{ \left\| \frac{1}{\alpha} \mathbf{h}_1 - \mathbf{h}_2 \right\|_2^2 + \|\alpha \mathbf{x}_1 - \mathbf{x}_2\|_2^2 \right\}.$$

To simplify notations, define  $\tilde{\mathbf{z}}^t$  as

$$\tilde{\mathbf{z}}^t = \begin{bmatrix} \tilde{\mathbf{h}}^t \\ \tilde{\mathbf{x}}^t \end{bmatrix} := \begin{bmatrix} \frac{1}{\alpha^t} \mathbf{h}^t \\ \alpha^t \mathbf{x}^t \end{bmatrix} \tag{81}$$

with the alignment parameter  $\alpha^t$  given in (38). Then we can see that  $\tilde{\mathbf{z}}^t$  is aligned with  $\mathbf{z}^*$  and

$$\text{dist}(\mathbf{z}^t, \mathbf{z}^*) = \text{dist}(\tilde{\mathbf{z}}^t, \mathbf{z}^*) = \|\tilde{\mathbf{z}}^t - \mathbf{z}^*\|_2.$$

### 8.1 Step 1: Characterizing Local Geometry in the RIC

#### 8.1.1 Local Geometry

The first step is to characterize the region of incoherence and contraction (RIC), where the empirical loss function enjoys restricted strong convexity and smoothness properties. To this end, we have the following lemma.

**Lemma 14** (Restricted strong convexity and smoothness for blind deconvolution) *Let  $c > 0$  be a sufficiently small constant and*

$$\delta = c / \log^2 m.$$

*Suppose the sample size satisfies  $m \geq c_0 \mu^2 K \log^9 m$  for some sufficiently large constant  $c_0 > 0$ . Then with probability  $1 - O(m^{-10} + e^{-K} \log m)$ , the Wirtinger Hessian  $\nabla^2 f(\mathbf{z})$  obeys*

$$\mathbf{u}^H \left[ \mathbf{D} \nabla^2 f(\mathbf{z}) + \nabla^2 f(\mathbf{z}) \mathbf{D} \right] \mathbf{u} \geq (1/4) \cdot \|\mathbf{u}\|_2^2 \quad \text{and} \quad \left\| \nabla^2 f(\mathbf{z}) \right\| \leq 3$$

*simultaneously for all*

$$\mathbf{z} = \begin{bmatrix} \mathbf{h} \\ \mathbf{x} \end{bmatrix} \quad \text{and} \quad \mathbf{u} = \begin{bmatrix} \mathbf{h}_1 - \mathbf{h}_2 \\ \frac{\mathbf{x}_1 - \mathbf{x}_2}{\mathbf{h}_1 - \mathbf{h}_2} \\ \mathbf{x}_1 - \mathbf{x}_2 \end{bmatrix} \quad \text{and} \quad \mathbf{D} = \begin{bmatrix} \gamma_1 \mathbf{I}_K & & & \\ & \gamma_2 \mathbf{I}_K & & \\ & & \gamma_1 \mathbf{I}_K & \\ & & & \gamma_2 \mathbf{I}_K \end{bmatrix},$$

*where  $\mathbf{z}$  satisfies*

$$\max \{ \|\mathbf{h} - \mathbf{h}^*\|_2, \|\mathbf{x} - \mathbf{x}^*\|_2 \} \leq \delta; \tag{82a}$$

$$\max_{1 \leq j \leq m} \left| \mathbf{a}_j^H (\mathbf{x} - \mathbf{x}^*) \right| \leq 2C_3 \frac{1}{\log^{3/2} m}; \tag{82b}$$

$$\max_{1 \leq j \leq m} \left| \mathbf{b}_j^H \mathbf{h} \right| \leq 2C_4 \frac{\mu}{\sqrt{m}} \log^2 m; \tag{82c}$$

$(\mathbf{h}_1, \mathbf{x}_1)$  is aligned with  $(\mathbf{h}_2, \mathbf{x}_2)$ , and they satisfy

$$\max \{ \|\mathbf{h}_1 - \mathbf{h}^*\|_2, \|\mathbf{h}_2 - \mathbf{h}^*\|_2, \|\mathbf{x}_1 - \mathbf{x}^*\|_2, \|\mathbf{x}_2 - \mathbf{x}^*\|_2 \} \leq \delta; \tag{83}$$

and finally,  $\mathbf{D}$  satisfies for  $\gamma_1, \gamma_2 \in \mathbb{R}$ ,

$$\max \{ |\gamma_1 - 1|, |\gamma_2 - 1| \} \leq \delta. \tag{84}$$

Here,  $C_3, C_4 > 0$  are numerical constants.

**Proof** See Appendix C.1. □

Lemma 14 characterizes the restricted strong convexity and smoothness of the loss function used in blind deconvolution. To the best of our knowledge, this provides the first characterization regarding geometric properties of the Hessian matrix for blind deconvolution. A few interpretations are in order.

- Conditions (82) specify the region of incoherence and contraction (RIC). In particular, (82a) specifies a neighborhood that is close to the ground truth in  $\ell_2$  norm, and (82b) and (82c) specify the incoherence region with respect to the sensing vectors  $\{\mathbf{a}_j\}$  and  $\{\mathbf{b}_j\}$ , respectively.
- Similar to matrix completion, the Hessian matrix is rank-deficient even at the population level. Consequently, we resort to a restricted form of strong convexity by focusing on certain directions. More specifically, these directions can be viewed as the difference between two pre-aligned points that are not far from the truth, which is characterized by (83).
- Finally, the diagonal matrix  $\mathbf{D}$  accounts for scaling factors that are not too far from 1 (see (84)), which allows us to account for different step sizes employed for  $\mathbf{h}$  and  $\mathbf{x}$ .

### 8.1.2 Error Contraction

The restricted strong convexity and smoothness allow us to establish the contraction of the error measured in terms of  $\text{dist}(\cdot, \mathbf{z}^*)$  as defined in (34) as long as the iterates stay in the RIC.

**Lemma 15** *Suppose the number of measurements satisfies  $m \geq C\mu^2 K \log^9 m$  for some sufficiently large constant  $C > 0$ , and the step size  $\eta > 0$  is some sufficiently small constant. There exists an event that does not depend on  $t$  and has probability  $1 - O(m^{-10} + e^{-K} \log m)$ , such that when it happens and*

$$\text{dist}(\mathbf{z}^t, \mathbf{z}^*) \leq \xi, \tag{85a}$$

$$\max_{1 \leq j \leq m} \left| \mathbf{a}_j^H (\tilde{\mathbf{x}}^t - \mathbf{x}^*) \right| \leq C_3 \frac{1}{\log^{1.5} m}, \tag{85b}$$

$$\max_{1 \leq j \leq m} \left| \mathbf{b}_j^H \tilde{\mathbf{h}}^t \right| \leq C_4 \frac{\mu}{\sqrt{m}} \log^2 m \tag{85c}$$

hold for some constants  $C_3, C_4 > 0$ , one has

$$\text{dist}(\mathbf{z}^{t+1}, \mathbf{z}^*) \leq (1 - \eta/16) \text{dist}(\mathbf{z}^t, \mathbf{z}^*).$$

Here,  $\tilde{\mathbf{h}}^t$  and  $\tilde{\mathbf{x}}^t$  are defined in (81), and  $\xi \ll 1/\log^2 m$ .

**Proof** See Appendix C.2. □

As a result, if  $\mathbf{z}^t$  satisfies condition (85) for all  $0 \leq t \leq T$ , then

$$\text{dist}(\mathbf{z}^t, \mathbf{z}^*) \leq \rho \text{dist}(\mathbf{z}^{t-1}, \mathbf{z}^*) \leq \rho^t \text{dist}(\mathbf{z}^0, \mathbf{z}^*) \leq \rho^t c_1, \quad 0 < t \leq T,$$

where  $\rho := 1 - \eta/16$ . Furthermore, similar to the case of phase retrieval (i.e., Lemma 3), as soon as we demonstrate that conditions (85) hold for all  $0 \leq t \leq m$ , then Theorem 3 holds true. The proof of this claim is exactly the same as for Lemma 3 and is thus omitted for conciseness. In what follows, we focus on establishing (85) for all  $0 \leq t \leq m$ .

Before concluding this subsection, we make note of another important result that concerns the alignment parameter  $\alpha^t$ , which will be useful in the subsequent analysis. Specifically, the alignment parameter sequence  $\{\alpha^t\}$  converges linearly to a constant whose magnitude is fairly close to 1, as long as the two initial vectors  $\mathbf{h}^0$  and  $\mathbf{x}^0$  have similar  $\ell_2$  norms and are close to the truth. Given that  $\alpha^t$  determines the global scaling of the iterates, this reveals rapid convergence of both  $\|\mathbf{h}^t\|_2$  and  $\|\mathbf{x}^t\|_2$ , which explains why there is no need to impose extra terms to regularize the  $\ell_2$  norm as employed in [58,76].

**Lemma 16** *When  $m > 1$  is sufficiently large, the following two claims hold true.*

- If  $|\alpha^t| - 1| \leq 1/2$  and  $\text{dist}(\mathbf{z}^t, \mathbf{z}^*) \leq C_1/\log^2 m$ , then

$$\left| \frac{\alpha^{t+1}}{\alpha^t} - 1 \right| \leq c \text{dist}(\mathbf{z}^t, \mathbf{z}^*) \leq \frac{cC_1}{\log^2 m}$$

for some absolute constant  $c > 0$ ;

- If  $|\alpha^0| - 1| \leq 1/4$  and  $\text{dist}(\mathbf{z}^s, \mathbf{z}^*) \leq C_1(1 - \eta/16)^s / \log^2 m$  for all  $0 \leq s \leq t$ , then one has

$$|\alpha^{s+1}| - 1| \leq 1/2, \quad 0 \leq s \leq t.$$

**Proof** See Appendix C.2. □

The initial condition  $|\alpha^0| - 1| < 1/4$  will be guaranteed to hold with high probability by Lemma 19.

### 8.2 Step 2: Introducing the Leave-One-Out Sequences

As demonstrated by the assumptions in Lemma 15, the key is to show that the whole trajectory lies in the region specified by (85a)–(85c). Once again, the difficulty lies in the statistical dependency between the iterates  $\{z^t\}$  and the measurement vectors  $\{a_j\}$ . We follow the general recipe and introduce the *leave-one-out* sequences, denoted by  $\{\mathbf{h}^{t,(l)}, \mathbf{x}^{t,(l)}\}_{t \geq 0}$  for each  $1 \leq l \leq m$ . Specifically,  $\{\mathbf{h}^{t,(l)}, \mathbf{x}^{t,(l)}\}_{t \geq 0}$  is the gradient sequence operating on the loss function

$$f^{(l)}(\mathbf{h}, \mathbf{x}) := \sum_{j:j \neq l} \left| \mathbf{b}_j^H (\mathbf{h}\mathbf{x}^H - \mathbf{h}^* \mathbf{x}^{*H}) \mathbf{a}_j \right|^2. \tag{86}$$

The whole sequence is constructed by running gradient descent with spectral initialization on the leave-one-out loss (86). The precise description is supplied in Algorithm 6.

For notational simplicity, we denote  $z^{t,(l)} = \begin{bmatrix} \mathbf{h}^{t,(l)} \\ \mathbf{x}^{t,(l)} \end{bmatrix}$  and use  $f(z^{t,(l)}) = f(\mathbf{h}^{t,(l)}, \mathbf{x}^{t,(l)})$  interchangeably. Define similarly the alignment parameters

$$\alpha^{t,(l)} := \arg \min_{\alpha \in \mathbb{C}} \left\| \frac{1}{\alpha} \mathbf{h}^{t,(l)} - \mathbf{h}^* \right\|_2^2 + \left\| \alpha \mathbf{x}^{t,(l)} - \mathbf{x}^* \right\|_2^2, \tag{87}$$

and denote  $\tilde{z}^{t,(l)} = \begin{bmatrix} \tilde{\mathbf{h}}^{t,(l)} \\ \tilde{\mathbf{x}}^{t,(l)} \end{bmatrix}$  where

$$\tilde{\mathbf{h}}^{t,(l)} = \frac{1}{\alpha^{t,(l)}} \mathbf{h}^{t,(l)} \quad \text{and} \quad \tilde{\mathbf{x}}^{t,(l)} = \alpha^{t,(l)} \mathbf{x}^{t,(l)}. \tag{88}$$

**Algorithm 6** The  $l$ th leave-one-out sequence for blind deconvolution

**Input:**  $\{\mathbf{a}_j\}_{1 \leq j \leq m, j \neq l}, \{\mathbf{b}_j\}_{1 \leq j \leq m, j \neq l}$  and  $\{y_j\}_{1 \leq j \leq m, j \neq l}$ .

**Spectral initialization:** Let  $\sigma_1(\mathbf{M}^{(l)})$ ,  $\check{\mathbf{h}}^{0,(l)}$  and  $\check{\mathbf{x}}^{0,(l)}$  be the leading singular value, left and right singular vectors of

$$\mathbf{M}^{(l)} := \sum_{j:j \neq l} y_j \mathbf{b}_j \mathbf{a}_j^H,$$

respectively. Set  $\mathbf{h}^{0,(l)} = \sqrt{\sigma_1(\mathbf{M}^{(l)})} \check{\mathbf{h}}^{0,(l)}$  and  $\mathbf{x}^{0,(l)} = \sqrt{\sigma_1(\mathbf{M}^{(l)})} \check{\mathbf{x}}^{0,(l)}$ .

**Gradient updates:** for  $t = 0, 1, 2, \dots, T - 1$  do

$$\begin{bmatrix} \mathbf{h}^{t+1,(l)} \\ \mathbf{x}^{t+1,(l)} \end{bmatrix} = \begin{bmatrix} \mathbf{h}^{t,(l)} \\ \mathbf{x}^{t,(l)} \end{bmatrix} - \eta \begin{bmatrix} \frac{1}{\|\mathbf{x}^{t,(l)}\|_2^2} \nabla_{\mathbf{h}} f^{(l)}(\mathbf{h}^{t,(l)}, \mathbf{x}^{t,(l)}) \\ \frac{1}{\|\mathbf{h}^{t,(l)}\|_2^2} \nabla_{\mathbf{x}} f^{(l)}(\mathbf{h}^{t,(l)}, \mathbf{x}^{t,(l)}) \end{bmatrix}. \tag{89}$$

### 8.3 Step 3: Establishing the Incoherence Condition by Induction

As usual, we continue the proof in an inductive manner. For clarity of presentation, we list below the set of induction hypotheses underlying our analysis:

$$\text{dist}(\mathbf{z}^t, \mathbf{z}^*) \leq C_1 \frac{1}{\log^2 m}, \tag{90a}$$

$$\max_{1 \leq l \leq m} \text{dist}(\mathbf{z}^{t,(l)}, \tilde{\mathbf{z}}^t) \leq C_2 \frac{\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 K \log^9 m}{m}}, \tag{90b}$$

$$\max_{1 \leq l \leq m} |\mathbf{a}_l^H(\tilde{\mathbf{x}}^t - \mathbf{x}^*)| \leq C_3 \frac{1}{\log^{1.5} m}, \tag{90c}$$

$$\max_{1 \leq l \leq m} |\mathbf{b}_l^H \tilde{\mathbf{h}}^t| \leq C_4 \frac{\mu}{\sqrt{m}} \log^2 m, \tag{90d}$$

where  $\tilde{\mathbf{h}}^t$ ,  $\tilde{\mathbf{x}}^t$ , and  $\tilde{\mathbf{z}}^t$  are defined in (81). Here,  $C_1, C_3 > 0$  are some sufficiently small constants, while  $C_2, C_4 > 0$  are some sufficiently large constants. We aim to show that if these hypotheses (90) hold up to the  $t$ th iteration, then the same would hold for the  $(t + 1)$ th iteration with exceedingly high probability (e.g.,  $1 - O(m^{-10})$ ). The first hypothesis (90a) has already been established in Lemma 15, and hence, the rest of this section focuses on establishing the remaining three. To justify the incoherence hypotheses (90c) and (90d) for the  $(t + 1)$ th iteration, we need to leverage the nice properties of the leave-one-out sequences and establish (90b) first. In the sequel, we follow the steps suggested in the general recipe.

- *Step 3(a): proximity between the original and the leave-one-out iterates* We first justify hypothesis (90b) for the  $(t + 1)$ th iteration via the following lemma.

**Lemma 17** *Suppose the sample complexity obeys  $m \geq C \mu^2 K \log^9 m$  for some sufficiently large constant  $C > 0$ . Let  $\mathcal{E}_t$  be the event where hypotheses (90a)–(90d) hold for the  $t$ th iteration. Then on an event  $\mathcal{E}_{t+1,1} \subseteq \mathcal{E}_t$  obeying  $\mathbb{P}(\mathcal{E}_t \cap \mathcal{E}_{t+1,1}^c) = O(m^{-10} + m e^{-cK})$  for some constant  $c > 0$ , one has*

$$\begin{aligned} \max_{1 \leq l \leq m} \text{dist}(\mathbf{z}^{t+1,(l)}, \tilde{\mathbf{z}}^{t+1}) &\leq C_2 \frac{\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 K \log^9 m}{m}} \\ \text{and} \quad \max_{1 \leq l \leq m} \|\tilde{\mathbf{z}}^{t+1,(l)} - \tilde{\mathbf{z}}^{t+1}\|_2 &\lesssim C_2 \frac{\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 K \log^9 m}{m}}, \end{aligned}$$

provided that the step size  $\eta > 0$  is some sufficiently small constant.

**Proof** As usual, this result follows from the restricted strong convexity, which forces the distance between the two sequences of interest to be contractive. See Appendix C.3. □

- *Step 3(b): incoherence of the leave-one-out iterate  $\mathbf{x}^{t+1,(l)}$  w.r.t.  $\mathbf{a}_l$*  Next, we show that the leave-one-out iterate  $\tilde{\mathbf{x}}^{t+1,(l)}$ —which is independent of  $\mathbf{a}_l$ —is incoherent w.r.t.  $\mathbf{a}_l$  in the sense that

$$\left| \mathbf{a}_l^H (\tilde{\mathbf{x}}^{t+1,(l)} - \mathbf{x}^*) \right| \leq 10C_1 \frac{1}{\log^{3/2} m} \tag{91}$$

with probability exceeding  $1 - O(m^{-10} + e^{-K} \log m)$ . To see why, use the statistical independence and the standard Gaussian concentration inequality to show that

$$\max_{1 \leq l \leq m} \left| \mathbf{a}_l^H (\tilde{\mathbf{x}}^{t+1,(l)} - \mathbf{x}^*) \right| \leq 5\sqrt{\log m} \max_{1 \leq l \leq m} \|\tilde{\mathbf{x}}^{t+1,(l)} - \mathbf{x}^*\|_2$$

with probability exceeding  $1 - O(m^{-10})$ . It then follows from the triangle inequality that

$$\begin{aligned} \|\tilde{\mathbf{x}}^{t+1,(l)} - \mathbf{x}^*\|_2 &\leq \|\tilde{\mathbf{x}}^{t+1,(l)} - \tilde{\mathbf{x}}^{t+1}\|_2 + \|\tilde{\mathbf{x}}^{t+1} - \mathbf{x}^*\|_2 \\ &\stackrel{(i)}{\leq} CC_2 \frac{\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 K \log^9 m}{m}} + C_1 \frac{1}{\log^2 m} \\ &\stackrel{(ii)}{\leq} 2C_1 \frac{1}{\log^2 m}, \end{aligned}$$

where (i) follows from Lemmas 15 and 17 and (ii) holds as soon as  $m/(\mu^2 \sqrt{K} \log^{13/2} m)$  is sufficiently large. Combining the preceding two bounds establishes (91).

- *Step 3(c): combining the bounds to show incoherence of  $\mathbf{x}^{t+1}$  w.r.t.  $\{\mathbf{a}_l\}$*  The above bounds immediately allow us to conclude that

$$\max_{1 \leq l \leq m} \left| \mathbf{a}_l^H (\tilde{\mathbf{x}}^{t+1} - \mathbf{x}^*) \right| \leq C_3 \frac{1}{\log^{3/2} m}$$

with probability at least  $1 - O(m^{-10} + e^{-K} \log m)$ , which is exactly hypothesis (90c) for the  $(t + 1)$ th iteration. Specifically, for each  $1 \leq l \leq m$ , the triangle inequality yields

$$\begin{aligned} \left| \mathbf{a}_l^H (\tilde{\mathbf{x}}^{t+1} - \mathbf{x}^*) \right| &\leq \left| \mathbf{a}_l^H (\tilde{\mathbf{x}}^{t+1} - \tilde{\mathbf{x}}^{t+1,(l)}) \right| + \left| \mathbf{a}_l^H (\tilde{\mathbf{x}}^{t+1,(l)} - \mathbf{x}^*) \right| \\ &\stackrel{(i)}{\leq} \|\mathbf{a}_l\|_2 \left\| \tilde{\mathbf{x}}^{t+1} - \tilde{\mathbf{x}}^{t+1,(l)} \right\|_2 + \left| \mathbf{a}_l^H (\tilde{\mathbf{x}}^{t+1,(l)} - \mathbf{x}^*) \right| \\ &\stackrel{(ii)}{\leq} 3\sqrt{K} \cdot CC_2 \frac{\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 K \log^9 m}{m}} + 10C_1 \frac{1}{\log^{3/2} m} \\ &\stackrel{(iii)}{\leq} C_3 \frac{1}{\log^{3/2} m}. \end{aligned}$$

Here (i) follows from Cauchy–Schwarz; (ii) is a consequence of (190), Lemma 17, and bound (91); and the last inequality holds as long as  $m/(\mu^2 K \log^6 m)$  is sufficiently large and  $C_3 \geq 11C_1$ .

- *Step 3(d): incoherence of  $\mathbf{h}^{t+1}$  w.r.t.  $\{\mathbf{b}_l\}$*  It remains to justify that  $\mathbf{h}^{t+1}$  is also incoherent w.r.t. its associated design vectors  $\{\mathbf{b}_l\}$ . This proof of this step, however, is much more involved and challenging, due to the deterministic nature of the  $\mathbf{b}_l$ 's. As a result, we would need to “propagate” the randomness brought about by  $\{\mathbf{a}_l\}$  to  $\mathbf{h}^{t+1}$  in order to facilitate the analysis. The result is summarized as follows.

**Lemma 18** *Suppose that the sample complexity obeys  $m \geq C\mu^2 K \log^9 m$  for some sufficiently large constant  $C > 0$ . Let  $\mathcal{E}_t$  be the event where hypotheses (90a)–(90d) hold for the  $t$ th iteration. Then on an event  $\mathcal{E}_{t+1,2} \subseteq \mathcal{E}_t$  obeying  $\mathbb{P}(\mathcal{E}_t \cap \mathcal{E}_{t+1,2}^c) = O(m^{-10})$ , one has*

$$\max_{1 \leq l \leq m} \left| \mathbf{b}_l^H \tilde{\mathbf{h}}^{t+1} \right| \leq C_4 \frac{\mu}{\sqrt{m}} \log^2 m$$

as long as  $C_4$  is sufficiently large and  $\eta > 0$  is taken to be some sufficiently small constant.

**Proof** The key idea is to divide  $\{1, \dots, m\}$  into consecutive bins each of size  $\text{poly}(\log(m))$ , and to exploit the randomness (namely, the randomness from  $\mathbf{a}_l$ ) within each bin. This binning idea is crucial in ensuring that the incoherence measure of interest does not blow up as  $t$  increases. See Appendix C.4. □

With these steps in place, we conclude the proof of Theorem 3 via induction and the union bound.

### 8.4 The Base Case: Spectral Initialization

In order to finish the induction steps, we still need to justify the induction hypotheses for the base cases; namely, we need to show that the spectral initializations  $\mathbf{z}^0$  and  $\{\mathbf{z}^{0,(l)}\}_{1 \leq l \leq m}$  satisfy the induction hypotheses (90) at  $t = 0$ .

To start with, the initializations are sufficiently close to the truth when measured by the  $\ell_2$  norm, as summarized by the following lemma.

**Lemma 19** *Fix any small constant  $\xi > 0$ . Suppose the sample size obeys  $m \geq C\mu^2 K \log^2 m / \xi^2$  for some sufficiently large constant  $C > 0$ . Then with probability at least  $1 - O(m^{-10})$ , we have*

$$\min_{\alpha \in \mathbb{C}, |\alpha|=1} \left\{ \|\alpha \mathbf{h}^0 - \mathbf{h}^*\|_2 + \|\alpha \mathbf{x}^0 - \mathbf{x}^*\|_2 \right\} \leq \xi \quad \text{and} \tag{92}$$

$$\min_{\alpha \in \mathbb{C}, |\alpha|=1} \left\{ \|\alpha \mathbf{h}^{0,(l)} - \mathbf{h}^*\|_2 + \|\alpha \mathbf{x}^{0,(l)} - \mathbf{x}^*\|_2 \right\} \leq \xi, \quad 1 \leq l \leq m, \tag{93}$$

and  $||\alpha_0| - 1| \leq 1/4$ .

**Proof** This follows from Wedin’s  $\sin\Theta$  theorem [121] and [76, Lemma 5.20]. See Appendix C.5. □



From the definition of  $\text{dist}(\cdot, \cdot)$  (cf. (34)), we immediately have

$$\begin{aligned} \text{dist}(\mathbf{z}^0, \mathbf{z}^*) &= \min_{\alpha \in \mathbb{C}} \sqrt{\left\| \frac{1}{\alpha} \mathbf{h} - \mathbf{h}^* \right\|_2^2 + \|\alpha \mathbf{x} - \mathbf{x}^*\|_2^2} \\ &\stackrel{(i)}{\leq} \min_{\alpha \in \mathbb{C}} \left\{ \left\| \frac{1}{\alpha} \mathbf{h} - \mathbf{h}^* \right\|_2 + \|\alpha \mathbf{x} - \mathbf{x}^*\|_2 \right\} \\ &\stackrel{(ii)}{\leq} \min_{\alpha \in \mathbb{C}, |\alpha|=1} \left\{ \|\alpha \mathbf{h}^0 - \mathbf{h}^*\|_2 + \|\alpha \mathbf{x}^0 - \mathbf{x}^*\|_2 \right\} \stackrel{(iii)}{\leq} C_1 \frac{1}{\log^2 m}, \end{aligned} \tag{94}$$

as long as  $m \geq C \mu^2 K \log^6 m$  for some sufficiently large constant  $C > 0$ . Here (i) follows from the elementary inequality that  $a^2 + b^2 \leq (a + b)^2$  for positive  $a$  and  $b$ , (ii) holds since the feasible set of the latter one is strictly smaller, and (iii) follows directly from Lemma 19. This finishes the proof of (90a) for  $t = 0$ . Similarly, with high probability we have

$$\begin{aligned} \text{dist}(\mathbf{z}^{0,(l)}, \mathbf{z}^*) &\leq \min_{\alpha \in \mathbb{C}, |\alpha|=1} \left\{ \|\alpha \mathbf{h}^{0,(l)} - \mathbf{h}^*\|_2 + \|\alpha \mathbf{x}^{0,(l)} - \mathbf{x}^*\|_2 \right\} \lesssim \frac{1}{\log^2 m}, \\ 1 \leq l \leq m. \end{aligned} \tag{95}$$

Next, when properly aligned, the true initial estimate  $\mathbf{z}^0$  and the leave-one-out estimate  $\mathbf{z}^{0,(l)}$  are expected to be sufficiently close, as claimed by the following lemma. Along the way, we show that  $\mathbf{h}^0$  is incoherent w.r.t. the sampling vectors  $\{\mathbf{b}_l\}$ . This establishes (90b) and (90d) for  $t = 0$ .

**Lemma 20** *Suppose that  $m \geq C \mu^2 K \log^3 m$  for some sufficiently large constant  $C > 0$ . Then with probability at least  $1 - O(m^{-10})$ , one has*

$$\max_{1 \leq l \leq m} \text{dist}(\mathbf{z}^{0,(l)}, \tilde{\mathbf{z}}^0) \leq C_2 \frac{\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 K \log^5 m}{m}} \tag{96}$$

and

$$\max_{1 \leq l \leq m} |\mathbf{b}_l^H \tilde{\mathbf{h}}^0| \leq C_4 \frac{\mu \log^2 m}{\sqrt{m}}. \tag{97}$$

**Proof** The key is to establish that  $\text{dist}(\mathbf{z}^{0,(l)}, \tilde{\mathbf{z}}^0)$  can be upper bounded by some linear scaling of  $|\mathbf{b}_l^H \tilde{\mathbf{h}}^0|$ , and vice versa. This allows us to derive bounds simultaneously for both quantities. See Appendix C.6.  $\square$

Finally, we establish (90c) regarding the incoherence of  $\mathbf{x}^0$  with respect to the design vectors  $\{\mathbf{a}_l\}$ .

**Lemma 21** *Suppose that  $m \geq C \mu^2 K \log^6 m$  for some sufficiently large constant  $C > 0$ . Then with probability exceeding  $1 - O(m^{-10})$ , we have*

$$\max_{1 \leq l \leq m} |a_l^H(\tilde{x}^0 - x^*)| \leq C_3 \frac{1}{\log^{1.5} m}.$$

**Proof** See Appendix C.7. □

## 9 Discussion

This paper showcases an important phenomenon in nonconvex optimization: even without explicit enforcement of regularization, the vanilla form of gradient descent effectively achieves implicit regularization for a large family of statistical estimation problems. We believe this phenomenon arises in problems far beyond the three cases studied herein, and our results are initial steps toward understanding this fundamental phenomenon. There are numerous avenues open for future investigation, and we point out a few of them.

- *Improving sample complexity* In the current paper, the required sample complexity  $O(\mu^3 r^3 n \log^3 n)$  for matrix completion is suboptimal when the rank  $r$  of the underlying matrix is large. While this allows us to achieve a dimension-free iteration complexity, it is slightly higher than the sample complexity derived for regularized gradient descent in [32]. We expect our results continue to hold under lower sample complexity  $O(\mu^2 r^2 n \log n)$ , but it calls for a more refined analysis (e.g., a generic chaining argument).
- *Leave-one-out tricks for more general designs* So far, our focus is on independent designs, including the i.i.d. Gaussian design adopted in phase retrieval and partially in blind deconvolution, as well as the independent sampling mechanism in matrix completion. Such independence property creates some sort of “statistical homogeneity,” for which the leave-one-out argument works beautifully. It remains unclear how to generalize such leave-one-out tricks for more general designs (e.g., more general sampling patterns in matrix completion and more structured Fourier designs in phase retrieval and blind deconvolution). In fact, the readers can already get a flavor of this issue in the analysis of blind deconvolution, where the Fourier design vectors require much more delicate treatments than purely Gaussian designs.
- *Uniform stability* The leave-one-out perturbation argument is established upon a basic fact: when we exclude one sample from consideration, the resulting estimates/predictions do not deviate much from the original ones. This leave-one-out stability bears similarity to the notion of uniform stability studied in statistical learning theory [8]. We expect our analysis framework to be helpful for analyzing other learning algorithms that are uniformly stable.
- *Other iterative methods and other loss functions* The focus of the current paper has been the analysis of vanilla GD tailored to the natural squared loss. This is by no means to advocate GD as the top-performing algorithm in practice; rather, we are using this simple algorithm to isolate some seemingly pervasive phenomena (i.e., implicit regularization) that generic optimization theory fails to account for. The simplicity of vanilla GD makes it an ideal object to initiate such discussions. That being said, practitioners should definitely explore as many algorithmic alter-

natives as possible before settling on a particular algorithm. Take phase retrieval for example: iterative methods other than GD and/or algorithms tailored to other loss functions have been proposed in the nonconvex optimization literature, including but not limited to alternating minimization, block coordinate descent, and sub-gradient methods and prox-linear methods tailed to nonsmooth losses. It would be interesting to develop a full theoretical understanding of a broader class of iterative algorithms, and to conduct a careful comparison regarding which loss functions lead to the most desirable practical performance.

- *Connections to deep learning?* We have focused on nonlinear systems that are bilinear or quadratic in this paper. Deep learning formulations/architectures, highly nonlinear, are notorious for their daunting nonconvex geometry. However, iterative methods including stochastic gradient descent have enjoyed enormous practical success in learning neural networks (e.g., [46, 103, 132]), even when the architecture is significantly over-parameterized without explicit regularization. We hope the message conveyed in this paper for several simple statistical models can shed light on why simple forms of gradient descent and variants work so well in learning complicated neural networks.

Finally, while the present paper provides a general recipe for problem-specific analyses of nonconvex algorithms, we acknowledge that a unified theory of this kind has yet to be developed. As a consequence, each problem requires delicate and somewhat lengthy analyses of its own. It would certainly be helpful if one could single out a few stylized structural properties/elements (like sparsity and incoherence in compressed sensing [13]) that enable near-optimal performance guarantees through an overarching method of analysis; with this in place, one would not need to start each problem from scratch. Having said that, we believe that our current theory elucidates a few ingredients (e.g., the region of incoherence and leave-one-out stability) that might serve as crucial building blocks for such a general theory. We invite the interested readers to contribute toward this path forward.

**Acknowledgements** Y. Chen is supported in part by the AFOSR YIP Award FA9550-19-1-0030, ONR Grant N00014-19-1-2120, ARO Grant W911NF-18-1-0303, NSF Grant CCF-1907661, and the Princeton SEAS innovation award. Y. Chi is supported in part by the Grants AFOSR FA9550-15-1-0205, ONR N00014-18-1-2142 and N00014-19-1-2404, ARO W911NF-18-1-0303, NSF CCF-1826519, ECCS-1818571, CCF-1806154. Y. Chen thanks Yudong Chen for inspiring discussions about matrix completion.

**Open Access** This article is distributed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made.

## A Proofs for Phase Retrieval

Before proceeding, we gather a few simple facts. The standard concentration inequality for  $\chi^2$  random variables together with the union bound reveals that the sampling vectors  $\{\mathbf{a}_j\}$  obey

$$\max_{1 \leq j \leq m} \|\mathbf{a}_j\|_2 \leq \sqrt{6n} \tag{98}$$

with probability at least  $1 - O(me^{-1.5n})$ . In addition, standard Gaussian concentration inequalities give

$$\max_{1 \leq j \leq m} \left| \mathbf{a}_j^\top \mathbf{x}^\star \right| \leq 5\sqrt{\log n} \tag{99}$$

with probability exceeding  $1 - O(mn^{-10})$ .

### A.1 Proof of Lemma 1

We start with the smoothness bound, namely  $\nabla^2 f(\mathbf{x}) \preceq O(\log n) \cdot \mathbf{I}_n$ . It suffices to prove the upper bound  $\|\nabla^2 f(\mathbf{x})\| \lesssim \log n$ . To this end, we first decompose the Hessian (cf. (44)) into three components as follows:

$$\begin{aligned} \nabla^2 f(\mathbf{x}) &= \underbrace{\frac{3}{m} \sum_{j=1}^m \left[ \left( \mathbf{a}_j^\top \mathbf{x} \right)^2 - \left( \mathbf{a}_j^\top \mathbf{x}^\star \right)^2 \right] \mathbf{a}_j \mathbf{a}_j^\top}_{:=\Lambda_1} \\ &\quad + \underbrace{\frac{2}{m} \sum_{j=1}^m \left( \mathbf{a}_j^\top \mathbf{x}^\star \right)^2 \mathbf{a}_j \mathbf{a}_j^\top - 2 \left( \mathbf{I}_n + 2\mathbf{x}^\star \mathbf{x}^{\star\top} \right)}_{:=\Lambda_2} + \underbrace{2 \left( \mathbf{I}_n + 2\mathbf{x}^\star \mathbf{x}^{\star\top} \right)}_{:=\Lambda_3}, \end{aligned}$$

where we have used  $y_j = (\mathbf{a}_j^\top \mathbf{x}^\star)^2$ . In the sequel, we control the three terms  $\Lambda_1$ ,  $\Lambda_2$ , and  $\Lambda_3$  in reverse order.

- The third term  $\Lambda_3$  can be easily bounded by

$$\|\Lambda_3\| \leq 2 \left( \|\mathbf{I}_n\| + 2 \|\mathbf{x}^\star \mathbf{x}^{\star\top}\| \right) = 6.$$

- The second term  $\Lambda_2$  can be controlled by means of Lemma 32:

$$\|\Lambda_2\| \leq 2\delta$$

for an arbitrarily small constant  $\delta > 0$ , as long as  $m \geq c_0 n \log n$  for  $c_0$  sufficiently large.

- It thus remains to control  $\Lambda_1$ . Toward this, we discover that

$$\|\Lambda_1\| \leq \left\| \frac{3}{m} \sum_{j=1}^m \left| \mathbf{a}_j^\top (\mathbf{x} - \mathbf{x}^\star) \right| \left| \mathbf{a}_j^\top (\mathbf{x} + \mathbf{x}^\star) \right| \mathbf{a}_j \mathbf{a}_j^\top \right\|. \tag{100}$$

Under the assumption  $\max_{1 \leq j \leq m} \left| \mathbf{a}_j^\top (\mathbf{x} - \mathbf{x}^\star) \right| \leq C_2 \sqrt{\log n}$  and fact (99), we can also obtain

$$\begin{aligned} \max_{1 \leq j \leq m} \left| \mathbf{a}_j^\top (\mathbf{x} + \mathbf{x}^\star) \right| &\leq 2 \max_{1 \leq j \leq m} \left| \mathbf{a}_j^\top \mathbf{x}^\star \right| + \max_{1 \leq j \leq m} \left| \mathbf{a}_j^\top (\mathbf{x} - \mathbf{x}^\star) \right| \\ &\leq (10 + C_2) \sqrt{\log n}. \end{aligned}$$

Substitution into (100) leads to

$$\|\mathbf{\Lambda}_1\| \leq 3C_2(10 + C_2) \log n \cdot \left\| \frac{1}{m} \sum_{j=1}^m \mathbf{a}_j \mathbf{a}_j^\top \right\| \leq 4C_2(10 + C_2) \log n,$$

where the last inequality is a direct consequence of Lemma 31.

Combining the above bounds on  $\mathbf{\Lambda}_1$ ,  $\mathbf{\Lambda}_2$ , and  $\mathbf{\Lambda}_3$  yields

$$\begin{aligned} \left\| \nabla^2 f(\mathbf{x}) \right\| &\leq \|\mathbf{\Lambda}_1\| + \|\mathbf{\Lambda}_2\| + \|\mathbf{\Lambda}_3\| \leq 4C_2(10 + C_2) \log n + 2\delta + 6 \\ &\leq 5C_2(10 + C_2) \log n, \end{aligned}$$

as long as  $n$  is sufficiently large. This establishes the claimed smoothness property.

Next, we move on to the strong convexity lower bound. Picking a constant  $C > 0$  and enforcing proper truncation, we get

$$\begin{aligned} \nabla^2 f(\mathbf{x}) &= \frac{1}{m} \sum_{j=1}^m \left[ 3 \left( \mathbf{a}_j^\top \mathbf{x} \right)^2 - y_j \right] \mathbf{a}_j \mathbf{a}_j^\top \\ &\geq \underbrace{\frac{3}{m} \sum_{j=1}^m \left( \mathbf{a}_j^\top \mathbf{x} \right)^2 \mathbb{1}_{\left\{ \left| \mathbf{a}_j^\top \mathbf{x} \right| \leq C \right\}} \mathbf{a}_j \mathbf{a}_j^\top}_{:=\mathbf{\Lambda}_4} - \underbrace{\frac{1}{m} \sum_{j=1}^m \left( \mathbf{a}_j^\top \mathbf{x}^* \right)^2 \mathbf{a}_j \mathbf{a}_j^\top}_{:=\mathbf{\Lambda}_5}. \end{aligned}$$

We begin with the simpler term  $\mathbf{\Lambda}_5$ . Lemma 32 implies that with probability at least  $1 - O(n^{-10})$ ,

$$\left\| \mathbf{\Lambda}_5 - \left( \mathbf{I}_n + 2\mathbf{x}^* \mathbf{x}^{*\top} \right) \right\| \leq \delta$$

holds for any small constant  $\delta > 0$ , as long as  $m/(n \log n)$  is sufficiently large. This reveals that

$$\mathbf{\Lambda}_5 \leq (1 + \delta) \cdot \mathbf{I}_n + 2\mathbf{x}^* \mathbf{x}^{*\top}.$$

To bound  $\mathbf{\Lambda}_4$ , invoke Lemma 33 to conclude that with probability at least  $1 - c_3 e^{-c_2 m}$  (for some constants  $c_2, c_3 > 0$ ),

$$\left\| \mathbf{\Lambda}_4 - 3 \left( \beta_1 \mathbf{x} \mathbf{x}^\top + \beta_2 \|\mathbf{x}\|_2^2 \mathbf{I}_n \right) \right\| \leq \delta \|\mathbf{x}\|_2^2$$

for any small constant  $\delta > 0$ , provided that  $m/n$  is sufficiently large. Here,

$$\beta_1 := \mathbb{E} \left[ \xi^4 \mathbb{1}_{\left\{ |\xi| \leq C \right\}} \right] - \mathbb{E} \left[ \xi^2 \mathbb{1}_{\left\{ |\xi| \leq C \right\}} \right] \quad \text{and} \quad \beta_2 := \mathbb{E} \left[ \xi^2 \mathbb{1}_{\left\{ |\xi| \leq C \right\}} \right],$$

where the expectation is taken with respect to  $\xi \sim \mathcal{N}(0, 1)$ . By the assumption  $\|\mathbf{x} - \mathbf{x}^*\|_2 \leq 2C_1$ , one has

$$\begin{aligned} \|\mathbf{x}\|_2 \leq 1 + 2C_1, \quad \left| \|\mathbf{x}\|_2^2 - \|\mathbf{x}^*\|_2^2 \right| \leq 2C_1(4C_1 + 1), \quad \left\| \mathbf{x}^*\mathbf{x}^{*\top} - \mathbf{x}\mathbf{x}^\top \right\| \\ \leq 6C_1(4C_1 + 1), \end{aligned}$$

which leads to

$$\begin{aligned} & \left\| \Lambda_4 - 3 \left( \beta_1 \mathbf{x}^*\mathbf{x}^{*\top} + \beta_2 \mathbf{I}_n \right) \right\| \\ & \leq \left\| \Lambda_4 - 3 \left( \beta_1 \mathbf{x}\mathbf{x}^\top + \beta_2 \|\mathbf{x}\|_2^2 \mathbf{I}_n \right) \right\| \\ & \quad + 3 \left\| \left( \beta_1 \mathbf{x}^*\mathbf{x}^{*\top} + \beta_2 \mathbf{I}_n \right) - \left( \beta_1 \mathbf{x}\mathbf{x}^\top + \beta_2 \|\mathbf{x}\|_2^2 \mathbf{I}_n \right) \right\| \\ & \leq \delta \|\mathbf{x}\|_2^2 + 3\beta_1 \left\| \mathbf{x}^*\mathbf{x}^{*\top} - \mathbf{x}\mathbf{x}^\top \right\| + 3\beta_2 \left\| \mathbf{I}_n - \|\mathbf{x}\|_2^2 \mathbf{I}_n \right\| \\ & \leq \delta(1 + 2C_1)^2 + 18\beta_1 C_1(4C_1 + 1) + 6\beta_2 C_1(4C_1 + 1). \end{aligned}$$

This further implies

$$\begin{aligned} \Lambda_4 \geq 3 \left( \beta_1 \mathbf{x}^*\mathbf{x}^{*\top} + \beta_2 \mathbf{I}_n \right) \\ - \left[ \delta(1 + 2C_1)^2 + 18\beta_1 C_1(4C_1 + 1) + 6\beta_2 C_1(4C_1 + 1) \right] \mathbf{I}_n. \end{aligned}$$

Recognizing that  $\beta_1$  (resp.  $\beta_2$ ) approaches 2 (resp. 1) as  $C$  grows, we can thus take  $C_1$  small enough and  $C$  large enough to guarantee that

$$\Lambda_4 \geq 5\mathbf{x}^*\mathbf{x}^{*\top} + 2\mathbf{I}_n.$$

Putting the preceding two bounds on  $\Lambda_4$  and  $\Lambda_5$  together yields

$$\nabla^2 f(\mathbf{x}) \geq 5\mathbf{x}^*\mathbf{x}^{*\top} + 2\mathbf{I}_n - \left[ (1 + \delta) \cdot \mathbf{I}_n + 2\mathbf{x}^*\mathbf{x}^{*\top} \right] \geq (1/2) \cdot \mathbf{I}_n$$

as claimed.

### A.2 Proof of Lemma 2

Using the update rule (cf. (17)) as well as the fundamental theorem of calculus [70, Chapter XIII, Theorem 4.2], we get

$$\begin{aligned} \mathbf{x}^{t+1} - \mathbf{x}^* &= \mathbf{x}^t - \eta \nabla f(\mathbf{x}^t) - [\mathbf{x}^* - \eta \nabla f(\mathbf{x}^*)] \\ &= \left[ \mathbf{I}_n - \eta \int_0^1 \nabla^2 f(\mathbf{x}(\tau)) \, d\tau \right] (\mathbf{x}^t - \mathbf{x}^*), \end{aligned}$$

where we denote  $\mathbf{x}(\tau) = \mathbf{x}^* + \tau(\mathbf{x}^t - \mathbf{x}^*)$ ,  $0 \leq \tau \leq 1$ . Here, the first equality makes use of the fact that  $\nabla f(\mathbf{x}^*) = \mathbf{0}$ . Under condition (45), it is self-evident that for all  $0 \leq \tau \leq 1$ ,

$$\begin{aligned} \|\mathbf{x}(\tau) - \mathbf{x}^*\|_2 &= \|\tau(\mathbf{x}^t - \mathbf{x}^*)\|_2 \leq 2C_1 \quad \text{and} \\ \max_{1 \leq l \leq m} \left| \mathbf{a}_l^\top (\mathbf{x}(\tau) - \mathbf{x}^*) \right| &\leq \max_{1 \leq l \leq m} \left| \mathbf{a}_l^\top \tau (\mathbf{x}^t - \mathbf{x}^*) \right| \leq C_2 \sqrt{\log n}. \end{aligned}$$

This means that for all  $0 \leq \tau \leq 1$ ,

$$(1/2) \cdot \mathbf{I}_n \leq \nabla^2 f(\mathbf{x}(\tau)) \leq [5C_2(10 + C_2) \log n] \cdot \mathbf{I}_n$$

in view of Lemma 1. Picking  $\eta \leq 1/[5C_2(10 + C_2) \log n]$  (and hence  $\|\eta \nabla^2 f(\mathbf{x}(\tau))\| \leq 1$ ), one sees that

$$\mathbf{0} \leq \mathbf{I}_n - \eta \int_0^1 \nabla^2 f(\mathbf{x}(\tau)) \, d\tau \leq (1 - \eta/2) \cdot \mathbf{I}_n,$$

which immediately yields

$$\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2 \leq \left\| \mathbf{I}_n - \eta \int_0^1 \nabla^2 f(\mathbf{x}(\tau)) \, d\tau \right\| \cdot \|\mathbf{x}^t - \mathbf{x}^*\|_2 \leq (1 - \eta/2) \|\mathbf{x}^t - \mathbf{x}^*\|_2.$$

### A.3 Proof of Lemma 3

We start with proving (19a). For all  $0 \leq t \leq T_0$ , invoke Lemma 2 recursively with conditions (47) to reach

$$\|\mathbf{x}^t - \mathbf{x}^*\|_2 \leq (1 - \eta/2)^t \|\mathbf{x}^0 - \mathbf{x}^*\|_2 \leq C_1(1 - \eta/2)^t \|\mathbf{x}^*\|_2. \tag{101}$$

This finishes the proof of (19a) for  $0 \leq t \leq T_0$  and also reveals that

$$\|\mathbf{x}^{T_0} - \mathbf{x}^*\|_2 \leq C_1(1 - \eta/2)^{T_0} \|\mathbf{x}^*\|_2 \ll \frac{1}{n} \|\mathbf{x}^*\|_2, \tag{102}$$

provided that  $\eta \asymp 1/\log n$ . Applying the Cauchy–Schwarz inequality and fact (98) indicate that

$$\max_{1 \leq l \leq m} \left| \mathbf{a}_l^\top (\mathbf{x}^{T_0} - \mathbf{x}^*) \right| \leq \max_{1 \leq l \leq m} \|\mathbf{a}_l\|_2 \|\mathbf{x}^{T_0} - \mathbf{x}^*\|_2 \leq \sqrt{6n} \cdot \frac{1}{n} \|\mathbf{x}^*\|_2 \ll C_2 \sqrt{\log n},$$

leading to the satisfaction of (45). Therefore, invoking Lemma 2 yields

$$\|\mathbf{x}^{T_0+1} - \mathbf{x}^*\|_2 \leq (1 - \eta/2) \|\mathbf{x}^{T_0} - \mathbf{x}^*\|_2 \ll \frac{1}{n} \|\mathbf{x}^*\|_2.$$

One can then repeat this argument to arrive at for all  $t > T_0$

$$\|\mathbf{x}^t - \mathbf{x}^*\|_2 \leq (1 - \eta/2)^t \|\mathbf{x}^0 - \mathbf{x}^*\|_2 \leq C_1(1 - \eta/2)^t \|\mathbf{x}^*\|_2 \ll \frac{1}{n} \|\mathbf{x}^*\|_2. \tag{103}$$

We are left with (19b). It is self-evident that the iterates from  $0 \leq t \leq T_0$  satisfy (19b) by assumptions. For  $t > T_0$ , we can use the Cauchy–Schwarz inequality to obtain

$$\max_{1 \leq j \leq m} |\mathbf{a}_j^\top (\mathbf{x}^t - \mathbf{x}^*)| \leq \max_{1 \leq j \leq m} \|\mathbf{a}_j\|_2 \|\mathbf{x}^t - \mathbf{x}^*\|_2 \ll \sqrt{n} \cdot \frac{1}{n} \leq C_2 \sqrt{\log n},$$

where the penultimate relation uses conditions (98) and (103).

### A.4 Proof of Lemma 4

First, going through the same derivation as in (54) and (55) will result in

$$\max_{1 \leq l \leq m} |\mathbf{a}_l^\top (\mathbf{x}^{t,(l)} - \mathbf{x}^*)| \leq C_4 \sqrt{\log n} \tag{104}$$

for some  $C_4 < C_2$ , which will be helpful for our analysis.

We use the gradient update rules once again to decompose

$$\begin{aligned} \mathbf{x}^{t+1} - \mathbf{x}^{t+1,(l)} &= \mathbf{x}^t - \eta \nabla f(\mathbf{x}^t) - \left[ \mathbf{x}^{t,(l)} - \eta \nabla f^{(l)}(\mathbf{x}^{t,(l)}) \right] \\ &= \mathbf{x}^t - \eta \nabla f(\mathbf{x}^t) - \left[ \mathbf{x}^{t,(l)} - \eta \nabla f(\mathbf{x}^{t,(l)}) \right] \\ &\quad - \eta \left[ \nabla f(\mathbf{x}^{t,(l)}) - \nabla f^{(l)}(\mathbf{x}^{t,(l)}) \right] \\ &= \underbrace{\mathbf{x}^t - \mathbf{x}^{t,(l)} - \eta \left[ \nabla f(\mathbf{x}^t) - \nabla f(\mathbf{x}^{t,(l)}) \right]}_{:= \mathbf{v}_1^{(l)}} \\ &\quad - \underbrace{\eta \frac{1}{m} \left[ (\mathbf{a}_l^\top \mathbf{x}^{t,(l)})^2 - (\mathbf{a}_l^\top \mathbf{x}^*)^2 \right] (\mathbf{a}_l^\top \mathbf{x}^{t,(l)}) \mathbf{a}_l}_{:= \mathbf{v}_2^{(l)}}, \end{aligned}$$

where the last line comes from the definition of  $\nabla f(\cdot)$  and  $\nabla f^{(l)}(\cdot)$ .

1. We first control the term  $\mathbf{v}_2^{(l)}$ , which is easier to deal with. Specifically,

$$\begin{aligned} \|\mathbf{v}_2^{(l)}\|_2 &\leq \eta \frac{\|\mathbf{a}_l\|_2}{m} \left| (\mathbf{a}_l^\top \mathbf{x}^{t,(l)})^2 - (\mathbf{a}_l^\top \mathbf{x}^*)^2 \right| |\mathbf{a}_l^\top \mathbf{x}^{t,(l)}| \\ &\stackrel{(i)}{\lesssim} C_4(C_4 + 5)(C_4 + 10) \eta \frac{n \log n}{m} \sqrt{\frac{\log n}{n}} \stackrel{(ii)}{\leq} c \eta \sqrt{\frac{\log n}{n}}, \end{aligned}$$

for any small constant  $c > 0$ . Here (i) follows since (98) and, in view of (99) and (104),

$$\begin{aligned} \left| (\mathbf{a}_l^\top \mathbf{x}^{t,(l)})^2 - (\mathbf{a}_l^\top \mathbf{x}^*)^2 \right| &\leq |\mathbf{a}_l^\top (\mathbf{x}^{t,(l)} - \mathbf{x}^*)| \left( |\mathbf{a}_l^\top (\mathbf{x}^{t,(l)} - \mathbf{x}^*)| + 2 |\mathbf{a}_l^\top \mathbf{x}^*| \right) \\ &\leq C_4(C_4 + 10) \log n, \end{aligned}$$



$$\text{and } \left| \mathbf{a}_l^\top \mathbf{x}^{t,(l)} \right| \leq \left| \mathbf{a}_l^\top (\mathbf{x}^{t,(l)} - \mathbf{x}^*) \right| + \left| \mathbf{a}_l^\top \mathbf{x}^* \right| \leq (C_4 + 5)\sqrt{\log n}.$$

And (ii) holds as long as  $m \gg n \log n$ .

- For the term  $\mathbf{v}_1^{(l)}$ , the fundamental theorem of calculus [70, Chapter XIII, Theorem 4.2] tells us that

$$\mathbf{v}_1^{(l)} = \left[ \mathbf{I}_n - \eta \int_0^1 \nabla^2 f(\mathbf{x}(\tau)) \, d\tau \right] (\mathbf{x}^t - \mathbf{x}^{t,(l)}),$$

where we abuse the notation and denote  $\mathbf{x}(\tau) = \mathbf{x}^{t,(l)} + \tau(\mathbf{x}^t - \mathbf{x}^{t,(l)})$ . By the induction hypotheses (51) and condition (104), one can verify that

$$\begin{aligned} \|\mathbf{x}(\tau) - \mathbf{x}^*\|_2 &\leq \tau \|\mathbf{x}^t - \mathbf{x}^*\|_2 + (1 - \tau) \|\mathbf{x}^{t,(l)} - \mathbf{x}^*\|_2 \leq 2C_1 \quad \text{and} \\ \max_{1 \leq l \leq m} \left| \mathbf{a}_l^\top (\mathbf{x}(\tau) - \mathbf{x}^*) \right| &\leq \tau \max_{1 \leq l \leq m} \left| \mathbf{a}_l^\top (\mathbf{x}^t - \mathbf{x}^*) \right| + (1 - \tau) \\ &\quad \max_{1 \leq l \leq m} \left| \mathbf{a}_l^\top (\mathbf{x}^{t,(l)} - \mathbf{x}^*) \right| \leq C_2 \sqrt{\log n} \end{aligned} \tag{105}$$

for all  $0 \leq \tau \leq 1$ , as long as  $C_4 \leq C_2$ . The second line follows directly from (104). To see why (105) holds, we note that

$$\|\mathbf{x}^{t,(l)} - \mathbf{x}^*\|_2 \leq \|\mathbf{x}^{t,(l)} - \mathbf{x}^t\|_2 + \|\mathbf{x}^t - \mathbf{x}^*\|_2 \leq C_3 \sqrt{\frac{\log n}{n}} + C_1,$$

where the second inequality follows from the induction hypotheses (51b) and (51a). This combined with (51a) gives

$$\|\mathbf{x}(\tau) - \mathbf{x}^*\|_2 \leq \tau C_1 + (1 - \tau) \left( C_3 \sqrt{\frac{\log n}{n}} + C_1 \right) \leq 2C_1$$

as long as  $n$  is large enough, thus justifying (105). Hence, by Lemma 1,  $\nabla^2 f(\mathbf{x}(\tau))$  is positive definite and almost well conditioned. By choosing  $0 < \eta \leq 1/[5C_2(10 + C_2)\log n]$ , we get

$$\|\mathbf{v}_1^{(l)}\|_2 \leq (1 - \eta/2) \|\mathbf{x}^t - \mathbf{x}^{t,(l)}\|_2.$$

- Combine the preceding bounds on  $\mathbf{v}_1^{(l)}$  and  $\mathbf{v}_2^{(l)}$  as well as the induction bound (51b) to arrive at

$$\|\mathbf{x}^{t+1} - \mathbf{x}^{t+1,(l)}\|_2 \leq (1 - \eta/2) \|\mathbf{x}^t - \mathbf{x}^{t,(l)}\|_2 + c\eta \sqrt{\frac{\log n}{n}} \leq C_3 \sqrt{\frac{\log n}{n}}. \tag{106}$$

This establishes (53) for the  $(t + 1)$ th iteration.

**A.5 Proof of Lemma 5**

In view of the assumption (42) that  $\|x^0 - x^*\|_2 \leq \|x^0 + x^*\|_2$  and the fact that  $x^0 = \sqrt{\lambda_1(Y)}/3 \tilde{x}^0$  for some  $\lambda_1(Y) > 0$  (which we will verify below), it is straightforward to see that

$$\|\tilde{x}^0 - x^*\|_2 \leq \|\tilde{x}^0 + x^*\|_2.$$

One can then invoke the Davis–Kahan  $\sin\Theta$  theorem [124, Corollary 1] to obtain

$$\|\tilde{x}^0 - x^*\|_2 \leq 2\sqrt{2} \frac{\|Y - \mathbb{E}[Y]\|}{\lambda_1(\mathbb{E}[Y]) - \lambda_2(\mathbb{E}[Y])}.$$

Note that (56)— $\|Y - \mathbb{E}[Y]\| \leq \delta$ —is a direct consequence of Lemma 32. Additionally, the fact that  $\mathbb{E}[Y] = I + 2x^*x^{*\top}$  gives  $\lambda_1(\mathbb{E}[Y]) = 3$ ,  $\lambda_2(\mathbb{E}[Y]) = 1$ , and  $\lambda_1(\mathbb{E}[Y]) - \lambda_2(\mathbb{E}[Y]) = 2$ . Combining this spectral gap and the inequality  $\|Y - \mathbb{E}[Y]\| \leq \delta$ , we arrive at

$$\|\tilde{x}^0 - x^*\|_2 \leq \sqrt{2}\delta.$$

To connect this bound with  $x^0$ , we need to take into account the scaling factor  $\sqrt{\lambda_1(Y)}/3$ . To this end, it follows from Weyl’s inequality and (56) that

$$|\lambda_1(Y) - 3| = |\lambda_1(Y) - \lambda_1(\mathbb{E}[Y])| \leq \|Y - \mathbb{E}[Y]\| \leq \delta$$

and, as a consequence,  $\lambda_1(Y) \geq 3 - \delta > 0$  when  $\delta \leq 1$ . This further implies that

$$\left| \sqrt{\frac{\lambda_1(Y)}{3}} - 1 \right| = \left| \frac{\frac{\lambda_1(Y)}{3} - 1}{\sqrt{\frac{\lambda_1(Y)}{3} + 1}} \right| \leq \left| \frac{\lambda_1(Y)}{3} - 1 \right| \leq \frac{1}{3}\delta, \tag{107}$$

where we have used the elementary identity  $\sqrt{a} - \sqrt{b} = (a - b) / (\sqrt{a} + \sqrt{b})$ . With these bounds in place, we can use the triangle inequality to get

$$\begin{aligned} \|x^0 - x^*\|_2 &= \left\| \sqrt{\frac{\lambda_1(Y)}{3}} \tilde{x}^0 - x^* \right\|_2 = \left\| \sqrt{\frac{\lambda_1(Y)}{3}} \tilde{x}^0 - \tilde{x}^0 + \tilde{x}^0 - x^* \right\|_2 \\ &\leq \left| \sqrt{\frac{\lambda_1(Y)}{3}} - 1 \right| + \|\tilde{x}^0 - x^*\|_2 \\ &\leq \frac{1}{3}\delta + \sqrt{2}\delta \leq 2\delta. \end{aligned}$$

### A.6 Proof of Lemma 6

To begin with, repeating the same argument as in Lemma 5 (which we omit here for conciseness), we see that for any fixed constant  $\delta > 0$ ,

$$\|Y^{(l)} - \mathbb{E}[Y^{(l)}]\| \leq \delta, \quad \|x^{0,(l)} - x^*\|_2 \leq 2\delta, \quad \|\tilde{x}^{0,(l)} - x^*\|_2 \leq \sqrt{2}\delta, \quad 1 \leq l \leq m \tag{108}$$

holds with probability at least  $1 - O(mn^{-10})$  as long as  $m \gg n \log n$ . The  $\ell_2$  bound on  $\|x^0 - x^{0,(l)}\|_2$  is derived as follows.

1. We start by controlling  $\|x^0 - \tilde{x}^{0,(l)}\|_2$ . Combining (57) and (108) yields

$$\|x^0 - \tilde{x}^{0,(l)}\|_2 \leq \|x^0 - x^*\|_2 + \|\tilde{x}^{0,(l)} - x^*\|_2 \leq 2\sqrt{2}\delta.$$

For  $\delta$  sufficiently small, this implies that  $\|x^0 - \tilde{x}^{0,(l)}\|_2 \leq \|x^0 + \tilde{x}^{0,(l)}\|_2$ , and hence, the Davis–Kahan  $\sin\Theta$  theorem [39] gives

$$\|\tilde{x}^0 - \tilde{x}^{0,(l)}\|_2 \leq \frac{\|(Y - Y^{(l)})\tilde{x}^{0,(l)}\|_2}{\lambda_1(Y) - \lambda_2(Y^{(l)})} \leq \|(Y - Y^{(l)})\tilde{x}^{0,(l)}\|_2. \tag{109}$$

Here, the second inequality uses Weyl’s inequality:

$$\begin{aligned} \lambda_1(Y) - \lambda_2(Y^{(l)}) &\geq \lambda_1(\mathbb{E}[Y]) - \|Y - \mathbb{E}[Y]\| - \lambda_2(\mathbb{E}[Y^{(l)}]) - \|Y^{(l)} - \mathbb{E}[Y^{(l)}]\| \\ &\geq 3 - \delta - 1 - \delta \geq 1, \end{aligned}$$

with the proviso that  $\delta \leq 1/2$ .

2. We now connect  $\|x^0 - x^{0,(l)}\|_2$  with  $\|\tilde{x}^0 - \tilde{x}^{0,(l)}\|_2$ . Applying the Weyl’s inequality and (56) yields

$$|\lambda_1(Y) - 3| \leq \|Y - \mathbb{E}[Y]\| \leq \delta \implies \lambda_1(Y) \in [3 - \delta, 3 + \delta] \subseteq [2, 4] \tag{110}$$

and, similarly,  $\lambda_1(Y^{(l)}), \|Y\|, \|Y^{(l)}\| \in [2, 4]$ . Invoke Lemma 34 to arrive at

$$\begin{aligned} \frac{1}{\sqrt{3}}\|x^0 - x^{0,(l)}\|_2 &\leq \frac{\|(Y - Y^{(l)})\tilde{x}^{0,(l)}\|_2}{2\sqrt{2}} + \left(2 + \frac{4}{\sqrt{2}}\right)\|\tilde{x}^0 - \tilde{x}^{0,(l)}\|_2 \\ &\leq 6\|(Y - Y^{(l)})\tilde{x}^{0,(l)}\|_2, \end{aligned} \tag{111}$$

where the last inequality comes from (109).

3. Everything then boils down to controlling  $\|(Y - Y^{(l)})\tilde{x}^{0,(l)}\|_2$ . Toward this, we observe that

$$\begin{aligned} \max_{1 \leq l \leq m} \|(Y - Y^{(l)})\tilde{x}^{0,(l)}\|_2 &= \max_{1 \leq l \leq m} \frac{1}{m} \left\| \left( a_l^\top x^* \right)^2 a_l a_l^\top \tilde{x}^{0,(l)} \right\|_2 \\ &\leq \max_{1 \leq l \leq m} \frac{\left( a_l^\top x^* \right)^2 \|a_l^\top \tilde{x}^{0,(l)}\| \|a_l\|_2}{m} \end{aligned}$$

$$\begin{aligned} &\stackrel{(i)}{\lesssim} \frac{\log n \cdot \sqrt{\log n} \cdot \sqrt{n}}{m} \\ &\asymp \sqrt{\frac{\log n}{n}} \cdot \frac{n \log n}{m}. \end{aligned} \tag{112}$$

Inequality (i) makes use of the fact  $\max_l |\mathbf{a}_l^\top \mathbf{x}^*| \leq 5\sqrt{\log n}$  (cf. (99)), the bound  $\max_l \|\mathbf{a}_l\|_2 \leq 6\sqrt{n}$  (cf. (98)), and  $\max_l |\mathbf{a}_l^\top \tilde{\mathbf{x}}^{0,(l)}| \leq 5\sqrt{\log n}$  (due to statistical independence and standard Gaussian concentration). As long as  $m/(n \log n)$  is sufficiently large, substituting the above bound (112) into (111) leads us to conclude that

$$\max_{1 \leq l \leq m} \|\mathbf{x}^0 - \mathbf{x}^{0,(l)}\|_2 \leq C_3 \sqrt{\frac{\log n}{n}} \tag{113}$$

for any constant  $C_3 > 0$ .

### B Proofs for Matrix Completion

Before proceeding to the proofs, let us record an immediate consequence of the incoherence property (25):

$$\|\mathbf{X}^*\|_{2,\infty} \leq \sqrt{\frac{\kappa\mu}{n}} \|\mathbf{X}^*\|_F \leq \sqrt{\frac{\kappa\mu r}{n}} \|\mathbf{X}^*\|. \tag{114}$$

where  $\kappa = \sigma_{\max}/\sigma_{\min}$  is the condition number of  $\mathbf{M}^*$ . This follows since

$$\begin{aligned} \|\mathbf{X}^*\|_{2,\infty} &= \|\mathbf{U}^*(\boldsymbol{\Sigma}^*)^{1/2}\|_{2,\infty} \leq \|\mathbf{U}^*\|_{2,\infty} \|(\boldsymbol{\Sigma}^*)^{1/2}\| \\ &\leq \sqrt{\frac{\mu}{n}} \|\mathbf{U}^*\|_F \|(\boldsymbol{\Sigma}^*)^{1/2}\| \leq \sqrt{\frac{\mu}{n}} \|\mathbf{U}^*\|_F \sqrt{\kappa\sigma_{\min}} \\ &\leq \sqrt{\frac{\kappa\mu}{n}} \|\mathbf{X}^*\|_F \leq \sqrt{\frac{\kappa\mu r}{n}} \|\mathbf{X}^*\|. \end{aligned}$$

Unless otherwise specified, we use the indicator variable  $\delta_{j,k}$  to denote whether the entry in the location  $(j, k)$  is included in  $\Omega$ . Under our model,  $\delta_{j,k}$  is a Bernoulli random variable with mean  $p$ .

#### B.1 Proof of Lemma 7

By the expression of the Hessian in (61), one can decompose

$$\begin{aligned} \text{vec}(\mathbf{V})^\top \nabla^2 f_{\text{clean}}(\mathbf{X}) \text{vec}(\mathbf{V}) &= \frac{1}{2p} \left\| \mathcal{P}_\Omega \left( \mathbf{V} \mathbf{X}^\top + \mathbf{X} \mathbf{V}^\top \right) \right\|_F^2 \\ &\quad + \frac{1}{p} \left\langle \mathcal{P}_\Omega \left( \mathbf{X} \mathbf{X}^\top - \mathbf{M}^* \right), \mathbf{V} \mathbf{V}^\top \right\rangle \end{aligned}$$

$$\begin{aligned}
 &= \underbrace{\frac{1}{2p} \left\| \mathcal{P}_\Omega \left( \mathbf{V} \mathbf{X}^\top + \mathbf{X} \mathbf{V}^\top \right) \right\|_F^2 - \frac{1}{2p} \left\| \mathcal{P}_\Omega \left( \mathbf{V} \mathbf{X}^{\star\top} + \mathbf{X}^* \mathbf{V}^\top \right) \right\|_F^2}_{:=\alpha_1} \\
 &\quad + \underbrace{\frac{1}{p} \left\langle \mathcal{P}_\Omega \left( \mathbf{X} \mathbf{X}^\top - \mathbf{M}^* \right), \mathbf{V} \mathbf{V}^\top \right\rangle}_{:=\alpha_2} \\
 &\quad + \underbrace{\frac{1}{2p} \left\| \mathcal{P}_\Omega \left( \mathbf{V} \mathbf{X}^{\star\top} + \mathbf{X}^* \mathbf{V}^\top \right) \right\|_F^2 - \frac{1}{2} \left\| \mathbf{V} \mathbf{X}^{\star\top} + \mathbf{X}^* \mathbf{V}^\top \right\|_F^2}_{:=\alpha_3} \\
 &\quad + \underbrace{\frac{1}{2} \left\| \mathbf{V} \mathbf{X}^{\star\top} + \mathbf{X}^* \mathbf{V}^\top \right\|_F^2}_{:=\alpha_4}.
 \end{aligned}$$

The basic idea is to demonstrate that: (1)  $\alpha_4$  is bounded both from above and from below, and (2) the first three terms are sufficiently small in size compared to  $\alpha_4$ .

1. We start by controlling  $\alpha_4$ . It is immediate to derive the following upper bound

$$\alpha_4 \leq \left\| \mathbf{V} \mathbf{X}^{\star\top} \right\|_F^2 + \left\| \mathbf{X}^* \mathbf{V}^\top \right\|_F^2 \leq 2 \left\| \mathbf{X}^* \right\|^2 \left\| \mathbf{V} \right\|_F^2 = 2 \sigma_{\max} \left\| \mathbf{V} \right\|_F^2.$$

When it comes to the lower bound, one discovers that

$$\begin{aligned}
 \alpha_4 &= \frac{1}{2} \left\{ \left\| \mathbf{V} \mathbf{X}^{\star\top} \right\|_F^2 + \left\| \mathbf{X}^* \mathbf{V}^\top \right\|_F^2 + 2 \text{Tr} \left( \mathbf{X}^{\star\top} \mathbf{V} \mathbf{X}^{\star\top} \mathbf{V} \right) \right\} \\
 &\geq \sigma_{\min} \left\| \mathbf{V} \right\|_F^2 + \text{Tr} \left[ \left( \mathbf{Z} + \mathbf{X}^* - \mathbf{Z} \right)^\top \mathbf{V} \left( \mathbf{Z} + \mathbf{X}^* - \mathbf{Z} \right)^\top \mathbf{V} \right] \\
 &\geq \sigma_{\min} \left\| \mathbf{V} \right\|_F^2 + \text{Tr} \left( \mathbf{Z}^\top \mathbf{V} \mathbf{Z}^\top \mathbf{V} \right) - 2 \left\| \mathbf{Z} - \mathbf{X}^* \right\| \left\| \mathbf{Z} \right\| \left\| \mathbf{V} \right\|_F^2 - \left\| \mathbf{Z} - \mathbf{X}^* \right\|^2 \left\| \mathbf{V} \right\|_F^2 \\
 &\geq \left( \sigma_{\min} - 5 \delta \sigma_{\max} \right) \left\| \mathbf{V} \right\|_F^2 + \text{Tr} \left( \mathbf{Z}^\top \mathbf{V} \mathbf{Z}^\top \mathbf{V} \right), \tag{115}
 \end{aligned}$$

where the last line comes from the assumptions that

$$\left\| \mathbf{Z} - \mathbf{X}^* \right\| \leq \delta \left\| \mathbf{X}^* \right\| \leq \left\| \mathbf{X}^* \right\| \quad \text{and} \quad \left\| \mathbf{Z} \right\| \leq \left\| \mathbf{Z} - \mathbf{X}^* \right\| + \left\| \mathbf{X}^* \right\| \leq 2 \left\| \mathbf{X}^* \right\|.$$

With our assumption  $\mathbf{V} = \mathbf{Y} \mathbf{H}_Y - \mathbf{Z}$  in mind, it comes down to controlling

$$\text{Tr} \left( \mathbf{Z}^\top \mathbf{V} \mathbf{Z}^\top \mathbf{V} \right) = \text{Tr} \left[ \mathbf{Z}^\top \left( \mathbf{Y} \mathbf{H}_Y - \mathbf{Z} \right) \mathbf{Z}^\top \left( \mathbf{Y} \mathbf{H}_Y - \mathbf{Z} \right) \right].$$

From the definition of  $\mathbf{H}_Y$ , we see from Lemma 35 that  $\mathbf{Z}^\top \mathbf{Y} \mathbf{H}_Y$  (and hence  $\mathbf{Z}^\top \left( \mathbf{Y} \mathbf{H}_Y - \mathbf{Z} \right)$ ) is a symmetric matrix, which implies that

$$\text{Tr} \left[ \mathbf{Z}^\top \left( \mathbf{Y} \mathbf{H}_Y - \mathbf{Z} \right) \mathbf{Z}^\top \left( \mathbf{Y} \mathbf{H}_Y - \mathbf{Z} \right) \right] \geq 0.$$

Substitution into (115) gives

$$\alpha_4 \geq (\sigma_{\min} - 5\delta\sigma_{\max}) \|\mathbf{V}\|_F^2 \geq \frac{9}{10}\sigma_{\min} \|\mathbf{V}\|_F^2,$$

provided that  $\kappa\delta \leq 1/50$ .

2. For  $\alpha_1$ , we consider the following quantity

$$\begin{aligned} & \|\mathcal{P}_\Omega(\mathbf{V}\mathbf{X}^\top + \mathbf{X}\mathbf{V}^\top)\|_F^2 \\ &= \langle \mathcal{P}_\Omega(\mathbf{V}\mathbf{X}^\top), \mathcal{P}_\Omega(\mathbf{V}\mathbf{X}^\top) \rangle + \langle \mathcal{P}_\Omega(\mathbf{V}\mathbf{X}^\top), \mathcal{P}_\Omega(\mathbf{X}\mathbf{V}^\top) \rangle \\ & \quad + \langle \mathcal{P}_\Omega(\mathbf{X}\mathbf{V}^\top), \mathcal{P}_\Omega(\mathbf{V}\mathbf{X}^\top) \rangle + \langle \mathcal{P}_\Omega(\mathbf{X}\mathbf{V}^\top), \mathcal{P}_\Omega(\mathbf{X}\mathbf{V}^\top) \rangle \\ &= 2\langle \mathcal{P}_\Omega(\mathbf{V}\mathbf{X}^\top), \mathcal{P}_\Omega(\mathbf{V}\mathbf{X}^\top) \rangle + 2\langle \mathcal{P}_\Omega(\mathbf{V}\mathbf{X}^\top), \mathcal{P}_\Omega(\mathbf{X}\mathbf{V}^\top) \rangle. \end{aligned}$$

Similar decomposition can be performed on  $\|\mathcal{P}_\Omega(\mathbf{V}\mathbf{X}^{*\top} + \mathbf{X}^*\mathbf{V}^\top)\|_F^2$  as well. These identities yield

$$\begin{aligned} \alpha_1 &= \frac{1}{p} \underbrace{\left[ \langle \mathcal{P}_\Omega(\mathbf{V}\mathbf{X}^\top), \mathcal{P}_\Omega(\mathbf{V}\mathbf{X}^\top) \rangle - \langle \mathcal{P}_\Omega(\mathbf{V}\mathbf{X}^{*\top}), \mathcal{P}_\Omega(\mathbf{V}\mathbf{X}^{*\top}) \rangle \right]}_{:=\beta_1} \\ & \quad + \frac{1}{p} \underbrace{\left[ \langle \mathcal{P}_\Omega(\mathbf{V}\mathbf{X}^\top), \mathcal{P}_\Omega(\mathbf{X}\mathbf{V}^\top) \rangle - \langle \mathcal{P}_\Omega(\mathbf{V}\mathbf{X}^{*\top}), \mathcal{P}_\Omega(\mathbf{X}^*\mathbf{V}^\top) \rangle \right]}_{:=\beta_2}. \end{aligned}$$

For  $\beta_2$ , one has

$$\begin{aligned} \beta_2 &= \frac{1}{p} \left\langle \mathcal{P}_\Omega(\mathbf{V}(\mathbf{X} - \mathbf{X}^*)^\top), \mathcal{P}_\Omega((\mathbf{X} - \mathbf{X}^*)\mathbf{V}^\top) \right\rangle \\ & \quad + \frac{1}{p} \left\langle \mathcal{P}_\Omega(\mathbf{V}(\mathbf{X} - \mathbf{X}^*)^\top), \mathcal{P}_\Omega(\mathbf{X}^*\mathbf{V}^\top) \right\rangle \\ & \quad + \frac{1}{p} \left\langle \mathcal{P}_\Omega(\mathbf{V}\mathbf{X}^{*\top}), \mathcal{P}_\Omega((\mathbf{X} - \mathbf{X}^*)\mathbf{V}^\top) \right\rangle \end{aligned}$$

which together with the inequality  $|\langle \mathbf{A}, \mathbf{B} \rangle| \leq \|\mathbf{A}\|_F \|\mathbf{B}\|_F$  gives

$$|\beta_2| \leq \frac{1}{p} \|\mathcal{P}_\Omega(\mathbf{V}(\mathbf{X} - \mathbf{X}^*)^\top)\|_F^2 + \frac{2}{p} \|\mathcal{P}_\Omega(\mathbf{V}(\mathbf{X} - \mathbf{X}^*)^\top)\|_F \|\mathcal{P}_\Omega(\mathbf{X}^*\mathbf{V}^\top)\|_F. \tag{116}$$

This then calls for upper bounds on the following two terms

$$\frac{1}{\sqrt{p}} \|\mathcal{P}_\Omega(\mathbf{V}(\mathbf{X} - \mathbf{X}^*)^\top)\|_F \quad \text{and} \quad \frac{1}{\sqrt{p}} \|\mathcal{P}_\Omega(\mathbf{X}^*\mathbf{V}^\top)\|_F.$$

The injectivity of  $\mathcal{P}_\Omega$  (cf. [19, Section 4.2] or Lemma 38)—when restricted to the tangent space of  $\mathbf{M}^*$ —gives: for any fixed constant  $\gamma > 0$ ,

$$\frac{1}{\sqrt{p}} \left\| \mathcal{P}_\Omega \left( \mathbf{X}^* \mathbf{V}^\top \right) \right\|_F \leq (1 + \gamma) \left\| \mathbf{X}^* \mathbf{V}^\top \right\|_F \leq (1 + \gamma) \left\| \mathbf{X}^* \right\| \left\| \mathbf{V} \right\|_F$$

with probability at least  $1 - O(n^{-10})$ , provided that  $n^2 p / (\mu n r \log n)$  is sufficiently large. In addition,

$$\begin{aligned} & \frac{1}{p} \left\| \mathcal{P}_\Omega \left( \mathbf{V} \left( \mathbf{X} - \mathbf{X}^* \right)^\top \right) \right\|_F^2 \\ &= \frac{1}{p} \sum_{1 \leq j, k \leq n} \delta_{j, k} \left[ \mathbf{V}_{j, \cdot} \left( \mathbf{X}_{k, \cdot} - \mathbf{X}_{k, \cdot}^* \right)^\top \right]^2 \\ &= \sum_{1 \leq j \leq n} \mathbf{V}_{j, \cdot} \left[ \frac{1}{p} \sum_{1 \leq k \leq n} \delta_{j, k} \left( \mathbf{X}_{k, \cdot} - \mathbf{X}_{k, \cdot}^* \right)^\top \left( \mathbf{X}_{k, \cdot} - \mathbf{X}_{k, \cdot}^* \right) \right] \mathbf{V}_{j, \cdot}^\top \\ &\leq \max_{1 \leq j \leq n} \left\| \frac{1}{p} \sum_{1 \leq k \leq n} \delta_{j, k} \left( \mathbf{X}_{k, \cdot} - \mathbf{X}_{k, \cdot}^* \right)^\top \left( \mathbf{X}_{k, \cdot} - \mathbf{X}_{k, \cdot}^* \right) \right\| \left\| \mathbf{V} \right\|_F^2 \\ &\leq \left\{ \frac{1}{p} \max_{1 \leq j \leq n} \sum_{1 \leq k \leq n} \delta_{j, k} \right\} \left\{ \max_{1 \leq k \leq n} \left\| \mathbf{X}_{k, \cdot} - \mathbf{X}_{k, \cdot}^* \right\|_2^2 \right\} \left\| \mathbf{V} \right\|_F^2 \\ &\leq (1 + \gamma) n \left\| \mathbf{X} - \mathbf{X}^* \right\|_{2, \infty}^2 \left\| \mathbf{V} \right\|_F^2, \end{aligned}$$

with probability exceeding  $1 - O(n^{-10})$ , which holds as long as  $np / \log n$  is sufficiently large. Taken collectively, the above bounds yield that for any small constant  $\gamma > 0$ ,

$$\begin{aligned} |\beta_2| &\leq (1 + \gamma) n \left\| \mathbf{X} - \mathbf{X}^* \right\|_{2, \infty}^2 \left\| \mathbf{V} \right\|_F^2 \\ &\quad + 2 \sqrt{(1 + \gamma) n \left\| \mathbf{X} - \mathbf{X}^* \right\|_{2, \infty}^2 \left\| \mathbf{V} \right\|_F^2 \cdot (1 + \gamma)^2 \left\| \mathbf{X}^* \right\|^2 \left\| \mathbf{V} \right\|_F^2} \\ &\lesssim \left( \epsilon^2 n \left\| \mathbf{X}^* \right\|_{2, \infty}^2 + \epsilon \sqrt{n} \left\| \mathbf{X}^* \right\|_{2, \infty} \left\| \mathbf{X}^* \right\| \right) \left\| \mathbf{V} \right\|_F^2, \end{aligned}$$

where the last inequality makes use of the assumption  $\left\| \mathbf{X} - \mathbf{X}^* \right\|_{2, \infty} \leq \epsilon \left\| \mathbf{X}^* \right\|_{2, \infty}$ . The same analysis can be repeated to control  $\beta_1$ . Altogether, we obtain

$$\begin{aligned} |\alpha_1| \leq |\beta_1| + |\beta_2| &\lesssim \left( n \epsilon^2 \left\| \mathbf{X}^* \right\|_{2, \infty}^2 + \sqrt{n} \epsilon \left\| \mathbf{X}^* \right\|_{2, \infty} \left\| \mathbf{X}^* \right\| \right) \left\| \mathbf{V} \right\|_F^2 \\ &\stackrel{(i)}{\leq} \left( n \epsilon^2 \frac{\kappa \mu r}{n} + \sqrt{n} \epsilon \sqrt{\frac{\kappa \mu r}{n}} \right) \sigma_{\max} \left\| \mathbf{V} \right\|_F^2 \stackrel{(ii)}{\leq} \frac{1}{10} \sigma_{\min} \left\| \mathbf{V} \right\|_F^2, \end{aligned}$$

where (i) utilizes the incoherence condition (114) and (ii) holds with the proviso that  $\epsilon \sqrt{\kappa^3 \mu r} \ll 1$ .

3. To bound  $\alpha_2$ , apply the Cauchy–Schwarz inequality to get

$$|\alpha_2| = \left| \left\langle \mathbf{V}, \frac{1}{p} \mathcal{P}_\Omega \left( \mathbf{X} \mathbf{X}^\top - \mathbf{M}^\star \right) \mathbf{V} \right\rangle \right| \leq \left\| \frac{1}{p} \mathcal{P}_\Omega \left( \mathbf{X} \mathbf{X}^\top - \mathbf{M}^\star \right) \right\| \|\mathbf{V}\|_{\mathbb{F}}^2.$$

In view of Lemma 43, with probability at least  $1 - O(n^{-10})$ ,

$$\begin{aligned} \left\| \frac{1}{p} \mathcal{P}_\Omega \left( \mathbf{X} \mathbf{X}^\top - \mathbf{M}^\star \right) \right\| &\leq 2n\epsilon^2 \|\mathbf{X}^\star\|_{2,\infty}^2 + 4\epsilon\sqrt{n} \log n \|\mathbf{X}^\star\|_{2,\infty} \|\mathbf{X}^\star\| \\ &\leq \left( 2n\epsilon^2 \frac{\kappa\mu r}{n} + 4\epsilon\sqrt{n} \log n \sqrt{\frac{\kappa\mu r}{n}} \right) \sigma_{\max} \leq \frac{1}{10} \sigma_{\min} \end{aligned}$$

as soon as  $\epsilon\sqrt{\kappa^3\mu r} \log n \ll 1$ , where we utilize the incoherence condition (114). This in turn implies that

$$|\alpha_2| \leq \frac{1}{10} \sigma_{\min} \|\mathbf{V}\|_{\mathbb{F}}^2.$$

Notably, this bound holds uniformly over all  $\mathbf{X}$  satisfying the condition in Lemma 7, regardless of the statistical dependence between  $\mathbf{X}$  and the sampling set  $\Omega$ .

4. The last term  $\alpha_3$  can also be controlled using the injectivity of  $\mathcal{P}_\Omega$  when restricted to the tangent space of  $\mathbf{M}^\star$ . Specifically, it follows from the bounds in [19, Section 4.2] or Lemma 38 that

$$|\alpha_3| \leq \gamma \left\| \mathbf{V} \mathbf{X}^{\star\top} + \mathbf{X}^\star \mathbf{V}^\top \right\|_{\mathbb{F}}^2 \leq 4\gamma \sigma_{\max} \|\mathbf{V}\|_{\mathbb{F}}^2 \leq \frac{1}{10} \sigma_{\min} \|\mathbf{V}\|_{\mathbb{F}}^2$$

for any  $\gamma > 0$  such that  $\kappa\gamma$  is a small constant, as soon as  $n^2 p \gg \kappa^2 \mu r n \log n$ .

5. Taking all the preceding bounds collectively yields

$$\begin{aligned} \text{vec}(\mathbf{V})^\top \nabla^2 f_{\text{clean}}(\mathbf{X}) \text{vec}(\mathbf{V}) &\geq \alpha_4 - |\alpha_1| - |\alpha_2| - |\alpha_3| \\ &\geq \left( \frac{9}{10} - \frac{3}{10} \right) \sigma_{\min} \|\mathbf{V}\|_{\mathbb{F}}^2 \geq \frac{1}{2} \sigma_{\min} \|\mathbf{V}\|_{\mathbb{F}}^2 \end{aligned}$$

for all  $\mathbf{V}$  satisfying our assumptions, and

$$\begin{aligned} \left| \text{vec}(\mathbf{V})^\top \nabla^2 f_{\text{clean}}(\mathbf{X}) \text{vec}(\mathbf{V}) \right| &\leq \alpha_4 + |\alpha_1| + |\alpha_2| + |\alpha_3| \\ &\leq \left( 2\sigma_{\max} + \frac{3}{10} \sigma_{\min} \right) \|\mathbf{V}\|_{\mathbb{F}}^2 \leq \frac{5}{2} \sigma_{\max} \|\mathbf{V}\|_{\mathbb{F}}^2 \end{aligned}$$

for all  $\mathbf{V}$ . Since this upper bound holds uniformly over all  $\mathbf{V}$ , we conclude that

$$\left\| \nabla^2 f_{\text{clean}}(\mathbf{X}) \right\| \leq \frac{5}{2} \sigma_{\max}$$

as claimed.



**B.2 Proof of Lemma 8**

Given that  $\widehat{\mathbf{H}}^{t+1}$  is chosen to minimize the error in terms of the Frobenius norm (cf. (26)), we have

$$\begin{aligned} \left\| \mathbf{X}^{t+1} \widehat{\mathbf{H}}^{t+1} - \mathbf{X}^* \right\|_{\mathbf{F}} &\leq \left\| \mathbf{X}^{t+1} \widehat{\mathbf{H}}^t - \mathbf{X}^* \right\|_{\mathbf{F}} = \left\| [\mathbf{X}^t - \eta \nabla f(\mathbf{X}^t)] \widehat{\mathbf{H}}^t - \mathbf{X}^* \right\|_{\mathbf{F}} \\ &\stackrel{(i)}{=} \left\| \mathbf{X}^t \widehat{\mathbf{H}}^t - \eta \nabla f(\mathbf{X}^t \widehat{\mathbf{H}}^t) - \mathbf{X}^* \right\|_{\mathbf{F}} \\ &\stackrel{(ii)}{=} \left\| \mathbf{X}^t \widehat{\mathbf{H}}^t - \eta \left[ \nabla f_{\text{clean}}(\mathbf{X}^t \widehat{\mathbf{H}}^t) - \frac{1}{p} \mathcal{P}_{\Omega}(\mathbf{E}) \mathbf{X}^t \widehat{\mathbf{H}}^t \right] - \mathbf{X}^* \right\|_{\mathbf{F}} \\ &\leq \underbrace{\left\| \mathbf{X}^t \widehat{\mathbf{H}}^t - \eta \nabla f_{\text{clean}}(\mathbf{X}^t \widehat{\mathbf{H}}^t) - (\mathbf{X}^* - \eta \nabla f_{\text{clean}}(\mathbf{X}^*)) \right\|_{\mathbf{F}}}_{:=\alpha_1} + \eta \underbrace{\left\| \frac{1}{p} \mathcal{P}_{\Omega}(\mathbf{E}) \mathbf{X}^t \widehat{\mathbf{H}}^t \right\|_{\mathbf{F}}}_{:=\alpha_2}, \end{aligned} \tag{117}$$

where (i) follows from the identity  $\nabla f(\mathbf{X}^t \mathbf{R}) = \nabla f(\mathbf{X}^t) \mathbf{R}$  for any orthonormal matrix  $\mathbf{R} \in \mathcal{O}^{r \times r}$ , (ii) arises from the definitions of  $\nabla f(\mathbf{X})$  and  $\nabla f_{\text{clean}}(\mathbf{X})$  (see (59) and (60), respectively), and the last inequality (117) utilizes the triangle inequality and the fact that  $\nabla f_{\text{clean}}(\mathbf{X}^*) = \mathbf{0}$ . It thus suffices to control  $\alpha_1$  and  $\alpha_2$ .

1. For the second term  $\alpha_2$  in (117), it is easy to see that with probability at least  $1 - O(n^{-10})$ ,

$$\alpha_2 \leq \eta \left\| \frac{1}{p} \mathcal{P}_{\Omega}(\mathbf{E}) \right\| \left\| \mathbf{X}^t \widehat{\mathbf{H}}^t \right\|_{\mathbf{F}} \leq 2\eta \left\| \frac{1}{p} \mathcal{P}_{\Omega}(\mathbf{E}) \right\| \left\| \mathbf{X}^* \right\|_{\mathbf{F}} \leq 2\eta C \sigma \sqrt{\frac{n}{p}} \left\| \mathbf{X}^* \right\|_{\mathbf{F}}$$

for some absolute constant  $C > 0$ . Here, the second inequality holds because  $\left\| \mathbf{X}^t \widehat{\mathbf{H}}^t \right\|_{\mathbf{F}} \leq \left\| \mathbf{X}^t \widehat{\mathbf{H}}^t - \mathbf{X}^* \right\|_{\mathbf{F}} + \left\| \mathbf{X}^* \right\|_{\mathbf{F}} \leq 2 \left\| \mathbf{X}^* \right\|_{\mathbf{F}}$ , following hypothesis (28a) together with our assumptions on the noise and the sample complexity. The last inequality makes use of Lemma 40.

2. For the first term  $\alpha_1$  in (117), the fundamental theorem of calculus [70, Chapter XIII, Theorem 4.2] reveals

$$\begin{aligned} &\text{vec} \left[ \mathbf{X}^t \widehat{\mathbf{H}}^t - \eta \nabla f_{\text{clean}}(\mathbf{X}^t \widehat{\mathbf{H}}^t) - (\mathbf{X}^* - \eta \nabla f_{\text{clean}}(\mathbf{X}^*)) \right] \\ &= \text{vec} \left[ \mathbf{X}^t \widehat{\mathbf{H}}^t - \mathbf{X}^* \right] - \eta \cdot \text{vec} \left[ \nabla f_{\text{clean}}(\mathbf{X}^t \widehat{\mathbf{H}}^t) - \nabla f_{\text{clean}}(\mathbf{X}^*) \right] \\ &= \left( \mathbf{I}_{nr} - \eta \underbrace{\int_0^1 \nabla^2 f_{\text{clean}}(\mathbf{X}(\tau)) \, d\tau}_{:=\mathbf{A}} \right) \text{vec} \left( \mathbf{X}^t \widehat{\mathbf{H}}^t - \mathbf{X}^* \right), \end{aligned} \tag{118}$$

where we denote  $\mathbf{X}(\tau) := \mathbf{X}^* + \tau(\mathbf{X}^t \widehat{\mathbf{H}}^t - \mathbf{X}^*)$ . Taking the squared Euclidean norm of both sides of equality (118) leads to

$$\begin{aligned} (\alpha_1)^2 &= \text{vec}(\mathbf{X}^t \widehat{\mathbf{H}}^t - \mathbf{X}^*)^\top (\mathbf{I}_{nr} - \eta \mathbf{A})^2 \text{vec}(\mathbf{X}^t \widehat{\mathbf{H}}^t - \mathbf{X}^*) \\ &= \text{vec}(\mathbf{X}^t \widehat{\mathbf{H}}^t - \mathbf{X}^*)^\top \left( \mathbf{I}_{nr} - 2\eta \mathbf{A} + \eta^2 \mathbf{A}^2 \right) \text{vec}(\mathbf{X}^t \widehat{\mathbf{H}}^t - \mathbf{X}^*) \end{aligned}$$

$$\begin{aligned} &\leq \left\| X^t \widehat{H}^t - X^* \right\|_F^2 + \eta^2 \|A\|^2 \left\| X^t \widehat{H}^t - X^* \right\|_F^2 \\ &\quad - 2\eta \operatorname{vec}(X^t \widehat{H}^t - X^*)^\top A \operatorname{vec}(X^t \widehat{H}^t - X^*), \end{aligned} \tag{119}$$

where in (119) we have used the fact that

$$\begin{aligned} \operatorname{vec}(X^t \widehat{H}^t - X^*)^\top A^2 \operatorname{vec}(X^t \widehat{H}^t - X^*) &\leq \|A\|^2 \left\| \operatorname{vec}(X^t \widehat{H}^t - X^*) \right\|_2^2 \\ &= \|A\|^2 \left\| X^t \widehat{H}^t - X^* \right\|_F^2. \end{aligned}$$

Based on condition (28b), it is easily seen that  $\forall \tau \in [0, 1]$ ,

$$\|X(\tau) - X^*\|_{2,\infty} \leq \left( C_5 \mu r \sqrt{\frac{\log n}{np}} + \frac{C_8}{\sigma_{\min}} \sigma \sqrt{\frac{n \log n}{p}} \right) \|X^*\|_{2,\infty}.$$

Taking  $X = X(\tau)$ ,  $Y = X^t$ , and  $Z = X^*$  in Lemma 7, one can easily verify the assumptions therein given our sample size condition  $n^2 p \gg \kappa^3 \mu^3 r^3 n \log^3 n$  and the noise condition (27). As a result,

$$\begin{aligned} \operatorname{vec}(X^t \widehat{H}^t - X^*)^\top A \operatorname{vec}(X^t \widehat{H}^t - X^*) &\geq \frac{\sigma_{\min}}{2} \|X^t \widehat{H}^t - X^*\|_F^2 \\ \text{and } \|A\| &\leq \frac{5}{2} \sigma_{\max}. \end{aligned}$$

Substituting these two inequalities into (119) yields

$$\begin{aligned} (\alpha_1)^2 &\leq \left( 1 + \frac{25}{4} \eta^2 \sigma_{\max}^2 - \sigma_{\min} \eta \right) \|X^t \widehat{H}^t - X^*\|_F^2 \\ &\leq \left( 1 - \frac{\sigma_{\min}}{2} \eta \right) \|X^t \widehat{H}^t - X^*\|_F^2 \end{aligned}$$

as long as  $0 < \eta \leq (2\sigma_{\min}) / (25\sigma_{\max}^2)$ , which further implies that

$$\alpha_1 \leq \left( 1 - \frac{\sigma_{\min}}{4} \eta \right) \|X^t \widehat{H}^t - X^*\|_F.$$

3. Combining the preceding bounds on both  $\alpha_1$  and  $\alpha_2$  and making use of hypothesis (28a), we have

$$\begin{aligned} \|X^{t+1} \widehat{H}^{t+1} - X^*\|_F &\leq \left( 1 - \frac{\sigma_{\min}}{4} \eta \right) \|X^t \widehat{H}^t - X^*\|_F + 2\eta C \sigma \sqrt{\frac{n}{p}} \|X^*\|_F \\ &\leq \left( 1 - \frac{\sigma_{\min}}{4} \eta \right) \left( C_4 \rho^t \mu r \frac{1}{\sqrt{np}} \|X^*\|_F + C_1 \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \|X^*\|_F \right) + 2\eta C \sigma \sqrt{\frac{n}{p}} \|X^*\|_F \\ &\leq \left( 1 - \frac{\sigma_{\min}}{4} \eta \right) C_4 \rho^t \mu r \frac{1}{\sqrt{np}} \|X^*\|_F + \left[ \left( 1 - \frac{\sigma_{\min}}{4} \eta \right) \frac{C_1}{\sigma_{\min}} + 2\eta C \right] \sigma \sqrt{\frac{n}{p}} \|X^*\|_F \\ &\leq C_4 \rho^{t+1} \mu r \frac{1}{\sqrt{np}} \|X^*\|_F + C_1 \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \|X^*\|_F \end{aligned}$$

as long as  $0 < \eta \leq (2\sigma_{\min})/(25\sigma_{\max}^2)$ ,  $1 - (\sigma_{\min}/4) \cdot \eta \leq \rho < 1$ , and  $C_1$  is sufficiently large. This completes the proof of the contraction with respect to the Frobenius norm.

### B.3 Proof of Lemma 9

To facilitate analysis, we construct an auxiliary matrix defined as follows

$$\tilde{X}^{t+1} := X^t \hat{H}^t - \eta \frac{1}{p} \mathcal{P}_\Omega \left[ X^t X^{t\top} - (M^* + E) \right] X^*. \tag{120}$$

With this auxiliary matrix in place, we invoke the triangle inequality to bound

$$\|X^{t+1} \hat{H}^{t+1} - X^*\| \leq \underbrace{\|X^{t+1} \hat{H}^{t+1} - \tilde{X}^{t+1}\|}_{:=\alpha_1} + \underbrace{\|\tilde{X}^{t+1} - X^*\|}_{:=\alpha_2}. \tag{121}$$

1. We start with the second term  $\alpha_2$  and show that the auxiliary matrix  $\tilde{X}^{t+1}$  is also not far from the truth. The definition of  $\tilde{X}^{t+1}$  allows one to express

$$\begin{aligned} \alpha_2 &= \left\| X^t \hat{H}^t - \eta \frac{1}{p} \mathcal{P}_\Omega \left[ X^t X^{t\top} - (M^* + E) \right] X^* - X^* \right\| \\ &\leq \eta \left\| \frac{1}{p} \mathcal{P}_\Omega (E) \right\| \|X^*\| + \left\| X^t \hat{H}^t - \eta \frac{1}{p} \mathcal{P}_\Omega \left( X^t X^{t\top} - X^* X^{*\top} \right) X^* - X^* \right\| \\ &\leq \eta \left\| \frac{1}{p} \mathcal{P}_\Omega (E) \right\| \|X^*\| + \underbrace{\left\| X^t \hat{H}^t - \eta \left( X^t X^{t\top} - X^* X^{*\top} \right) X^* - X^* \right\|}_{:=\beta_1} \\ &\quad + \underbrace{\eta \left\| \frac{1}{p} \mathcal{P}_\Omega \left( X^t X^{t\top} - X^* X^{*\top} \right) X^* - \left( X^t X^{t\top} - X^* X^{*\top} \right) X^* \right\|}_{:=\beta_2}, \end{aligned} \tag{122}$$

where we have used the triangle inequality to separate the population-level component (i.e.,  $\beta_1$ ), the perturbation (i.e.,  $\beta_2$ ), and the noise component. In what follows, we will denote

$$\Delta^t := X^t \hat{H}^t - X^*$$

which, by Lemma 35, satisfies the following symmetry property

$$\hat{H}^{t\top} X^{t\top} X^* = X^{*\top} X^t \hat{H}^t \implies \Delta^{t\top} X^* = X^{*\top} \Delta^t. \tag{124}$$

- (a) The population-level component  $\beta_1$  is easier to control. Specifically, we first simplify its expression as

$$\begin{aligned} \beta_1 &= \left\| \Delta^t - \eta \left( \Delta^t \Delta^{t\top} + \Delta^t X^{*\top} + X^* \Delta^{t\top} \right) X^* \right\| \\ &\leq \underbrace{\left\| \Delta^t - \eta \left( \Delta^t X^{*\top} + X^* \Delta^{t\top} \right) X^* \right\|}_{:=\gamma_1} + \eta \underbrace{\left\| \Delta^t \Delta^{t\top} X^* \right\|}_{:=\gamma_2}. \end{aligned}$$

The leading term  $\gamma_1$  can be upper bounded by

$$\begin{aligned} \gamma_1 &= \left\| \Delta^t - \eta \Delta^t \Sigma^* - \eta X^* \Delta^{t\top} X^* \right\| = \left\| \Delta^t - \eta \Delta^t \Sigma^* - \eta X^* X^{*\top} \Delta^t \right\| \\ &= \left\| \frac{1}{2} \Delta^t (I_r - 2\eta \Sigma^*) + \frac{1}{2} (I_n - 2\eta M^*) \Delta^t \right\| \\ &\leq \frac{1}{2} (\|I_r - 2\eta \Sigma^*\| + \|I_n - 2\eta M^*\|) \|\Delta^t\| \end{aligned}$$

where the second identity follows from the symmetry property (124). By choosing  $\eta \leq 1/(2\sigma_{\max})$ , one has  $\mathbf{0} \leq I_r - 2\eta \Sigma^* \leq (1 - 2\eta\sigma_{\min}) I_r$  and  $\mathbf{0} \leq I_n - 2\eta M^* \leq I_n$ , and further one can ensure

$$\gamma_1 \leq \frac{1}{2} [(1 - 2\eta\sigma_{\min}) + 1] \|\Delta^t\| = (1 - \eta\sigma_{\min}) \|\Delta^t\|. \tag{125}$$

Next, regarding the higher-order term  $\gamma_2$ , we can easily obtain

$$\gamma_2 \leq \eta \|\Delta^t\|^2 \|X^*\|. \tag{126}$$

Bounds (125) and (126) taken collectively give

$$\beta_1 \leq (1 - \eta\sigma_{\min}) \|\Delta^t\| + \eta \|\Delta^t\|^2 \|X^*\|. \tag{127}$$

- (b) We now turn to the perturbation part  $\beta_2$  by showing that

$$\begin{aligned} \frac{1}{\eta} \beta_2 &= \left\| \frac{1}{p} \mathcal{P}_\Omega \left( \Delta^t \Delta^{t\top} + \Delta^t X^{*\top} + X^* \Delta^{t\top} \right) X^* \right. \\ &\quad \left. - \left[ \Delta^t \Delta^{t\top} + \Delta^t X^{*\top} + X^* \Delta^{t\top} \right] X^* \right\| \\ &\leq \underbrace{\left\| \frac{1}{p} \mathcal{P}_\Omega \left( \Delta^t X^{*\top} \right) X^* - \left( \Delta^t X^{*\top} \right) X^* \right\|_F}_{:=\theta_1} \end{aligned}$$

$$\begin{aligned}
 &+ \underbrace{\left\| \frac{1}{p} \mathcal{P}_\Omega \left( \mathbf{X}^* \Delta^{t\top} \right) \mathbf{X}^* - \left( \mathbf{X}^* \Delta^{t\top} \right) \mathbf{X}^* \right\|_F}_{:=\theta_2} \\
 &+ \underbrace{\left\| \frac{1}{p} \mathcal{P}_\Omega \left( \Delta^t \Delta^{t\top} \right) \mathbf{X}^* - \left( \Delta^t \Delta^{t\top} \right) \mathbf{X}^* \right\|_F}_{:=\theta_3}, \tag{128}
 \end{aligned}$$

where the last inequality holds due to the triangle inequality as well as the fact that  $\|A\| \leq \|A\|_F$ . In the sequel, we shall bound the three terms separately.

- For the first term  $\theta_1$  in (128), the  $l$ th row of  $\frac{1}{p} \mathcal{P}_\Omega \left( \Delta^t \mathbf{X}^{*\top} \right) \mathbf{X}^* - \left( \Delta^t \mathbf{X}^{*\top} \right) \mathbf{X}^*$  is given by

$$\frac{1}{p} \sum_{j=1}^n (\delta_{l,j} - p) \Delta_{l,\cdot}^t \mathbf{X}_{j,\cdot}^{*\top} \mathbf{X}_{j,\cdot}^* = \Delta_{l,\cdot}^t \cdot \left[ \frac{1}{p} \sum_{j=1}^n (\delta_{l,j} - p) \mathbf{X}_{j,\cdot}^{*\top} \mathbf{X}_{j,\cdot}^* \right]$$

where, as usual,  $\delta_{l,j} = \mathbb{1}_{\{(l,j) \in \Omega\}}$ . Lemma 41 together with the union bound reveals that

$$\begin{aligned}
 &\left\| \frac{1}{p} \sum_{j=1}^n (\delta_{l,j} - p) \mathbf{X}_{j,\cdot}^{*\top} \mathbf{X}_{j,\cdot}^* \right\| \\
 &\lesssim \frac{1}{p} \left( \sqrt{p \|\mathbf{X}^*\|_{2,\infty}^2 \|\mathbf{X}^*\|^2 \log n} + \|\mathbf{X}^*\|_{2,\infty}^2 \log n \right) \\
 &\asymp \sqrt{\frac{\|\mathbf{X}^*\|_{2,\infty}^2 \sigma_{\max} \log n}{p}} + \frac{\|\mathbf{X}^*\|_{2,\infty}^2 \log n}{p}
 \end{aligned}$$

for all  $1 \leq l \leq n$  with high probability. This gives

$$\begin{aligned}
 &\left\| \Delta_{l,\cdot}^t \cdot \left[ \frac{1}{p} \sum_{j=1}^n (\delta_{l,j} - p) \mathbf{X}_{j,\cdot}^{*\top} \mathbf{X}_{j,\cdot}^* \right] \right\|_2 \\
 &\leq \|\Delta_{l,\cdot}^t\|_2 \left\| \frac{1}{p} \sum_j (\delta_{l,j} - p) \mathbf{X}_{j,\cdot}^{*\top} \mathbf{X}_{j,\cdot}^* \right\| \\
 &\lesssim \|\Delta_{l,\cdot}^t\|_2 \left\{ \sqrt{\frac{\|\mathbf{X}^*\|_{2,\infty}^2 \sigma_{\max} \log n}{p}} + \frac{\|\mathbf{X}^*\|_{2,\infty}^2 \log n}{p} \right\},
 \end{aligned}$$

which further reveals that

$$\begin{aligned} \theta_1 &= \sqrt{\sum_{l=1}^n \left\| \frac{1}{p} \sum_j (\delta_{l,j} - p) \Delta_{l,\cdot}^t \mathbf{X}_{j,\cdot}^{\star\top} \mathbf{X}_{j,\cdot}^* \right\|_2^2} \\ &\lesssim \|\Delta^t\|_F \left\{ \sqrt{\frac{\|\mathbf{X}^*\|_{2,\infty}^2 \sigma_{\max} \log n}{p}} + \frac{\|\mathbf{X}^*\|_{2,\infty}^2 \log n}{p} \right\} \\ &\stackrel{(i)}{\lesssim} \|\Delta^t\| \left\{ \sqrt{\frac{\|\mathbf{X}^*\|_{2,\infty}^2 r \sigma_{\max} \log n}{p}} + \frac{\sqrt{r} \|\mathbf{X}^*\|_{2,\infty}^2 \log n}{p} \right\} \\ &\stackrel{(ii)}{\lesssim} \|\Delta^t\| \left\{ \sqrt{\frac{\kappa \mu r^2 \log n}{np}} + \frac{\kappa \mu r^{3/2} \log n}{np} \right\} \sigma_{\max} \\ &\stackrel{(iii)}{\leq} \gamma \sigma_{\min} \|\Delta^t\|, \end{aligned}$$

for arbitrarily small  $\gamma > 0$ . Here, (i) follows from  $\|\Delta^t\|_F \leq \sqrt{r} \|\Delta^t\|$ , (ii) holds owing to the incoherence condition (114), and (iii) follows as long as  $n^2 p \gg \kappa^3 \mu r^2 n \log n$ .

- For the second term  $\theta_2$  in (128), denote

$$A = \mathcal{P}_\Omega \left( \mathbf{X}^* \Delta^{t\top} \right) \mathbf{X}^* - p \left( \mathbf{X}^* \Delta^{t\top} \right) \mathbf{X}^*,$$

whose  $l$ th row is given by

$$A_{l,\cdot} = \mathbf{X}_{l,\cdot}^* \sum_{j=1}^n (\delta_{l,j} - p) \Delta_{j,\cdot}^{\top} \mathbf{X}_{j,\cdot}^*. \tag{129}$$

Recalling the induction hypotheses (28b) and (28c), we define

$$\|\Delta^t\|_{2,\infty} \leq C_5 \rho^t \mu r \sqrt{\frac{\log n}{np}} \|\mathbf{X}^*\|_{2,\infty} + C_8 \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} \|\mathbf{X}^*\|_{2,\infty} := \xi \tag{130}$$

$$\|\Delta^t\| \leq C_9 \rho^t \mu r \frac{1}{\sqrt{np}} \|\mathbf{X}^*\| + C_{10} \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \|\mathbf{X}^*\| := \psi. \tag{131}$$

With these two definitions in place, we now introduce a ‘‘truncation level’’

$$\omega := 2p\xi\sigma_{\max} \tag{132}$$

that allows us to bound  $\theta_2$  in terms of the following two terms

$$\theta_2 = \frac{1}{p} \sqrt{\sum_{l=1}^n \|A_{l,\cdot}\|_2^2} \leq \frac{1}{p} \underbrace{\sqrt{\sum_{l=1}^n \|A_{l,\cdot}\|_2^2 \mathbb{1}_{\{\|A_{l,\cdot}\|_2 \leq \omega\}}}}_{:=\phi_1} + \frac{1}{p} \underbrace{\sqrt{\sum_{l=1}^n \|A_{l,\cdot}\|_2^2 \mathbb{1}_{\{\|A_{l,\cdot}\|_2 \geq \omega\}}}}_{:=\phi_2}.$$

We will apply different strategies when upper bounding the terms  $\phi_1$  and  $\phi_2$ , with their bounds given in the following two lemmas under the induction hypotheses (28b) and (28c).

**Lemma 22** *Under the conditions in Lemma 9, there exist some constants  $c, C > 0$  such that with probability exceeding  $1 - c \exp(-Cnr \log n)$ ,*

$$\phi_1 \lesssim \xi \sqrt{p \sigma_{\max} \|X^*\|_{2,\infty}^2 nr \log^2 n} \tag{133}$$

holds simultaneously for all  $\Delta^t$  obeying (130) and (131). Here,  $\xi$  is defined in (130).

**Lemma 23** *Under the conditions in Lemma 9, with probability at least  $1 - O(n^{-10})$ ,*

$$\phi_2 \lesssim \xi \sqrt{\kappa \mu r^2 p \log^2 n} \|X^*\|^2 \tag{134}$$

holds simultaneously for all  $\Delta^t$  obeying (130) and (131). Here,  $\xi$  is defined in (130).

Bounds (133) and (134) together with the incoherence condition (114) yield

$$\begin{aligned} \theta_2 &\lesssim \frac{1}{p} \xi \sqrt{p \sigma_{\max} \|X^*\|_{2,\infty}^2 nr \log^2 n} + \frac{1}{p} \xi \sqrt{\kappa \mu r^2 p \log^2 n} \|X^*\|^2 \\ &\lesssim \sqrt{\frac{\kappa \mu r^2 \log^2 n}{p}} \xi \sigma_{\max}. \end{aligned}$$

- Next, we assert that the third term  $\theta_3$  in (128) has the same upper bound as  $\theta_2$ . The proof follows by repeating the same argument used in bounding  $\theta_2$  and is hence omitted.

Take the previous three bounds on  $\theta_1, \theta_2,$  and  $\theta_3$  together to arrive at

$$\beta_2 \leq \eta (|\theta_1| + |\theta_2| + |\theta_3|) \leq \eta \gamma \sigma_{\min} \|\Delta^t\| + \tilde{C} \eta \sqrt{\frac{\kappa \mu r^2 \log^2 n}{p}} \xi \sigma_{\max}$$

for some constant  $\tilde{C} > 0$ .

(c) Substituting the preceding bounds on  $\beta_1$  and  $\beta_2$  into (123), we reach

$$\begin{aligned}
 \alpha_2 &\stackrel{(i)}{\leq} \left(1 - \eta\sigma_{\min} + \eta\gamma\sigma_{\min} + \eta \|\Delta^t\| \|X^*\| \right) \|\Delta^t\| + \eta \left\| \frac{1}{p} \mathcal{P}_\Omega(E) \right\| \|X^*\| \\
 &\quad + \tilde{C}\eta \sqrt{\frac{\kappa\mu r^2 \log^2 n}{p}} \sigma_{\max} \left( C_5 \rho^t \mu r \sqrt{\frac{\log n}{np}} \|X^*\|_{2,\infty} + C_8 \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} \|X^*\|_{2,\infty} \right) \\
 &\stackrel{(ii)}{\leq} \left(1 - \frac{\sigma_{\min}}{2} \eta\right) \|\Delta^t\| + \eta \left\| \frac{1}{p} \mathcal{P}_\Omega(E) \right\| \|X^*\| \\
 &\quad + \tilde{C}\eta \sqrt{\frac{\kappa\mu r^2 \log^2 n}{p}} \sigma_{\max} \left( C_5 \rho^t \mu r \sqrt{\frac{\log n}{np}} \|X^*\|_{2,\infty} + C_8 \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} \|X^*\|_{2,\infty} \right) \\
 &\stackrel{(iii)}{\leq} \left(1 - \frac{\sigma_{\min}}{2} \eta\right) \|\Delta^t\| + C\eta\sigma \sqrt{\frac{n}{p}} \|X^*\| \\
 &\quad + \tilde{C}\eta \sqrt{\frac{\kappa^2 \mu^2 r^3 \log^3 n}{np}} \sigma_{\max} \left( C_5 \rho^t \mu r \sqrt{\frac{1}{np}} + C_8 \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \right) \|X^*\| \tag{135}
 \end{aligned}$$

for some constant  $C > 0$ . Here, (i) uses the definition of  $\xi$  (cf. (130)), (ii) holds if  $\gamma$  is small enough and  $\|\Delta^t\| \|X^*\| \ll \sigma_{\min}$ , and (iii) follows from Lemma 40 as well as the incoherence condition (114). An immediate consequence of (135) is that under the sample size condition and the noise condition of this lemma, one has

$$\|\tilde{X}^{t+1} - X^*\| \|X^*\| \leq \sigma_{\min}/2 \tag{136}$$

if  $0 < \eta \leq 1/\sigma_{\max}$ .

2. We then move on to the first term  $\alpha_1$  in (121), which can be rewritten as

$$\alpha_1 = \|X^{t+1} \widehat{H}^t R_1 - \tilde{X}^{t+1}\|,$$

with

$$R_1 = (\widehat{H}^t)^{-1} \widehat{H}^{t+1} := \arg \min_{R \in \mathcal{O}^{r \times r}} \|X^{t+1} \widehat{H}^t R - X^*\|_F. \tag{137}$$

(a) First, we claim that  $\tilde{X}^{t+1}$  satisfies

$$I_r = \arg \min_{R \in \mathcal{O}^{r \times r}} \|\tilde{X}^{t+1} R - X^*\|_F, \tag{138}$$

meaning that  $\tilde{X}^{t+1}$  is already rotated to the direction that is most “aligned” with  $X^*$ . This important property eases the analysis. In fact, in view of Lemma 35, (138) follows if one can show that  $X^{*\top} \tilde{X}^{t+1}$  is symmetric and positive semidefinite. First of all, it follows from Lemma 35 that  $X^{*\top} X^t \widehat{H}^t$  is symmetric and, hence, by definition,

$$X^{*\top} \tilde{X}^{t+1} = X^{*\top} X^t \widehat{H}^t - \frac{\eta}{p} X^{*\top} \mathcal{P}_\Omega \left[ X^t X^{t\top} - (M^* + E) \right] X^*$$



is also symmetric. Additionally,

$$\|X^{*\top} \tilde{X}^{t+1} - M^*\| \leq \|\tilde{X}^{t+1} - X^*\| \|X^*\| \leq \sigma_{\min}/2,$$

where the second inequality holds according to (136). Weyl’s inequality guarantees that

$$X^{*\top} \tilde{X}^{t+1} \succeq \frac{1}{2} \sigma_{\min} I_r,$$

thus justifying (138) via Lemma 35.

- (b) With (137) and (138) in place, we resort to Lemma 37 to establish the bound. Specifically, take  $X_1 = \tilde{X}^{t+1}$  and  $X_2 = X^{t+1} \hat{H}^t$ , and it comes from (136) that

$$\|X_1 - X^*\| \|X^*\| \leq \sigma_{\min}/2.$$

Moreover, we have

$$\|X_1 - X_2\| \|X^*\| = \|X^{t+1} \hat{H}^t - \tilde{X}^{t+1}\| \|X^*\|,$$

in which

$$\begin{aligned} X^{t+1} \hat{H}^t - \tilde{X}^{t+1} &= \left( X^t - \eta \frac{1}{p} \mathcal{P}_\Omega \left[ X^t X^{t\top} - (M^* + E) \right] X^t \right) \hat{H}^t \\ &\quad - \left[ X^t \hat{H}^t - \eta \frac{1}{p} \mathcal{P}_\Omega \left[ X^t X^{t\top} - (M^* + E) \right] X^* \right] \\ &= -\eta \frac{1}{p} \mathcal{P}_\Omega \left[ X^t X^{t\top} - (M^* + E) \right] (X^t \hat{H}^t - X^*). \end{aligned}$$

This allows one to derive

$$\begin{aligned} &\|X^{t+1} \hat{H}^t - \tilde{X}^{t+1}\| \\ &\leq \eta \left\| \frac{1}{p} \mathcal{P}_\Omega \left[ X^t X^{t\top} - M^* \right] (X^t \hat{H}^t - X^*) \right\| \\ &\quad + \eta \left\| \frac{1}{p} \mathcal{P}_\Omega(E) (X^t \hat{H}^t - X^*) \right\| \\ &\leq \eta \left( 2n \|\Delta^t\|_{2,\infty}^2 + 4\sqrt{n} \log n \|\Delta^t\|_{2,\infty} \|X^*\| + C\sigma \sqrt{\frac{n}{p}} \right) \|\Delta^t\| \quad (139) \end{aligned}$$

for some absolute constant  $C > 0$ . Here the last inequality follows from Lemmas 40 and 43. As a consequence,

$$\begin{aligned} & \|X_1 - X_2\| \|X^*\| \\ & \leq \eta \left( 2n \|\Delta^t\|_{2,\infty}^2 + 4\sqrt{n} \log n \|\Delta^t\|_{2,\infty} \|X^*\| + C\sigma \sqrt{\frac{n}{p}} \right) \|\Delta^t\| \|X^*\|. \end{aligned}$$

Under our sample size condition and the noise condition (27) and the induction hypotheses (28), one can show

$$\|X_1 - X_2\| \|X^*\| \leq \sigma_{\min}/4.$$

Apply Lemma 37 and (139) to reach

$$\begin{aligned} \alpha_1 & \leq 5\kappa \|X^{t+1} \widehat{H}^t - \widetilde{X}^{t+1}\| \\ & \leq 5\kappa \eta \left( 2n \|\Delta^t\|_{2,\infty}^2 + 2\sqrt{n} \log n \|\Delta^t\|_{2,\infty} \|X^*\| + C\sigma \sqrt{\frac{n}{p}} \right) \|\Delta^t\|. \end{aligned}$$

3. Combining the above bounds on  $\alpha_1$  and  $\alpha_2$ , we arrive at

$$\begin{aligned} & \|X^{t+1} \widehat{H}^{t+1} - X^*\| \\ & \leq \left( 1 - \frac{\sigma_{\min}}{2} \eta \right) \|\Delta^t\| + \eta C\sigma \sqrt{\frac{n}{p}} \|X^*\| \\ & \quad + \widetilde{C} \eta \sqrt{\frac{\kappa^2 \mu^2 r^3 \log^3 n}{np}} \sigma_{\max} \left( C_5 \rho^t \mu r \sqrt{\frac{1}{np}} + \frac{C_8}{\sigma_{\min}} \sigma \sqrt{\frac{n}{p}} \right) \|X^*\| \\ & \quad + 5\eta \kappa \left( 2n \|\Delta^t\|_{2,\infty}^2 + 2\sqrt{n} \log n \|\Delta^t\|_{2,\infty} \|X^*\| + C\sigma \sqrt{\frac{n}{p}} \right) \|\Delta^t\| \\ & \leq C_9 \rho^{t+1} \mu r \frac{1}{\sqrt{np}} \|X^*\| + C_{10} \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \|X^*\|, \end{aligned}$$

with the proviso that  $\rho \geq 1 - (\sigma_{\min}/3) \cdot \eta$ ,  $\kappa$  is a constant, and  $n^2 p \gg \mu^3 r^3 n \log^3 n$ .

### B.3.1 Proof of Lemma 22

In what follows, we first assume that the  $\delta_{j,k}$ 's are independent and then use the standard decoupling trick to extend the result to symmetric sampling case (i.e.,  $\delta_{j,k} = \delta_{k,j}$ ).

To begin with, we justify the concentration bound for any  $\Delta^t$  independent of  $\Omega$ , followed by the standard covering argument that extends the bound to all  $\Delta^t$ . For any  $\Delta^t$  independent of  $\Omega$ , one has

$$\begin{aligned} B & := \max_{1 \leq j \leq n} \left\| X_{l,\cdot}^* (\delta_{l,j} - p) \Delta_{j,\cdot}^{t\top} X_{j,\cdot}^* \right\|_2 \leq \|X^*\|_{2,\infty}^2 \xi \\ \text{and } V & := \left\| \mathbb{E} \left[ \sum_{j=1}^n (\delta_{l,j} - p)^2 X_{l,\cdot}^* \Delta_{j,\cdot}^{t\top} X_{j,\cdot}^* \left( X_{l,\cdot}^* \Delta_{j,\cdot}^{t\top} X_{j,\cdot}^* \right)^\top \right] \right\| \end{aligned}$$

$$\begin{aligned} &\leq p \|X_{l,\cdot}^*\|_2^2 \|X^*\|_{2,\infty}^2 \left\| \sum_{j=1}^n \Delta_{j,\cdot}^{t\top} \Delta_{j,\cdot}^t \right\| \\ &\leq p \|X_{l,\cdot}^*\|_2^2 \|X^*\|_{2,\infty}^2 \psi^2 \\ &\leq 2p \|X^*\|_{2,\infty}^2 \xi^2 \sigma_{\max}, \end{aligned}$$

where  $\xi$  and  $\psi$  are defined, respectively, in (130) and (131). Here, the last line makes use of the fact that

$$\|X^*\|_{2,\infty} \psi \ll \xi \|X^*\| = \xi \sqrt{\sigma_{\max}}, \tag{140}$$

as long as  $n$  is sufficiently large. Apply the matrix Bernstein inequality [114, Theorem 6.1.1] to get

$$\begin{aligned} \mathbb{P} \{ \|A_{l,\cdot}\|_2 \geq t \} &\leq 2r \exp \left( - \frac{ct^2}{2p\xi^2\sigma_{\max} \|X^*\|_{2,\infty}^2 + t \cdot \|X^*\|_{2,\infty}^2 \xi} \right) \\ &\leq 2r \exp \left( - \frac{ct^2}{4p\xi^2\sigma_{\max} \|X^*\|_{2,\infty}^2} \right) \end{aligned}$$

for some constant  $c > 0$ , provided that

$$t \leq 2p\sigma_{\max}\xi.$$

This upper bound on  $t$  is exactly the truncation level  $\omega$  we introduce in (132). With this in mind, we can easily verify that

$$\|A_{l,\cdot}\|_2 \mathbb{1}_{\{\|A_{l,\cdot}\|_2 \leq \omega\}}$$

is a sub-Gaussian random variable with variance proxy not exceeding  $O(p\xi^2\sigma_{\max} \|X^*\|_{2,\infty}^2 \log r)$ . Therefore, invoking the concentration bounds for quadratic functions [57, Theorem 2.1] yields that for some constants  $C_0, C > 0$ , with probability at least  $1 - C_0e^{-Cnr \log n}$ ,

$$\phi_1^2 = \sum_{l=1}^n \|A_{l,\cdot}\|_2^2 \mathbb{1}_{\{\|A_{l,\cdot}\|_2 \leq \omega\}} \lesssim p\xi^2\sigma_{\max} \|X^*\|_{2,\infty}^2 nr \log^2 n.$$

Now that we have established an upper bound on any fixed matrix  $\Delta^l$  (which holds with exponentially high probability), we can proceed to invoke the standard epsilon-net argument to establish a uniform bound over all feasible  $\Delta^l$ . This argument is fairly standard and is thus omitted; see [111, Section 2.3.1] or the proof of Lemma 42. In conclusion, we have that with probability exceeding  $1 - C_0e^{-\frac{1}{2}Cnr \log n}$ ,

$$\phi_1 = \sqrt{\sum_{l=1}^n \|A_{l,\cdot}\|_2^2 \mathbb{1}_{\{\|A_{l,\cdot}\|_2 \leq \omega\}}} \lesssim \sqrt{p\xi^2 \sigma_{\max} \|X^*\|_{2,\infty}^2 nr \log^2 n} \tag{141}$$

holds simultaneously for all  $\Delta^t \in \mathbb{R}^{n \times r}$  obeying the conditions of the lemma.

In the end, we comment on how to extend the bound to the symmetric sampling pattern where  $\delta_{j,k} = \delta_{k,j}$ . Recall from (129) that the diagonal element  $\delta_{l,l}$  cannot change the  $\ell_2$  norm of  $A_{l,\cdot}$  by more than  $\|X^*\|_{2,\infty}^2 \xi$ . As a result, changing all the diagonals  $\{\delta_{l,l}\}$  cannot change the quantity of interest (i.e.,  $\phi_1$ ) by more than  $\sqrt{n} \|X^*\|_{2,\infty}^2 \xi$ . This is smaller than the right-hand side of (141) under our incoherence and sample size conditions. Hence, from now on, we ignore the effect of  $\{\delta_{l,l}\}$  and focus on off-diagonal terms. The proof then follows from the same argument as in [48, Theorem D.2]. More specifically, we can employ the construction of Bernoulli random variables introduced therein to demonstrate that the upper bound in (141) still holds if the indicator  $\delta_{i,j}$  is replaced by  $(\tau_{i,j} + \tau'_{i,j})/2$ , where  $\tau_{i,j}$  and  $\tau'_{i,j}$  are independent copies of the symmetric Bernoulli random variables. Recognizing that  $\sup_{\Delta^t} \phi_1$  is a norm of the Bernoulli random variables  $\tau_{i,j}$ , one can repeat the decoupling argument in [48, Claim D.3] to finish the proof. We omit the details here for brevity.

### B.3.2 Proof of Lemma 23

Observe from (129) that

$$\begin{aligned} \|A_{l,\cdot}\|_2 &\leq \|X^*\|_{2,\infty} \left\| \sum_{j=1}^n (\delta_{l,j} - p) \Delta_{j,\cdot}^{t\top} X_{j,\cdot}^* \right\| \tag{142} \\ &\leq \|X^*\|_{2,\infty} \left( \left\| \sum_{j=1}^n \delta_{l,j} \Delta_{j,\cdot}^{t\top} X_{j,\cdot}^* \right\| + p \|\Delta^t\| \|X^*\| \right) \\ &\leq \|X^*\|_{2,\infty} \left( \left\| [\delta_{l,1} \Delta_{1,\cdot}^{t\top}, \dots, \delta_{l,n} \Delta_{n,\cdot}^{t\top}] \right\| \left\| \begin{bmatrix} \delta_{l,1} X_{1,\cdot}^* \\ \vdots \\ \delta_{l,n} X_{n,\cdot}^* \end{bmatrix} \right\| + p\psi \|X^*\| \right) \\ &\leq \|X^*\|_{2,\infty} (\|G_l(\Delta^t)\| \cdot 1.2\sqrt{p} \|X^*\| + p\psi \|X^*\|), \tag{143} \end{aligned}$$

where  $\psi$  is as defined in (131) and  $G_l(\cdot)$  is as defined in Lemma 41. Here, the last inequality follows from Lemma 41; namely, for some constant  $C > 0$ , the following holds with probability at least  $1 - O(n^{-10})$

$$\left\| \begin{bmatrix} \delta_{l,1} X_{1,\cdot}^* \\ \vdots \\ \delta_{l,n} X_{n,\cdot}^* \end{bmatrix} \right\| \leq \left( p \|X^*\|^2 + C \sqrt{p \|X^*\|_{2,\infty}^2 \|X^*\|^2 \log n} + C \|X^*\|_{2,\infty}^2 \log n \right)^{\frac{1}{2}}$$

$$\leq \left( p + C\sqrt{p\frac{\kappa\mu r}{n}\log n} + C\frac{\kappa\mu r\log n}{n} \right)^{\frac{1}{2}} \|X^*\| \leq 1.2\sqrt{p} \|X^*\|, \tag{144}$$

where we also use the incoherence condition (114) and the sample complexity condition  $n^2 p \gg \kappa\mu r n \log n$ . Hence, the event

$$\|A_{l,\cdot}\|_2 \geq \omega = 2p\sigma_{\max}\xi$$

together with (142) and (143) necessarily implies that

$$\begin{aligned} \left\| \sum_{j=1}^n (\delta_{l,j} - p) \Delta_{j,\cdot}^t \top X_{j,\cdot}^* \right\| &\geq 2p\sigma_{\max} \frac{\xi}{\|X^*\|_{2,\infty}} \quad \text{and} \\ \|G_l(\Delta^t)\| &\geq \frac{2p\sigma_{\max}\xi}{1.2\sqrt{p}} - p\psi \geq \frac{2\sqrt{p}\|X^*\|_{2,\infty}\xi - \sqrt{p}\psi}{1.2} \geq 1.5\sqrt{p} \frac{\xi}{\|X^*\|_{2,\infty}} \|X^*\|, \end{aligned}$$

where the last inequality follows from bound (140). As a result, with probability at least  $1 - O(n^{-10})$  (i.e., when (144) holds for all  $l$ 's) we can upper bound  $\phi_2$  by

$$\phi_2 = \sqrt{\sum_{l=1}^n \|A_{l,\cdot}\|_2^2 \mathbb{1}_{\{\|A_{l,\cdot}\|_2 \geq \omega\}}} \leq \sqrt{\sum_{l=1}^n \|A_{l,\cdot}\|_2^2 \mathbb{1}_{\{\|G_l(\Delta^t)\| \geq \frac{1.5\sqrt{p}\xi\sqrt{\sigma_{\max}}}{\|X^*\|_{2,\infty}}\}}},$$

where the indicator functions are now specified with respect to  $\|G_l(\Delta^t)\|$ .

Next, we divide into multiple cases based on the size of  $\|G_l(\Delta^t)\|$ . By Lemma 42, for some constants  $c_1, c_2 > 0$ , with probability at least  $1 - c_1 \exp(-c_2 nr \log n)$ ,

$$\sum_{l=1}^n \mathbb{1}_{\{\|G_l(\Delta^t)\| \geq 4\sqrt{p}\psi + \sqrt{2^k r \xi}\}} \leq \frac{\alpha n}{2^{k-3}} \tag{145}$$

for any  $k \geq 0$  and any  $\alpha \gtrsim \log n$ . We claim that it suffices to consider the set of sufficiently large  $k$  obeying

$$\sqrt{2^k r \xi} \geq 4\sqrt{p}\psi \quad \text{or equivalently} \quad k \geq \log \frac{16p\psi^2}{r\xi^2}; \tag{146}$$

otherwise, we can use (140) to obtain

$$4\sqrt{p}\psi + \sqrt{2^k r \xi} \leq 8\sqrt{p}\psi \ll 1.5\sqrt{p} \frac{\xi}{\|X^*\|_{2,\infty}} \|X^*\|,$$

which contradicts the event  $\|A_{l \cdot}\|_2 \geq \omega$ . Consequently, we divide all indices into the following sets

$$S_k = \left\{ 1 \leq l \leq n : \|G_l(\Delta^t)\| \in (\sqrt{2^k r \xi}, \sqrt{2^{k+1} r \xi}] \right\} \tag{147}$$

defined for each integer  $k$  obeying (146). Under condition (146), it follows from (145) that

$$\sum_{l=1}^n \mathbb{1} \left\{ \|G_l(\Delta^t)\| \geq \sqrt{2^{k+2} r \xi} \right\} \leq \sum_{l=1}^n \mathbb{1} \left\{ \|G_l(\Delta^t)\| \geq 4\sqrt{p}\psi + \sqrt{2^k r \xi} \right\} \leq \frac{\alpha n}{2^{k-3}},$$

meaning that the cardinality of  $S_k$  satisfies

$$|S_{k+2}| \leq \frac{\alpha n}{2^{k-3}} \quad \text{or} \quad |S_k| \leq \frac{\alpha n}{2^{k-5}}$$

which decays exponentially fast as  $k$  increases. Therefore, when restricting attention to the set of indices within  $S_k$ , we can obtain

$$\begin{aligned} \sqrt{\sum_{l \in S_k} \|A_{l \cdot}\|_2^2} &\stackrel{(i)}{\leq} \sqrt{|S_k| \cdot \|X^*\|_{2,\infty}^2 \left( 1.2\sqrt{2^{k+1} r \xi} \sqrt{p} \|X^*\| + p\psi \|X^*\| \right)^2} \\ &\leq \sqrt{\frac{\alpha n}{2^{k-5}} \|X^*\|_{2,\infty}^2 \left( 2\sqrt{2^{k+1} r \xi} \sqrt{p} \|X^*\| + p\psi \|X^*\| \right)} \\ &\stackrel{(ii)}{\leq} 4\sqrt{\frac{\alpha n}{2^{k-5}} \|X^*\|_{2,\infty}^2 \sqrt{2^{k+1} r \xi} \sqrt{p} \|X^*\|} \\ &\stackrel{(iii)}{\leq} 32\sqrt{\alpha \kappa \mu r^2 p \xi} \|X^*\|^2, \end{aligned}$$

where (i) follows from bound (143) and constraint (147) in  $S_k$ , (ii) is a consequence of (146), and (iii) uses the incoherence condition (114).

Now that we have developed an upper bound with respect to each  $S_k$ , we can add them up to yield the final upper bound. Note that there are in total no more than  $O(\log n)$  different sets, i.e.,  $S_k = \emptyset$  if  $k \geq c_1 \log n$  for  $c_1$  sufficiently large. This arises since

$$\|G_l(\Delta^t)\| \leq \|\Delta^t\|_F \leq \sqrt{n} \|\Delta^t\|_{2,\infty} \leq \sqrt{n} \xi \leq \sqrt{n} \sqrt{r} \xi$$

and hence

$$\mathbb{1} \left\{ \|G_l(\Delta^t)\| \geq 4\sqrt{p}\psi + \sqrt{2^k r \xi} \right\} = 0 \quad \text{and} \quad S_k = \emptyset$$

if  $k / \log n$  is sufficiently large. One can thus conclude that

$$\phi_2^2 \leq \sum_{k=\log \frac{16p\psi^2}{r\xi^2}}^{c_1 \log n} \sum_{l \in S_k} \|A_{l,\cdot}\|_2^2 \lesssim \left( \sqrt{\alpha\kappa\mu r^2 p \xi} \|X^*\|^2 \right)^2 \cdot \log n,$$

leading to  $\phi_2 \lesssim \xi \sqrt{\alpha\kappa\mu r^2 p \log n} \|X^*\|^2$ . The proof is finished by taking  $\alpha = c \log n$  for some sufficiently large constant  $c > 0$ .

**B.4 Proof of Lemma 10**

1. To obtain (73a), we invoke Lemma 37. Setting  $X_1 = X^t \widehat{H}^t$  and  $X_2 = X^{t,(l)} R^{t,(l)}$ , we get

$$\|X_1 - X^*\| \|X^*\| \stackrel{(i)}{\leq} C_9 \rho^t \mu r \frac{1}{\sqrt{np}} \sigma_{\max} + \frac{C_{10}}{\sigma_{\min}} \sigma \sqrt{\frac{n \log n}{p}} \sigma_{\max} \stackrel{(ii)}{\leq} \frac{1}{2} \sigma_{\min},$$

where (i) follows from (70c) and (ii) holds as long as  $n^2 p \gg \kappa^2 \mu^2 r^2 n$  and the noise satisfies (27). In addition,

$$\begin{aligned} \|X_1 - X_2\| \|X^*\| &\leq \|X_1 - X_2\|_F \|X^*\| \\ &\stackrel{(i)}{\leq} \left( C_3 \rho^t \mu r \sqrt{\frac{\log n}{np}} \|X^*\|_{2,\infty} + \frac{C_7}{\sigma_{\min}} \sigma \sqrt{\frac{n \log n}{p}} \|X^*\|_{2,\infty} \right) \|X^*\| \\ &\stackrel{(ii)}{\leq} C_3 \rho^t \mu r \sqrt{\frac{\log n}{np}} \sigma_{\max} + \frac{C_7}{\sigma_{\min}} \sigma \sqrt{\frac{n \log n}{p}} \sigma_{\max} \\ &\stackrel{(iii)}{\leq} \frac{1}{2} \sigma_{\min}, \end{aligned}$$

where (i) utilizes (70d), (ii) follows since  $\|X^*\|_{2,\infty} \leq \|X^*\|$ , and (iii) holds if  $n^2 p \gg \kappa^2 \mu^2 r^2 n \log n$  and the noise satisfies (27). With these in place, Lemma 37 immediately yields (73a).

2. The first inequality in (73b) follows directly from the definition of  $\widehat{H}^{t,(l)}$ . The second inequality is concerned with the estimation error of  $X^{t,(l)} R^{t,(l)}$  with respect to the Frobenius norm. Combining (70a), (70d), and the triangle inequality yields

$$\begin{aligned} \|X^{t,(l)} R^{t,(l)} - X^*\|_F &\leq \|X^t \widehat{H}^t - X^*\|_F + \|X^t \widehat{H}^t - X^{t,(l)} R^{t,(l)}\|_F \\ &\leq C_4 \rho^t \mu r \frac{1}{\sqrt{np}} \|X^*\|_F + \frac{C_{10}}{\sigma_{\min}} \sigma \sqrt{\frac{n}{p}} \|X^*\|_F + C_3 \rho^t \mu r \sqrt{\frac{\log n}{np}} \|X^*\|_{2,\infty} \\ &\quad + \frac{C_7 \sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} \|X^*\|_{2,\infty} \end{aligned}$$

$$\begin{aligned}
&\leq C_4 \rho^t \mu r \frac{1}{\sqrt{np}} \|X^*\|_F + \frac{C_1 \sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \|X^*\|_F + C_3 \rho^t \mu r \sqrt{\frac{\log n}{np}} \sqrt{\frac{\kappa \mu}{n}} \|X^*\|_F \\
&\quad + \frac{C_7 \sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} \sqrt{\frac{\kappa \mu}{n}} \|X^*\|_F \\
&\leq 2C_4 \rho^t \mu r \frac{1}{\sqrt{np}} \|X^*\|_F + \frac{2C_1 \sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \|X^*\|_F, \tag{148}
\end{aligned}$$

where the last step holds true as long as  $n \gg \kappa \mu \log n$ .

3. To obtain (73c), we use (70d) and (70b) to get

$$\begin{aligned}
\|X^{t,(l)} R^{t,(l)} - X^*\|_{2,\infty} &\leq \|X^t \widehat{H}^t - X^*\|_{2,\infty} + \|X^t \widehat{H}^t - X^{t,(l)} R^{t,(l)}\|_F \\
&\leq C_5 \rho^t \mu r \sqrt{\frac{\log n}{np}} \|X^*\|_{2,\infty} + \frac{C_8 \sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} \|X^*\|_{2,\infty} \\
&\quad + C_3 \rho^t \mu r \sqrt{\frac{\log n}{np}} \|X^*\|_{2,\infty} + \frac{C_7 \sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} \|X^*\|_{2,\infty} \\
&\leq (C_3 + C_5) \rho^t \mu r \sqrt{\frac{\log n}{np}} \|X^*\|_{2,\infty} + \frac{C_8 + C_7}{\sigma_{\min}} \sigma \sqrt{\frac{n \log n}{p}} \|X^*\|_{2,\infty}.
\end{aligned}$$

4. Finally, to obtain (73d), one can take the triangle inequality

$$\begin{aligned}
\|X^{t,(l)} \widehat{H}^{t,(l)} - X^*\| &\leq \|X^{t,(l)} \widehat{H}^{t,(l)} - X^t \widehat{H}^t\|_F + \|X^t \widehat{H}^t - X^*\| \\
&\leq 5\kappa \|X^t \widehat{H}^t - X^{t,(l)} R^{t,(l)}\|_F + \|X^t \widehat{H}^t - X^*\|,
\end{aligned}$$

where the second line follows from (73a). Combine (70d) and (70c) to yield

$$\begin{aligned}
&\|X^{t,(l)} \widehat{H}^{t,(l)} - X^*\| \\
&\leq 5\kappa \left( C_3 \rho^t \mu r \sqrt{\frac{\log n}{np}} \|X^*\|_{2,\infty} + \frac{C_7}{\sigma_{\min}} \sigma \sqrt{\frac{n \log n}{p}} \|X^*\|_{2,\infty} \right) \\
&\quad + C_9 \rho^t \mu r \frac{1}{\sqrt{np}} \|X^*\| + \frac{C_{10} \sigma}{\sigma_{\min}} \sigma \sqrt{\frac{n}{p}} \|X^*\| \\
&\leq 5\kappa \sqrt{\frac{\kappa \mu r}{n}} \|X^*\| \left( C_3 \rho^t \mu r \sqrt{\frac{\log n}{np}} + \frac{C_7}{\sigma_{\min}} \sigma \sqrt{\frac{n \log n}{p}} \right) \\
&\quad + C_9 \rho^t \mu r \frac{1}{\sqrt{np}} \|X^*\| + \frac{C_{10} \sigma}{\sigma_{\min}} \sigma \sqrt{\frac{n}{p}} \|X^*\| \\
&\leq 2C_9 \rho^t \mu r \frac{1}{\sqrt{np}} \|X^*\| + \frac{2C_{10} \sigma}{\sigma_{\min}} \sigma \sqrt{\frac{n}{p}} \|X^*\|,
\end{aligned}$$



where the second inequality uses the incoherence of  $X^*$  (cf. (114)) and the last inequality holds as long as  $n \gg \kappa^3 \mu r \log n$ .

**B.5 Proof of Lemma 11**

From the definition of  $R^{t+1,(l)}$  (see (72)), we must have

$$\left\| X^{t+1} \widehat{H}^{t+1} - X^{t+1,(l)} R^{t+1,(l)} \right\|_F \leq \left\| X^{t+1} \widehat{H}^t - X^{t+1,(l)} R^{t,(l)} \right\|_F.$$

The gradient update rules in (24) and (69) allow one to express

$$\begin{aligned} & X^{t+1} \widehat{H}^t - X^{t+1,(l)} R^{t,(l)} \\ &= [X^t - \eta \nabla f(X^t)] \widehat{H}^t - [X^{t,(l)} - \eta \nabla f^{(l)}(X^{t,(l)})] R^{t,(l)} \\ &= X^t \widehat{H}^t - \eta \nabla f(X^t \widehat{H}^t) - [X^{t,(l)} R^{t,(l)} - \eta \nabla f^{(l)}(X^{t,(l)} R^{t,(l)})] \\ &= (X^t \widehat{H}^t - X^{t,(l)} R^{t,(l)}) - \eta [\nabla f(X^t \widehat{H}^t) - \nabla f(X^{t,(l)} R^{t,(l)})] \\ &\quad - \eta [\nabla f(X^{t,(l)} R^{t,(l)}) - \nabla f^{(l)}(X^{t,(l)} R^{t,(l)})], \end{aligned}$$

where we have again used the fact that  $\nabla f(X^t) R = \nabla f(X^t R)$  for any orthonormal matrix  $R \in \mathcal{O}^{r \times r}$  (similarly for  $\nabla f^{(l)}(X^{t,(l)})$ ). Relate the right-hand side of the above equation with  $\nabla f_{\text{clean}}(X)$  to reach

$$\begin{aligned} & X^{t+1} \widehat{H}^t - X^{t+1,(l)} R^{t,(l)} \\ &= \underbrace{(X^t \widehat{H}^t - X^{t,(l)} R^{t,(l)}) - \eta [\nabla f_{\text{clean}}(X^t \widehat{H}^t) - \nabla f_{\text{clean}}(X^{t,(l)} R^{t,(l)})]}_{:=B_1^{(l)}} \\ &\quad - \underbrace{\eta \left[ \frac{1}{p} \mathcal{P}_{\Omega_l} \left( X^{t,(l)} X^{t,(l)\top} - M^* \right) - \mathcal{P}_l \left( X^{t,(l)} X^{t,(l)\top} - M^* \right) \right]}_{:=B_2^{(l)}} X^{t,(l)} R^{t,(l)} \\ &\quad + \underbrace{\eta \frac{1}{p} \mathcal{P}_{\Omega} (E) \left( X^t \widehat{H}^t - X^{t,(l)} R^{t,(l)} \right)}_{:=B_3^{(l)}} + \underbrace{\eta \frac{1}{p} \mathcal{P}_{\Omega_l} (E) X^{t,(l)} R^{t,(l)}}_{:=B_4^{(l)}}, \end{aligned} \tag{149}$$

where we have used the following relationship between  $\nabla f^{(l)}(X)$  and  $\nabla f(X)$ :

$$\nabla f^{(l)}(X) = \nabla f(X) - \frac{1}{p} \mathcal{P}_{\Omega_l} [X X^\top - (M^* + E)] X + \mathcal{P}_l (X X^\top - M^*) X \tag{150}$$

for all  $X \in \mathbb{R}^{n \times r}$  with  $\mathcal{P}_{\Omega_l}$  and  $\mathcal{P}_l$  defined, respectively, in (66) and (67). In the sequel, we control the four terms in reverse order.

1. The last term  $\mathbf{B}_4^{(l)}$  is controlled via the following lemma.

**Lemma 24** *Suppose that the sample size obeys  $n^2 p > C \mu^2 r^2 n \log^2 n$  for some sufficiently large constant  $C > 0$ . Then with probability at least  $1 - O(n^{-10})$ , the matrix  $\mathbf{B}_4^{(l)}$  as defined in (149) satisfies*

$$\|\mathbf{B}_4^{(l)}\|_F \lesssim \eta \sigma \sqrt{\frac{n \log n}{p}} \|\mathbf{X}^*\|_{2,\infty}.$$

2. The third term  $\mathbf{B}_3^{(l)}$  can be bounded as follows

$$\|\mathbf{B}_3^{(l)}\|_F \leq \eta \left\| \frac{1}{p} \mathcal{P}_\Omega(\mathbf{E}) \right\| \left\| \mathbf{X}^t \widehat{\mathbf{H}}^t - \mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)} \right\|_F \lesssim \eta \sigma \sqrt{\frac{n}{p}} \left\| \mathbf{X}^t \widehat{\mathbf{H}}^t - \mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)} \right\|_F,$$

where the second inequality comes from Lemma 40.

3. For the second term  $\mathbf{B}_2^{(l)}$ , we have the following lemma.

**Lemma 25** *Suppose that the sample size obeys  $n^2 p \gg \mu^2 r^2 n \log n$ . Then with probability exceeding  $1 - O(n^{-10})$ , the matrix  $\mathbf{B}_2^{(l)}$  as defined in (149) satisfies*

$$\|\mathbf{B}_2^{(l)}\|_F \lesssim \eta \sqrt{\frac{\kappa^2 \mu^2 r^2 \log n}{np}} \left\| \mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)} - \mathbf{X}^* \right\|_{2,\infty} \sigma_{\max}. \tag{151}$$

4. Regarding the first term  $\mathbf{B}_1^{(l)}$ , apply the fundamental theorem of calculus [70, Chapter XIII, Theorem 4.2] to get

$$\text{vec}(\mathbf{B}_1^{(l)}) = \left( \mathbf{I}_{nr} - \eta \int_0^1 \nabla^2 f_{\text{clean}}(\mathbf{X}(\tau)) d\tau \right) \text{vec} \left( \mathbf{X}^t \widehat{\mathbf{H}}^t - \mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)} \right), \tag{152}$$

where we abuse the notation and denote  $\mathbf{X}(\tau) := \mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)} + \tau \left( \mathbf{X}^t \widehat{\mathbf{H}}^t - \mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)} \right)$ . Going through the same derivations as in the proof of Lemma 8 (see Appendix B.2), we get

$$\|\mathbf{B}_1^{(l)}\|_F \leq \left( 1 - \frac{\sigma_{\min}}{4} \eta \right) \left\| \mathbf{X}^t \widehat{\mathbf{H}}^t - \mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)} \right\|_F \tag{153}$$

with the proviso that  $0 < \eta \leq (2\sigma_{\min}) / (25\sigma_{\max}^2)$ .

Applying the triangle inequality to (149) and invoking the preceding four bounds, we arrive at

$$\begin{aligned} & \left\| \mathbf{X}^{t+1} \widehat{\mathbf{H}}^{t+1} - \mathbf{X}^{t+1,(l)} \mathbf{R}^{t+1,(l)} \right\|_F \\ & \leq \left( 1 - \frac{\sigma_{\min}}{4} \eta \right) \left\| \mathbf{X}^t \widehat{\mathbf{H}}^t - \mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)} \right\|_F \end{aligned}$$

$$\begin{aligned}
 &+ \tilde{C}\eta\sqrt{\frac{\kappa^2\mu^2r^2\log n}{np}}\|X^{t,(l)}R^{t,(l)} - X^*\|_{2,\infty}\sigma_{\max} \\
 &+ \tilde{C}\eta\sigma\sqrt{\frac{n}{p}}\|X^t\hat{H}^t - X^{t,(l)}R^{t,(l)}\|_F + \tilde{C}\eta\sigma\sqrt{\frac{n\log n}{p}}\|X^*\|_{2,\infty} \\
 = &\left(1 - \frac{\sigma_{\min}}{4}\eta + \tilde{C}\eta\sigma\sqrt{\frac{n}{p}}\right)\|X^t\hat{H}^t - X^{t,(l)}R^{t,(l)}\|_F \\
 &+ \tilde{C}\eta\sqrt{\frac{\kappa^2\mu^2r^2\log n}{np}}\|X^{t,(l)}R^{t,(l)} - X^*\|_{2,\infty}\sigma_{\max} \\
 &+ \tilde{C}\eta\sigma\sqrt{\frac{n\log n}{p}}\|X^*\|_{2,\infty} \\
 \leq &\left(1 - \frac{2\sigma_{\min}}{9}\eta\right)\|X^t\hat{H}^t - X^{t,(l)}R^{t,(l)}\|_F \\
 &+ \tilde{C}\eta\sqrt{\frac{\kappa^2\mu^2r^2\log n}{np}}\|X^{t,(l)}R^{t,(l)} - X^*\|_{2,\infty}\sigma_{\max} \\
 &+ \tilde{C}\eta\sigma\sqrt{\frac{n\log n}{p}}\|X^*\|_{2,\infty}
 \end{aligned}$$

for some absolute constant  $\tilde{C} > 0$ . Here the last inequality holds as long as  $\sigma\sqrt{n/p} \ll \sigma_{\min}$ , which is satisfied under our noise condition (27). This taken collectively with hypotheses (70d) and (73c) leads to

$$\begin{aligned}
 &\|X^{t+1}\hat{H}^{t+1} - X^{t+1,(l)}R^{t+1,(l)}\|_F \\
 \leq &\left(1 - \frac{2\sigma_{\min}}{9}\eta\right)\left(C_3\rho^t\mu r\sqrt{\frac{\log n}{np}}\|X^*\|_{2,\infty} + C_7\frac{\sigma}{\sigma_{\min}}\sqrt{\frac{n\log n}{p}}\|X^*\|_{2,\infty}\right) \\
 &+ \tilde{C}\eta\sqrt{\frac{\kappa^2\mu^2r^2\log n}{np}}\left[(C_3 + C_5)\rho^t\mu r\sqrt{\frac{\log n}{np}} + (C_8 + C_7)\frac{\sigma}{\sigma_{\min}}\sqrt{\frac{n\log n}{p}}\right]\|X^*\|_{2,\infty}\sigma_{\max} \\
 &+ \tilde{C}\eta\sigma\sqrt{\frac{n\log n}{p}}\|X^*\|_{2,\infty} \\
 \leq &\left(1 - \frac{\sigma_{\min}}{5}\eta\right)C_3\rho^t\mu r\sqrt{\frac{\log n}{np}}\|X^*\|_{2,\infty} + C_7\frac{\sigma}{\sigma_{\min}}\sqrt{\frac{n\log n}{p}}\|X^*\|_{2,\infty}
 \end{aligned}$$

as long as  $C_7 > 0$  is sufficiently large, where we have used the sample complexity assumption  $n^2p \gg \kappa^4\mu^2r^2n\log n$  and the step size  $0 < \eta \leq 1/(2\sigma_{\max}) \leq 1/(2\sigma_{\min})$ . This finishes the proof.

### B.5.1 Proof of Lemma 24

By the unitary invariance of the Frobenius norm, one has

$$\left\| \mathbf{B}_4^{(l)} \right\|_F = \frac{\eta}{p} \left\| \mathcal{P}_{\Omega_l}(\mathbf{E}) \mathbf{X}^{t,(l)} \right\|_F,$$

where all nonzero entries of the matrix  $\mathcal{P}_{\Omega_l}(\mathbf{E})$  reside in the  $l$ th row/column. Decouple the effects of the  $l$ th row and the  $l$ th column of  $\mathcal{P}_{\Omega_l}(\mathbf{E})$  to reach

$$\frac{p}{\eta} \left\| \mathbf{B}_4^{(l)} \right\|_F \leq \left\| \underbrace{\sum_{j=1}^n \delta_{l,j} E_{l,j} \mathbf{X}_{j,\cdot}^{t,(l)}}_{:=\mathbf{u}_j} \right\|_2 + \left\| \underbrace{\sum_{j:j \neq l} \delta_{l,j} E_{l,j} \mathbf{X}_{l,\cdot}^{t,(l)}}_{:=\alpha} \right\|_2, \tag{154}$$

where  $\delta_{l,j} := \mathbb{1}_{\{(l,j) \in \Omega\}}$  indicates whether the  $(l, j)$ th entry is observed. Since  $\mathbf{X}^{t,(l)}$  is independent of  $\{\delta_{l,j}\}_{1 \leq j \leq n}$  and  $\{E_{l,j}\}_{1 \leq j \leq n}$ , we can treat the first term as a sum of independent vectors  $\{\mathbf{u}_j\}$ . It is easy to verify that

$$\left\| \mathbf{u}_j \right\|_{\psi_1} \leq \left\| \mathbf{X}^{t,(l)} \right\|_{2,\infty} \left\| \delta_{l,j} E_{l,j} \right\|_{\psi_1} \lesssim \sigma \left\| \mathbf{X}^{t,(l)} \right\|_{2,\infty},$$

where  $\|\cdot\|_{\psi_1}$  denotes the sub-exponential norm [66, Section A.1]. Further, one can calculate

$$\begin{aligned} V &:= \left\| \mathbb{E} \left[ \sum_{j=1}^n (\delta_{l,j} E_{l,j})^2 \mathbf{X}_{j,\cdot}^{t,(l)} \mathbf{X}_{j,\cdot}^{t,(l)\top} \right] \right\| \lesssim p\sigma^2 \left\| \mathbb{E} \left[ \sum_{j=1}^n \mathbf{X}_{j,\cdot}^{t,(l)} \mathbf{X}_{j,\cdot}^{t,(l)\top} \right] \right\| \\ &= p\sigma^2 \left\| \mathbf{X}^{t,(l)} \right\|_F^2. \end{aligned}$$

Invoke the matrix Bernstein inequality [66, Theorem 2.7] to discover that with probability at least  $1 - O(n^{-10})$ ,

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

where the third inequality follows from  $\left\| \mathbf{X}^{t,(l)} \right\|_F^2 \leq n \left\| \mathbf{X}^{t,(l)} \right\|_{2,\infty}^2$  and the last inequality holds as long as  $np \gg \log^2 n$ .

Additionally, the remaining term  $\alpha$  in (154) can be controlled using the same argument, giving rise to

$$\alpha \lesssim \sigma \sqrt{np \log n} \left\| \mathbf{X}^{t,(l)} \right\|_{2,\infty}.$$

We then complete the proof by observing that

$$\|X^{t,(l)}\|_{2,\infty} = \|X^{t,(l)} R^{t,(l)}\|_{2,\infty} \leq \|X^{t,(l)} R^{t,(l)} - X^*\|_{2,\infty} + \|X^*\|_{2,\infty} \leq 2\|X^*\|_{2,\infty}, \tag{155}$$

where the last inequality follows by combining (73c), the sample complexity condition  $n^2 p \gg \mu^2 r^2 n \log n$ , and the noise condition (27).

### B.5.2 Proof of Lemma 25

For notational simplicity, we denote

$$C := X^{t,(l)} X^{t,(l)\top} - M^* = X^{t,(l)} X^{t,(l)\top} - X^* X^{*\top}. \tag{156}$$

Since the Frobenius norm is unitarily invariant, we have

$$\|B_2^{(l)}\|_F = \eta \left\| \underbrace{\left[ \frac{1}{p} \mathcal{P}_{\Omega_l}(C) - \mathcal{P}_l(C) \right]}_{:=W} X^{t,(l)} \right\|_F.$$

Again, all nonzero entries of the matrix  $W$  reside in its  $l$ th row/column. We can deal with the  $l$ th row and the  $l$ th column of  $W$  separately as follows

$$\begin{aligned} \frac{p}{\eta} \|B_2^{(l)}\|_F &\leq \left\| \sum_{j=1}^n (\delta_{l,j} - p) C_{l,j} X_{j,\cdot}^{t,(l)} \right\|_2 + \sqrt{\sum_{j:j \neq l} (\delta_{l,j} - p)^2} \|C\|_\infty \|X_{l,\cdot}^{t,(l)}\|_2 \\ &\lesssim \left\| \sum_{j=1}^n (\delta_{l,j} - p) C_{l,j} X_{j,\cdot}^{t,(l)} \right\|_2 + \sqrt{np} \|C\|_\infty \|X_{l,\cdot}^{t,(l)}\|_2, \end{aligned}$$

where  $\delta_{l,j} := \mathbb{1}_{\{(l,j) \in \Omega\}}$  and the second line relies on the fact that  $\sum_{j:j \neq l} (\delta_{l,j} - p)^2 \asymp np$ . It follows that

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

Here, (i) is a consequence of (155). In addition, (ii) follows from

$$\|X^{t,(l)}\|_F = \|X^{t,(l)} R^{t,(l)}\|_F \leq \|X^{t,(l)} R^{t,(l)} - X^*\|_F + \|X^*\|_F \leq 2\|X^*\|_F,$$

where the last inequality comes from (73b), the sample complexity condition  $n^2 p \gg \mu^2 r^2 n \log n$ , and the noise condition (27). The matrix Bernstein inequality [114, Theorem 6.1.1] reveals that

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

with probability exceeding  $1 - O(n^{-10})$ , and as a result,

$$\frac{p}{\eta} \|B_2^{(l)}\|_F \lesssim \sqrt{p \log n} \|C\|_\infty \|X^*\|_F + \sqrt{np} \|C\|_\infty \|X^*\|_{2,\infty} \tag{157}$$

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

To finish up, we make the observation that

$$\begin{aligned} \|C\|_\infty &= \left\| \mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)} (\mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)})^\top - X^* X^{*\top} \right\|_\infty \\ &\leq \left\| (\mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)} - X^*) (\mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)})^\top \right\|_\infty + \left\| X^* (\mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)} - X^*)^\top - X^* X^{*\top} \right\|_\infty \\ &\leq \left\| \mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)} - X^* \right\|_{2,\infty} \left\| \mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)} \right\|_{2,\infty} + \|X^*\|_{2,\infty} \left\| \mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)} - X^* \right\|_{2,\infty} \\ &\leq 3 \left\| \mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)} - X^* \right\|_{2,\infty} \|X^*\|_{2,\infty}, \end{aligned} \tag{158}$$

where the last line arises from (155). This combined with (157) gives

$$\begin{aligned} \|B_2^{(l)}\|_F &\lesssim \eta \sqrt{\frac{\log n}{p}} \|C\|_\infty \|X^*\|_F + \eta \sqrt{\frac{n}{p}} \|C\|_\infty \|X^*\|_{2,\infty} \\ &\stackrel{(i)}{\lesssim} \eta \sqrt{\frac{\log n}{p}} \left\| \mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)} - X^* \right\|_{2,\infty} \|X^*\|_{2,\infty} \|X^*\|_F \\ &\quad + \eta \sqrt{\frac{n}{p}} \left\| \mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)} - X^* \right\|_{2,\infty} \|X^*\|_{2,\infty}^2 \\ &\stackrel{(ii)}{\lesssim} \eta \sqrt{\frac{\log n}{p}} \left\| \mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)} - X^* \right\|_{2,\infty} \sqrt{\frac{\kappa \mu r^2}{n}} \sigma_{\max} \\ &\quad + \eta \sqrt{\frac{n}{p}} \left\| \mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)} - X^* \right\|_{2,\infty} \frac{\kappa \mu r}{n} \sigma_{\max} \\ &\lesssim \eta \sqrt{\frac{\kappa^2 \mu^2 r^2 \log n}{np}} \left\| \mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)} - X^* \right\|_{2,\infty} \sigma_{\max}, \end{aligned}$$

where (i) comes from (158) and (ii) makes use of the incoherence condition (114).

**B.6 Proof of Lemma 12**

We first introduce an auxiliary matrix

$$\begin{aligned} \tilde{X}^{t+1,(l)} := & X^{t,(l)} \widehat{H}^{t,(l)} - \eta \left[ \frac{1}{p} \mathcal{P}_{\Omega^c} \left[ X^{t,(l)} X^{t,(l)\top} - (M^* + E) \right] \right. \\ & \left. + \mathcal{P}_l \left( X^{t,(l)} X^{t,(l)\top} - M^* \right) \right] X^*. \end{aligned} \tag{159}$$

With this in place, we can use the triangle inequality to obtain

$$\begin{aligned} \left\| (X^{t+1,(l)} \widehat{H}^{t+1,(l)} - X^*)_{l,\cdot} \right\|_2 &\leq \underbrace{\left\| (X^{t+1,(l)} \widehat{H}^{t+1,(l)} - \tilde{X}^{t+1,(l)})_{l,\cdot} \right\|_2}_{:=\alpha_1} \\ &+ \underbrace{\left\| (\tilde{X}^{t+1,(l)} - X^*)_{l,\cdot} \right\|_2}_{:=\alpha_2}. \end{aligned} \tag{160}$$

In what follows, we bound the two terms  $\alpha_1$  and  $\alpha_2$  separately.

1. Regarding the second term  $\alpha_2$  of (160), we see from the definition of  $\tilde{X}^{t+1,(l)}$  (see (159)) that

$$\left( \tilde{X}^{t+1,(l)} - X^* \right)_{l,\cdot} = \left[ X^{t,(l)} \widehat{H}^{t,(l)} - \eta (X^{t,(l)} X^{t,(l)\top} - X^* X^{*\top}) X^* - X^* \right]_{l,\cdot}, \tag{161}$$

where we also utilize the definitions of  $\mathcal{P}_{\Omega^c}$  and  $\mathcal{P}_l$  in (67). For notational convenience, we denote

$$\Delta^{t,(l)} := X^{t,(l)} \widehat{H}^{t,(l)} - X^*. \tag{162}$$

This allows us to rewrite (161) as

$$\begin{aligned} \left( \tilde{X}^{t+1,(l)} - X^* \right)_{l,\cdot} &= \Delta_{l,\cdot}^{t,(l)} - \eta \left[ \left( \Delta^{t,(l)} X^{*\top} + X^* \Delta^{t,(l)\top} \right) X^* \right]_{l,\cdot} \\ &\quad - \eta \left[ \Delta^{t,(l)} \Delta^{t,(l)\top} X^* \right]_{l,\cdot} \\ &= \Delta_{l,\cdot}^{t,(l)} - \eta \Delta_{l,\cdot}^{t,(l)} \Sigma^* - \eta X_{l,\cdot}^* \Delta^{t,(l)\top} X^* - \eta \Delta_{l,\cdot}^{t,(l)} \Delta^{t,(l)\top} X^*, \end{aligned}$$

which further implies that

$$\begin{aligned} \alpha_2 &\leq \left\| \Delta_{l,\cdot}^{t,(l)} - \eta \Delta_{l,\cdot}^{t,(l)} \Sigma^* \right\|_2 + \eta \left\| X_{l,\cdot}^* \Delta^{t,(l)\top} X^* \right\|_2 + \eta \left\| \Delta_{l,\cdot}^{t,(l)} \Delta^{t,(l)\top} X^* \right\|_2 \\ &\leq \left\| \Delta_{l,\cdot}^{t,(l)} \right\|_2 \left\| I_r - \eta \Sigma^* \right\| + \eta \left\| X^* \right\|_{2,\infty} \left\| \Delta^{t,(l)} \right\| \left\| X^* \right\| + \eta \left\| \Delta_{l,\cdot}^{t,(l)} \right\|_2 \left\| \Delta^{t,(l)} \right\| \left\| X^* \right\| \\ &\leq \left\| \Delta_{l,\cdot}^{t,(l)} \right\|_2 \left\| I_r - \eta \Sigma^* \right\| + 2\eta \left\| X^* \right\|_{2,\infty} \left\| \Delta^{t,(l)} \right\| \left\| X^* \right\|. \end{aligned}$$

Here, the last line follows from the fact that  $\|\Delta_{l,\cdot}^{t,(l)}\|_2 \leq \|X^*\|_{2,\infty}$ . To see this, one can use the induction hypothesis (70e) to get

$$\|\Delta_{l,\cdot}^{t,(l)}\|_2 \leq C_2 \rho^t \mu r \frac{1}{\sqrt{np}} \|X^*\|_{2,\infty} + C_6 \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} \|X^*\|_{2,\infty} \ll \|X^*\|_{2,\infty} \tag{163}$$

as long as  $np \gg \mu^2 r^2$  and  $\sigma \sqrt{(n \log n)/p} \ll \sigma_{\min}$ . By taking  $0 < \eta \leq 1/\sigma_{\max}$ , we have  $\mathbf{0} \leq I_r - \eta \Sigma^* \leq (1 - \eta \sigma_{\min}) I_r$  and hence can obtain

$$\alpha_2 \leq (1 - \eta \sigma_{\min}) \|\Delta_{l,\cdot}^{t,(l)}\|_2 + 2\eta \|X^*\|_{2,\infty} \|\Delta_{l,\cdot}^{t,(l)}\| \|X^*\|. \tag{164}$$

An immediate consequence of the above two inequalities and (73d) is

$$\alpha_2 \leq \|X^*\|_{2,\infty}. \tag{165}$$

2. The first term  $\alpha_1$  of (160) can be equivalently written as

$$\alpha_1 = \left\| (X^{t+1,(l)} \widehat{H}^{t,(l)} R_1 - \widetilde{X}^{t+1,(l)})_{l,\cdot} \right\|_2,$$

where

$$R_1 = (\widehat{H}^{t,(l)})^{-1} \widehat{H}^{t+1,(l)} := \arg \min_{R \in \mathcal{O}^{r \times r}} \|X^{t+1,(l)} \widehat{H}^{t,(l)} R - X^*\|_F.$$

Simple algebra yields

$$\begin{aligned} \alpha_1 &\leq \left\| (X^{t+1,(l)} \widehat{H}^{t,(l)} - \widetilde{X}^{t+1,(l)})_{l,\cdot} R_1 \right\|_2 + \|\widetilde{X}^{t+1,(l)}\|_2 \|R_1 - I_r\| \\ &\leq \underbrace{\left\| (X^{t+1,(l)} \widehat{H}^{t,(l)} - \widetilde{X}^{t+1,(l)})_{l,\cdot} \right\|_2}_{:=\beta_1} + 2 \|X^*\|_{2,\infty} \underbrace{\|R_1 - I_r\|}_{:=\beta_2}. \end{aligned}$$

Here, to bound the second term, we have used

$$\|\widetilde{X}_{l,\cdot}^{t+1,(l)}\|_2 \leq \|\widetilde{X}_{l,\cdot}^{t+1,(l)} - X_{l,\cdot}^*\|_2 + \|X_{l,\cdot}^*\|_2 = \alpha_2 + \|X_{l,\cdot}^*\|_2 \leq 2 \|X^*\|_{2,\infty},$$

where the last inequality follows from (165). It remains to upper bound  $\beta_1$  and  $\beta_2$ . For both  $\beta_1$  and  $\beta_2$ , a central quantity to control is  $X^{t+1,(l)} \widehat{H}^{t,(l)} - \widetilde{X}^{t+1,(l)}$ . By the definition of  $\widetilde{X}^{t+1,(l)}$  in (159) and the gradient update rule for  $X^{t+1,(l)}$  (see (69)), one has

$$\begin{aligned} &X^{t+1,(l)} \widehat{H}^{t,(l)} - \widetilde{X}^{t+1,(l)} \\ &= \left\{ X^{t,(l)} \widehat{H}^{t,(l)} - \eta \left[ \frac{1}{p} \mathcal{P}_{\Omega^t} \left[ X^{t,(l)} X^{t,(l)\top} - (M^* + E) \right] \right. \right. \end{aligned}$$



$$\begin{aligned}
 & + \mathcal{P}_l \left( \mathbf{X}^{t,(l)} \mathbf{X}^{t,(l)\top} - \mathbf{M}^* \right) \mathbf{X}^{t,(l)} \widehat{\mathbf{H}}^{t,(l)} \} \\
 & - \left\{ \mathbf{X}^{t,(l)} \widehat{\mathbf{H}}^{t,(l)} - \eta \left[ \frac{1}{p} \mathcal{P}_{\Omega^{-l}} \left[ \mathbf{X}^{t,(l)} \mathbf{X}^{t,(l)\top} - (\mathbf{M}^* + \mathbf{E}) \right] + \mathcal{P}_l \left( \mathbf{X}^{t,(l)} \mathbf{X}^{t,(l)\top} - \mathbf{M}^* \right) \right] \mathbf{X}^* \right\} \\
 = & -\eta \left[ \frac{1}{p} \mathcal{P}_{\Omega^{-l}} \left( \mathbf{X}^{t,(l)} \mathbf{X}^{t,(l)\top} - \mathbf{X}^* \mathbf{X}^{*\top} \right) + \mathcal{P}_l \left( \mathbf{X}^{t,(l)} \mathbf{X}^{t,(l)\top} - \mathbf{X}^* \mathbf{X}^{*\top} \right) \right] \Delta^{t,(l)} \\
 & + \frac{\eta}{p} \mathcal{P}_{\Omega^{-l}} (\mathbf{E}) \Delta^{t,(l)}. \tag{166}
 \end{aligned}$$

It is easy to verify that

$$\frac{1}{p} \left\| \mathcal{P}_{\Omega^{-l}} (\mathbf{E}) \right\| \stackrel{(i)}{\leq} \frac{1}{p} \left\| \mathcal{P}_{\Omega} (\mathbf{E}) \right\| \stackrel{(ii)}{\lesssim} \sigma \sqrt{\frac{n}{p}} \stackrel{(iii)}{\leq} \frac{\delta}{2} \sigma_{\min}$$

for  $\delta > 0$  sufficiently small. Here, (i) uses the elementary fact that the spectral norm of a submatrix is no more than that of the matrix itself, (ii) arises from Lemma 40, and (iii) is a consequence of the noise condition (27). Therefore, in order to control (166), we need to upper bound the following quantity

$$\gamma := \left\| \frac{1}{p} \mathcal{P}_{\Omega^{-l}} \left( \mathbf{X}^{t,(l)} \mathbf{X}^{t,(l)\top} - \mathbf{X}^* \mathbf{X}^{*\top} \right) + \mathcal{P}_l \left( \mathbf{X}^{t,(l)} \mathbf{X}^{t,(l)\top} - \mathbf{X}^* \mathbf{X}^{*\top} \right) \right\|. \tag{167}$$

To this end, we make the observation that

$$\begin{aligned}
 \gamma & \leq \underbrace{\left\| \frac{1}{p} \mathcal{P}_{\Omega} \left( \mathbf{X}^{t,(l)} \mathbf{X}^{t,(l)\top} - \mathbf{X}^* \mathbf{X}^{*\top} \right) \right\|}_{:=\gamma_1} \\
 & + \underbrace{\left\| \frac{1}{p} \mathcal{P}_{\Omega_l} \left( \mathbf{X}^{t,(l)} \mathbf{X}^{t,(l)\top} - \mathbf{X}^* \mathbf{X}^{*\top} \right) - \mathcal{P}_l \left( \mathbf{X}^{t,(l)} \mathbf{X}^{t,(l)\top} - \mathbf{X}^* \mathbf{X}^{*\top} \right) \right\|}_{:=\gamma_2}, \tag{168}
 \end{aligned}$$

where  $\mathcal{P}_{\Omega_l}$  is defined in (66). An application of Lemma 43 reveals that

$$\gamma_1 \leq 2n \left\| \mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)} - \mathbf{X}^* \right\|_{2,\infty}^2 + 4\sqrt{n} \log n \left\| \mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)} - \mathbf{X}^* \right\|_{2,\infty} \left\| \mathbf{X}^* \right\|,$$

where  $\mathbf{R}^{t,(l)} \in \mathcal{O}^{r \times r}$  is defined in (72). Let  $\mathbf{C} = \mathbf{X}^{t,(l)} \mathbf{X}^{t,(l)\top} - \mathbf{X}^* \mathbf{X}^{*\top}$  as in (156), and one can bound the other term  $\gamma_2$  by taking advantage of the triangle inequality and the symmetry property:

$$\gamma_2 \leq \frac{2}{p} \sqrt{\sum_{j=1}^n (\delta_{l,j} - p)^2 C_{l,j}^2} \stackrel{(i)}{\lesssim} \sqrt{\frac{n}{p}} \left\| \mathbf{C} \right\|_{\infty} \stackrel{(ii)}{\lesssim} \sqrt{\frac{n}{p}} \left\| \mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)} - \mathbf{X}^* \right\|_{2,\infty} \left\| \mathbf{X}^* \right\|_{2,\infty},$$

where (i) comes from the standard Chernoff bound  $\sum_{j=1}^n (\delta_{l,j} - p)^2 \asymp np$ , and in (ii) we utilize the bound established in (158). The previous two bounds taken collectively give

$$\begin{aligned} \gamma &\leq 2n \left\| \mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)} - \mathbf{X}^\star \right\|_{2,\infty}^2 + 4\sqrt{n} \log n \left\| \mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)} - \mathbf{X}^\star \right\|_{2,\infty} \|\mathbf{X}^\star\| \\ &\quad + \tilde{C} \sqrt{\frac{n}{p}} \left\| \mathbf{X}^{t,(l)} \mathbf{R}^{t,(l)} - \mathbf{X}^\star \right\|_{2,\infty} \|\mathbf{X}^\star\|_{2,\infty} \leq \frac{\delta}{2} \sigma_{\min} \end{aligned} \tag{169}$$

for some constant  $\tilde{C} > 0$  and  $\delta > 0$  sufficiently small. The last inequality follows from (73c), the incoherence condition (114), and our sample size condition. In summary, we obtain

$$\left\| \mathbf{X}^{t+1,(l)} \widehat{\mathbf{H}}^{t,(l)} - \tilde{\mathbf{X}}^{t+1,(l)} \right\| \leq \eta \left( \gamma + \left\| \frac{1}{p} \mathcal{P}_{\Omega^{-l}}(\mathbf{E}) \right\| \right) \|\Delta^{t,(l)}\| \leq \eta \delta \sigma_{\min} \|\Delta^{t,(l)}\|, \tag{170}$$

for  $\delta > 0$  sufficiently small. With the estimate (170) in place, we can continue our derivation on  $\beta_1$  and  $\beta_2$ .

(a) With regard to  $\beta_1$ , in view of (166) we can obtain

$$\begin{aligned} \beta_1 &\stackrel{(i)}{=} \eta \left\| \left( \mathbf{X}^{t,(l)} \mathbf{X}^{t,(l)\top} - \mathbf{X}^\star \mathbf{X}^{\star\top} \right)_{l,\cdot} \Delta^{t,(l)} \right\|_2 \\ &\leq \eta \left\| \left( \mathbf{X}^{t,(l)} \mathbf{X}^{t,(l)\top} - \mathbf{X}^\star \mathbf{X}^{\star\top} \right)_{l,\cdot} \right\|_2 \|\Delta^{t,(l)}\| \\ &\stackrel{(ii)}{=} \eta \left\| \left[ \Delta^{t,(l)} \left( \mathbf{X}^{t,(l)} \widehat{\mathbf{H}}^{t,(l)} \right)^\top + \mathbf{X}^\star \Delta^{t,(l)\top} \right]_{l,\cdot} \right\|_2 \|\Delta^{t,(l)}\| \\ &\leq \eta \left( \|\Delta_{l,\cdot}^{t,(l)}\|_2 \|\mathbf{X}^{t,(l)}\| + \|\mathbf{X}_{l,\cdot}^\star\|_2 \|\Delta^{t,(l)}\| \right) \|\Delta^{t,(l)}\| \\ &\leq \eta \|\Delta_{l,\cdot}^{t,(l)}\|_2 \|\mathbf{X}^{t,(l)}\| \|\Delta^{t,(l)}\| + \eta \|\mathbf{X}_{l,\cdot}^\star\|_2 \|\Delta^{t,(l)}\|^2, \end{aligned} \tag{171}$$

where (i) follows from the definitions of  $\mathcal{P}_{\Omega^{-l}}$  and  $\mathcal{P}_l$  (see (67) and note that all entries in the  $l$ th row of  $\mathcal{P}_{\Omega^{-l}}(\cdot)$  are identically zero) and identity (ii) is due to the definition of  $\Delta^{t,(l)}$  in (162).

(b) For  $\beta_2$ , we first claim that

$$\mathbf{I}_r := \arg \min_{\mathbf{R} \in \mathcal{O}^{r \times r}} \left\| \tilde{\mathbf{X}}^{t+1,(l)} \mathbf{R} - \mathbf{X}^\star \right\|_{\mathbf{F}}, \tag{172}$$

whose justification follows similar reasonings as that of (138) and is therefore omitted. In particular, it gives rise to the facts that  $\mathbf{X}^{\star\top} \tilde{\mathbf{X}}^{t+1,(l)}$  is symmetric and

$$\left( \tilde{\mathbf{X}}^{t+1,(l)} \right)^\top \mathbf{X}^\star \succeq \frac{1}{2} \sigma_{\min} \mathbf{I}_r. \tag{173}$$

We are now ready to invoke Lemma 36 to bound  $\beta_2$ . We abuse the notation and denote  $\mathbf{C} := (\tilde{\mathbf{X}}^{t+1,(l)})^\top \mathbf{X}^\star$  and  $\mathbf{E} := (\mathbf{X}^{t+1,(l)} \widehat{\mathbf{H}}^{t,(l)} - \tilde{\mathbf{X}}^{t+1,(l)})^\top \mathbf{X}^\star$ . We have

$$\|\mathbf{E}\| \leq \frac{1}{2} \sigma_{\min} \leq \sigma_r(\mathbf{C}).$$

The first inequality arises from (170), namely

$$\begin{aligned} \|\mathbf{E}\| &\leq \left\| \mathbf{X}^{t+1,(l)} \widehat{\mathbf{H}}^{t,(l)} - \tilde{\mathbf{X}}^{t+1,(l)} \right\| \|\mathbf{X}^\star\| \leq \eta \delta \sigma_{\min} \left\| \Delta^{t,(l)} \right\| \|\mathbf{X}^\star\| \\ &\stackrel{(i)}{\leq} \eta \delta \sigma_{\min} \|\mathbf{X}^\star\|^2 \stackrel{(ii)}{\leq} \frac{1}{2} \sigma_{\min}, \end{aligned}$$

where (i) holds since  $\|\Delta^{t,(l)}\| \leq \|\mathbf{X}^\star\|$  and (ii) holds true for  $\delta$  sufficiently small and  $\eta \leq 1/\sigma_{\max}$ . Invoke Lemma 36 to obtain

$$\begin{aligned} \beta_2 = \|\mathbf{R}_1 - \mathbf{I}_r\| &\leq \frac{2}{\sigma_{r-1}(\mathbf{C}) + \sigma_r(\mathbf{C})} \|\mathbf{E}\| \\ &\leq \frac{2}{\sigma_{\min}} \left\| \mathbf{X}^{t+1,(l)} \widehat{\mathbf{H}}^{t,(l)} - \tilde{\mathbf{X}}^{t+1,(l)} \right\| \|\mathbf{X}^\star\| \quad (174) \end{aligned}$$

$$\leq 2\delta\eta \left\| \Delta^{t,(l)} \right\| \|\mathbf{X}^\star\|, \quad (175)$$

where (174) follows since  $\sigma_{r-1}(\mathbf{C}) \geq \sigma_r(\mathbf{C}) \geq \sigma_{\min}/2$  from (173), and the last line comes from (170).

(c) Putting the previous bounds (171) and (175) together yields

$$\begin{aligned} \alpha_1 &\leq \eta \left\| \Delta_{l,\cdot}^{t,(l)} \right\|_2 \left\| \mathbf{X}^{t,(l)} \right\| \left\| \Delta^{t,(l)} \right\| + \eta \left\| \mathbf{X}_{l,\cdot}^\star \right\|_2 \left\| \Delta^{t,(l)} \right\|^2 \\ &\quad + 4\delta\eta \|\mathbf{X}^\star\|_{2,\infty} \left\| \Delta^{t,(l)} \right\| \|\mathbf{X}^\star\|. \quad (176) \end{aligned}$$

3. Combine (160), (164), and (176) to reach

$$\begin{aligned} \left\| (\mathbf{X}^{t+1,(l)} \widehat{\mathbf{H}}^{t+1,(l)} - \mathbf{X}^\star)_{l,\cdot} \right\|_2 &\leq (1 - \eta \sigma_{\min}) \left\| \Delta_{l,\cdot}^{t,(l)} \right\|_2 + 2\eta \|\mathbf{X}^\star\|_{2,\infty} \left\| \Delta^{t,(l)} \right\| \|\mathbf{X}^\star\| \\ &\quad + \eta \left\| \Delta_{l,\cdot}^{t,(l)} \right\|_2 \left\| \mathbf{X}^{t,(l)} \right\| \left\| \Delta^{t,(l)} \right\| + \eta \left\| \mathbf{X}_{l,\cdot}^\star \right\|_2 \left\| \Delta^{t,(l)} \right\|^2 + 4\delta\eta \|\mathbf{X}^\star\|_{2,\infty} \left\| \Delta^{t,(l)} \right\| \|\mathbf{X}^\star\| \\ &\stackrel{(i)}{\leq} \left( 1 - \eta \sigma_{\min} + \eta \left\| \mathbf{X}^{t,(l)} \right\| \left\| \Delta^{t,(l)} \right\| \right) \left\| \Delta_{l,\cdot}^{t,(l)} \right\|_2 + 4\eta \|\mathbf{X}^\star\|_{2,\infty} \left\| \Delta^{t,(l)} \right\| \|\mathbf{X}^\star\| \\ &\stackrel{(ii)}{\leq} \left( 1 - \frac{\sigma_{\min}}{2} \eta \right) \left( C_2 \rho^t \mu r \frac{1}{\sqrt{np}} + \frac{C_6}{\sigma_{\min}} \sigma \sqrt{\frac{n \log n}{p}} \right) \|\mathbf{X}^\star\|_{2,\infty} \\ &\quad + 4\eta \|\mathbf{X}^\star\| \|\mathbf{X}^\star\|_{2,\infty} \left( 2C_9 \rho^t \mu r \frac{1}{\sqrt{np}} \|\mathbf{X}^\star\| + \frac{2C_{10}}{\sigma_{\min}} \sigma \sqrt{\frac{n}{p}} \|\mathbf{X}^\star\| \right) \\ &\stackrel{(iii)}{\leq} C_2 \rho^{t+1} \mu r \frac{1}{\sqrt{np}} \|\mathbf{X}^\star\|_{2,\infty} + \frac{C_6}{\sigma_{\min}} \sigma \sqrt{\frac{n \log n}{p}} \|\mathbf{X}^\star\|_{2,\infty}. \end{aligned}$$

Here, (i) follows since  $\|\Delta^{t,(l)}\| \leq \|X^*\|$  and  $\delta$  is sufficiently small, (ii) invokes hypotheses (70e) and (73d) and recognizes that

$$\|X^{t,(l)}\| \|\Delta^{t,(l)}\| \leq 2 \|X^*\| \left( 2C_9\mu r \frac{1}{\sqrt{np}} \|X^*\| + \frac{2C_{10}}{\sigma_{\min}} \sigma \sqrt{\frac{n \log n}{np}} \|X^*\| \right) \leq \frac{\sigma_{\min}}{2}$$

holds under the sample size and noise condition, while (iii) is valid as long as  $1 - (\sigma_{\min}/3) \cdot \eta \leq \rho < 1$ ,  $C_2 \gg \kappa C_9$ , and  $C_6 \gg \kappa C_{10}/\sqrt{\log n}$ .

**B.7 Proof of Lemma 13**

For notational convenience, we define the following two orthonormal matrices

$$Q := \arg \min_{R \in \mathcal{O}^{r \times r}} \|U^0 R - U^*\|_F \quad \text{and} \quad Q^{(l)} := \arg \min_{R \in \mathcal{O}^{r \times r}} \|U^{0,(l)} R - U^*\|_F.$$

The problem of finding  $\hat{H}^t$  (see (26)) is called the *orthogonal Procrustes problem* [112]. It is well known that the minimizer  $\hat{H}^t$  always exists and is given by

$$\hat{H}^t = \text{sgn}(X^{t\top} X^*).$$

Here, the sign matrix  $\text{sgn}(B)$  is defined as

$$\text{sgn}(B) := UV^\top \tag{177}$$

for any matrix  $B$  with singular value decomposition  $B = U\Sigma V^\top$ , where the columns of  $U$  and  $V$  are left and right singular vectors, respectively.

Before proceeding, we make note of the following perturbation bounds on  $M^0$  and  $M^{(l)}$  (as defined in Algorithms 2 and 5, respectively):

$$\begin{aligned} \|M^0 - M^*\| &\stackrel{(i)}{\leq} \left\| \frac{1}{p} \mathcal{P}_\Omega(M^*) - M^* \right\| + \left\| \frac{1}{p} \mathcal{P}_\Omega(E) \right\| \\ &\stackrel{(ii)}{\leq} C \sqrt{\frac{n}{p}} \|M^*\|_\infty + C\sigma \sqrt{\frac{n}{p}} = C \sqrt{\frac{n}{p}} \|X^*\|_{2,\infty}^2 + C \frac{\sigma}{\sqrt{\sigma_{\min}}} \sqrt{\frac{n}{p}} \sqrt{\sigma_{\min}} \\ &\stackrel{(iii)}{\leq} C \left\{ \mu r \sqrt{\frac{1}{np}} \sqrt{\sigma_{\max}} + \frac{\sigma}{\sqrt{\sigma_{\min}}} \sqrt{\frac{n}{p}} \right\} \|X^*\| \stackrel{(iv)}{\ll} \sigma_{\min}, \end{aligned} \tag{178}$$

for some universal constant  $C > 0$ . Here, (i) arises from the triangle inequality, (ii) utilizes Lemmas 39 and 40, (iii) follows from the incoherence condition (114), and (iv) holds under our sample complexity assumption that  $n^2 p \gg \mu^2 r^2 n$  and the noise

condition (27). Similarly, we have

$$\|M^{(l)} - M^*\| \lesssim \left\{ \mu r \sqrt{\frac{1}{np}} \sqrt{\sigma_{\max}} + \frac{\sigma}{\sqrt{\sigma_{\min}}} \sqrt{\frac{n}{p}} \right\} \|X^*\| \ll \sigma_{\min}. \tag{179}$$

Combine Weyl’s inequality, (178), and (179) to obtain

$$\|\Sigma^0 - \Sigma^*\| \leq \|M^0 - M^*\| \ll \sigma_{\min} \quad \text{and} \quad \|\Sigma^{(l)} - \Sigma^*\| \leq \|M^{(l)} - M^*\| \ll \sigma_{\min}, \tag{180}$$

which further implies

$$\frac{1}{2} \sigma_{\min} \leq \sigma_r(\Sigma^0) \leq \sigma_1(\Sigma^0) \leq 2\sigma_{\max} \quad \text{and} \quad \frac{1}{2} \sigma_{\min} \leq \sigma_r(\Sigma^{(l)}) \leq \sigma_1(\Sigma^{(l)}) \leq 2\sigma_{\max}. \tag{181}$$

We start by proving (70a), (70b), and (70c). The key decomposition we need is the following

$$\begin{aligned} X^0 \widehat{H}^0 - X^* &= U^0 (\Sigma^0)^{1/2} (\widehat{H}^0 - Q) + U^0 \left[ (\Sigma^0)^{1/2} Q - Q (\Sigma^*)^{1/2} \right] \\ &\quad + (U^0 Q - U^*) (\Sigma^*)^{1/2}. \end{aligned} \tag{182}$$

1. For the spectral norm error bound in (70c), the triangle inequality together with (182) yields

$$\begin{aligned} \|X^0 \widehat{H}^0 - X^*\| &\leq \|(\Sigma^0)^{1/2}\| \|\widehat{H}^0 - Q\| + \|(\Sigma^0)^{1/2} Q - Q (\Sigma^*)^{1/2}\| \\ &\quad + \sqrt{\sigma_{\max}} \|U^0 Q - U^*\|, \end{aligned}$$

where we have also used the fact that  $\|U^0\| = 1$ . Recognizing that  $\|M^0 - M^*\| \ll \sigma_{\min}$  (see (178)) and the assumption  $\sigma_{\max}/\sigma_{\min} \lesssim 1$ , we can apply Lemmas 47, 46, and 45 to obtain

$$\|\widehat{H}^0 - Q\| \lesssim \frac{1}{\sigma_{\min}} \|M^0 - M^*\|, \tag{183a}$$

$$\|(\Sigma^0)^{1/2} Q - Q (\Sigma^*)^{1/2}\| \lesssim \frac{1}{\sqrt{\sigma_{\min}}} \|M^0 - M^*\|, \tag{183b}$$

$$\|U^0 Q - U^*\| \lesssim \frac{1}{\sigma_{\min}} \|M^0 - M^*\|. \tag{183c}$$

These taken collectively imply the advertised upper bound

$$\begin{aligned} \|X^0 \widehat{H}^0 - X^*\| &\lesssim \sqrt{\sigma_{\max}} \frac{1}{\sigma_{\min}} \|M^0 - M^*\| + \frac{1}{\sqrt{\sigma_{\min}}} \|M^0 - M^*\| \\ &\lesssim \frac{1}{\sqrt{\sigma_{\min}}} \|M^0 - M^*\| \\ &\lesssim \left\{ \mu r \sqrt{\frac{1}{np}} \sqrt{\frac{\sigma_{\max}}{\sigma_{\min}}} + \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \right\} \|X^*\|, \end{aligned}$$

where we also utilize the fact that  $\|(\Sigma^0)^{1/2}\| \leq \sqrt{2\sigma_{\max}}$  (see (181)) and the bounded condition number assumption, i.e.,  $\sigma_{\max}/\sigma_{\min} \lesssim 1$ . This finishes the proof of (70c).

2. With regard to the Frobenius norm bound in (70a), one has

$$\begin{aligned} \|X^0 \widehat{H}^0 - X^*\|_F &\leq \sqrt{r} \|X^0 \widehat{H}^0 - X^*\| \\ &\stackrel{(i)}{\lesssim} \left\{ \mu r \sqrt{\frac{1}{np}} + \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \right\} \sqrt{r} \|X^*\| \\ &= \left\{ \mu r \sqrt{\frac{1}{np}} + \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \right\} \sqrt{r} \frac{\sqrt{\sigma_{\max}}}{\sqrt{\sigma_{\min}}} \sqrt{\sigma_{\min}} \\ &\stackrel{(ii)}{\lesssim} \left\{ \mu r \sqrt{\frac{1}{np}} + \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \right\} \sqrt{r} \|X^*\|_F. \end{aligned}$$

Here (i) arises from (70c) and (ii) holds true since  $\sigma_{\max}/\sigma_{\min} \asymp 1$  and  $\sqrt{r} \sqrt{\sigma_{\min}} \leq \|X^*\|_F$ , thus completing the proof of (70a).

3. The proof of (70b) follows from similar arguments as used in proving (70c). Combine (182) and the triangle inequality to reach

$$\begin{aligned} &\|X^0 \widehat{H}^0 - X^*\|_{2,\infty} \\ &\leq \|U^0\|_{2,\infty} \left\{ \|(\Sigma^0)^{1/2}\| \|\widehat{H}^0 - Q\| + \|(\Sigma^0)^{1/2} Q - Q(\Sigma^*)^{1/2}\| \right\} \\ &\quad + \sqrt{\sigma_{\max}} \|U^0 Q - U^*\|_{2,\infty}. \end{aligned}$$

Plugging in the estimates (178), (181), (183a), and (183b) results in

$$\begin{aligned} \|X^0 \widehat{H}^0 - X^*\|_{2,\infty} &\lesssim \left\{ \mu r \sqrt{\frac{1}{np}} + \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \right\} \|X^*\| \|U^0\|_{2,\infty} \\ &\quad + \sqrt{\sigma_{\max}} \|U^0 Q - U^*\|_{2,\infty}. \end{aligned}$$

It remains to study the component-wise error of  $U^0$ . To this end, it has already been shown in [1, Lemma 14] that

$$\|U^0 Q - U^*\|_{2,\infty} \lesssim \left( \mu r \sqrt{\frac{1}{np}} + \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \right) \|U^*\|_{2,\infty} \quad \text{and} \quad \|U^0\|_{2,\infty} \lesssim \|U^*\|_{2,\infty} \tag{184}$$

under our assumptions. These combined with the previous inequality give

$$\begin{aligned} \|X^0 \widehat{H}^0 - X^*\|_{2,\infty} &\lesssim \left\{ \mu r \sqrt{\frac{1}{np}} + \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \right\} \sqrt{\sigma_{\max}} \|U^*\|_{2,\infty} \\ &\lesssim \left\{ \mu r \sqrt{\frac{1}{np}} + \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \right\} \|X^*\|_{2,\infty}, \end{aligned}$$

where the last relation is due to the observation that

$$\sqrt{\sigma_{\max}} \|U^*\|_{2,\infty} \lesssim \sqrt{\sigma_{\min}} \|U^*\|_{2,\infty} \leq \|X^*\|_{2,\infty}.$$

4. We now move on to proving (70e). Recall that  $Q^{(l)} = \arg \min_{R \in \mathcal{O}^{r \times r}} \|U^{0,(l)} R - U^*\|_F$ . By the triangle inequality,

$$\begin{aligned} \|(X^{0,(l)} \widehat{H}^{0,(l)} - X^*)_{l,\cdot}\|_2 &\leq \|X_{l,\cdot}^{0,(l)} (\widehat{H}^{0,(l)} - Q^{(l)})\|_2 + \|(X^{0,(l)} Q^{(l)} - X^*)_{l,\cdot}\|_2 \\ &\leq \|X_{l,\cdot}^{0,(l)}\|_2 \|\widehat{H}^{0,(l)} - Q^{(l)}\| + \|(X^{0,(l)} Q^{(l)} - X^*)_{l,\cdot}\|_2. \end{aligned} \tag{185}$$

Note that  $X_{l,\cdot}^* = M_{l,\cdot}^* (U^*)^{-1/2}$  and, by construction of  $M^{(l)}$ ,

$$X_{l,\cdot}^{0,(l)} = M_{l,\cdot}^{(l)} U^{0,(l)} (\Sigma^{(l)})^{-1/2} = M_{l,\cdot}^* U^{0,(l)} (\Sigma^{(l)})^{-1/2}.$$

We can thus decompose

$$\begin{aligned} (X^{0,(l)} Q^{(l)} - X^*)_{l,\cdot} &= M_{l,\cdot}^* \left\{ U^{0,(l)} \left[ (\Sigma^{(l)})^{-1/2} Q^{(l)} - Q^{(l)} (\Sigma^*)^{-1/2} \right] \right. \\ &\quad \left. + (U^{0,(l)} Q^{(l)} - U^*) (\Sigma^*)^{-1/2} \right\}, \end{aligned}$$

which further implies that

$$\begin{aligned} \|(X^{0,(l)} Q^{(l)} - X^*)_{l,\cdot}\|_2 &\leq \|M^*\|_{2,\infty} \left\{ \left\| (\Sigma^{(l)})^{-1/2} Q^{(l)} - Q^{(l)} (\Sigma^*)^{-1/2} \right\| \right. \\ &\quad \left. + \frac{1}{\sqrt{\sigma_{\min}}} \|U^{0,(l)} Q^{(l)} - U^*\| \right\}. \end{aligned} \tag{186}$$

In order to control this, we first see that

$$\begin{aligned} & \left\| (\boldsymbol{\Sigma}^{(l)})^{-1/2} \boldsymbol{Q}^{(l)} - \boldsymbol{Q}^{(l)} (\boldsymbol{\Sigma}^*)^{-1/2} \right\| \\ &= \left\| (\boldsymbol{\Sigma}^{(l)})^{-1/2} \left[ \boldsymbol{Q}^{(l)} (\boldsymbol{\Sigma}^*)^{1/2} - (\boldsymbol{\Sigma}^{(l)})^{1/2} \boldsymbol{Q}^{(l)} \right] (\boldsymbol{\Sigma}^*)^{-1/2} \right\| \\ &\lesssim \frac{1}{\sigma_{\min}} \left\| \boldsymbol{Q}^{(l)} (\boldsymbol{\Sigma}^*)^{1/2} - (\boldsymbol{\Sigma}^{(l)})^{1/2} \boldsymbol{Q}^{(l)} \right\| \\ &\lesssim \frac{1}{\sigma_{\min}^{3/2}} \left\| \boldsymbol{M}^{(l)} - \boldsymbol{M}^* \right\|, \end{aligned}$$

where the penultimate inequality uses (181) and the last inequality arises from Lemma 46. Additionally, Lemma 45 gives

$$\left\| \boldsymbol{U}^{0,(l)} \boldsymbol{Q}^{(l)} - \boldsymbol{U}^* \right\| \lesssim \frac{1}{\sigma_{\min}} \left\| \boldsymbol{M}^{(l)} - \boldsymbol{M}^* \right\|.$$

Plugging the previous two bounds into (186), we reach

$$\begin{aligned} \left\| (\boldsymbol{X}^{0,(l)} \boldsymbol{Q}^{(l)} - \boldsymbol{X}^*)_{l,\cdot} \right\|_2 &\lesssim \frac{1}{\sigma_{\min}^{3/2}} \left\| \boldsymbol{M}^{(l)} - \boldsymbol{M}^* \right\| \left\| \boldsymbol{M}^* \right\|_{2,\infty} \\ &\lesssim \left\{ \mu r \sqrt{\frac{1}{np}} + \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \right\} \left\| \boldsymbol{X}^* \right\|_{2,\infty}. \end{aligned}$$

where the last relation follows from  $\left\| \boldsymbol{M}^* \right\|_{2,\infty} = \left\| \boldsymbol{X}^* \boldsymbol{X}^{*\top} \right\|_{2,\infty} \leq \sqrt{\sigma_{\max}} \left\| \boldsymbol{X}^* \right\|_{2,\infty}$  and estimate (179). Note that this also implies that  $\left\| \boldsymbol{X}_{l,\cdot}^{0,(l)} \right\|_2 \leq 2 \left\| \boldsymbol{X}^* \right\|_{2,\infty}$ . To see this, one has by the unitary invariance of  $\left\| (\cdot)_{l,\cdot} \right\|_2$ ,

$$\left\| \boldsymbol{X}_{l,\cdot}^{0,(l)} \right\|_2 = \left\| \boldsymbol{X}_{l,\cdot}^{0,(l)} \boldsymbol{Q}^{(l)} \right\|_2 \leq \left\| (\boldsymbol{X}^{0,(l)} \boldsymbol{Q}^{(l)} - \boldsymbol{X}^*)_{l,\cdot} \right\|_2 + \left\| \boldsymbol{X}_{l,\cdot}^* \right\|_2 \leq 2 \left\| \boldsymbol{X}^* \right\|_{2,\infty}.$$

Substituting the above bounds back to (185) yields in

$$\begin{aligned} \left\| (\boldsymbol{X}^{0,(l)} \widehat{\boldsymbol{H}}^{0,(l)} - \boldsymbol{X}^*)_{l,\cdot} \right\|_2 &\lesssim \left\| \boldsymbol{X}^* \right\|_{2,\infty} \left\| \widehat{\boldsymbol{H}}^{0,(l)} - \boldsymbol{Q}^{(l)} \right\| \\ &\quad + \left\{ \mu r \sqrt{\frac{1}{np}} + \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \right\} \left\| \boldsymbol{X}^* \right\|_{2,\infty} \\ &\lesssim \left\{ \mu r \sqrt{\frac{1}{np}} + \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \right\} \left\| \boldsymbol{X}^* \right\|_{2,\infty}, \end{aligned}$$

where the second line relies on Lemma 47, bound (179), and condition  $\sigma_{\max}/\sigma_{\min} \asymp 1$ . This establishes (70e).



5. Our final step is to justify (70d). Define  $\mathbf{B} := \arg \min_{\mathbf{R} \in \mathcal{O}^{r \times r}} \|\mathbf{U}^{0,(l)} \mathbf{R} - \mathbf{U}^0\|_F$ . From the definition of  $\mathbf{R}^{0,(l)}$  (cf. (72)), one has

$$\|\mathbf{X}^0 \widehat{\mathbf{H}}^0 - \mathbf{X}^{0,(l)} \mathbf{R}^{0,(l)}\|_F \leq \|\mathbf{X}^{0,(l)} \mathbf{B} - \mathbf{X}^0\|_F.$$

Recognizing that

$$\mathbf{X}^{0,(l)} \mathbf{B} - \mathbf{X}^0 = \mathbf{U}^{0,(l)} \left[ (\boldsymbol{\Sigma}^{(l)})^{1/2} \mathbf{B} - \mathbf{B} (\boldsymbol{\Sigma}^0)^{1/2} \right] + (\mathbf{U}^{0,(l)} \mathbf{B} - \mathbf{U}^0) (\boldsymbol{\Sigma}^0)^{1/2},$$

we can use the triangle inequality to bound

$$\|\mathbf{X}^{0,(l)} \mathbf{B} - \mathbf{X}^0\|_F \leq \left\| (\boldsymbol{\Sigma}^{(l)})^{1/2} \mathbf{B} - \mathbf{B} (\boldsymbol{\Sigma}^0)^{1/2} \right\|_F + \|\mathbf{U}^{0,(l)} \mathbf{B} - \mathbf{U}^0\|_F \left\| (\boldsymbol{\Sigma}^0)^{1/2} \right\|.$$

In view of Lemma 46 and bounds (178) and (179), one has

$$\left\| (\boldsymbol{\Sigma}^{(l)})^{-1/2} \mathbf{B} - \mathbf{B} \boldsymbol{\Sigma}^{1/2} \right\|_F \lesssim \frac{1}{\sqrt{\sigma_{\min}}} \left\| (\mathbf{M}^0 - \mathbf{M}^{(l)}) \mathbf{U}^{0,(l)} \right\|_F.$$

From Davis–Kahan’s  $\sin \Theta$  theorem [39] we see that

$$\|\mathbf{U}^{0,(l)} \mathbf{B} - \mathbf{U}^0\|_F \lesssim \frac{1}{\sigma_{\min}} \left\| (\mathbf{M}^0 - \mathbf{M}^{(l)}) \mathbf{U}^{0,(l)} \right\|_F.$$

These estimates taken together with (181) give

$$\|\mathbf{X}^{0,(l)} \mathbf{B} - \mathbf{X}^0\|_F \lesssim \frac{1}{\sqrt{\sigma_{\min}}} \left\| (\mathbf{M}^0 - \mathbf{M}^{(l)}) \mathbf{U}^{0,(l)} \right\|_F.$$

It then boils down to controlling  $\left\| (\mathbf{M}^0 - \mathbf{M}^{(l)}) \mathbf{U}^{0,(l)} \right\|_F$ . Quantities of this type have shown up multiple times already, and hence, we omit the proof details for conciseness (see Appendix B.5). With probability at least  $1 - O(n^{-10})$ ,

$$\left\| (\mathbf{M}^0 - \mathbf{M}^{(l)}) \mathbf{U}^{0,(l)} \right\|_F \lesssim \left\{ \mu r \sqrt{\frac{\log n}{np}} \sigma_{\max} + \sigma \sqrt{\frac{n \log n}{p}} \right\} \|\mathbf{U}^{0,(l)}\|_{2,\infty}.$$

If one further has

$$\|\mathbf{U}^{0,(l)}\|_{2,\infty} \lesssim \|\mathbf{U}^*\|_{2,\infty} \lesssim \frac{1}{\sqrt{\sigma_{\min}}} \|\mathbf{X}^*\|_{2,\infty}, \tag{187}$$

then taking the previous bounds collectively establishes the desired bound

$$\|\mathbf{X}^0 \widehat{\mathbf{H}}^0 - \mathbf{X}^{0,(l)} \mathbf{R}^{0,(l)}\|_F \lesssim \left\{ \mu r \sqrt{\frac{\log n}{np}} + \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} \right\} \|\mathbf{X}^*\|_{2,\infty}.$$

**Proof of Claim (187)** Denote by  $M^{(l),zero}$  the matrix derived by zeroing out the  $l$ th row/column of  $M^{(l)}$ , and  $U^{(l),zero} \in \mathbb{R}^{n \times r}$  containing the leading  $r$  eigenvectors of  $M^{(l),zero}$ . On the one hand, [1, Lemma 4 and Lemma 14] demonstrate that

$$\max_{1 \leq l \leq n} \|U^{(l),zero}\|_{2,\infty} \lesssim \|U^*\|_{2,\infty}.$$

On the other hand, by the Davis–Kahan  $\sin \Theta$  theorem [39] we obtain

$$\|U^{0,(l)} \operatorname{sgn}(U^{0,(l)\top} U^{(l),zero}) - U^{(l),zero}\|_F \lesssim \frac{1}{\sigma_{\min}} \left\| (M^{(l)} - M^{(l),zero}) U^{(l),zero} \right\|_F, \tag{188}$$

where  $\operatorname{sgn}(A)$  denotes the sign matrix of  $A$ . For any  $j \neq l$ , one has

$$(M^{(l)} - M^{(l),zero})_{j,\cdot} U^{(l),zero} = (M^{(l)} - M^{(l),zero})_{j,l} U_{l,\cdot}^{(l),zero} = \mathbf{0}_{1 \times r},$$

since the  $l$ th row of  $U_{l,\cdot}^{(l),zero}$  is identically zero by construction. In addition,

$$\left\| (M^{(l)} - M^{(l),zero})_{l,\cdot} U^{(l),zero} \right\|_2 = \left\| M_{l,\cdot}^* U^{(l),zero} \right\|_2 \leq \|M^*\|_{2,\infty} \leq \sigma_{\max} \|U^*\|_{2,\infty}.$$

As a consequence, one has

$$\left\| (M^{(l)} - M^{(l),zero}) U^{(l),zero} \right\|_F = \left\| (M^{(l)} - M^{(l),zero})_{l,\cdot} U^{(l),zero} \right\|_2 \leq \sigma_{\max} \|U^*\|_{2,\infty},$$

which combined with (188) and the assumption  $\sigma_{\max}/\sigma_{\min} \asymp 1$  yields

$$\|U^{0,(l)} \operatorname{sgn}(U^{0,(l)\top} U^{(l),zero}) - U^{(l),zero}\|_F \lesssim \|U^*\|_{2,\infty}$$

Claim (187) then follows by combining the above estimates:

$$\begin{aligned} \|U^{0,(l)}\|_{2,\infty} &= \|U^{0,(l)} \operatorname{sgn}(U^{0,(l)\top} U^{(l),zero})\|_{2,\infty} \\ &\leq \|U^{(l),zero}\|_{2,\infty} + \|U^{0,(l)} \operatorname{sgn}(U^{0,(l)\top} U^{(l),zero}) - U^{(l),zero}\|_F \lesssim \|U^*\|_{2,\infty}, \end{aligned}$$

where we have utilized the unitary invariance of  $\|\cdot\|_{2,\infty}$ . □

### C Proofs for Blind Deconvolution

Before proceeding to the proofs, we make note of the following concentration results. The standard Gaussian concentration inequality and the union bound give

$$\max_{1 \leq l \leq m} |a_l^H x^*| \leq 5\sqrt{\log m} \tag{189}$$

with probability at least  $1 - O(m^{-10})$ . In addition, with probability exceeding  $1 - Cm \exp(-cK)$  for some constants  $c, C > 0$ ,

$$\max_{1 \leq l \leq m} \|a_l\|_2 \leq 3\sqrt{K}. \tag{190}$$

In addition, the population/expected Wirtinger Hessian at the truth  $z^*$  is given by

$$\nabla^2 F(z^*) = \begin{bmatrix} I_K & \mathbf{0} & \mathbf{0} & h^* x^{*\top} \\ \mathbf{0} & I_K & x^* h^{*\top} & \mathbf{0} \\ \mathbf{0} & (x^* h^{*\top})^H & I_K & \mathbf{0} \\ (h^* x^{*\top})^H & \mathbf{0} & \mathbf{0} & I_K \end{bmatrix}. \tag{191}$$

#### C.1 Proof of Lemma 14

First, we find it convenient to decompose the Wirtinger Hessian (cf. (80)) into the expected Wirtinger Hessian at the truth (cf. (191)) and the perturbation part as follows:

$$\nabla^2 f(z) = \nabla^2 F(z^*) + \left( \nabla^2 f(z) - \nabla^2 F(z^*) \right). \tag{192}$$

The proof then proceeds by showing that (i) the population Hessian  $\nabla^2 F(z^*)$  satisfies the restricted strong convexity and smoothness properties as advertised and (ii) the perturbation  $\nabla^2 f(z) - \nabla^2 F(z^*)$  is well controlled under our assumptions. We start by controlling the population Hessian in the following lemma.

**Lemma 26** *Instate the notation and the conditions of Lemma 14. We have*

$$\|\nabla^2 F(z^*)\| = 2 \quad \text{and} \quad u^H \left[ D \nabla^2 F(z^*) + \nabla^2 F(z^*) D \right] u \geq \|u\|_2^2.$$

The next step is to bound the perturbation. To this end, we define the set

$$\mathcal{S} := \{z : z \text{ satisfies (82)}\},$$

and derive the following lemma.

**Lemma 27** *Suppose the sample complexity satisfies  $m \gg \mu^2 K \log^9 m$ ,  $c > 0$  is a sufficiently small constant, and  $\delta = c/\log^2 m$ . Then with probability at least  $1 - O(m^{-10} + e^{-K} \log m)$ , one has*

$$\sup_{z \in S} \left\| \nabla^2 f(z) - \nabla^2 F(z^*) \right\| \leq 1/4.$$

Combining the two lemmas, we can easily see that for  $z \in S$ ,

$$\left\| \nabla^2 f(z) \right\| \leq \left\| \nabla^2 F(z^*) \right\| + \left\| \nabla^2 f(z) - \nabla^2 F(z^*) \right\| \leq 2 + 1/4 \leq 3,$$

which verifies the smoothness upper bound. In addition,

$$\begin{aligned} & \mathbf{u}^H \left[ \mathbf{D} \nabla^2 f(z) + \nabla^2 f(z) \mathbf{D} \right] \mathbf{u} \\ &= \mathbf{u}^H \left[ \mathbf{D} \nabla^2 F(z^*) + \nabla^2 F(z^*) \mathbf{D} \right] \mathbf{u} + \mathbf{u}^H \mathbf{D} \left[ \nabla^2 f(z) - \nabla^2 F(z^*) \right] \mathbf{u} \\ & \quad + \mathbf{u}^H \left[ \nabla^2 f(z) - \nabla^2 F(z^*) \right] \mathbf{D} \mathbf{u} \\ & \stackrel{(i)}{\geq} \mathbf{u}^H \left[ \mathbf{D} \nabla^2 F(z^*) + \nabla^2 F(z^*) \mathbf{D} \right] \mathbf{u} - 2 \|\mathbf{D}\| \left\| \nabla^2 f(z) - \nabla^2 F(z^*) \right\| \|\mathbf{u}\|_2^2 \\ & \stackrel{(ii)}{\geq} \|\mathbf{u}\|_2^2 - 2(1 + \delta) \cdot \frac{1}{4} \|\mathbf{u}\|_2^2 \\ & \stackrel{(iii)}{\geq} \frac{1}{4} \|\mathbf{u}\|_2^2, \end{aligned}$$

where (i) uses the triangle inequality, (ii) holds because of Lemma 27 and the fact that  $\|\mathbf{D}\| \leq 1 + \delta$ , and (iii) follows if  $\delta \leq 1/2$ . This establishes the claim on the restricted strong convexity.

### C.1.1 Proof of Lemma 26

We start by proving the identity  $\left\| \nabla^2 F(z^*) \right\| = 2$ . Let

$$\begin{aligned} \mathbf{u}_1 &= \frac{1}{\sqrt{2}} \begin{bmatrix} \mathbf{h}^* \\ \mathbf{0} \\ \mathbf{0} \\ \mathbf{x}^* \end{bmatrix}, & \mathbf{u}_2 &= \frac{1}{\sqrt{2}} \begin{bmatrix} \mathbf{0} \\ \mathbf{x}^* \\ \mathbf{h}^* \\ \mathbf{0} \end{bmatrix}, & \mathbf{u}_3 &= \frac{1}{\sqrt{2}} \begin{bmatrix} \mathbf{h}^* \\ \mathbf{0} \\ \mathbf{0} \\ -\mathbf{x}^* \end{bmatrix}, \\ \mathbf{u}_4 &= \frac{1}{\sqrt{2}} \begin{bmatrix} \mathbf{0} \\ \mathbf{x}^* \\ -\mathbf{h}^* \\ \mathbf{0} \end{bmatrix}. \end{aligned}$$

Recalling that  $\|\mathbf{h}^*\|_2 = \|\mathbf{x}^*\|_2 = 1$ , we can easily check that these four vectors form an orthonormal set. A little algebra reveals that

$$\nabla^2 F(\mathbf{z}^*) = \mathbf{I}_{4K} + \mathbf{u}_1 \mathbf{u}_1^H + \mathbf{u}_2 \mathbf{u}_2^H - \mathbf{u}_3 \mathbf{u}_3^H - \mathbf{u}_4 \mathbf{u}_4^H,$$

which immediately implies

$$\|\nabla^2 F(\mathbf{z}^*)\| = 2.$$

We now turn attention to the restricted strong convexity. Since  $\mathbf{u}^H \mathbf{D} \nabla^2 F(\mathbf{z}^*) \mathbf{u}$  is the complex conjugate of  $\mathbf{u}^H \nabla^2 F(\mathbf{z}^*) \mathbf{D} \mathbf{u}$  as both  $\nabla^2 F(\mathbf{z}^*)$  and  $\mathbf{D}$  are Hermitian, we will focus on the first term  $\mathbf{u}^H \mathbf{D} \nabla^2 F(\mathbf{z}^*) \mathbf{u}$ . This term can be rewritten as

$$\begin{aligned} & \mathbf{u}^H \mathbf{D} \nabla^2 F(\mathbf{z}^*) \mathbf{u} \\ & \stackrel{(i)}{=} \left[ (\mathbf{h}_1 - \mathbf{h}_2)^H, (\mathbf{x}_1 - \mathbf{x}_2)^H, (\overline{\mathbf{h}_1 - \mathbf{h}_2})^H, (\overline{\mathbf{x}_1 - \mathbf{x}_2})^H \right] \mathbf{D} \\ & \quad \begin{bmatrix} \mathbf{I}_K & \mathbf{0} & \mathbf{0} & \mathbf{h}^* \mathbf{x}^{*\top} \\ \mathbf{0} & \mathbf{I}_K & \mathbf{x}^* \mathbf{h}^{*\top} & \mathbf{0} \\ \mathbf{0} & (\mathbf{x}^* \mathbf{h}^{*\top})^H & \mathbf{I}_K & \mathbf{0} \\ (\mathbf{h}^* \mathbf{x}^{*\top})^H & \mathbf{0} & \mathbf{0} & \mathbf{I}_K \end{bmatrix} \begin{bmatrix} \mathbf{h}_1 - \mathbf{h}_2 \\ \mathbf{x}_1 - \mathbf{x}_2 \\ \overline{\mathbf{h}_1 - \mathbf{h}_2} \\ \overline{\mathbf{x}_1 - \mathbf{x}_2} \end{bmatrix} \\ & \stackrel{(ii)}{=} \left[ \gamma_1 (\mathbf{h}_1 - \mathbf{h}_2)^H, \gamma_2 (\mathbf{x}_1 - \mathbf{x}_2)^H, \gamma_1 (\overline{\mathbf{h}_1 - \mathbf{h}_2})^H, \gamma_2 (\overline{\mathbf{x}_1 - \mathbf{x}_2})^H \right] \\ & \quad \begin{bmatrix} \mathbf{h}_1 - \mathbf{h}_2 + \mathbf{h}^* \mathbf{x}^{*\top} \overline{(\mathbf{x}_1 - \mathbf{x}_2)} \\ \mathbf{x}_1 - \mathbf{x}_2 + \mathbf{x}^* \mathbf{h}^{*\top} (\overline{\mathbf{h}_1 - \mathbf{h}_2}) \\ (\mathbf{x}^* \mathbf{h}^{*\top})^H (\mathbf{x}_1 - \mathbf{x}_2) + \overline{(\mathbf{h}_1 - \mathbf{h}_2)} \\ (\mathbf{h}^* \mathbf{x}^{*\top})^H (\mathbf{h}_1 - \mathbf{h}_2) + \overline{(\mathbf{x}_1 - \mathbf{x}_2)} \end{bmatrix} \\ & = \left[ \gamma_1 (\mathbf{h}_1 - \mathbf{h}_2)^H, \gamma_2 (\mathbf{x}_1 - \mathbf{x}_2)^H, \gamma_1 (\overline{\mathbf{h}_1 - \mathbf{h}_2})^H, \gamma_2 (\overline{\mathbf{x}_1 - \mathbf{x}_2})^H \right] \\ & \quad \begin{bmatrix} \mathbf{h}_1 - \mathbf{h}_2 + \mathbf{h}^* (\mathbf{x}_1 - \mathbf{x}_2)^H \mathbf{x}^* \\ \mathbf{x}_1 - \mathbf{x}_2 + \mathbf{x}^* (\mathbf{h}_1 - \mathbf{h}_2)^H \mathbf{h}^* \\ \overline{\mathbf{h}_1 - \mathbf{h}_2} + \overline{\mathbf{h}^* (\mathbf{x}_1 - \mathbf{x}_2)^H \mathbf{x}^*} \\ \overline{\mathbf{x}_1 - \mathbf{x}_2} + \overline{\mathbf{x}^* (\mathbf{h}_1 - \mathbf{h}_2)^H \mathbf{h}^*} \end{bmatrix} \\ & = 2\gamma_1 \|\mathbf{h}_1 - \mathbf{h}_2\|_2^2 + 2\gamma_2 \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 \\ & \quad + (\gamma_1 + \gamma_2) \underbrace{(\mathbf{h}_1 - \mathbf{h}_2)^H \mathbf{h}^* (\mathbf{x}_1 - \mathbf{x}_2)^H \mathbf{x}^*}_{:=\beta} + (\gamma_1 + \gamma_2) \underbrace{(\overline{\mathbf{h}_1 - \mathbf{h}_2})^H \overline{\mathbf{h}^* (\mathbf{x}_1 - \mathbf{x}_2)^H \mathbf{x}^*}}_{=\bar{\beta}}, \end{aligned} \tag{193}$$

where (i) uses the definitions of  $\mathbf{u}$  and  $\nabla^2 F(\mathbf{z}^*)$ , and (ii) follows from the definition of  $\mathbf{D}$ . In view of the assumption (84), we can obtain

$$\begin{aligned} 2\gamma_1 \|\mathbf{h}_1 - \mathbf{h}_2\|_2^2 + 2\gamma_2 \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 & \geq 2 \min\{\gamma_1, \gamma_2\} \left( \|\mathbf{h}_1 - \mathbf{h}_2\|_2^2 + \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 \right) \\ & \geq (1 - \delta) \|\mathbf{u}\|_2^2, \end{aligned}$$

where the last inequality utilizes the identity

$$2 \|\mathbf{h}_1 - \mathbf{h}_2\|_2^2 + 2 \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 = \|\mathbf{u}\|_2^2. \tag{194}$$

It then boils down to controlling  $\beta$ . Toward this goal, we decompose  $\beta$  into the following four terms

$$\begin{aligned} \beta = & \underbrace{(\mathbf{h}_1 - \mathbf{h}_2)^H \mathbf{h}_2 (\mathbf{x}_1 - \mathbf{x}_2)^H \mathbf{x}_2}_{:=\beta_1} + \underbrace{(\mathbf{h}_1 - \mathbf{h}_2)^H (\mathbf{h}^* - \mathbf{h}_2) (\mathbf{x}_1 - \mathbf{x}_2)^H (\mathbf{x}^* - \mathbf{x}_2)}_{:=\beta_2} \\ & + \underbrace{(\mathbf{h}_1 - \mathbf{h}_2)^H (\mathbf{h}^* - \mathbf{h}_2) (\mathbf{x}_1 - \mathbf{x}_2)^H \mathbf{x}_2}_{:=\beta_3} + \underbrace{(\mathbf{h}_1 - \mathbf{h}_2)^H \mathbf{h}_2 (\mathbf{x}_1 - \mathbf{x}_2)^H (\mathbf{x}^* - \mathbf{x}_2)}_{:=\beta_4}. \end{aligned}$$

Since  $\|\mathbf{h}_2 - \mathbf{h}^*\|_2$  and  $\|\mathbf{x}_2 - \mathbf{x}^*\|_2$  are both small by (83),  $\beta_2, \beta_3,$  and  $\beta_4$  are well bounded. Specifically, regarding  $\beta_2$ , we discover that

$$\begin{aligned} |\beta_2| & \leq \|\mathbf{h}^* - \mathbf{h}_2\|_2 \|\mathbf{x}^* - \mathbf{x}_2\|_2 \|\mathbf{h}_1 - \mathbf{h}_2\|_2 \|\mathbf{x}_1 - \mathbf{x}_2\|_2 \\ & \leq \delta^2 \|\mathbf{h}_1 - \mathbf{h}_2\|_2 \|\mathbf{x}_1 - \mathbf{x}_2\|_2 \leq \delta \|\mathbf{h}_1 - \mathbf{h}_2\|_2 \|\mathbf{x}_1 - \mathbf{x}_2\|_2, \end{aligned}$$

where the second inequality is due to (83) and the last one holds since  $\delta < 1$ . Similarly, we can obtain

$$\begin{aligned} |\beta_3| & \leq \delta \|\mathbf{x}_2\|_2 \|\mathbf{h}_1 - \mathbf{h}_2\|_2 \|\mathbf{x}_1 - \mathbf{x}_2\|_2 \leq 2\delta \|\mathbf{h}_1 - \mathbf{h}_2\|_2 \|\mathbf{x}_1 - \mathbf{x}_2\|_2, \text{ and} \\ |\beta_4| & \leq \delta \|\mathbf{h}_2\|_2 \|\mathbf{h}_1 - \mathbf{h}_2\|_2 \|\mathbf{x}_1 - \mathbf{x}_2\|_2 \leq 2\delta \|\mathbf{h}_1 - \mathbf{h}_2\|_2 \|\mathbf{x}_1 - \mathbf{x}_2\|_2, \end{aligned}$$

where both lines make use of the facts that

$$\begin{aligned} \|\mathbf{x}_2\|_2 & \leq \|\mathbf{x}_2 - \mathbf{x}^*\|_2 + \|\mathbf{x}^*\|_2 \leq 1 + \delta \leq 2 \quad \text{and} \\ \|\mathbf{h}_2\|_2 & \leq \|\mathbf{h}_2 - \mathbf{h}^*\|_2 + \|\mathbf{h}^*\|_2 \leq 1 + \delta \leq 2. \end{aligned} \tag{195}$$

Combine the previous three bounds to reach

$$\begin{aligned} |\beta_2| + |\beta_3| + |\beta_4| & \leq 5\delta \|\mathbf{h}_1 - \mathbf{h}_2\|_2 \|\mathbf{x}_1 - \mathbf{x}_2\|_2 \\ & \leq 5\delta \frac{\|\mathbf{h}_1 - \mathbf{h}_2\|_2^2 + \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2}{2} = \frac{5}{4} \delta \|\mathbf{u}\|_2^2, \end{aligned}$$

where we utilize the elementary inequality  $ab \leq (a^2 + b^2)/2$  and identity (194).

The only remaining term is thus  $\beta_1$ . Recalling that  $(\mathbf{h}_1, \mathbf{x}_1)$  and  $(\mathbf{h}_2, \mathbf{x}_2)$  are aligned by our assumption, we can invoke Lemma 56 to obtain

$$(\mathbf{h}_1 - \mathbf{h}_2)^H \mathbf{h}_2 = \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 + \mathbf{x}_2^H (\mathbf{x}_1 - \mathbf{x}_2) - \|\mathbf{h}_1 - \mathbf{h}_2\|_2^2,$$

which allows one to rewrite  $\beta_1$  as

$$\begin{aligned} \beta_1 & = \left\{ \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 + \mathbf{x}_2^H (\mathbf{x}_1 - \mathbf{x}_2) - \|\mathbf{h}_1 - \mathbf{h}_2\|_2^2 \right\} \cdot (\mathbf{x}_1 - \mathbf{x}_2)^H \mathbf{x}_2 \\ & = (\mathbf{x}_1 - \mathbf{x}_2)^H \mathbf{x}_2 \left( \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 - \|\mathbf{h}_1 - \mathbf{h}_2\|_2^2 \right) + \left| (\mathbf{x}_1 - \mathbf{x}_2)^H \mathbf{x}_2 \right|^2. \end{aligned}$$

Consequently,

$$\begin{aligned} \beta_1 + \overline{\beta_1} &= 2 \left| (\mathbf{x}_1 - \mathbf{x}_2)^H \mathbf{x}_2 \right|_2^2 + 2 \operatorname{Re} \left[ (\mathbf{x}_1 - \mathbf{x}_2)^H \mathbf{x}_2 \right] \left( \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 - \|\mathbf{h}_1 - \mathbf{h}_2\|_2^2 \right) \\ &\geq 2 \operatorname{Re} \left[ (\mathbf{x}_1 - \mathbf{x}_2)^H \mathbf{x}_2 \right] \left( \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 - \|\mathbf{h}_1 - \mathbf{h}_2\|_2^2 \right) \\ &\stackrel{(i)}{\geq} - \left| (\mathbf{x}_1 - \mathbf{x}_2)^H \mathbf{x}_2 \right| \|\mathbf{u}\|_2^2 \\ &\stackrel{(ii)}{\geq} -4\delta \|\mathbf{u}\|_2^2. \end{aligned}$$

Here, (i) arises from the triangle inequality that

$$\left| \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 - \|\mathbf{h}_1 - \mathbf{h}_2\|_2^2 \right| \leq \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 + \|\mathbf{h}_1 - \mathbf{h}_2\|_2^2 = \frac{1}{2} \|\mathbf{u}\|_2^2,$$

and (ii) occurs since  $\|\mathbf{x}_1 - \mathbf{x}_2\|_2 \leq \|\mathbf{x}_1 - \mathbf{x}^*\|_2 + \|\mathbf{x}_2 - \mathbf{x}^*\|_2 \leq 2\delta$  and  $\|\mathbf{x}_2\|_2 \leq 2$  (see (195)).

To finish up, note that  $\gamma_1 + \gamma_2 \leq 2(1 + \delta) \leq 3$  for  $\delta < 1/2$ . Substitute these bounds into (193) to obtain

$$\begin{aligned} \mathbf{u}^H \mathbf{D} \nabla^2 F(\mathbf{z}^*) \mathbf{u} &\geq (1 - \delta) \|\mathbf{u}\|_2^2 + (\gamma_1 + \gamma_2) (\beta + \overline{\beta}) \\ &\geq (1 - \delta) \|\mathbf{u}\|_2^2 + (\gamma_1 + \gamma_2) (\beta_1 + \overline{\beta_1}) - 2(\gamma_1 + \gamma_2) (|\beta_2| + |\beta_3| + |\beta_4|) \\ &\geq (1 - \delta) \|\mathbf{u}\|_2^2 - 12\delta \|\mathbf{u}\|_2^2 - 6 \cdot \frac{5}{4} \delta \|\mathbf{u}\|_2^2 \\ &\geq (1 - 20.5\delta) \|\mathbf{u}\|_2^2 \\ &\geq \frac{1}{2} \|\mathbf{u}\|_2^2 \end{aligned}$$

as long as  $\delta$  is small enough.

### C.1.2 Proof of Lemma 27

In view of the expressions of  $\nabla^2 f(\mathbf{z})$  and  $\nabla^2 F(\mathbf{z}^*)$  (cf. (80) and (191)) and the triangle inequality, we get

$$\left\| \nabla^2 f(\mathbf{z}) - \nabla^2 F(\mathbf{z}^*) \right\| \leq 2\alpha_1 + 2\alpha_2 + 4\alpha_3 + 4\alpha_4, \tag{196}$$

where the four terms on the right-hand side are defined as follows

$$\begin{aligned} \alpha_1 &= \left\| \sum_{j=1}^m \left| \mathbf{a}_j^H \mathbf{x} \right|^2 \mathbf{b}_j \mathbf{b}_j^H - \mathbf{I}_K \right\|, & \alpha_2 &= \left\| \sum_{j=1}^m \left| \mathbf{b}_j^H \mathbf{h} \right|^2 \mathbf{a}_j \mathbf{a}_j^H - \mathbf{I}_K \right\|, \\ \alpha_3 &= \left\| \sum_{j=1}^m \left( \mathbf{b}_j^H \mathbf{h} \mathbf{x}^H \mathbf{a}_j - y_j \right) \mathbf{b}_j \mathbf{a}_j^H \right\|, & \alpha_4 &= \left\| \sum_{j=1}^m \mathbf{b}_j \mathbf{b}_j^H \mathbf{h} \left( \mathbf{a}_j \mathbf{a}_j^H \mathbf{x} \right)^T - \mathbf{h}^* \mathbf{x}^{*T} \right\|. \end{aligned}$$

In what follows, we shall control  $\sup_{z \in \mathcal{S}} \alpha_j$  for  $j = 1, 2, 3, 4$  separately.

1. Regarding the first term  $\alpha_1$ , the triangle inequality gives

$$\alpha_1 \leq \underbrace{\left\| \sum_{j=1}^m |a_j^H x|^2 b_j b_j^H - \sum_{j=1}^m |a_j^H x^*|^2 b_j b_j^H \right\|}_{:=\beta_1} + \underbrace{\left\| \sum_{j=1}^m |a_j^H x^*|^2 b_j b_j^H - I_K \right\|}_{:=\beta_2}.$$

- To control  $\beta_1$ , the key observation is that  $a_j^H x$  and  $a_j^H x^*$  are extremely close. We can rewrite  $\beta_1$  as

$$\beta_1 = \left\| \sum_{j=1}^m \left( |a_j^H x|^2 - |a_j^H x^*|^2 \right) b_j b_j^H \right\| \leq \left\| \sum_{j=1}^m \left| |a_j^H x|^2 - |a_j^H x^*|^2 \right| b_j b_j^H \right\|, \tag{197}$$

where

$$\begin{aligned} & \left| |a_j^H x|^2 - |a_j^H x^*|^2 \right| \\ & \stackrel{(i)}{=} \left| \left[ a_j^H (x - x^*) \right]^H a_j^H (x - x^*) + \left[ a_j^H (x - x^*) \right]^H a_j^H x^* + \left( a_j^H x^* \right)^H a_j^H (x - x^*) \right| \\ & \stackrel{(ii)}{\leq} \left| a_j^H (x - x^*) \right|^2 + 2 \left| a_j^H (x - x^*) \right| \left| a_j^H x^* \right| \\ & \stackrel{(iii)}{\leq} 4C_3^2 \frac{1}{\log^3 m} + 4C_3 \frac{1}{\log^{3/2} m} \cdot 5\sqrt{\log m} \\ & \lesssim C_3 \frac{1}{\log m}. \end{aligned}$$

Here, the first line (i) uses the identity for  $u, v \in \mathbb{C}$ ,

$$|u|^2 - |v|^2 = u^H u - v^H v = (u - v)^H (u - v) + (u - v)^H v + v^H (u - v),$$

the second relation (ii) comes from the triangle inequality, and the third line (iii) follows from (189) and assumption (82b). Substitution into (197) gives

$$\beta_1 \leq \max_{1 \leq j \leq m} \left| |a_j^H x|^2 - |a_j^H x^*|^2 \right| \left\| \sum_{j=1}^m b_j b_j^H \right\| \lesssim C_3 \frac{1}{\log m},$$

where the last inequality comes from the fact that  $\sum_{j=1}^m b_j b_j^H = I_K$ .

- The other term  $\beta_2$  can be bounded through Lemma 59, which reveals that with probability  $1 - O(m^{-10})$ ,

$$\beta_2 \lesssim \sqrt{\frac{K}{m} \log m}.$$



Taken collectively, the preceding two bounds give

$$\sup_{z \in \mathcal{S}} \alpha_1 \lesssim \sqrt{\frac{K}{m} \log m} + C_3 \frac{1}{\log m}.$$

Hence,  $\mathbb{P}(\sup_{z \in \mathcal{S}} \alpha_1 \leq 1/32) = 1 - O(m^{-10})$ .

2. We are going to prove that  $\mathbb{P}(\sup_{z \in \mathcal{S}} \alpha_2 \leq 1/32) = 1 - O(m^{-10})$ . The triangle inequality allows us to bound  $\alpha_2$  as

$$\alpha_2 \leq \underbrace{\left\| \sum_{j=1}^m \left| \mathbf{b}_j^H \mathbf{h} \right|^2 \mathbf{a}_j \mathbf{a}_j^H - \|\mathbf{h}\|_2^2 \mathbf{I}_K \right\|}_{:=\theta_1(\mathbf{h})} + \underbrace{\left\| \|\mathbf{h}\|_2^2 \mathbf{I}_K - \mathbf{I}_K \right\|}_{:=\theta_2(\mathbf{h})}.$$

The second term  $\theta_2(\mathbf{h})$  is easy to control. To see this, we have

$$\theta_2(\mathbf{h}) = \left| \|\mathbf{h}\|_2^2 - 1 \right| = \left| \|\mathbf{h}\|_2 - 1 \right| (\|\mathbf{h}\|_2 + 1) \leq 3\delta < 1/64,$$

where the penultimate relation uses the assumption that  $\|\mathbf{h} - \mathbf{h}^*\|_2 \leq \delta$  and hence

$$\left| \|\mathbf{h}\|_2 - 1 \right| \leq \delta, \quad \|\mathbf{h}\|_2 \leq 1 + \delta \leq 2.$$

For the first term  $\theta_1(\mathbf{h})$ , we define a new set

$$\mathcal{H} := \left\{ \mathbf{h} \in \mathbb{C}^K : \|\mathbf{h} - \mathbf{h}^*\|_2 \leq \delta \text{ and } \max_{1 \leq j \leq m} \left| \mathbf{b}_j^H \mathbf{h} \right| \leq \frac{2C_4\mu \log^2 m}{\sqrt{m}} \right\}.$$

It is easily seen that  $\sup_{z \in \mathcal{S}} \theta_1 \leq \sup_{\mathbf{h} \in \mathcal{H}} \theta_1$ . We plan to use the standard covering argument to show that

$$\mathbb{P} \left( \sup_{\mathbf{h} \in \mathcal{H}} \theta_1(\mathbf{h}) \leq 1/64 \right) = 1 - O(m^{-10}). \tag{198}$$

To this end, we define  $c_j(\mathbf{h}) = \left| \mathbf{b}_j^H \mathbf{h} \right|^2$  for every  $1 \leq j \leq m$ . It is straightforward to check that

$$\theta_1(\mathbf{h}) = \left\| \sum_{j=1}^m c_j(\mathbf{h}) \left( \mathbf{a}_j \mathbf{a}_j^H - \mathbf{I}_K \right) \right\|, \quad \max_{1 \leq j \leq m} |c_j| \leq \left( \frac{2C_4\mu \log^2 m}{\sqrt{m}} \right)^2, \tag{199}$$

$$\begin{aligned} \sum_{j=1}^m c_j^2 &= \sum_{j=1}^m |\mathbf{b}_j^H \mathbf{h}|^4 \leq \left\{ \max_{1 \leq j \leq m} |\mathbf{b}_j^H \mathbf{h}|^2 \right\} \sum_{j=1}^m |\mathbf{b}_j^H \mathbf{h}|^2 = \left\{ \max_{1 \leq j \leq m} |\mathbf{b}_j^H \mathbf{h}|^2 \right\} \|\mathbf{h}\|_2^2 \\ &\leq 4 \left( \frac{2C_4 \mu \log^2 m}{\sqrt{m}} \right)^2 \end{aligned} \tag{200}$$

for  $\mathbf{h} \in \mathcal{H}$ . In the above argument, we have used the facts that  $\sum_{j=1}^m \mathbf{b}_j \mathbf{b}_j^H = \mathbf{I}_K$  and

$$\sum_{j=1}^m |\mathbf{b}_j^H \mathbf{h}|^2 = \mathbf{h}^H \left( \sum_{j=1}^m \mathbf{b}_j \mathbf{b}_j^H \right) \mathbf{h} = \|\mathbf{h}\|_2^2 \leq (1 + \delta)^2 \leq 4,$$

together with the definition of  $\mathcal{H}$ . Lemma 57 combined with (199) and (200) readily yields that for any fixed  $\mathbf{h} \in \mathcal{H}$  and any  $t \geq 0$ ,

$$\begin{aligned} \mathbb{P}(\theta_1(\mathbf{h}) \geq t) &\leq 2 \exp \left( \tilde{C}_1 K - \tilde{C}_2 \min \left\{ \frac{t}{\max_{1 \leq j \leq m} |c_j|}, \frac{t^2}{\sum_{j=1}^m c_j^2} \right\} \right) \\ &\leq 2 \exp \left( \tilde{C}_1 K - \tilde{C}_2 \frac{mt \min \{1, t/4\}}{4C_4^2 \mu^2 \log^4 m} \right), \end{aligned} \tag{201}$$

where  $\tilde{C}_1, \tilde{C}_2 > 0$  are some universal constants.

Now we are in a position to strengthen this bound to obtain uniform control of  $\theta_1$  over  $\mathcal{H}$ . Note that for any  $\mathbf{h}_1, \mathbf{h}_2 \in \mathcal{H}$ ,

$$\begin{aligned} |\theta_1(\mathbf{h}_1) - \theta_1(\mathbf{h}_2)| &\leq \left\| \sum_{j=1}^m \left( |\mathbf{b}_j^H \mathbf{h}_1|^2 - |\mathbf{b}_j^H \mathbf{h}_2|^2 \right) \mathbf{a}_j \mathbf{a}_j^H \right\| + \left| \|\mathbf{h}_1\|_2^2 - \|\mathbf{h}_2\|_2^2 \right| \\ &= \max_{1 \leq j \leq m} \left| |\mathbf{b}_j^H \mathbf{h}_1|^2 - |\mathbf{b}_j^H \mathbf{h}_2|^2 \right| \left\| \sum_{j=1}^m \mathbf{a}_j \mathbf{a}_j^H \right\| + \left| \|\mathbf{h}_1\|_2^2 - \|\mathbf{h}_2\|_2^2 \right|, \end{aligned}$$

where

$$\begin{aligned} \left| |\mathbf{b}_j^H \mathbf{h}_2|^2 - |\mathbf{b}_j^H \mathbf{h}_1|^2 \right| &= \left| (\mathbf{h}_2 - \mathbf{h}_1)^H \mathbf{b}_j \mathbf{b}_j^H \mathbf{h}_2 + \mathbf{h}_1^H \mathbf{b}_j \mathbf{b}_j^H (\mathbf{h}_2 - \mathbf{h}_1) \right| \\ &\leq 2 \max\{\|\mathbf{h}_1\|_2, \|\mathbf{h}_2\|_2\} \|\mathbf{h}_2 - \mathbf{h}_1\|_2 \|\mathbf{b}_j\|_2^2 \\ &\leq 4 \|\mathbf{h}_2 - \mathbf{h}_1\|_2 \|\mathbf{b}_j\|_2^2 \leq \frac{4K}{m} \|\mathbf{h}_2 - \mathbf{h}_1\|_2 \end{aligned}$$

and

$$\begin{aligned} \left| \|\mathbf{h}_1\|_2^2 - \|\mathbf{h}_2\|_2^2 \right| &= \left| \mathbf{h}_1^H (\mathbf{h}_1 - \mathbf{h}_2) - (\mathbf{h}_1 - \mathbf{h}_2)^H \mathbf{h}_2 \right| \\ &\leq 2 \max\{\|\mathbf{h}_1\|_2, \|\mathbf{h}_2\|_2\} \|\mathbf{h}_2 - \mathbf{h}_1\|_2 \leq 4 \|\mathbf{h}_1 - \mathbf{h}_2\|_2. \end{aligned}$$

Define an event  $\mathcal{E}_0 = \left\{ \left\| \sum_{j=1}^m \mathbf{a}_j \mathbf{a}_j^H \right\| \leq 2m \right\}$ . When  $\mathcal{E}_0$  happens, the previous estimates give

$$|\theta_1(\mathbf{h}_1) - \theta_1(\mathbf{h}_2)| \leq (8K + 4)\|\mathbf{h}_1 - \mathbf{h}_2\|_2 \leq 10K\|\mathbf{h}_1 - \mathbf{h}_2\|_2, \quad \forall \mathbf{h}_1, \mathbf{h}_2 \in \mathcal{H}.$$

Let  $\varepsilon = 1/(1280K)$ , and  $\tilde{\mathcal{H}}$  be an  $\varepsilon$ -net covering  $\mathcal{H}$  (see [116, Definition 5.1]). We have

$$\left( \left\{ \sup_{\mathbf{h} \in \tilde{\mathcal{H}}} \theta_1(\mathbf{h}) \leq \frac{1}{128} \right\} \cap \mathcal{E}_0 \right) \subseteq \left\{ \sup_{\mathbf{h} \in \mathcal{H}} \theta_1 \leq \frac{1}{64} \right\}$$

and, as a result,

$$\begin{aligned} \mathbb{P} \left( \sup_{\mathbf{h} \in \mathcal{H}} \theta_1(\mathbf{h}) \geq \frac{1}{64} \right) &\leq \mathbb{P} \left( \sup_{\mathbf{h} \in \tilde{\mathcal{H}}} \theta_1(\mathbf{h}) \geq \frac{1}{128} \right) + \mathbb{P}(\mathcal{E}_0^c) \\ &\leq |\tilde{\mathcal{H}}| \cdot \max_{\mathbf{h} \in \tilde{\mathcal{H}}} \mathbb{P} \left( \theta_1(\mathbf{h}) \geq \frac{1}{128} \right) + \mathbb{P}(\mathcal{E}_0^c). \end{aligned}$$

Lemma 57 forces that  $\mathbb{P}(\mathcal{E}_0^c) = O(m^{-10})$ . Additionally, we have  $\log |\tilde{\mathcal{H}}| \leq \tilde{C}_3 K \log K$  for some absolute constant  $\tilde{C}_3 > 0$  according to [116, Lemma 5.2]. Hence, (201) leads to

$$\begin{aligned} &|\tilde{\mathcal{H}}| \cdot \max_{\mathbf{h} \in \tilde{\mathcal{H}}} \mathbb{P} \left( \theta_1(\mathbf{h}) \geq \frac{1}{128} \right) \\ &\leq 2 \exp \left( \tilde{C}_3 K \log K + \tilde{C}_1 K - \tilde{C}_2 \frac{m(1/128) \min \{1, (1/128)/4\}}{4C_4^2 \mu^2 \log^4 m} \right) \\ &\leq 2 \exp \left( 2\tilde{C}_3 K \log m - \frac{\tilde{C}_4 m}{\mu^2 \log^4 m} \right) \end{aligned}$$

for some constant  $\tilde{C}_4 > 0$ . Under the sample complexity  $m \gg \mu^2 K \log^5 m$ , the right-hand side of the above display is at most  $O(m^{-10})$ . Combine the estimates above to establish the desired high-probability bound for  $\sup_{z \in \mathcal{S}} \alpha_2$ .

3. Next, we will demonstrate that

$$\mathbb{P}(\sup_{z \in \mathcal{S}} \alpha_3 \leq 1/96) = 1 - O \left( m^{-10} + e^{-K} \log m \right).$$

To this end, we let

$$\mathbf{A} = \begin{bmatrix} \mathbf{a}_1^H \\ \vdots \\ \mathbf{a}_m^H \end{bmatrix} \in \mathbb{C}^{m \times K}, \quad \mathbf{B} = \begin{bmatrix} \mathbf{b}_1^H \\ \vdots \\ \mathbf{b}_m^H \end{bmatrix} \in \mathbb{C}^{m \times K},$$

$$C = \begin{bmatrix} c_1(z) & & & \\ & c_2(z) & & \\ & & \dots & \\ & & & c_m(z) \end{bmatrix} \in \mathbb{C}^{m \times m},$$

where for each  $1 \leq j \leq m$ ,

$$c_j(z) := \mathbf{b}_j^H \mathbf{h} \mathbf{x}^H \mathbf{a}_j - y_j = \mathbf{b}_j^H (\mathbf{h} \mathbf{x}^H - \mathbf{h}^* \mathbf{x}^{*H}) \mathbf{a}_j.$$

As a consequence, we can write  $\alpha_3 = \|\mathbf{B}^H \mathbf{C} \mathbf{A}\|$ .

The key observation is that both the  $\ell_\infty$  norm and the Frobenius norm of  $C$  are well controlled. Specifically, we claim for the moment that with probability at least  $1 - O(m^{-10})$ ,

$$\|C\|_\infty = \max_{1 \leq j \leq m} |c_j| \leq C \frac{\mu \log^{5/2} m}{\sqrt{m}}; \tag{202a}$$

$$\|C\|_F^2 = \sum_{j=1}^m |c_j|^2 \leq 12\delta^2, \tag{202b}$$

where  $C > 0$  is some absolute constant. This motivates us to divide the entries in  $C$  into multiple groups based on their magnitudes.

To be precise, introduce  $R := 1 + \lceil \log_2(C\mu \log^{7/2} m) \rceil$  sets  $\{\mathcal{I}_r\}_{1 \leq r \leq R}$ , where

$$\mathcal{I}_r = \left\{ j \in [m] : \frac{C\mu \log^{5/2} m}{2^r \sqrt{m}} < |c_j| \leq \frac{C\mu \log^{5/2} m}{2^{r-1} \sqrt{m}} \right\}, \quad 1 \leq r \leq R - 1$$

and  $\mathcal{I}_R = \{1, \dots, m\} \setminus (\bigcup_{r=1}^{R-1} \mathcal{I}_r)$ . An immediate consequence of the definition of  $\mathcal{I}_r$  and the norm constraints in (202) is the following cardinality bound

$$|\mathcal{I}_r| \leq \frac{\|C\|_F^2}{\min_{j \in \mathcal{I}_r} |c_j|^2} \leq \frac{12\delta^2}{\left(\frac{C\mu \log^{5/2} m}{2^r \sqrt{m}}\right)^2} = \underbrace{\frac{12\delta^2 4^r}{C^2 \mu^2 \log^5 m}}_{\delta_r} m \tag{203}$$

for  $1 \leq r \leq R - 1$ . Since  $\{\mathcal{I}_r\}_{1 \leq r \leq R}$  form a partition of the index set  $\{1, \dots, m\}$ , it is easy to see that

$$\mathbf{B}^H \mathbf{C} \mathbf{A} = \sum_{r=1}^R (\mathbf{B}_{\mathcal{I}_r, \cdot})^H \mathbf{C}_{\mathcal{I}_r, \mathcal{I}_r} \mathbf{A}_{\mathcal{I}_r, \cdot},$$

where  $\mathbf{D}_{\mathcal{I}, \mathcal{J}}$  denotes the submatrix of  $\mathbf{D}$  induced by the rows and columns of  $\mathbf{D}$  having indices from  $\mathcal{I}$  and  $\mathcal{J}$ , respectively, and  $\mathbf{D}_{\mathcal{I}, \cdot}$  refers to the submatrix

formed by the rows from the index set  $\mathcal{I}$ . As a result, one can invoke the triangle inequality to derive

$$\alpha_3 \leq \sum_{r=1}^{R-1} \| \mathbf{B}_{\mathcal{I}_r, \cdot} \| \cdot \| \mathbf{C}_{\mathcal{I}_r, \mathcal{I}_r} \| \cdot \| \mathbf{A}_{\mathcal{I}_r, \cdot} \| + \| \mathbf{B}_{\mathcal{I}_R, \cdot} \| \cdot \| \mathbf{C}_{\mathcal{I}_R, \mathcal{I}_R} \| \cdot \| \mathbf{A}_{\mathcal{I}_R, \cdot} \| . \tag{204}$$

Recognizing that  $\mathbf{B}^H \mathbf{B} = \mathbf{I}_K$ , we obtain

$$\| \mathbf{B}_{\mathcal{I}_r, \cdot} \| \leq \| \mathbf{B} \| = 1$$

for every  $1 \leq r \leq R$ . In addition, by construction of  $\mathcal{I}_r$ , we have

$$\| \mathbf{C}_{\mathcal{I}_r, \mathcal{I}_r} \| = \max_{j \in \mathcal{I}_r} |c_j| \leq \frac{C \mu \log^{5/2} m}{2^{r-1} \sqrt{m}}$$

for  $1 \leq r \leq R$ , and specifically for  $R$ , one has

$$\| \mathbf{C}_{\mathcal{I}_R, \mathcal{I}_R} \| = \max_{j \in \mathcal{I}_R} |c_j| \leq \frac{C \mu \log^{5/2} m}{2^{R-1} \sqrt{m}} \leq \frac{1}{\sqrt{m} \log m} ,$$

which follows from the definition of  $R$ , i.e.,  $R = 1 + \lceil \log_2(C \mu \log^{7/2} m) \rceil$ . Regarding  $\| \mathbf{A}_{\mathcal{I}_r, \cdot} \|$ , we discover that  $\| \mathbf{A}_{\mathcal{I}_R, \cdot} \| \leq \| \mathbf{A} \|$  and, in view of (203),

$$\| \mathbf{A}_{\mathcal{I}_r, \cdot} \| \leq \sup_{\mathcal{I}: |\mathcal{I}| \leq \delta_r m} \| \mathbf{A}_{\mathcal{I}, \cdot} \| , \quad 1 \leq r \leq R - 1 .$$

Substitute the above estimates into (204) to get

$$\alpha_3 \leq \sum_{r=1}^{R-1} \frac{C \mu \log^{5/2} m}{2^{r-1} \sqrt{m}} \sup_{\mathcal{I}: |\mathcal{I}| \leq \delta_r m} \| \mathbf{A}_{\mathcal{I}, \cdot} \| + \frac{\| \mathbf{A} \|}{\sqrt{m} \log m} . \tag{205}$$

It remains to upper bound  $\| \mathbf{A} \|$  and  $\sup_{\mathcal{I}: |\mathcal{I}| \leq \delta_r m} \| \mathbf{A}_{\mathcal{I}, \cdot} \|$ . Lemma 57 tells us that  $\| \mathbf{A} \| \leq 2\sqrt{m}$  with probability at least  $1 - O(m^{-10})$ . Furthermore, we can invoke Lemma 58 to bound  $\sup_{\mathcal{I}: |\mathcal{I}| \leq \delta_r m} \| \mathbf{A}_{\mathcal{I}, \cdot} \|$  for each  $1 \leq r \leq R - 1$ . It is easily seen from our assumptions  $m \gg \mu^2 K \log^9 m$  and  $\delta = c / \log^2 m$  that  $\delta_r \gg K / m$ . In addition,

$$\delta_r \leq \frac{12 \delta^2 4^{R-1}}{C^2 \mu^2 \log^5 m} \leq \frac{12 \delta^2 4^{1 + \log_2(C \mu \log^{7/2} m)}}{C^2 \mu^2 \log^5 m} = 48 \delta^2 \log^2 m = \frac{48c}{\log^2 m} \ll 1 .$$

By Lemma 58, we obtain that for some constants  $\tilde{C}_2, \tilde{C}_3 > 0$

$$\begin{aligned} \mathbb{P} \left( \sup_{\mathcal{I}: |\mathcal{I}| \leq \delta_r m} \|A_{\mathcal{I}, \cdot}\| \geq \sqrt{4\tilde{C}_3 \delta_r m \log(e/\delta_r)} \right) &\leq 2 \exp \left( -\frac{\tilde{C}_2 \tilde{C}_3}{3} \delta_r m \log(e/\delta_r) \right) \\ &\leq 2 \exp \left( -\frac{\tilde{C}_2 \tilde{C}_3}{3} \delta_r m \right) \leq 2e^{-K}. \end{aligned}$$

Taking the union bound and substituting the estimates above into (205), we see that with probability at least  $1 - O(m^{-10}) - O((R - 1)e^{-K})$ ,

$$\begin{aligned} \alpha_3 &\leq \sum_{r=1}^{R-1} \frac{C\mu \log^{5/2} m}{2^{r-1} \sqrt{m}} \cdot \sqrt{4\tilde{C}_3 \delta_r m \log(e/\delta_r)} + \frac{2\sqrt{m}}{\sqrt{m} \log m} \\ &\leq \sum_{r=1}^{R-1} 4\delta \sqrt{12\tilde{C}_3 \log(e/\delta_r)} + \frac{2}{\log m} \\ &\lesssim (R - 1)\delta \sqrt{\log(e/\delta_1)} + \frac{1}{\log m}. \end{aligned}$$

Note that  $\mu \leq \sqrt{m}$ ,  $R - 1 = \lceil \log_2(C\mu \log^{7/2} m) \rceil \lesssim \log m$ , and

$$\sqrt{\log \frac{e}{\delta_1}} = \sqrt{\log \left( \frac{eC^2\mu^2 \log^5 m}{48\delta^2} \right)} \lesssim \log m.$$

Therefore, with probability exceeding  $1 - O(m^{-10}) - O(e^{-K} \log m)$ ,

$$\sup_{z \in \mathcal{S}} \alpha_3 \lesssim \delta \log^2 m + \frac{1}{\log m}.$$

By taking  $c$  to be small enough in  $\delta = c/\log^2 m$ , we get

$$\mathbb{P} \left( \sup_{z \in \mathcal{S}} \alpha_3 \geq 1/96 \right) \leq O(m^{-10}) + O(e^{-K} \log m)$$

as claimed.

Finally, it remains to justify (202). For all  $z \in \mathcal{S}$ , the triangle inequality tells us that

$$\begin{aligned} |c_j| &\leq \left| \mathbf{b}_j^H \mathbf{h} (\mathbf{x} - \mathbf{x}^*)^H \mathbf{a}_j \right| + \left| \mathbf{b}_j^H (\mathbf{h} - \mathbf{h}^*) \mathbf{x}^{*H} \mathbf{a}_j \right| \\ &\leq \left| \mathbf{b}_j^H \mathbf{h} \right| \cdot \left| \mathbf{a}_j^H (\mathbf{x} - \mathbf{x}^*) \right| + \left( \left| \mathbf{b}_j^H \mathbf{h} \right| + \left| \mathbf{b}_j^H \mathbf{h}^* \right| \right) \cdot \left| \mathbf{a}_j^H \mathbf{x}^* \right| \\ &\leq \frac{2C_4\mu \log^2 m}{\sqrt{m}} \cdot \frac{2C_3}{\log^{3/2} m} + \left( \frac{2C_4\mu \log^2 m}{\sqrt{m}} + \frac{\mu}{\sqrt{m}} \right) 5\sqrt{\log m} \end{aligned}$$

$$\leq C \frac{\mu \log^{5/2} m}{\sqrt{m}},$$

for some large constant  $C > 0$ , where we have used the definition of  $\mathcal{S}$  and fact (189). Claim (202b) follows directly from [76, Lemma 5.14]. To avoid confusion, we use  $\mu_1$  to refer to the parameter  $\mu$  therein. Let  $L = m$ ,  $N = K$ ,  $d_0 = 1$ ,  $\mu_1 = C_4 \mu \log^2 m/2$ , and  $\varepsilon = 1/15$ . Then

$$\mathcal{S} \subseteq \mathcal{N}_{d_0} \cap \mathcal{N}_{\mu_1} \cap \mathcal{N}_\varepsilon,$$

and the sample complexity condition  $L \gg \mu_1^2 (K + N) \log^2 L$  is satisfied because we have assumed  $m \gg \mu^2 K \log^6 m$ . Therefore, with probability exceeding  $1 - O(m^{-10} + e^{-K})$ , we obtain that for all  $z \in \mathcal{S}$ ,

$$\|C\|_F^2 \leq \frac{5}{4} \|\mathbf{h}\mathbf{x}^H - \mathbf{h}^* \mathbf{x}^{*H}\|_F^2.$$

Claim (202b) can then be justified by observing that

$$\begin{aligned} \|\mathbf{h}\mathbf{x}^H - \mathbf{h}^* \mathbf{x}^{*H}\|_F &= \|\mathbf{h}(\mathbf{x} - \mathbf{x}^*)^H + (\mathbf{h} - \mathbf{h}^*) \mathbf{x}^{*H}\|_F \\ &\leq \|\mathbf{h}\|_2 \|\mathbf{x} - \mathbf{x}^*\|_2 + \|\mathbf{h} - \mathbf{h}^*\|_2 \|\mathbf{x}^*\|_2 \leq 3\delta. \end{aligned}$$

4. It remains to control  $\alpha_4$ , for which we make note of the following inequality

$$\alpha_4 \leq \underbrace{\left\| \sum_{j=1}^m \mathbf{b}_j \mathbf{b}_j^H (\mathbf{h}\mathbf{x}^T - \mathbf{h}^* \mathbf{x}^{*T}) \overline{\mathbf{a}_j} \overline{\mathbf{a}_j}^H \right\|}_{\theta_3} + \underbrace{\left\| \sum_{j=1}^m \mathbf{b}_j \mathbf{b}_j^H \mathbf{h}^* \mathbf{x}^{*T} (\overline{\mathbf{a}_j} \overline{\mathbf{a}_j}^H - \mathbf{I}_K) \right\|}_{\theta_4}$$

with  $\overline{\mathbf{a}_j}$  denoting the entrywise conjugate of  $\mathbf{a}_j$ . Since  $\{\overline{\mathbf{a}_j}\}$  has the same joint distribution as  $\{\mathbf{a}_j\}$ , by the same argument used for bounding  $\alpha_3$  we obtain control of the first term, namely

$$\mathbb{P}\left(\sup_{z \in \mathcal{S}} \theta_3 \geq 1/96\right) = O(m^{-10} + e^{-K} \log m).$$

Note that  $m \gg \mu^2 K \log m / \delta^2$  and  $\delta \ll 1$ . According to [76, Lemma 5.20],

$$\mathbb{P}\left(\sup_{z \in \mathcal{S}} \theta_4 \geq 1/96\right) \leq \mathbb{P}\left(\sup_{z \in \mathcal{S}} \theta_4 \geq \delta\right) = O(m^{-10}).$$

Putting together the above bounds, we reach  $\mathbb{P}(\sup_{z \in \mathcal{S}} \alpha_4 \leq 1/48) = 1 - O(m^{-10} + e^{-K} \log m)$ .

5. Combining all the previous bounds for  $\sup_{z \in \mathcal{S}} \alpha_j$  and (196), we deduce that with probability  $1 - O(m^{-10} + e^{-K} \log m)$ ,

$$\left\| \nabla^2 f(z) - \nabla^2 F(z^*) \right\| \leq 2 \cdot \frac{1}{32} + 2 \cdot \frac{1}{32} + 4 \cdot \frac{1}{96} + 4 \cdot \frac{1}{48} = \frac{1}{4}.$$

**C.2 Proofs of Lemmas 15 and 16**

*Proof of Lemma 15* In view of the definition of  $\alpha^{t+1}$  (see (38)), one has

$$\begin{aligned} \text{dist}(z^{t+1}, z^*)^2 &= \left\| \frac{1}{\alpha^{t+1}} \mathbf{h}^{t+1} - \mathbf{h}^* \right\|_2^2 + \left\| \alpha^{t+1} \mathbf{x}^{t+1} - \mathbf{x}^* \right\|_2^2 \\ &\leq \left\| \frac{1}{\alpha^t} \mathbf{h}^{t+1} - \mathbf{h}^* \right\|_2^2 + \left\| \alpha^t \mathbf{x}^{t+1} - \mathbf{x}^* \right\|_2^2. \end{aligned}$$

The gradient update rules (79) imply that

$$\begin{aligned} \frac{1}{\alpha^t} \mathbf{h}^{t+1} &= \frac{1}{\alpha^t} \left( \mathbf{h}^t - \frac{\eta}{\|\mathbf{x}^t\|_2^2} \nabla_{\mathbf{h}} f(z^t) \right) = \tilde{\mathbf{h}}^t - \frac{\eta}{\|\tilde{\mathbf{x}}^t\|_2^2} \nabla_{\mathbf{h}} f(\tilde{z}^t), \\ \alpha^t \mathbf{x}^{t+1} &= \alpha^t \left( \mathbf{x}^t - \frac{\eta}{\|\mathbf{h}^t\|_2^2} \nabla_{\mathbf{x}} f(z^t) \right) = \tilde{\mathbf{x}}^t - \frac{\eta}{\|\tilde{\mathbf{h}}^t\|_2^2} \nabla_{\mathbf{x}} f(\tilde{z}^t), \end{aligned}$$

where we denote  $\tilde{\mathbf{h}}^t = \frac{1}{\alpha^t} \mathbf{h}^t$  and  $\tilde{\mathbf{x}}^t = \alpha^t \mathbf{x}^t$  as in (81). Let  $\hat{\mathbf{h}}^{t+1} = \frac{1}{\alpha^t} \mathbf{h}^{t+1}$  and  $\hat{\mathbf{x}}^{t+1} = \alpha^t \mathbf{x}^{t+1}$ . We further get

$$\begin{aligned} \begin{bmatrix} \hat{\mathbf{h}}^{t+1} - \mathbf{h}^* \\ \hat{\mathbf{x}}^{t+1} - \mathbf{x}^* \\ \tilde{\mathbf{h}}^{t+1} - \mathbf{h}^* \\ \tilde{\mathbf{x}}^{t+1} - \mathbf{x}^* \end{bmatrix} &= \begin{bmatrix} \tilde{\mathbf{h}}^t - \mathbf{h}^* \\ \tilde{\mathbf{x}}^t - \mathbf{x}^* \\ \tilde{\mathbf{h}}^t - \mathbf{h}^* \\ \tilde{\mathbf{x}}^t - \mathbf{x}^* \end{bmatrix} - \eta \underbrace{\begin{bmatrix} \|\tilde{\mathbf{x}}^t\|_2^{-2} \mathbf{I}_K & & & \\ & \|\tilde{\mathbf{h}}^t\|_2^{-2} \mathbf{I}_K & & \\ & & \|\tilde{\mathbf{x}}^t\|_2^{-2} \mathbf{I}_K & \\ & & & \|\tilde{\mathbf{h}}^t\|_2^{-2} \mathbf{I}_K \end{bmatrix}}_{:=D} \\ &\quad \begin{bmatrix} \nabla_{\mathbf{h}} f(\tilde{z}^t) \\ \nabla_{\mathbf{x}} f(\tilde{z}^t) \\ \nabla_{\mathbf{h}} f(\tilde{z}^t) \\ \nabla_{\mathbf{x}} f(\tilde{z}^t) \end{bmatrix}. \end{aligned} \tag{206}$$



The fundamental theorem of calculus (see Appendix D.3.1) together with the fact that  $\nabla f(z^*) = \mathbf{0}$  tells us

$$\begin{bmatrix} \nabla_{\mathbf{h}} f(\tilde{\mathbf{z}}^t) \\ \nabla_{\mathbf{x}} f(\tilde{\mathbf{z}}^t) \\ \nabla_{\mathbf{h}} f(\tilde{\mathbf{z}}^t) \\ \nabla_{\mathbf{x}} f(\tilde{\mathbf{z}}^t) \end{bmatrix} = \begin{bmatrix} \nabla_{\mathbf{h}} f(\tilde{\mathbf{z}}^t) - \nabla_{\mathbf{h}} f(\mathbf{z}^*) \\ \nabla_{\mathbf{x}} f(\tilde{\mathbf{z}}^t) - \nabla_{\mathbf{x}} f(\mathbf{z}^*) \\ \nabla_{\mathbf{h}} f(\tilde{\mathbf{z}}^t) - \nabla_{\mathbf{h}} f(\mathbf{z}^*) \\ \nabla_{\mathbf{x}} f(\tilde{\mathbf{z}}^t) - \nabla_{\mathbf{x}} f(\mathbf{z}^*) \end{bmatrix} = \underbrace{\int_0^1 \nabla^2 f(\mathbf{z}(\tau)) \, d\tau}_{:=\mathbf{A}} \begin{bmatrix} \tilde{\mathbf{h}}^t - \mathbf{h}^* \\ \tilde{\mathbf{x}}^t - \mathbf{x}^* \\ \tilde{\mathbf{h}}^t - \mathbf{h}^* \\ \tilde{\mathbf{x}}^t - \mathbf{x}^* \end{bmatrix}, \tag{207}$$

where we denote  $\mathbf{z}(\tau) := \mathbf{z}^* + \tau(\tilde{\mathbf{z}}^t - \mathbf{z}^*)$  and  $\nabla^2 f$  is the Wirtinger Hessian. To further simplify notation, denote  $\tilde{\mathbf{z}}^{t+1} = \begin{bmatrix} \tilde{\mathbf{h}}^{t+1} \\ \tilde{\mathbf{x}}^{t+1} \end{bmatrix}$ . Identity (207) allows us to rewrite (206) as

$$\begin{bmatrix} \tilde{\mathbf{z}}^{t+1} - \mathbf{z}^* \\ \tilde{\mathbf{z}}^{t+1} - \mathbf{z}^* \end{bmatrix} = (\mathbf{I} - \eta \mathbf{DA}) \begin{bmatrix} \tilde{\mathbf{z}}^t - \mathbf{z}^* \\ \tilde{\mathbf{z}}^t - \mathbf{z}^* \end{bmatrix}. \tag{208}$$

Take the squared Euclidean norm of both sides of (208) to reach

$$\begin{aligned} \|\tilde{\mathbf{z}}^{t+1} - \mathbf{z}^*\|_2^2 &= \frac{1}{2} \begin{bmatrix} \tilde{\mathbf{z}}^t - \mathbf{z}^* \\ \tilde{\mathbf{z}}^t - \mathbf{z}^* \end{bmatrix}^H (\mathbf{I} - \eta \mathbf{DA})^H (\mathbf{I} - \eta \mathbf{DA}) \begin{bmatrix} \tilde{\mathbf{z}}^t - \mathbf{z}^* \\ \tilde{\mathbf{z}}^t - \mathbf{z}^* \end{bmatrix} \\ &= \frac{1}{2} \begin{bmatrix} \tilde{\mathbf{z}}^t - \mathbf{z}^* \\ \tilde{\mathbf{z}}^t - \mathbf{z}^* \end{bmatrix}^H \left( \mathbf{I} + \eta^2 \mathbf{AD}^2 \mathbf{A} - \eta(\mathbf{DA} + \mathbf{AD}) \right) \begin{bmatrix} \tilde{\mathbf{z}}^t - \mathbf{z}^* \\ \tilde{\mathbf{z}}^t - \mathbf{z}^* \end{bmatrix} \\ &\leq (1 + \eta^2 \|\mathbf{A}\|^2 \|\mathbf{D}\|^2) \|\tilde{\mathbf{z}}^t - \mathbf{z}^*\|_2^2 - \frac{\eta}{2} \begin{bmatrix} \tilde{\mathbf{z}}^t - \mathbf{z}^* \\ \tilde{\mathbf{z}}^t - \mathbf{z}^* \end{bmatrix}^H (\mathbf{DA} + \mathbf{AD}) \begin{bmatrix} \tilde{\mathbf{z}}^t - \mathbf{z}^* \\ \tilde{\mathbf{z}}^t - \mathbf{z}^* \end{bmatrix}. \end{aligned} \tag{209}$$

Since  $\mathbf{z}(\tau)$  lies between  $\tilde{\mathbf{z}}^t$  and  $\mathbf{z}^*$ , we conclude from assumptions (85) that for all  $0 \leq \tau \leq 1$ ,

$$\begin{aligned} \max \{ \|\mathbf{h}(\tau) - \mathbf{h}^*\|_2, \|\mathbf{x}(\tau) - \mathbf{x}^*\|_2 \} &\leq \text{dist}(\mathbf{z}^t, \mathbf{z}^*) \leq \xi \leq \delta; \\ \max_{1 \leq j \leq m} \left| \mathbf{a}_j^H(\mathbf{x}(\tau) - \mathbf{x}^*) \right| &\leq C_3 \frac{1}{\log^{3/2} m}; \\ \max_{1 \leq j \leq m} \left| \mathbf{b}_j^H \mathbf{h}(\tau) \right| &\leq C_4 \frac{\mu}{\sqrt{m}} \log^2 m \end{aligned}$$

for  $\xi > 0$  sufficiently small. Moreover, it is straightforward to see that

$$\gamma_1 := \|\tilde{\mathbf{x}}^t\|_2^{-2} \quad \text{and} \quad \gamma_2 := \|\tilde{\mathbf{h}}^t\|_2^{-2}$$

satisfy

$$\max \{ |\gamma_1 - 1|, |\gamma_2 - 1| \} \lesssim \max \left\{ \|\tilde{\mathbf{h}}^t - \mathbf{h}^*\|_2, \|\tilde{\mathbf{x}}^t - \mathbf{x}^*\|_2 \right\} \leq \delta$$

as long as  $\xi > 0$  is sufficiently small. We can now readily invoke Lemma 14 to arrive at

$$\|A\| \|D\| \leq 3(1 + \delta) \leq 4 \quad \text{and}$$

$$\begin{bmatrix} \tilde{z}^t - z^* \\ \tilde{z}^t - z^* \end{bmatrix}^H (DA + AD) \begin{bmatrix} \tilde{z}^t - z^* \\ \tilde{z}^t - z^* \end{bmatrix} \geq \frac{1}{4} \left\| \begin{bmatrix} \tilde{z}^t - z^* \\ \tilde{z}^t - z^* \end{bmatrix} \right\|_2^2 = \frac{1}{2} \|\tilde{z}^t - z^*\|_2^2.$$

Substitution into (209) indicates that

$$\|\tilde{z}^{t+1} - z^*\|_2^2 \leq (1 + 16\eta^2 - \eta/4) \|\tilde{z}^t - z^*\|_2^2.$$

When  $0 < \eta \leq 1/128$ , this implies that

$$\|\tilde{z}^t - z^*\|_2^2 \leq (1 - \eta/8) \|\tilde{z}^t - z^*\|_2^2,$$

and hence

$$\|\tilde{z}^{t+1} - z^*\|_2 \leq \|\tilde{z}^{t+1} - z^*\|_2 \leq (1 - \eta/8)^{1/2} \|\tilde{z}^t - z^*\|_2 \leq (1 - \eta/16) \text{dist}(z^t, z^*). \tag{210}$$

This completes the proof of Lemma 15. □

**Proof of Lemma 16** Reuse the notation in this subsection, namely  $\hat{z}^{t+1} = \begin{bmatrix} \hat{h}^{t+1} \\ \hat{x}^{t+1} \end{bmatrix}$

with  $\hat{h}^{t+1} = \frac{1}{\alpha^t} h^{t+1}$  and  $\hat{x}^{t+1} = \alpha^t x^{t+1}$ . From (210), one can tell that

$$\|\hat{z}^{t+1} - z^*\|_2 \leq \|\hat{z}^{t+1} - z^*\|_2 \leq \text{dist}(z^t, z^*).$$

Invoke Lemma 52 with  $\beta = \alpha^t$  to get

$$|\alpha^{t+1} - \alpha^t| \lesssim \|\hat{z}^{t+1} - z^*\|_2 \leq \text{dist}(z^t, z^*).$$

This combined with the assumption  $||\alpha^t| - 1| \leq 1/2$  implies that

$$|\alpha^t| \geq \frac{1}{2} \quad \text{and} \quad \left| \frac{\alpha^{t+1}}{\alpha^t} - 1 \right| = \left| \frac{\alpha^{t+1} - \alpha^t}{\alpha^t} \right| \lesssim \text{dist}(z^t, z^*) \lesssim C_1 \frac{1}{\log^2 m}.$$

This finishes the proof of the first claim.

The second claim can be proved by induction. Suppose that  $||\alpha^s| - 1| \leq 1/2$  and  $\text{dist}(z^s, z^*) \leq C_1(1 - \eta/16)^s / \log^2 m$  hold for all  $0 \leq s \leq \tau \leq t$ , then using our result in the first part gives

$$\begin{aligned} \left| |\alpha^{\tau+1}| - 1 \right| &\leq \left| |\alpha^0| - 1 \right| + \sum_{s=0}^{\tau} \left| \alpha^{s+1} - \alpha^s \right| \leq \frac{1}{4} + c \sum_{s=0}^{\tau} \text{dist}(z^s, z^*) \\ &\leq \frac{1}{4} + \frac{cC_1}{\frac{\eta}{16} \log^2 m} \leq \frac{1}{2} \end{aligned}$$

for  $m$  sufficiently large. The proof is then complete by induction. □

### C.3 Proof of Lemma 17

Define the alignment parameter between  $z^{t,(l)}$  and  $\tilde{z}^t$  as

$$\alpha_{\text{mutual}}^{t,(l)} := \operatorname{argmin}_{\alpha \in \mathbb{C}} \left\| \frac{1}{\alpha} \mathbf{h}^{t,(l)} - \frac{1}{\alpha^t} \mathbf{h}^t \right\|_2^2 + \left\| \alpha \mathbf{x}^{t,(l)} - \alpha^t \mathbf{x}^t \right\|_2^2.$$

Further denote, for simplicity of presentation,  $\hat{z}^{t,(l)} = \begin{bmatrix} \hat{\mathbf{h}}^{t,(l)} \\ \hat{\mathbf{x}}^{t,(l)} \end{bmatrix}$  with

$$\hat{\mathbf{h}}^{t,(l)} := \frac{1}{\alpha_{\text{mutual}}^{t,(l)}} \mathbf{h}^{t,(l)} \quad \text{and} \quad \hat{\mathbf{x}}^{t,(l)} := \alpha_{\text{mutual}}^{t,(l)} \mathbf{x}^{t,(l)}.$$

Clearly,  $\hat{z}^{t,(l)}$  is aligned with  $\tilde{z}^t$ .

Armed with the above notation, we have

$$\begin{aligned} \text{dist}(z^{t+1,(l)}, \tilde{z}^{t+1}) &= \min_{\alpha} \sqrt{\left\| \frac{1}{\alpha} \mathbf{h}^{t+1,(l)} - \frac{1}{\alpha^{t+1}} \mathbf{h}^{t+1} \right\|_2^2 + \left\| \alpha \mathbf{x}^{t+1,(l)} - \alpha^{t+1} \mathbf{x}^{t+1} \right\|_2^2} \\ &= \min_{\alpha} \sqrt{\left\| \left( \frac{\bar{\alpha}^t}{\alpha^{t+1}} \right) \left( \frac{1}{\bar{\alpha}} \frac{\alpha^{t+1}}{\alpha^t} \mathbf{h}^{t+1,(l)} - \frac{1}{\alpha^t} \mathbf{h}^{t+1} \right) \right\|_2^2 + \left\| \left( \frac{\alpha^{t+1}}{\alpha^t} \right) \left( \alpha \frac{\alpha^t}{\alpha^{t+1}} \mathbf{x}^{t+1,(l)} - \alpha^t \mathbf{x}^{t+1} \right) \right\|_2^2} \\ &\leq \sqrt{\left\| \left( \frac{\bar{\alpha}^t}{\alpha^{t+1}} \right) \left( \frac{1}{\alpha_{\text{mutual}}^{t,(l)}} \mathbf{h}^{t+1,(l)} - \frac{1}{\alpha^t} \mathbf{h}^{t+1} \right) \right\|_2^2 + \left\| \left( \frac{\alpha^{t+1}}{\alpha^t} \right) \left( \alpha_{\text{mutual}}^{t,(l)} \mathbf{x}^{t+1,(l)} - \alpha^t \mathbf{x}^{t+1} \right) \right\|_2^2} \end{aligned} \tag{211}$$

$$\leq \max \left\{ \left| \frac{\alpha^{t+1}}{\alpha^t} \right|, \left| \frac{\alpha^t}{\alpha^{t+1}} \right| \right\} \sqrt{\left\| \frac{1}{\alpha_{\text{mutual}}^{t,(l)}} \mathbf{h}^{t+1,(l)} - \frac{1}{\alpha^t} \mathbf{h}^{t+1} \right\|_2^2}, \tag{212}$$

where (211) follows by taking  $\alpha = \frac{\alpha^{t+1}}{\alpha^t} \alpha_{\text{mutual}}^{t,(l)}$ . The latter bound is more convenient to work with when controlling the gap between  $z^{t,(l)}$  and  $z^t$ .

We can then apply the gradient update rules (79) and (89) to get

$$\begin{bmatrix} \frac{1}{\alpha_{\text{mutual}}^{t,(l)}} \mathbf{h}^{t+1,(l)} - \frac{1}{\alpha^t} \mathbf{h}^{t+1} \\ \alpha_{\text{mutual}}^{t,(l)} \mathbf{x}^{t+1,(l)} - \alpha^t \mathbf{x}^{t+1} \end{bmatrix}$$

$$\begin{aligned}
 &= \left[ \begin{aligned} &\frac{1}{\alpha_{\text{mutual}}^{t,(l)}} \left( \mathbf{h}^{t,(l)} - \frac{\eta}{\|\mathbf{x}^{t,(l)}\|_2^2} \nabla_{\mathbf{h}} f^{(l)} \left( \mathbf{h}^{t,(l)}, \mathbf{x}^{t,(l)} \right) \right) - \frac{1}{\alpha^t} \left( \mathbf{h}^t - \frac{\eta}{\|\mathbf{x}^t\|_2^2} \nabla_{\mathbf{h}} f \left( \mathbf{h}^t, \mathbf{x}^t \right) \right) \\ &\alpha_{\text{mutual}}^{t,(l)} \left( \mathbf{x}^{t,(l)} - \frac{\eta}{\|\mathbf{h}^{t,(l)}\|_2^2} \nabla_{\mathbf{x}} f^{(l)} \left( \mathbf{h}^{t,(l)}, \mathbf{x}^{t,(l)} \right) \right) - \alpha^t \left( \mathbf{x}^t - \frac{\eta}{\|\mathbf{h}^t\|_2^2} \nabla_{\mathbf{x}} f \left( \mathbf{h}^t, \mathbf{x}^t \right) \right) \end{aligned} \right] \\
 &= \left[ \begin{aligned} &\widehat{\mathbf{h}}^{t,(l)} - \frac{\eta}{\|\widehat{\mathbf{x}}^{t,(l)}\|_2^2} \nabla_{\mathbf{h}} f^{(l)} \left( \widehat{\mathbf{h}}^{t,(l)}, \widehat{\mathbf{x}}^{t,(l)} \right) - \left( \widetilde{\mathbf{h}}^t - \frac{\eta}{\|\widetilde{\mathbf{x}}^t\|_2^2} \nabla_{\mathbf{h}} f \left( \widetilde{\mathbf{h}}^t, \widetilde{\mathbf{x}}^t \right) \right) \\ &\widehat{\mathbf{x}}^{t,(l)} - \frac{\eta}{\|\widehat{\mathbf{h}}^{t,(l)}\|_2^2} \nabla_{\mathbf{x}} f^{(l)} \left( \widehat{\mathbf{h}}^{t,(l)}, \widehat{\mathbf{x}}^{t,(l)} \right) - \left( \widetilde{\mathbf{x}}^t - \frac{\eta}{\|\widetilde{\mathbf{h}}^t\|_2^2} \nabla_{\mathbf{x}} f \left( \widetilde{\mathbf{h}}^t, \widetilde{\mathbf{x}}^t \right) \right) \end{aligned} \right].
 \end{aligned}$$

By construction, we can write the leave-one-out gradients as

$$\begin{aligned}
 \nabla_{\mathbf{h}} f^{(l)}(\mathbf{h}, \mathbf{x}) &= \nabla_{\mathbf{h}} f(\mathbf{h}, \mathbf{x}) - \left( \mathbf{b}_l^H \mathbf{h} \mathbf{x}^H \mathbf{a}_l - y_l \right) \mathbf{b}_l \mathbf{a}_l^H \mathbf{x} \quad \text{and} \\
 \nabla_{\mathbf{x}} f^{(l)}(\mathbf{h}, \mathbf{x}) &= \nabla_{\mathbf{x}} f(\mathbf{h}, \mathbf{x}) - \overline{\left( \mathbf{b}_l^H \mathbf{h} \mathbf{x}^H \mathbf{a}_l - y_l \right)} \mathbf{a}_l \mathbf{b}_l^H \mathbf{h},
 \end{aligned}$$

which allow us to continue the derivation and obtain

$$\begin{aligned}
 &\left[ \begin{aligned} &\frac{1}{\alpha_{\text{mutual}}^{t,(l)}} \mathbf{h}^{t+1,(l)} - \frac{1}{\alpha^t} \mathbf{h}^{t+1} \\ &\alpha_{\text{mutual}}^{t,(l)} \mathbf{x}^{t+1,(l)} - \alpha^t \mathbf{x}^{t+1} \end{aligned} \right] \\
 &= \left[ \begin{aligned} &\widehat{\mathbf{h}}^{t,(l)} - \frac{\eta}{\|\widehat{\mathbf{x}}^{t,(l)}\|_2^2} \nabla_{\mathbf{h}} f \left( \widehat{\mathbf{h}}^{t,(l)}, \widehat{\mathbf{x}}^{t,(l)} \right) - \left( \widetilde{\mathbf{h}}^t - \frac{\eta}{\|\widetilde{\mathbf{x}}^t\|_2^2} \nabla_{\mathbf{h}} f \left( \widetilde{\mathbf{h}}^t, \widetilde{\mathbf{x}}^t \right) \right) \\ &\widehat{\mathbf{x}}^{t,(l)} - \frac{\eta}{\|\widehat{\mathbf{h}}^{t,(l)}\|_2^2} \nabla_{\mathbf{x}} f \left( \widehat{\mathbf{h}}^{t,(l)}, \widehat{\mathbf{x}}^{t,(l)} \right) - \left( \widetilde{\mathbf{x}}^t - \frac{\eta}{\|\widetilde{\mathbf{h}}^t\|_2^2} \nabla_{\mathbf{x}} f \left( \widetilde{\mathbf{h}}^t, \widetilde{\mathbf{x}}^t \right) \right) \end{aligned} \right] \\
 &\quad - \eta \underbrace{\left[ \begin{aligned} &\frac{1}{\|\widehat{\mathbf{x}}^{t,(l)}\|_2^2} \left( \mathbf{b}_l^H \widehat{\mathbf{h}}^{t,(l)} \widehat{\mathbf{x}}^{t,(l)H} \mathbf{a}_l - y_l \right) \mathbf{b}_l \mathbf{a}_l^H \widehat{\mathbf{x}}^{t,(l)} \\ &\frac{1}{\|\widehat{\mathbf{h}}^{t,(l)}\|_2^2} \left( \mathbf{b}_l^H \widehat{\mathbf{h}}^{t,(l)} \widehat{\mathbf{x}}^{t,(l)H} \mathbf{a}_l - y_l \right) \mathbf{a}_l \mathbf{b}_l^H \widehat{\mathbf{h}}^{t,(l)} \end{aligned} \right]}_{:=J_3}.
 \end{aligned}$$

This further gives

$$\begin{aligned}
 &\left[ \begin{aligned} &\frac{1}{\alpha_{\text{mutual}}^{t,(l)}} \mathbf{h}^{t+1,(l)} - \frac{1}{\alpha^t} \mathbf{h}^{t+1} \\ &\alpha_{\text{mutual}}^{t,(l)} \mathbf{x}^{t+1,(l)} - \alpha^t \mathbf{x}^{t+1} \end{aligned} \right] \\
 &= \underbrace{\left[ \begin{aligned} &\widehat{\mathbf{h}}^{t,(l)} - \frac{\eta}{\|\widehat{\mathbf{x}}^{t,(l)}\|_2^2} \nabla_{\mathbf{h}} f \left( \widehat{\mathbf{h}}^{t,(l)}, \widehat{\mathbf{x}}^{t,(l)} \right) - \left( \widetilde{\mathbf{h}}^t - \frac{\eta}{\|\widetilde{\mathbf{x}}^t\|_2^2} \nabla_{\mathbf{h}} f \left( \widetilde{\mathbf{h}}^t, \widetilde{\mathbf{x}}^t \right) \right) \\ &\widehat{\mathbf{x}}^{t,(l)} - \frac{\eta}{\|\widehat{\mathbf{h}}^{t,(l)}\|_2^2} \nabla_{\mathbf{x}} f \left( \widehat{\mathbf{h}}^{t,(l)}, \widehat{\mathbf{x}}^{t,(l)} \right) - \left( \widetilde{\mathbf{x}}^t - \frac{\eta}{\|\widetilde{\mathbf{h}}^t\|_2^2} \nabla_{\mathbf{x}} f \left( \widetilde{\mathbf{h}}^t, \widetilde{\mathbf{x}}^t \right) \right) \end{aligned} \right]}_{:=v_1}
 \end{aligned}$$

$$+ \eta \underbrace{\left[ \begin{array}{c} \left( \frac{1}{\|\tilde{\mathbf{x}}^t\|_2^2} - \frac{1}{\|\tilde{\mathbf{x}}^{t,(l)}\|_2^2} \right) \nabla_{\mathbf{h}} f(\tilde{\mathbf{h}}^t, \tilde{\mathbf{x}}^t) \\ \left( \frac{1}{\|\tilde{\mathbf{h}}^t\|_2^2} - \frac{1}{\|\tilde{\mathbf{h}}^{t,(l)}\|_2^2} \right) \nabla_{\mathbf{x}} f(\tilde{\mathbf{h}}^t, \tilde{\mathbf{x}}^t) \end{array} \right]}_{:=\mathbf{v}_2} - \eta \mathbf{v}_3. \tag{213}$$

In what follows, we bound the three terms  $\mathbf{v}_1$ ,  $\mathbf{v}_2$ , and  $\mathbf{v}_3$  separately.

1. Regarding the first term  $\mathbf{v}_1$ , one can adopt the same strategy as in Appendix C.2. Specifically, write

$$\begin{aligned}
 & \left[ \begin{array}{c} \widehat{\mathbf{h}}^{t,(l)} - \frac{\eta}{\|\widehat{\mathbf{x}}^{t,(l)}\|_2^2} \nabla_{\mathbf{h}} f(\widehat{\mathbf{z}}^{t,(l)}) - \left( \tilde{\mathbf{h}}^t - \frac{\eta}{\|\tilde{\mathbf{x}}^{t,(l)}\|_2^2} \nabla_{\mathbf{h}} f(\tilde{\mathbf{z}}^t) \right) \\ \widehat{\mathbf{x}}^{t,(l)} - \frac{\eta}{\|\widehat{\mathbf{h}}^{t,(l)}\|_2^2} \nabla_{\mathbf{x}} f(\widehat{\mathbf{z}}^{t,(l)}) - \left( \tilde{\mathbf{x}}^t - \frac{\eta}{\|\tilde{\mathbf{h}}^{t,(l)}\|_2^2} \nabla_{\mathbf{x}} f(\tilde{\mathbf{z}}^t) \right) \\ \widehat{\mathbf{h}}^{t,(l)} - \frac{\eta}{\|\widehat{\mathbf{x}}^{t,(l)}\|_2^2} \nabla_{\mathbf{h}} f(\widehat{\mathbf{z}}^{t,(l)}) - \left( \tilde{\mathbf{h}}^t - \frac{\eta}{\|\tilde{\mathbf{x}}^{t,(l)}\|_2^2} \nabla_{\mathbf{h}} f(\tilde{\mathbf{z}}^t) \right) \\ \widehat{\mathbf{x}}^{t,(l)} - \frac{\eta}{\|\widehat{\mathbf{h}}^{t,(l)}\|_2^2} \nabla_{\mathbf{x}} f(\widehat{\mathbf{z}}^{t,(l)}) - \left( \tilde{\mathbf{x}}^t - \frac{\eta}{\|\tilde{\mathbf{h}}^{t,(l)}\|_2^2} \nabla_{\mathbf{x}} f(\tilde{\mathbf{z}}^t) \right) \end{array} \right] = \left[ \begin{array}{c} \widehat{\mathbf{h}}^{t,(l)} - \tilde{\mathbf{h}}^t \\ \widehat{\mathbf{x}}^{t,(l)} - \tilde{\mathbf{x}}^t \\ \widehat{\mathbf{h}}^{t,(l)} - \tilde{\mathbf{h}}^t \\ \widehat{\mathbf{x}}^{t,(l)} - \tilde{\mathbf{x}}^t \end{array} \right] \\
 & - \eta \underbrace{\left[ \begin{array}{cccc} \|\widehat{\mathbf{x}}^{t,(l)}\|_2^{-2} \mathbf{I}_K & & & \\ & \|\widehat{\mathbf{h}}^{t,(l)}\|_2^{-2} \mathbf{I}_K & & \\ & & \|\widehat{\mathbf{x}}^{t,(l)}\|_2^{-2} \mathbf{I}_K & \\ & & & \|\widehat{\mathbf{h}}^{t,(l)}\|_2^{-2} \mathbf{I}_K \end{array} \right]}_{:=\mathbf{D}} \\
 & \left[ \begin{array}{c} \nabla_{\mathbf{h}} f(\widehat{\mathbf{z}}^{t,(l)}) - \nabla_{\mathbf{h}} f(\tilde{\mathbf{z}}^t) \\ \nabla_{\mathbf{x}} f(\widehat{\mathbf{z}}^{t,(l)}) - \nabla_{\mathbf{x}} f(\tilde{\mathbf{z}}^t) \\ \nabla_{\mathbf{h}} f(\widehat{\mathbf{z}}^{t,(l)}) - \nabla_{\mathbf{h}} f(\tilde{\mathbf{z}}^t) \\ \nabla_{\mathbf{x}} f(\widehat{\mathbf{z}}^{t,(l)}) - \nabla_{\mathbf{x}} f(\tilde{\mathbf{z}}^t) \end{array} \right].
 \end{aligned}$$

The fundamental theorem of calculus (see Appendix D.3.1) reveals that

$$\left[ \begin{array}{c} \nabla_{\mathbf{h}} f(\widehat{\mathbf{z}}^{t,(l)}) - \nabla_{\mathbf{h}} f(\tilde{\mathbf{z}}^t) \\ \nabla_{\mathbf{x}} f(\widehat{\mathbf{z}}^{t,(l)}) - \nabla_{\mathbf{x}} f(\tilde{\mathbf{z}}^t) \\ \nabla_{\mathbf{h}} f(\widehat{\mathbf{z}}^{t,(l)}) - \nabla_{\mathbf{h}} f(\tilde{\mathbf{z}}^t) \\ \nabla_{\mathbf{x}} f(\widehat{\mathbf{z}}^{t,(l)}) - \nabla_{\mathbf{x}} f(\tilde{\mathbf{z}}^t) \end{array} \right] = \underbrace{\int_0^1 \nabla^2 f(\mathbf{z}(\tau)) \, d\tau}_{:=\mathbf{A}} \left[ \begin{array}{c} \widehat{\mathbf{h}}^{t,(l)} - \tilde{\mathbf{h}}^t \\ \widehat{\mathbf{x}}^{t,(l)} - \tilde{\mathbf{x}}^t \\ \widehat{\mathbf{h}}^{t,(l)} - \tilde{\mathbf{h}}^t \\ \widehat{\mathbf{x}}^{t,(l)} - \tilde{\mathbf{x}}^t \end{array} \right],$$

where we abuse the notation and denote  $\mathbf{z}(\tau) = \tilde{\mathbf{z}}^t + \tau(\widehat{\mathbf{z}}^{t,(l)} - \tilde{\mathbf{z}}^t)$ . In order to invoke Lemma 14, we need to verify the conditions required therein. Recall the induction hypothesis (90b) that

$$\text{dist}(\mathbf{z}^{t,(l)}, \tilde{\mathbf{z}}^t) = \|\widehat{\mathbf{z}}^{t,(l)} - \tilde{\mathbf{z}}^t\|_2 \leq C_2 \frac{\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 K \log^9 m}{m}},$$

and the fact that  $\mathbf{z}(\tau)$  lies between  $\widehat{\mathbf{z}}^{t,(l)}$  and  $\widetilde{\mathbf{z}}^t$ . For all  $0 \leq \tau \leq 1$ :

(a) If  $m \gg \mu^2 \sqrt{K} \log^{13/2} m$ , then

$$\begin{aligned} \|\mathbf{z}(\tau) - \mathbf{z}^*\|_2 &\leq \max \left\{ \|\widehat{\mathbf{z}}^{t,(l)} - \mathbf{z}^*\|_2, \|\widetilde{\mathbf{z}}^t - \mathbf{z}^*\|_2 \right\} \leq \|\widetilde{\mathbf{z}}^t - \mathbf{z}^*\|_2 + \|\widehat{\mathbf{z}}^{t,(l)} - \widetilde{\mathbf{z}}^t\|_2 \\ &\leq C_1 \frac{1}{\log^2 m} + C_2 \frac{\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 K \log^9 m}{m}} \leq 2C_1 \frac{1}{\log^2 m}, \end{aligned}$$

where we have used the induction hypotheses (90a) and (90b);

(b) If  $m \gg \mu^2 K \log^6 m$ , then

$$\begin{aligned} &\max_{1 \leq j \leq m} \left| \mathbf{a}_j^H (\mathbf{x}(\tau) - \mathbf{x}^*) \right| \\ &= \max_{1 \leq j \leq m} \left| \tau \mathbf{a}_j^H (\widehat{\mathbf{x}}^{t,(l)} - \widetilde{\mathbf{x}}^t) + \mathbf{a}_j^H (\widetilde{\mathbf{x}}^t - \mathbf{x}^*) \right| \\ &\leq \max_{1 \leq j \leq m} \left| \mathbf{a}_j^H (\widehat{\mathbf{x}}^{t,(l)} - \widetilde{\mathbf{x}}^t) \right| + \max_{1 \leq j \leq m} \left| \mathbf{a}_j^H (\widetilde{\mathbf{x}}^t - \mathbf{x}^*) \right| \\ &\leq \max_{1 \leq j \leq m} \|\mathbf{a}_j\|_2 \|\widehat{\mathbf{x}}^{t,(l)} - \widetilde{\mathbf{x}}^t\|_2 + C_3 \frac{1}{\log^{3/2} m} \\ &\leq 3\sqrt{K} \cdot C_2 \frac{\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 K \log^9 m}{m}} + C_3 \frac{1}{\log^{3/2} m} \leq 2C_3 \frac{1}{\log^{3/2} m}, \end{aligned} \tag{214}$$

which follows from bound (190) and the induction hypotheses (90b) and (90c);

(c) If  $m \gg \mu K \log^{5/2} m$ , then

$$\begin{aligned} \max_{1 \leq j \leq m} \left| \mathbf{b}_j^H \mathbf{h}(\tau) \right| &= \max_{1 \leq j \leq m} \left| \tau \mathbf{b}_j^H (\widehat{\mathbf{h}}^{t,(l)} - \widetilde{\mathbf{h}}^t) + \mathbf{b}_j^H \widetilde{\mathbf{h}}^t \right| \\ &\leq \max_{1 \leq j \leq m} \left| \mathbf{b}_j^H (\widehat{\mathbf{h}}^{t,(l)} - \widetilde{\mathbf{h}}^t) \right| + \max_{1 \leq j \leq m} \left| \mathbf{b}_j^H \widetilde{\mathbf{h}}^t \right| \\ &\leq \max_{1 \leq j \leq m} \|\mathbf{b}_j\|_2 \|\widehat{\mathbf{h}}^{t,(l)} - \widetilde{\mathbf{h}}^t\|_2 + \max_{1 \leq j \leq m} \left| \mathbf{b}_j^H \widetilde{\mathbf{h}}^t \right| \\ &\leq \sqrt{\frac{K}{m}} \cdot C_2 \frac{\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 K \log^9 m}{m}} + C_4 \frac{\mu}{\sqrt{m}} \log^2 m \leq 2C_4 \frac{\mu}{\sqrt{m}} \log^2 m, \end{aligned} \tag{215}$$

which makes use of the fact  $\|\mathbf{b}_j\|_2 = \sqrt{K/m}$  as well as the induction hypotheses (90b) and (90d).

These properties satisfy condition (82) required in Lemma 14. The other two conditions (83) and (84) are also straightforward to check, and hence, we omit it. Thus, we can repeat the argument used in Appendix C.2 to obtain

$$\|\mathbf{v}_1\|_2 \leq (1 - \eta/16) \cdot \|\widehat{\mathbf{z}}^{t,(l)} - \widetilde{\mathbf{z}}^t\|_2.$$

2. In terms of the second term  $\mathbf{v}_2$ , it is easily seen that

$$\|\mathbf{v}_2\|_2 \leq \max \left\{ \left| \frac{1}{\|\hat{\mathbf{x}}^t\|_2^2} - \frac{1}{\|\hat{\mathbf{x}}^{t,(l)}\|_2^2} \right|, \left| \frac{1}{\|\hat{\mathbf{h}}^t\|_2^2} - \frac{1}{\|\hat{\mathbf{h}}^{t,(l)}\|_2^2} \right| \right\} \left\| \begin{bmatrix} \nabla_{\mathbf{h}} f(\tilde{\mathbf{z}}^t) \\ \nabla_{\mathbf{x}} f(\tilde{\mathbf{z}}^t) \end{bmatrix} \right\|_2.$$

We first note that the upper bound on  $\|\nabla^2 f(\cdot)\|$  (which essentially provides a Lipschitz constant on the gradient) in Lemma 14 forces

$$\left\| \begin{bmatrix} \nabla_{\mathbf{h}} f(\tilde{\mathbf{z}}^t) \\ \nabla_{\mathbf{x}} f(\tilde{\mathbf{z}}^t) \end{bmatrix} \right\|_2 = \left\| \begin{bmatrix} \nabla_{\mathbf{h}} f(\tilde{\mathbf{z}}^t) - \nabla_{\mathbf{h}} f(\mathbf{z}^*) \\ \nabla_{\mathbf{x}} f(\tilde{\mathbf{z}}^t) - \nabla_{\mathbf{x}} f(\mathbf{z}^*) \end{bmatrix} \right\|_2 \lesssim \|\tilde{\mathbf{z}}^t - \mathbf{z}^*\|_2 \lesssim C_1 \frac{1}{\log^2 m},$$

where the first identity follows since  $\nabla_{\mathbf{h}} f(\mathbf{z}^*) = \mathbf{0}$ , and the last inequality comes from the induction hypothesis (90a). Additionally, recognizing that  $\|\hat{\mathbf{x}}^t\|_2 \asymp \|\hat{\mathbf{x}}^{t,(l)}\|_2 \asymp 1$ , one can easily verify that

$$\left| \frac{1}{\|\hat{\mathbf{x}}^t\|_2^2} - \frac{1}{\|\hat{\mathbf{x}}^{t,(l)}\|_2^2} \right| = \left| \frac{\|\hat{\mathbf{x}}^{t,(l)}\|_2^2 - \|\hat{\mathbf{x}}^t\|_2^2}{\|\hat{\mathbf{x}}^t\|_2^2 \cdot \|\hat{\mathbf{x}}^{t,(l)}\|_2^2} \right| \lesssim \left| \|\hat{\mathbf{x}}^{t,(l)}\|_2 - \|\hat{\mathbf{x}}^t\|_2 \right| \lesssim \|\hat{\mathbf{x}}^{t,(l)} - \hat{\mathbf{x}}^t\|_2.$$

A similar bound holds for the other term involving  $\mathbf{h}$ . Combining the estimates above thus yields

$$\|\mathbf{v}_2\|_2 \lesssim C_1 \frac{1}{\log^2 m} \|\hat{\mathbf{z}}^{t,(l)} - \tilde{\mathbf{z}}^t\|_2.$$

3. When it comes to the last term  $\mathbf{v}_3$ , one first sees that

$$\left\| \left( \mathbf{b}_l^H \hat{\mathbf{h}}^{t,(l)} \hat{\mathbf{x}}^{t,(l)H} \mathbf{a}_l - y_l \right) \mathbf{b}_l \mathbf{a}_l^H \hat{\mathbf{x}}^{t,(l)} \right\|_2 \leq \left| \mathbf{b}_l^H \hat{\mathbf{h}}^{t,(l)} \hat{\mathbf{x}}^{t,(l)H} \mathbf{a}_l - y_l \right| \|\mathbf{b}_l\|_2 \left| \mathbf{a}_l^H \hat{\mathbf{x}}^{t,(l)} \right|. \tag{216}$$

Bounds (189) and (214) taken collectively yield

$$\left| \mathbf{a}_l^H \hat{\mathbf{x}}^{t,(l)} \right| \leq \left| \mathbf{a}_l^H \mathbf{x}^* \right| + \left| \mathbf{a}_l^H (\hat{\mathbf{x}}^{t,(l)} - \mathbf{x}^*) \right| \lesssim \sqrt{\log m} + C_3 \frac{1}{\log^{3/2} m} \asymp \sqrt{\log m}.$$

In addition, the same argument as in obtaining (215) tells us that

$$\left| \mathbf{b}_l^H (\hat{\mathbf{h}}^{t,(l)} - \mathbf{h}^*) \right| \lesssim C_4 \frac{\mu}{\sqrt{m}} \log^2 m.$$

Combine the previous two bounds to obtain

$$\begin{aligned} \left| \mathbf{b}_l^H \hat{\mathbf{h}}^{t,(l)} \hat{\mathbf{x}}^{t,(l)H} \mathbf{a}_l - y_l \right| &\leq \left| \mathbf{b}_l^H \hat{\mathbf{h}}^{t,(l)} (\hat{\mathbf{x}}^{t,(l)} - \mathbf{x}^*)^H \mathbf{a}_l \right| + \left| \mathbf{b}_l^H (\hat{\mathbf{h}}^{t,(l)} - \mathbf{h}^*) \mathbf{x}^{*H} \mathbf{a}_l \right| \\ &\leq \left| \mathbf{b}_l^H \hat{\mathbf{h}}^{t,(l)} \right| \cdot \left| \mathbf{a}_l^H (\hat{\mathbf{x}}^{t,(l)} - \mathbf{x}^*) \right| + \left| \mathbf{b}_l^H (\hat{\mathbf{h}}^{t,(l)} - \mathbf{h}^*) \right| \cdot \left| \mathbf{a}_l^H \mathbf{x}^* \right| \end{aligned}$$

$$\begin{aligned} &\leq \left( |\mathbf{b}_l^H \widehat{\mathbf{h}}^{t,(l)} - \mathbf{h}^*| + |\mathbf{b}_l^H \mathbf{h}^*| \right) \cdot |\mathbf{a}_l^H (\widehat{\mathbf{x}}^{t,(l)} - \mathbf{x}^*)| + |\mathbf{b}_l^H \widehat{\mathbf{h}}^{t,(l)} - \mathbf{h}^*| \cdot |\mathbf{a}_l^H \mathbf{x}^*| \\ &\lesssim \left( C_4 \mu \frac{\log^2 m}{\sqrt{m}} + \frac{\mu}{\sqrt{m}} \right) \cdot C_3 \frac{1}{\log^{3/2} m} + C_4 \mu \frac{\log^2 m}{\sqrt{m}} \cdot \sqrt{\log m} \lesssim C_4 \mu \frac{\log^{5/2} m}{\sqrt{m}}. \end{aligned}$$

Substitution into (216) gives

$$\left\| \left( \mathbf{b}_l^H \widehat{\mathbf{h}}^{t,(l)} \widehat{\mathbf{x}}^{t,(l)H} \mathbf{a}_l - y_l \right) \mathbf{b}_l \mathbf{a}_l^H \widehat{\mathbf{x}}^{t,(l)} \right\|_2 \lesssim C_4 \mu \frac{\log^{5/2} m}{\sqrt{m}} \cdot \sqrt{\frac{K}{m}} \cdot \sqrt{\log m}. \tag{217}$$

Similarly, we can also derive

$$\begin{aligned} \left\| \overline{\left( \mathbf{b}_l^H \widehat{\mathbf{h}}^{t,(l)} \widehat{\mathbf{x}}^{t,(l)H} \mathbf{a}_l - y_l \right) \mathbf{a}_l \mathbf{b}_l^H \widehat{\mathbf{h}}^{t,(l)}} \right\| &\leq \left| \mathbf{b}_l^H \widehat{\mathbf{h}}^{t,(l)} \widehat{\mathbf{x}}^{t,(l)H} \mathbf{a}_l - y_l \right| \|\mathbf{a}_l\|_2 \left| \mathbf{b}_l^H \widehat{\mathbf{h}}^{t,(l)} \right| \\ &\lesssim C_4 \mu \frac{\log^{5/2} m}{\sqrt{m}} \cdot \sqrt{K} \cdot C_4 \frac{\mu}{\sqrt{m}} \log^2 m \end{aligned}$$

Putting these bounds together indicates that

$$\|\mathbf{v}_3\|_2 \lesssim (C_4)^2 \frac{\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 K \log^9 m}{m}}.$$

The above bounds taken together with (212) and (213) ensure the existence of a constant  $C > 0$  such that

$$\begin{aligned} \text{dist}(\mathbf{z}^{t+1,(l)}, \widetilde{\mathbf{z}}^{t+1}) &\leq \max \left\{ \left| \frac{\alpha^{t+1}}{\alpha^t} \right|, \left| \frac{\alpha^t}{\alpha^{t+1}} \right| \right\} \left\{ \left( 1 - \frac{\eta}{16} + C C_1 \eta \frac{1}{\log^2 m} \right) \right. \\ &\quad \left. \|\widetilde{\mathbf{z}}^{t,(l)} - \widetilde{\mathbf{z}}^t\|_2 + C (C_4)^2 \eta \frac{\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 K \log^9 m}{m}} \right\} \\ &\stackrel{(i)}{\leq} \frac{1 - \eta/21}{1 - \eta/20} \left\{ \left( 1 - \frac{\eta}{20} \right) \|\widetilde{\mathbf{z}}^{t,(l)} - \widetilde{\mathbf{z}}^t\|_2 \right. \\ &\quad \left. + C (C_4)^2 \eta \frac{\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 K \log^9 m}{m}} \right\} \\ &\leq \left( 1 - \frac{\eta}{21} \right) \|\widetilde{\mathbf{z}}^{t,(l)} - \widetilde{\mathbf{z}}^t\|_2 + 2C (C_4)^2 \eta \frac{\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 K \log^9 m}{m}} \\ &= \left( 1 - \frac{\eta}{21} \right) \text{dist}(\mathbf{z}^{t,(l)}, \widetilde{\mathbf{z}}^t) + 2C (C_4)^2 \eta \frac{\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 K \log^9 m}{m}} \\ &\stackrel{(ii)}{\leq} C_2 \frac{\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 K \log^9 m}{m}}. \end{aligned}$$



Here, (i) holds as long as  $m$  is sufficiently large such that  $CC_1 1/\log^2 m \ll 1$  and

$$\max \left\{ \left| \frac{\alpha^{t+1}}{\alpha^t} \right|, \left| \frac{\alpha^t}{\alpha^{t+1}} \right| \right\} < \frac{1 - \eta/21}{1 - \eta/20}, \tag{218}$$

which is guaranteed by Lemma 16. Inequality (ii) arises from the induction hypothesis (90b) and taking  $C_2 > 0$  is sufficiently large.

Finally, we establish the second inequality claimed in the lemma. Take  $(\mathbf{h}_1, \mathbf{x}_1) = (\tilde{\mathbf{h}}^{t+1}, \tilde{\mathbf{x}}^{t+1})$  and  $(\mathbf{h}_2, \mathbf{x}_2) = (\hat{\mathbf{h}}^{t+1,(l)}, \hat{\mathbf{x}}^{t+1,(l)})$  in Lemma 55. Since both  $(\mathbf{h}_1, \mathbf{x}_1)$  and  $(\mathbf{h}_2, \mathbf{x}_2)$  are close enough to  $(\mathbf{h}^*, \mathbf{x}^*)$ , we deduce that

$$\|\tilde{\mathbf{z}}^{t+1,(l)} - \tilde{\mathbf{z}}^{t+1}\|_2 \lesssim \|\hat{\mathbf{z}}^{t+1,(l)} - \tilde{\mathbf{z}}^{t+1}\|_2 \lesssim C_2 \frac{\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 K \log^9 m}{m}}$$

as claimed.

### C.4 Proof of Lemma 18

Before going forward, we make note of the following inequality

$$\max_{1 \leq l \leq m} \left| \mathbf{b}_l^H \frac{1}{\alpha^{t+1}} \mathbf{h}^{t+1} \right| \leq \left| \frac{\alpha^t}{\alpha^{t+1}} \right| \max_{1 \leq l \leq m} \left| \mathbf{b}_l^H \frac{1}{\alpha^t} \mathbf{h}^{t+1} \right| \leq (1 + \delta) \max_{1 \leq l \leq m} \left| \mathbf{b}_l^H \frac{1}{\alpha^t} \mathbf{h}^{t+1} \right|$$

for some small  $\delta \asymp \log^{-2} m$ , where the last relation follows from Lemma 16 that

$$\left| \frac{\alpha^{t+1}}{\alpha^t} - 1 \right| \lesssim \frac{1}{\log^2 m} \leq \delta$$

for  $m$  sufficiently large. In view of the above inequality, the focus of our subsequent analysis will be to control  $\max_l \left| \mathbf{b}_l^H \frac{1}{\alpha^t} \mathbf{h}^{t+1} \right|$ .

The gradient update rule for  $\mathbf{h}^{t+1}$  (cf. (79a)) gives

$$\frac{1}{\alpha^t} \mathbf{h}^{t+1} = \tilde{\mathbf{h}}^t - \eta \xi \sum_{j=1}^m \mathbf{b}_j \mathbf{b}_j^H (\tilde{\mathbf{h}}^t \tilde{\mathbf{x}}^{tH} - \mathbf{h}^* \mathbf{x}^{*H}) \mathbf{a}_j \mathbf{a}_j^H \tilde{\mathbf{x}}^t,$$

where  $\tilde{\mathbf{h}}^t = \frac{1}{\alpha^t} \mathbf{h}^t$  and  $\tilde{\mathbf{x}}^t = \alpha^t \mathbf{x}^t$ . Here and below, we denote  $\xi = 1/\|\tilde{\mathbf{x}}^t\|_2^2$  for notational convenience. The above formula can be further decomposed into the following terms

$$\frac{1}{\alpha^t} \mathbf{h}^{t+1} = \tilde{\mathbf{h}}^t - \eta \xi \sum_{j=1}^m \mathbf{b}_j \mathbf{b}_j^H \tilde{\mathbf{h}}^t \left| \mathbf{a}_j^H \tilde{\mathbf{x}}^t \right|^2 + \eta \xi \sum_{j=1}^m \mathbf{b}_j \mathbf{b}_j^H \mathbf{h}^* \mathbf{x}^{*H} \mathbf{a}_j \mathbf{a}_j^H \tilde{\mathbf{x}}^t$$

$$\begin{aligned}
 &= \left(1 - \eta\xi \|x^*\|^2\right) \tilde{h}^t - \underbrace{\eta\xi \sum_{j=1}^m b_j b_j^H \tilde{h}^t (|a_j^H \tilde{x}^t|^2 - |a_j^H x^*|^2)}_{:=v_1} \\
 &\quad - \underbrace{\eta\xi \sum_{j=1}^m b_j b_j^H \tilde{h}^t (|a_j^H x^*|^2 - \|x^*\|^2)}_{:=v_2} + \underbrace{\eta\xi \sum_{j=1}^m b_j b_j^H h^* x^{*H} a_j a_j^H \tilde{x}^t}_{:=v_3},
 \end{aligned}$$

where we use the fact that  $\sum_{j=1}^m b_j b_j^H = I_K$ . In the sequel, we shall control each term separately.

1. We start with  $|b_l^H v_1|$  by making the observation that

$$\begin{aligned}
 \frac{1}{\eta\xi} |b_l^H v_1| &= \left| \sum_{j=1}^m b_l^H b_j b_j^H \tilde{h}^t \left[ a_j^H (\tilde{x}^t - x^*) (a_j^H \tilde{x}^t)^H + a_j^H x^* (a_j^H (\tilde{x}^t - x^*))^H \right] \right| \\
 &\leq \sum_{j=1}^m |b_l^H b_j| \left\{ \max_{1 \leq j \leq m} |b_j^H \tilde{h}^t| \right\} \left\{ \max_{1 \leq j \leq m} |a_j^H (\tilde{x}^t - x^*)| (|a_j^H \tilde{x}^t| + |a_j^H x^*|) \right\}.
 \end{aligned} \tag{219}$$

Combining the induction hypothesis (90c) and condition (189) yields

$$\begin{aligned}
 \max_{1 \leq j \leq m} |a_j^H \tilde{x}^t| &\leq \max_{1 \leq j \leq m} |a_j^H (\tilde{x}^t - x^*)| + \max_{1 \leq j \leq m} |a_j^H x^*| \\
 &\leq C_3 \frac{1}{\log^{3/2} m} + 5\sqrt{\log m} \leq 6\sqrt{\log m}
 \end{aligned}$$

as long as  $m$  is sufficiently large. This further implies

$$\max_{1 \leq j \leq m} |a_j^H (\tilde{x}^t - x^*)| (|a_j^H \tilde{x}^t| + |a_j^H x^*|) \leq C_3 \frac{1}{\log^{3/2} m} \cdot 11\sqrt{\log m} \leq 11C_3 \frac{1}{\log m}.$$

Substituting it into (219) and taking Lemma 48, we arrive at

$$\frac{1}{\eta\xi} |b_l^H v_1| \lesssim \log m \cdot \left\{ \max_{1 \leq j \leq m} |b_j^H \tilde{h}^t| \right\} \cdot C_3 \frac{1}{\log m} \lesssim C_3 \max_{1 \leq j \leq m} |b_j^H \tilde{h}^t| \leq 0.1 \max_{1 \leq j \leq m} |b_j^H \tilde{h}^t|,$$

with the proviso that  $C_3$  is sufficiently small.

2. We then move on to  $|b_l^H v_3|$ , which obeys

$$\frac{1}{\eta\xi} |b_l^H v_3| \leq \left| \sum_{j=1}^m b_l^H b_j b_j^H h^* x^{*H} a_j a_j^H x^* \right| + \left| \sum_{j=1}^m b_l^H b_j b_j^H h^* x^{*H} a_j a_j^H (\tilde{x}^t - x^*) \right|. \tag{220}$$

Regarding the first term, we have the following lemma, whose proof is given in Appendix C.4.1.

**Lemma 28** *Suppose  $m \geq CK \log^2 m$  for some sufficiently large constant  $C > 0$ . Then with probability at least  $1 - O(m^{-10})$ , one has*

$$\left| \sum_{j=1}^m \mathbf{b}_l^H \mathbf{b}_j \mathbf{b}_j^H \mathbf{h}^* \mathbf{x}^{*H} \mathbf{a}_j \mathbf{a}_j^H \mathbf{x}^* - \mathbf{b}_l^H \mathbf{h}^* \right| \lesssim \frac{\mu}{\sqrt{m}}.$$

For the remaining term, we apply the same strategy as in bounding  $|\mathbf{b}_l^H \mathbf{v}_1|$  to get

$$\begin{aligned} & \left| \sum_{j=1}^m \mathbf{b}_l^H \mathbf{b}_j \mathbf{b}_j^H \mathbf{h}^* \mathbf{x}^{*H} \mathbf{a}_j \mathbf{a}_j^H (\tilde{\mathbf{x}}^t - \mathbf{x}^*) \right| \\ & \leq \sum_{j=1}^m |\mathbf{b}_l^H \mathbf{b}_j| \left\{ \max_{1 \leq j \leq m} |\mathbf{b}_j^H \mathbf{h}^*| \right\} \left\{ \max_{1 \leq j \leq m} |\mathbf{a}_j^H (\tilde{\mathbf{x}}^t - \mathbf{x}^*)| \right\} \left\{ \max_{1 \leq j \leq m} |\mathbf{a}_j^H \mathbf{x}^*| \right\} \\ & \leq 4 \log m \cdot \frac{\mu}{\sqrt{m}} \cdot C_3 \frac{1}{\log^{3/2} m} \cdot 5\sqrt{\log m} \\ & \lesssim C_3 \frac{\mu}{\sqrt{m}}, \end{aligned}$$

where the second line follows from incoherence (36), the induction hypothesis (90c), condition (189), and Lemma 48. Combining the above three inequalities and incoherence (36) yields

$$\frac{1}{\eta\xi} |\mathbf{b}_l^H \mathbf{v}_3| \lesssim |\mathbf{b}_l^H \mathbf{h}^*| + \frac{\mu}{\sqrt{m}} + C_3 \frac{\mu}{\sqrt{m}} \lesssim (1 + C_3) \frac{\mu}{\sqrt{m}}.$$

3. Finally, we need to control  $|\mathbf{b}_l^H \mathbf{v}_2|$ . For convenience of presentation, we will only bound  $|\mathbf{b}_1^H \mathbf{v}_2|$  in the sequel, but the argument easily extends to all other  $\mathbf{b}_l$ 's. The idea is to group  $\{\mathbf{b}_j\}_{1 \leq j \leq m}$  into bins each containing  $\tau$  adjacent vectors, and to look at each bin separately. Here,  $\tau \asymp \text{poly} \log(m)$  is some integer to be specified later. For notational simplicity, we assume  $m/\tau$  to be an integer, although all arguments continue to hold when  $m/\tau$  is not an integer. For each  $0 \leq l \leq m - \tau$ , the following summation over  $\tau$  adjacent data obeys

$$\begin{aligned} & \mathbf{b}_1^H \sum_{j=1}^{\tau} \mathbf{b}_{l+j} \mathbf{b}_{l+j}^H \tilde{\mathbf{h}}^t \left( |\mathbf{a}_{l+j}^H \mathbf{x}^*|^2 - \|\mathbf{x}^*\|_2^2 \right) \\ & = \mathbf{b}_1^H \sum_{j=1}^{\tau} \mathbf{b}_{l+1} \mathbf{b}_{l+1}^H \tilde{\mathbf{h}}^t \left( |\mathbf{a}_{l+j}^H \mathbf{x}^*|^2 - \|\mathbf{x}^*\|_2^2 \right) \end{aligned}$$

$$\begin{aligned}
 & + \mathbf{b}_1^H \sum_{j=1}^{\tau} \left( \mathbf{b}_{l+j} \mathbf{b}_{l+j}^H - \mathbf{b}_{l+1} \mathbf{b}_{l+1}^H \right) \tilde{\mathbf{h}}^t \left( \left| \mathbf{a}_{l+j}^H \mathbf{x}^* \right|^2 - \|\mathbf{x}^*\|_2^2 \right) \\
 & = \left\{ \sum_{j=1}^{\tau} \left( \left| \mathbf{a}_{l+j}^H \mathbf{x}^* \right|^2 - \|\mathbf{x}^*\|_2^2 \right) \right\} \mathbf{b}_1^H \mathbf{b}_{l+1} \mathbf{b}_{l+1}^H \tilde{\mathbf{h}}^t \\
 & \quad + \mathbf{b}_1^H \sum_{j=1}^{\tau} (\mathbf{b}_{l+j} - \mathbf{b}_{l+1}) \mathbf{b}_{l+j}^H \tilde{\mathbf{h}}^t \left( \left| \mathbf{a}_{l+j}^H \mathbf{x}^* \right|^2 - \|\mathbf{x}^*\|_2^2 \right) \\
 & \quad + \mathbf{b}_1^H \sum_{j=1}^{\tau} \mathbf{b}_{l+1} (\mathbf{b}_{l+j} - \mathbf{b}_{l+1})^H \tilde{\mathbf{h}}^t \left( \left| \mathbf{a}_{l+j}^H \mathbf{x}^* \right|^2 - \|\mathbf{x}^*\|_2^2 \right). \tag{221}
 \end{aligned}$$

We will now bound each term in (221) separately.

- Before bounding the first term in (221), we first bound the pre-factor  $\left| \sum_{j=1}^{\tau} \left( \left| \mathbf{a}_{l+j}^H \mathbf{x}^* \right|^2 - \|\mathbf{x}^*\|_2^2 \right) \right|$ . Notably, the fluctuation of this quantity does not grow fast as it is the sum of i.i.d. random variables over a group of relatively large size, i.e.,  $\tau$ . Since  $2 \left| \mathbf{a}_j^H \mathbf{x}^* \right|^2$  follows the  $\chi_2^2$  distribution, by standard concentration results (e.g., [95, Theorem 1.1]), with probability exceeding  $1 - O(m^{-10})$ ,

$$\left| \sum_{j=1}^{\tau} \left( \left| \mathbf{a}_{l+j}^H \mathbf{x}^* \right|^2 - \|\mathbf{x}^*\|_2^2 \right) \right| \lesssim \sqrt{\tau \log m}.$$

With this result in place, we can bound the first term in (221) as

$$\left| \left\{ \sum_{j=1}^{\tau} \left( \left| \mathbf{a}_{l+j}^H \mathbf{x}^* \right|^2 - \|\mathbf{x}^*\|_2^2 \right) \right\} \mathbf{b}_1^H \mathbf{b}_{l+1} \mathbf{b}_{l+1}^H \tilde{\mathbf{h}}^t \right| \lesssim \sqrt{\tau \log m} \left| \mathbf{b}_1^H \mathbf{b}_{l+1} \right| \max_{1 \leq l \leq m} \left| \mathbf{b}_l^H \tilde{\mathbf{h}}^t \right|.$$

Taking the summation over all bins gives

$$\begin{aligned}
 & \sum_{k=0}^{\frac{m}{\tau}-1} \left| \left\{ \sum_{j=1}^{\tau} \left( \left| \mathbf{a}_{k\tau+j}^H \mathbf{x}^* \right|^2 - \|\mathbf{x}^*\|_2^2 \right) \right\} \mathbf{b}_1^H \mathbf{b}_{k\tau+1} \mathbf{b}_{k\tau+1}^H \tilde{\mathbf{h}}^t \right| \\
 & \lesssim \sqrt{\tau \log m} \sum_{k=0}^{\frac{m}{\tau}-1} \left| \mathbf{b}_1^H \mathbf{b}_{k\tau+1} \right| \max_{1 \leq l \leq m} \left| \mathbf{b}_l^H \tilde{\mathbf{h}}^t \right|. \tag{222}
 \end{aligned}$$

It is straightforward to see from the proof of Lemma 48 that

$$\sum_{k=0}^{\frac{m}{\tau}-1} \left| \mathbf{b}_1^H \mathbf{b}_{k\tau+1} \right| = \|\mathbf{b}_1\|_2^2 + \sum_{k=1}^{\frac{m}{\tau}-1} \left| \mathbf{b}_1^H \mathbf{b}_{k\tau+1} \right| \leq \frac{K}{m} + O\left(\frac{\log m}{\tau}\right). \tag{223}$$

Substituting (223) into the previous inequality (222) gives

$$\begin{aligned} & \sum_{k=0}^{\frac{m}{\tau}-1} \left| \left\{ \sum_{j=1}^{\tau} (|a_{k\tau+j}^H \mathbf{x}^*|^2 - \|\mathbf{x}^*\|_2^2) \right\} \mathbf{b}_1^H \mathbf{b}_{k\tau+1} \mathbf{b}_{k\tau+1}^H \tilde{\mathbf{h}}^t \right| \\ & \lesssim \left( \frac{K\sqrt{\tau \log m}}{m} + \sqrt{\frac{\log^3 m}{\tau}} \right) \max_{1 \leq l \leq m} |\mathbf{b}_l^H \tilde{\mathbf{h}}^t| \\ & \leq 0.1 \max_{1 \leq l \leq m} |\mathbf{b}_l^H \tilde{\mathbf{h}}^t|, \end{aligned}$$

as long as  $m \gg K\sqrt{\tau \log m}$  and  $\tau \gg \log^3 m$ .

- The second term of (221) obeys

$$\begin{aligned} & \left| \mathbf{b}_1^H \sum_{j=1}^{\tau} (\mathbf{b}_{l+j} - \mathbf{b}_{l+1}) \mathbf{b}_{l+j}^H \tilde{\mathbf{h}}^t \left( |a_{l+j}^H \mathbf{x}^*|^2 - \|\mathbf{x}^*\|_2^2 \right) \right| \\ & \leq \max_{1 \leq l \leq m} |\mathbf{b}_l^H \tilde{\mathbf{h}}^t| \sqrt{\sum_{j=1}^{\tau} |\mathbf{b}_1^H (\mathbf{b}_{l+j} - \mathbf{b}_{l+1})|^2} \sqrt{\sum_{j=1}^{\tau} (|a_{l+j}^H \mathbf{x}^*|^2 - \|\mathbf{x}^*\|_2^2)^2} \\ & \lesssim \sqrt{\tau} \max_{1 \leq l \leq m} |\mathbf{b}_l^H \tilde{\mathbf{h}}^t| \sqrt{\sum_{j=1}^{\tau} |\mathbf{b}_1^H (\mathbf{b}_{l+j} - \mathbf{b}_{l+1})|^2}, \end{aligned}$$

where the first inequality is due to Cauchy–Schwarz, and the second one holds because of the following lemma, whose proof can be found in Appendix C.4.2.

**Lemma 29** *Suppose  $\tau \geq C \log^4 m$  for some sufficiently large constant  $C > 0$ . Then with probability exceeding  $1 - O(m^{-10})$ ,*

$$\sum_{j=1}^{\tau} \left( |a_j^H \mathbf{x}^*|^2 - \|\mathbf{x}^*\|_2^2 \right)^2 \lesssim \tau.$$

With the above bound in mind, we can sum over all bins of size  $\tau$  to obtain

$$\begin{aligned} & \left| \mathbf{b}_1^H \sum_{k=0}^{\frac{m}{\tau}-1} \sum_{j=1}^{\tau} (\mathbf{b}_{k\tau+j} - \mathbf{b}_{k\tau+1}) \mathbf{b}_{k\tau+j}^H \tilde{\mathbf{h}}^t \left\{ |a_{l+j}^H \mathbf{x}^*|^2 - \|\mathbf{x}^*\|_2^2 \right\} \right| \\ & \lesssim \left\{ \sqrt{\tau} \sum_{k=0}^{\frac{m}{\tau}-1} \sqrt{\sum_{j=1}^{\tau} |\mathbf{b}_1^H (\mathbf{b}_{k\tau+j} - \mathbf{b}_{k\tau+1})|^2} \right\} \max_{1 \leq l \leq m} |\mathbf{b}_l^H \tilde{\mathbf{h}}^t| \\ & \leq 0.1 \max_{1 \leq l \leq m} |\mathbf{b}_l^H \tilde{\mathbf{h}}^t|. \end{aligned}$$

Here, the last line arises from Lemma 51, which says that for any small constant  $c > 0$ , as long as  $m \gg \tau K \log m$

$$\sum_{k=0}^{\frac{m}{\tau}-1} \sqrt{\sum_{j=1}^{\tau} \left| \mathbf{b}_1^H (\mathbf{b}_{k\tau+j} - \mathbf{b}_{k\tau+1}) \right|^2} \leq c \frac{1}{\sqrt{\tau}}.$$

- The third term of (221) obeys

$$\begin{aligned} & \left| \mathbf{b}_1^H \sum_{j=1}^{\tau} \mathbf{b}_{l+1} (\mathbf{b}_{l+j} - \mathbf{b}_{l+1})^H \tilde{\mathbf{h}}^t \left\{ \left| \mathbf{a}_{l+j}^H \mathbf{x}^* \right|^2 - \|\mathbf{x}^*\|_2^2 \right\} \right| \\ & \leq \left| \mathbf{b}_1^H \mathbf{b}_{l+1} \right| \left\{ \sum_{j=1}^{\tau} \left| \left| \mathbf{a}_{l+j}^H \mathbf{x}^* \right|^2 - \|\mathbf{x}^*\|_2^2 \right| \right\} \max_{0 \leq l \leq m-\tau, 1 \leq j \leq \tau} \left| (\mathbf{b}_{l+j} - \mathbf{b}_{l+1})^H \tilde{\mathbf{h}}^t \right| \\ & \lesssim \tau \left| \mathbf{b}_1^H \mathbf{b}_{l+1} \right| \max_{0 \leq l \leq m-\tau, 1 \leq j \leq \tau} \left| (\mathbf{b}_{l+j} - \mathbf{b}_{l+1})^H \tilde{\mathbf{h}}^t \right|, \end{aligned}$$

where the last line relies on the inequality

$$\sum_{j=1}^{\tau} \left| \left| \mathbf{a}_{l+j}^H \mathbf{x}^* \right|^2 - \|\mathbf{x}^*\|_2^2 \right| \leq \sqrt{\tau} \sqrt{\sum_{j=1}^{\tau} \left( \left| \mathbf{a}_{l+j}^H \mathbf{x}^* \right|^2 - \|\mathbf{x}^*\|_2^2 \right)^2} \lesssim \tau$$

owing to Lemma 29 and the Cauchy–Schwarz inequality. Summing over all bins gives

$$\begin{aligned} & \sum_{k=0}^{\frac{m}{\tau}-1} \left| \mathbf{b}_1^H \sum_{j=1}^{\tau} \mathbf{b}_{k\tau+1} (\mathbf{b}_{k\tau+j} - \mathbf{b}_{k\tau+1})^H \tilde{\mathbf{h}}^t \left\{ \left| \mathbf{a}_{k\tau+j}^H \mathbf{x}^* \right|^2 - \|\mathbf{x}^*\|_2^2 \right\} \right| \\ & \lesssim \tau \sum_{k=0}^{\frac{m}{\tau}-1} \left| \mathbf{b}_1^H \mathbf{b}_{k\tau+1} \right| \max_{0 \leq l \leq m-\tau, 1 \leq j \leq \tau} \left| (\mathbf{b}_{l+j} - \mathbf{b}_{l+1})^H \tilde{\mathbf{h}}^t \right| \\ & \lesssim \log m \max_{0 \leq l \leq m-\tau, 1 \leq j \leq \tau} \left| (\mathbf{b}_{l+j} - \mathbf{b}_{l+1})^H \tilde{\mathbf{h}}^t \right|, \end{aligned}$$

where the last relation makes use of (223) with the proviso that  $m \gg K\tau$ . It then boils down to bounding  $\max_{0 \leq l \leq m-\tau, 1 \leq j \leq \tau} \left| (\mathbf{b}_{l+j} - \mathbf{b}_{l+1})^H \tilde{\mathbf{h}}^t \right|$ . Without loss of generality, it suffices to look at  $\left| (\mathbf{b}_j - \mathbf{b}_1)^H \tilde{\mathbf{h}}^t \right|$  for all  $1 \leq j \leq \tau$ . Specifically, we claim for the moment that

$$\max_{1 \leq j \leq \tau} \left| (\mathbf{b}_j - \mathbf{b}_1)^H \tilde{\mathbf{h}}^t \right| \leq cC_4 \frac{\mu}{\sqrt{m}} \log m \quad (224)$$

for some sufficiently small constant  $c > 0$ , provided that  $m \gg \tau K \log^4 m$ . As a result,

$$\sum_{k=0}^{\frac{m}{\tau}-1} \left| \mathbf{b}_1^H \sum_{j=1}^{\tau} \mathbf{b}_{k\tau+1} (\mathbf{b}_{k\tau+j} - \mathbf{b}_{k\tau+1})^H \tilde{\mathbf{h}}^t \left\{ \left| \mathbf{a}_{k\tau+j}^H \mathbf{x}^* \right|^2 - \|\mathbf{x}^*\|_2^2 \right\} \right| \lesssim cC_4 \frac{\mu}{\sqrt{m}} \log^2 m.$$

- Putting the above results together, we get

$$\begin{aligned} \frac{1}{\eta\xi} \left| \mathbf{b}_1^H \mathbf{v}_2 \right| &\leq \sum_{k=0}^{\frac{m}{\tau}-1} \left| \mathbf{b}_1^H \sum_{j=1}^{\tau} \mathbf{b}_{k\tau+j} \mathbf{b}_{k\tau+j}^H \tilde{\mathbf{h}}^t \left\{ \left| \mathbf{a}_{k\tau+j}^H \mathbf{x}^* \right|^2 - \|\mathbf{x}^*\|_2^2 \right\} \right| \\ &\leq 0.2 \max_{1 \leq l \leq m} \left| \mathbf{b}_l^H \tilde{\mathbf{h}}^t \right| + O \left( cC_4 \frac{\mu}{\sqrt{m}} \log^2 m \right). \end{aligned}$$

4. Combining the preceding bounds guarantees the existence of some constant  $C_8 > 0$  such that

$$\begin{aligned} \left| \mathbf{b}_l^H \tilde{\mathbf{h}}^{t+1} \right| &\leq (1 + \delta) \left\{ (1 - \eta\xi) \left| \mathbf{b}_l^H \tilde{\mathbf{h}}^t \right| + 0.3\eta\xi \max_{1 \leq l \leq m} \left| \mathbf{b}_l^H \tilde{\mathbf{h}}^t \right| \right. \\ &\quad \left. + C_8(1 + C_3)\eta\xi \frac{\mu}{\sqrt{m}} + C_8\eta\xi cC_4 \frac{\mu}{\sqrt{m}} \log^2 m \right\} \\ &\stackrel{(i)}{\leq} \left( 1 + O \left( \frac{1}{\log^2 m} \right) \right) \left\{ (1 - 0.7\eta\xi) C_4 \frac{\mu}{\sqrt{m}} \log^2 m \right. \\ &\quad \left. + C_8(1 + C_3)\eta\xi \frac{\mu}{\sqrt{m}} + C_8\eta\xi cC_4 \frac{\mu}{\sqrt{m}} \log^2 m \right\} \\ &\stackrel{(ii)}{\leq} C_4 \frac{\mu}{\sqrt{m}} \log^2 m. \end{aligned}$$

Here, (i) uses the induction hypothesis (90d), and (ii) holds as long as  $c > 0$  is sufficiently small (so that  $(1 + \delta)C_8\eta\xi c \ll 1$ ) and  $\eta > 0$  is some sufficiently small constant. In order for the proof to go through, it suffices to pick

$$\tau = c_{10} \log^4 m$$

for some sufficiently large constant  $c_{10} > 0$ . Accordingly, we need the sample size to exceed

$$m \gg \mu^2 \tau K \log^4 m \asymp \mu^2 K \log^8 m.$$

Finally, it remains to verify claim (224), which we accomplish in Appendix C.4.3.

**C.4.1 Proof of Lemma 28**

Denote

$$w_j = \mathbf{b}_l^H \mathbf{b}_j \mathbf{b}_j^H \mathbf{h}^* \mathbf{x}^{*H} \mathbf{a}_j \mathbf{a}_j^H \mathbf{x}^*.$$

Recognizing that  $\mathbb{E}[\mathbf{a}_j \mathbf{a}_j^H] = \mathbf{I}_K$  and  $\sum_{j=1}^m \mathbf{b}_j \mathbf{b}_j^H = \mathbf{I}_K$ , we can write the quantity of interest as the sum of independent random variables, namely

$$\sum_{j=1}^m \mathbf{b}_l^H \mathbf{b}_j \mathbf{b}_j^H \mathbf{h}^* \mathbf{x}^{*H} \mathbf{a}_j \mathbf{a}_j^H \mathbf{x}^* - \mathbf{b}_l^H \mathbf{h}^* = \sum_{j=1}^m (w_j - \mathbb{E}[w_j]).$$

Further, the sub-exponential norm (see definition in [116]) of  $w_j - \mathbb{E}[w_j]$  obeys

$$\|w_j - \mathbb{E}[w_j]\|_{\psi_1} \stackrel{(i)}{\leq} 2 \|w_j\|_{\psi_1} \stackrel{(ii)}{\leq} 4 \left| \mathbf{b}_l^H \mathbf{b}_j \right| \left| \mathbf{b}_j^H \mathbf{h}^* \right| \left\| \mathbf{a}_j^H \mathbf{x}^* \right\|_{\psi_2}^2 \stackrel{(iii)}{\lesssim} \left| \mathbf{b}_l^H \mathbf{b}_j \right| \frac{\mu}{\sqrt{m}} \stackrel{(iv)}{\leq} \frac{\mu \sqrt{K}}{m},$$

where (i) arises from the centering property of the sub-exponential norm (see [116, Remark 5.18]), (ii) utilizes the relationship between the sub-exponential norm and the sub-Gaussian norm [116, Lemma 5.14], (iii) is a consequence of the incoherence condition (36) and the fact that  $\left\| \mathbf{a}_j^H \mathbf{x}^* \right\|_{\psi_2} \lesssim 1$ , and (iv) follows from  $\|\mathbf{b}_j\|_2 = \sqrt{K/m}$ .

Let  $M = \max_{j \in [m]} \|w_j - \mathbb{E}[w_j]\|_{\psi_1}$  and

$$V^2 = \sum_{j=1}^m \|w_j - \mathbb{E}[w_j]\|_{\psi_1}^2 \lesssim \sum_{j=1}^m \left( \left| \mathbf{b}_l^H \mathbf{b}_j \right| \frac{\mu}{\sqrt{m}} \right)^2 = \frac{\mu^2}{m} \|\mathbf{b}_l\|_2^2 = \frac{\mu^2 K}{m^2},$$

which follows since  $\sum_{j=1}^m \left| \mathbf{b}_l^H \mathbf{b}_j \right|^2 = \mathbf{b}_l^H \left( \sum_{j=1}^m \mathbf{b}_j \mathbf{b}_j^H \right) \mathbf{b}_l = \|\mathbf{b}_l\|_2^2 = K/m$ . Let  $a_j = \|w_j - \mathbb{E}[w_j]\|_{\psi_1}$  and  $X_j = (w_j - \mathbb{E}[w_j])/a_j$ . Since  $\|X_j\|_{\psi_1} = 1$ ,  $\sum_{j=1}^m a_j^2 = V^2$  and  $\max_{j \in [m]} |a_j| = M$ , we can invoke [116, Proposition 5.16] to obtain that

$$\mathbb{P} \left( \left| \sum_{j=1}^m a_j X_j \right| \geq t \right) \leq 2 \exp \left( -c \min \left\{ \frac{t}{M}, \frac{t^2}{V^2} \right\} \right),$$

where  $c > 0$  is some universal constant. By taking  $t = \mu/\sqrt{m}$ , we see there exists some constant  $c'$  such that

$$\begin{aligned} & \mathbb{P} \left( \left| \sum_{j=1}^m \mathbf{b}_l^H \mathbf{b}_j \mathbf{b}_j^H \mathbf{h}^* \mathbf{x}^{*H} \mathbf{a}_j \mathbf{a}_j^H \mathbf{x}^* - \mathbf{b}_l^H \mathbf{h}^* \right| \geq \frac{\mu}{\sqrt{m}} \right) \\ & \leq 2 \exp \left( -c \min \left\{ \frac{\mu/\sqrt{m}}{M}, \frac{\mu^2/m}{V^2} \right\} \right) \end{aligned}$$



$$\begin{aligned} &\leq 2 \exp \left( -c' \min \left\{ \frac{\mu/\sqrt{m}}{\mu\sqrt{K}/m}, \frac{\mu^2/m}{\mu^2 K/m^2} \right\} \right) \\ &= 2 \exp \left( -c' \min \left\{ \sqrt{m/K}, m/K \right\} \right). \end{aligned}$$

We conclude the proof by observing that  $m \gg K \log^2 m$  as stated in the assumption.

**C.4.2 Proof of Lemma 29**

From the elementary inequality  $(a - b)^2 \leq 2(a^2 + b^2)$ , we see that

$$\sum_{j=1}^{\tau} \left( \left| \mathbf{a}_j^H \mathbf{x}^* \right|^2 - \|\mathbf{x}^*\|_2^2 \right)^2 \leq 2 \sum_{j=1}^{\tau} \left( \left| \mathbf{a}_j^H \mathbf{x}^* \right|^4 + \|\mathbf{x}^*\|_2^4 \right) = 2 \sum_{j=1}^{\tau} \left| \mathbf{a}_j^H \mathbf{x}^* \right|^4 + 2\tau, \tag{225}$$

where the last identity holds true since  $\|\mathbf{x}^*\|_2 = 1$ . It thus suffices to control  $\sum_{j=1}^{\tau} \left| \mathbf{a}_j^H \mathbf{x}^* \right|^4$ . Let  $\xi_i = \mathbf{a}_i^H \mathbf{x}^*$ , which is a standard complex Gaussian random variable. Since the  $\xi_i$ 's are statistically independent, one has

$$\text{Var} \left( \sum_{i=1}^{\tau} |\xi_i|^4 \right) \leq C_4 \tau$$

for some constant  $C_4 > 0$ . It then follows from the hypercontractivity concentration result for Gaussian polynomials that [99, Theorem 1.9]

$$\begin{aligned} &\mathbb{P} \left\{ \sum_{i=1}^{\tau} \left( |\xi_i|^4 - \mathbb{E} \left[ |\xi_i|^4 \right] \right) \geq c\tau \right\} \\ &\leq C \exp \left( -c_2 \left( \frac{c^2 \tau^2}{\text{Var} \left( \sum_{i=1}^{\tau} |\xi_i|^4 \right)} \right)^{1/4} \right) \\ &\leq C \exp \left( -c_2 \left( \frac{c^2 \tau^2}{C_4 \tau} \right)^{1/4} \right) = C \exp \left( -c_2 \left( \frac{c^2}{C_4} \right)^{1/4} \tau^{1/4} \right) \\ &\leq O(m^{-10}), \end{aligned}$$

for some constants  $c, c_2, C > 0$ , with the proviso that  $\tau \gg \log^4 m$ . As a consequence, with probability at least  $1 - O(m^{-10})$ ,

$$\sum_{j=1}^{\tau} \left| \mathbf{a}_j^H \mathbf{x}^* \right|^4 \lesssim \tau + \sum_{j=1}^{\tau} \mathbb{E} \left[ \left| \mathbf{a}_j^H \mathbf{x}^* \right|^4 \right] \asymp \tau,$$

which together with (225) concludes the proof.

**C.4.3 Proof of Claim (224)**

We will prove the claim by induction. Again, observe that

$$\begin{aligned} \left| (\mathbf{b}_j - \mathbf{b}_1)^H \tilde{\mathbf{h}}^t \right| &= \left| (\mathbf{b}_j - \mathbf{b}_1)^H \frac{1}{\alpha^t} \mathbf{h}^t \right| = \left| \frac{\alpha^{t-1}}{\alpha^t} \right| \left| (\mathbf{b}_j - \mathbf{b}_1)^H \frac{1}{\alpha^{t-1}} \mathbf{h}^t \right| \\ &\leq (1 + \delta) \left| (\mathbf{b}_j - \mathbf{b}_1)^H \frac{1}{\alpha^{t-1}} \mathbf{h}^t \right| \end{aligned}$$

for some  $\delta \asymp \log^{-2} m$ , which allows us to look at  $(\mathbf{b}_j - \mathbf{b}_1)^H \frac{1}{\alpha^{t-1}} \mathbf{h}^t$  instead.

Use the gradient update rule for  $\mathbf{h}^t$  (cf. (79a)) once again to get

$$\begin{aligned} \frac{1}{\alpha^{t-1}} \mathbf{h}^t &= \frac{1}{\alpha^{t-1}} \left( \mathbf{h}^{t-1} - \frac{\eta}{\|\mathbf{x}^{t-1}\|_2^2} \sum_{l=1}^m \mathbf{b}_l \mathbf{b}_l^H \left( \mathbf{h}^{t-1} \mathbf{x}^{t-1H} - \mathbf{h}^* \mathbf{x}^{*H} \right) \mathbf{a}_l \mathbf{a}_l^H \mathbf{x}^{t-1} \right) \\ &= \tilde{\mathbf{h}}^{t-1} - \eta \theta \sum_{l=1}^m \mathbf{b}_l \mathbf{b}_l^H \left( \tilde{\mathbf{h}}^{t-1} \tilde{\mathbf{x}}^{t-1H} - \mathbf{h}^* \mathbf{x}^{*H} \right) \mathbf{a}_l \mathbf{a}_l^H \tilde{\mathbf{x}}^{t-1}, \end{aligned}$$

where we denote  $\theta := 1/\|\tilde{\mathbf{x}}^{t-1}\|_2^2$ . This further gives rise to

$$\begin{aligned} &(\mathbf{b}_j - \mathbf{b}_1)^H \frac{1}{\alpha^{t-1}} \mathbf{h}^t \\ &= (\mathbf{b}_j - \mathbf{b}_1)^H \tilde{\mathbf{h}}^{t-1} - \eta \theta (\mathbf{b}_j - \mathbf{b}_1)^H \sum_{l=1}^m \mathbf{b}_l \mathbf{b}_l^H \left( \tilde{\mathbf{h}}^{t-1} \tilde{\mathbf{x}}^{t-1H} - \mathbf{h}^* \mathbf{x}^{*H} \right) \mathbf{a}_l \mathbf{a}_l^H \tilde{\mathbf{x}}^{t-1} \\ &= (\mathbf{b}_j - \mathbf{b}_1)^H \tilde{\mathbf{h}}^{t-1} - \eta \theta (\mathbf{b}_j - \mathbf{b}_1)^H \sum_{l=1}^m \mathbf{b}_l \mathbf{b}_l^H \left( \tilde{\mathbf{h}}^{t-1} \tilde{\mathbf{x}}^{t-1H} - \mathbf{h}^* \mathbf{x}^{*H} \right) \tilde{\mathbf{x}}^{t-1} \\ &\quad - \eta \theta (\mathbf{b}_j - \mathbf{b}_1)^H \sum_{l=1}^m \mathbf{b}_l \mathbf{b}_l^H \left( \tilde{\mathbf{h}}^{t-1} \tilde{\mathbf{x}}^{t-1H} - \mathbf{h}^* \mathbf{x}^{*H} \right) \left( \mathbf{a}_l \mathbf{a}_l^H - \mathbf{I}_K \right) \tilde{\mathbf{x}}^{t-1} \\ &= \left( 1 - \eta \theta \|\tilde{\mathbf{x}}^{t-1}\|_2^2 \right) (\mathbf{b}_j - \mathbf{b}_1)^H \tilde{\mathbf{h}}^{t-1} + \underbrace{\eta \theta (\mathbf{b}_j - \mathbf{b}_1)^H \mathbf{h}^* \left( \mathbf{x}^{*H} \tilde{\mathbf{x}}^{t-1} \right)}_{:=\beta_1} \\ &\quad - \underbrace{\eta \theta (\mathbf{b}_j - \mathbf{b}_1)^H \sum_{l=1}^m \mathbf{b}_l \mathbf{b}_l^H \left( \tilde{\mathbf{h}}^{t-1} \tilde{\mathbf{x}}^{t-1H} - \mathbf{h}^* \mathbf{x}^{*H} \right) \left( \mathbf{a}_l \mathbf{a}_l^H - \mathbf{I}_K \right) \tilde{\mathbf{x}}^{t-1}}_{:=\beta_2}, \end{aligned}$$

where the last identity makes use of the fact that  $\sum_{l=1}^m \mathbf{b}_l \mathbf{b}_l^H = \mathbf{I}_K$ . For  $\beta_1$ , one can get

$$\frac{1}{\eta \theta} |\beta_1| \leq \left| (\mathbf{b}_j - \mathbf{b}_1)^H \mathbf{h}^* \right| \|\mathbf{x}^*\|_2 \|\tilde{\mathbf{x}}^{t-1}\|_2 \leq 4 \frac{\mu}{\sqrt{m}},$$

where we utilize the incoherence condition (36) and the fact that  $\tilde{\mathbf{x}}^{t-1}$  and  $\mathbf{x}^*$  are extremely close, i.e.,

$$\|\tilde{\mathbf{x}}^{t-1} - \mathbf{x}^*\|_2 \leq \text{dist}(\mathbf{z}^{t-1}, \mathbf{z}^*) \ll 1 \implies \|\tilde{\mathbf{x}}^{t-1}\|_2 \leq 2.$$

Regarding the second term  $\beta_2$ , we have

$$\frac{1}{\eta\theta} |\beta_2| \leq \left\{ \sum_{l=1}^m \left| (\mathbf{b}_j - \mathbf{b}_1)^H \mathbf{b}_l \right| \right\} \underbrace{\max_{1 \leq l \leq m} \left| \mathbf{b}_l^H \left( \tilde{\mathbf{h}}^{t-1} \tilde{\mathbf{x}}^{t-1H} - \mathbf{h}^* \mathbf{x}^{*H} \right) \left( \mathbf{a}_l \mathbf{a}_l^H - \mathbf{I}_K \right) \tilde{\mathbf{x}}^{t-1} \right|}_{:=\psi}.$$

The term  $\psi$  can be bounded as follows

$$\begin{aligned} \psi &\leq \max_{1 \leq l \leq m} \left| \mathbf{b}_l^H \tilde{\mathbf{h}}^{t-1} \tilde{\mathbf{x}}^{t-1H} \left( \mathbf{a}_l \mathbf{a}_l^H - \mathbf{I} \right) \tilde{\mathbf{x}}^{t-1} \right| + \max_{1 \leq l \leq m} \left| \mathbf{b}_l^H \mathbf{h}^* \mathbf{x}^{*H} \left( \mathbf{a}_l \mathbf{a}_l^H - \mathbf{I}_K \right) \tilde{\mathbf{x}}^{t-1} \right| \\ &\leq \max_{1 \leq l \leq m} \left| \mathbf{b}_l^H \tilde{\mathbf{h}}^{t-1} \right| \max_{1 \leq l \leq m} \left| \tilde{\mathbf{x}}^{t-1H} \left( \mathbf{a}_l \mathbf{a}_l^H - \mathbf{I}_K \right) \tilde{\mathbf{x}}^{t-1} \right| \\ &\quad + \max_{1 \leq l \leq m} \left| \mathbf{b}_l^H \mathbf{h}^* \right| \max_{1 \leq l \leq m} \left| \mathbf{x}^{*H} \left( \mathbf{a}_l \mathbf{a}_l^H - \mathbf{I}_K \right) \tilde{\mathbf{x}}^{t-1} \right| \\ &\lesssim \log m \left\{ \max_{1 \leq l \leq m} \left| \mathbf{b}_l^H \tilde{\mathbf{h}}^{t-1} \right| + \frac{\mu}{\sqrt{m}} \right\}. \end{aligned}$$

Here, we have used the incoherence condition (36) and the facts that

$$\begin{aligned} \left| (\tilde{\mathbf{x}}^{t-1})^H \left( \mathbf{a}_l \mathbf{a}_l^H - \mathbf{I} \right) \tilde{\mathbf{x}}^{t-1} \right| &\leq \|\mathbf{a}_l \tilde{\mathbf{x}}^{t-1}\|_2^2 + \|\tilde{\mathbf{x}}^{t-1}\|_2^2 \lesssim \log m, \\ \left| \mathbf{x}^{*H} \left( \mathbf{a}_l \mathbf{a}_l^H - \mathbf{I} \right) \tilde{\mathbf{x}}^{t-1} \right| &\leq \|\mathbf{a}_l \tilde{\mathbf{x}}^{t-1}\|_2 \|\mathbf{a}_l^H \mathbf{x}^*\|_2 + \|\tilde{\mathbf{x}}^{t-1}\|_2 \|\mathbf{x}^*\|_2 \lesssim \log m, \end{aligned}$$

which are immediate consequences of (90c) and (189). Combining this with Lemma 50, we see that for any small constant  $c > 0$

$$\frac{1}{\eta\theta} |\beta_2| \leq c \frac{1}{\log m} \left\{ \max_{1 \leq l \leq m} \left| \mathbf{b}_l^H \tilde{\mathbf{h}}^{t-1} \right| + \frac{\mu}{\sqrt{m}} \right\}$$

holds as long as  $m \gg \tau K \log^4 m$ .

To summarize, we arrive at

$$\begin{aligned} \left| (\mathbf{b}_j - \mathbf{b}_1)^H \tilde{\mathbf{h}}^t \right| &\leq (1 + \delta) \left\{ \left( 1 - \eta\theta \|\tilde{\mathbf{x}}^{t-1}\|_2^2 \right) \left| (\mathbf{b}_j - \mathbf{b}_1)^H \tilde{\mathbf{h}}^{t-1} \right| \right. \\ &\quad \left. + 4\eta\theta \frac{\mu}{\sqrt{m}} + c\eta\theta \frac{1}{\log m} \left[ \max_{1 \leq l \leq m} \left| \mathbf{b}_l^H \tilde{\mathbf{h}}^{t-1} \right| + \frac{\mu}{\sqrt{m}} \right] \right\}. \end{aligned}$$

Making use of the induction hypothesis (85c) and the fact that  $\|\tilde{\mathbf{x}}^{t-1}\|_2^2 \geq 0.9$ , we reach

$$\left| (\mathbf{b}_j - \mathbf{b}_1)^H \tilde{\mathbf{h}}^t \right| \leq (1 + \delta) \left\{ (1 - 0.9\eta\theta) \left| (\mathbf{b}_j - \mathbf{b}_1)^H \tilde{\mathbf{h}}^{t-1} \right| + cC_4\eta\theta \frac{\mu}{\sqrt{m}} \log m + \frac{c\mu\eta\theta}{\sqrt{m} \log m} \right\}.$$

Recall that  $\delta \asymp 1/\log^2 m$ . As a result, if  $\eta > 0$  is some sufficiently small constant and if

$$\left| (\mathbf{b}_j - \mathbf{b}_1)^H \tilde{\mathbf{h}}^{t-1} \right| \leq 10c \left( C_4 \frac{\mu}{\sqrt{m}} \log m + \frac{\mu}{\eta\theta\sqrt{m} \log m} \right) \leq 20cC_4 \frac{\mu}{\sqrt{m}} \log m$$

holds, then one has

$$\left| (\mathbf{b}_j - \mathbf{b}_1)^H \tilde{\mathbf{h}}^t \right| \leq 20cC_4 \frac{\mu}{\sqrt{m}} \log m.$$

Therefore, this concludes the proof of claim (224) by induction, provided that the base case is true, i.e., for some  $c > 0$  sufficiently small

$$\left| (\mathbf{b}_j - \mathbf{b}_1)^H \tilde{\mathbf{h}}^0 \right| \leq 20cC_4 \frac{\mu}{\sqrt{m}} \log m. \tag{226}$$

Claim (226) is proved in Appendix C.6 (see Lemma 30).

### C.5 Proof of Lemma 19

Recall that  $\check{\mathbf{h}}^0$  and  $\check{\mathbf{x}}^0$  are the leading left and right singular vectors of  $\mathbf{M}$ , respectively. Applying a variant of Wedin’s  $\sin\Theta$  theorem [42, Theorem 2.1], we derive that

$$\min_{\alpha \in \mathbb{C}, |\alpha|=1} \left\{ \|\alpha\check{\mathbf{h}}^0 - \mathbf{h}^*\|_2 + \|\alpha\check{\mathbf{x}}^0 - \mathbf{x}^*\|_2 \right\} \leq \frac{c_1 \|\mathbf{M} - \mathbb{E}[\mathbf{M}]\|}{\sigma_1(\mathbb{E}[\mathbf{M}]) - \sigma_2(\mathbf{M})}, \tag{227}$$

for some universal constant  $c_1 > 0$ . Regarding the numerator of (227), it has been shown in [76, Lemma 5.20] that for any  $\xi > 0$ ,

$$\|\mathbf{M} - \mathbb{E}[\mathbf{M}]\| \leq \xi \tag{228}$$

with probability exceeding  $1 - O(m^{-10})$ , provided that

$$m \geq \frac{c_2\mu^2 K \log^2 m}{\xi^2}$$

for some universal constant  $c_2 > 0$ . For the denominator of (227), we can take (228) together with Weyl’s inequality to demonstrate that

$$\sigma_1(\mathbb{E}[\mathbf{M}]) - \sigma_2(\mathbf{M}) \geq \sigma_1(\mathbb{E}[\mathbf{M}]) - \sigma_2(\mathbb{E}[\mathbf{M}]) - \|\mathbf{M} - \mathbb{E}[\mathbf{M}]\| \geq 1 - \xi,$$

where the last inequality utilizes the facts that  $\sigma_1(\mathbb{E}[\mathbf{M}]) = 1$  and  $\sigma_2(\mathbb{E}[\mathbf{M}]) = 0$ . These together with (227) reveal that

$$\min_{\alpha \in \mathbb{C}, |\alpha|=1} \left\{ \|\alpha \check{\mathbf{h}}^0 - \mathbf{h}^*\|_2 + \|\alpha \check{\mathbf{x}}^0 - \mathbf{x}^*\|_2 \right\} \leq \frac{c_1 \xi}{1 - \xi} \leq 2c_1 \xi \tag{229}$$

as long as  $\xi \leq 1/2$ .

Now we connect the preceding bound (229) with the scaled singular vectors  $\mathbf{h}^0 = \sqrt{\sigma_1(\mathbf{M})} \check{\mathbf{h}}^0$  and  $\mathbf{x}^0 = \sqrt{\sigma_1(\mathbf{M})} \check{\mathbf{x}}^0$ . For any  $\alpha \in \mathbb{C}$  with  $|\alpha| = 1$ , from the definition of  $\mathbf{h}^0$  and  $\mathbf{x}^0$  we have

$$\|\alpha \mathbf{h}^0 - \mathbf{h}^*\|_2 + \|\alpha \mathbf{x}^0 - \mathbf{x}^*\|_2 = \left\| \sqrt{\sigma_1(\mathbf{M})} (\alpha \check{\mathbf{h}}^0) - \mathbf{h}^* \right\|_2 + \left\| \sqrt{\sigma_1(\mathbf{M})} (\alpha \check{\mathbf{x}}^0) - \mathbf{x}^* \right\|_2.$$

Since  $\alpha \check{\mathbf{h}}^0, \alpha \check{\mathbf{x}}^0$  are also the leading left and right singular vectors of  $\mathbf{M}$ , we can invoke Lemma 60 to get

$$\begin{aligned} & \left\| \alpha \mathbf{h}^0 - \mathbf{h}^* \right\|_2 + \left\| \alpha \mathbf{x}^0 - \mathbf{x}^* \right\|_2 \\ & \leq \sqrt{\sigma_1(\mathbb{E}[\mathbf{M}])} \left( \left\| \alpha \check{\mathbf{h}}^0 - \mathbf{h}^* \right\|_2 + \left\| \alpha \check{\mathbf{x}}^0 - \mathbf{x}^* \right\|_2 \right) + \frac{2 |\sigma_1(\mathbf{M}) - \sigma_1(\mathbb{E}[\mathbf{M}])|}{\sqrt{\sigma_1(\mathbf{M})} + \sqrt{\sigma_1(\mathbb{E}[\mathbf{M}])}} \\ & = \left\| \alpha \check{\mathbf{h}}^0 - \mathbf{h}^* \right\|_2 + \left\| \alpha \check{\mathbf{x}}^0 - \mathbf{x}^* \right\|_2 + \frac{2 |\sigma_1(\mathbf{M}) - \sigma_1(\mathbb{E}[\mathbf{M}])|}{\sqrt{\sigma_1(\mathbf{M})} + 1}. \end{aligned} \tag{230}$$

In addition, we can apply Weyl’s inequality once again to deduce that

$$|\sigma_1(\mathbf{M}) - \sigma_1(\mathbb{E}[\mathbf{M}])| \leq \|\mathbf{M} - \mathbb{E}[\mathbf{M}]\| \leq \xi, \tag{231}$$

where the last inequality comes from (228). Substitute (231) into (230) to obtain

$$\left\| \alpha \mathbf{h}^0 - \mathbf{h}^* \right\|_2 + \left\| \alpha \mathbf{x}^0 - \mathbf{x}^* \right\|_2 \leq \left\| \alpha \check{\mathbf{h}}^0 - \mathbf{h}^* \right\|_2 + \left\| \alpha \check{\mathbf{x}}^0 - \mathbf{x}^* \right\|_2 + 2\xi. \tag{232}$$

Taking the minimum over  $\alpha$ , one can thus conclude that

$$\begin{aligned} & \min_{\alpha \in \mathbb{C}, |\alpha|=1} \left\{ \left\| \alpha \mathbf{h}^0 - \mathbf{h}^* \right\|_2 + \left\| \alpha \mathbf{x}^0 - \mathbf{x}^* \right\|_2 \right\} \\ & \leq \min_{\alpha \in \mathbb{C}, |\alpha|=1} \left\{ \left\| \alpha \check{\mathbf{h}}^0 - \mathbf{h}^* \right\|_2 + \left\| \alpha \check{\mathbf{x}}^0 - \mathbf{x}^* \right\|_2 \right\} + 2\xi \leq 2c_1 \xi + 2\xi, \end{aligned}$$

where the last inequality comes from (229). Since  $\xi$  is arbitrary, by taking  $m/(\mu^2 K \log^2 m)$  to be large enough, we finish the proof for (92). Carrying out similar arguments (which we omit here), we can also establish (93).

The last claim in Lemma 19 that  $|\alpha_0| - 1| \leq 1/4$  is a direct corollary of (92) and Lemma 52.

### C.6 Proof of Lemma 20

The proof is composed of three steps:

- In the first step, we show that the normalized singular vectors of  $M$  and  $M^{(l)}$  are close enough; see (240).
- We then proceed by passing this proximity result to the scaled singular vectors; see (243).
- Finally, we translate the usual  $\ell_2$  distance metric to the distance function we defined in (34); see (245). Along the way, we also prove the incoherence of  $\check{h}^0$  with respect to  $\{b_l\}$ .

Here comes the formal proof. Recall that  $\check{h}^0$  and  $\check{x}^0$  are, respectively, the leading left and right singular vectors of  $M$ , and  $\check{h}^{0,(l)}$  and  $\check{x}^{0,(l)}$  are, respectively, the leading left and right singular vectors of  $M^{(l)}$ . Invoke Wedin’s  $\sin\Theta$  theorem [42, Theorem 2.1] to obtain

$$\begin{aligned} & \min_{\alpha \in \mathbb{C}, |\alpha|=1} \left\{ \|\alpha \check{h}^0 - \check{h}^{0,(l)}\|_2 + \|\alpha \check{x}^0 - \check{x}^{0,(l)}\|_2 \right\} \\ & \leq c_1 \frac{\left\| (M - M^{(l)}) \check{x}^{0,(l)} \right\|_2 + \left\| \check{h}^{0,(l)H} (M - M^{(l)}) \right\|_2}{\sigma_1(M^{(l)}) - \sigma_2(M)} \end{aligned}$$

for some universal constant  $c_1 > 0$ . Using the Weyl’s inequality we get

$$\begin{aligned} \sigma_1(M^{(l)}) - \sigma_2(M) & \geq \sigma_1(\mathbb{E}[M^{(l)}]) - \|M^{(l)} - \mathbb{E}[M^{(l)}]\| - \sigma_2(\mathbb{E}[M]) - \|M - \mathbb{E}[M]\| \\ & \geq 3/4 - \|M^{(l)} - \mathbb{E}[M^{(l)}]\| - \|M - \mathbb{E}[M]\| \geq 1/2, \end{aligned}$$

where the penultimate inequality follows from

$$\sigma_1(\mathbb{E}[M^{(l)}]) \geq 3/4$$

for  $m$  sufficiently large, and the last inequality comes from [76, Lemma 5.20], provided that  $m \geq c_2 \mu^2 K \log^2 m$  for some sufficiently large constant  $c_2 > 0$ . As a result, denoting

$$\beta^{0,(l)} := \operatorname{argmin}_{\alpha \in \mathbb{C}, |\alpha|=1} \left\{ \|\alpha \check{h}^0 - \check{h}^{0,(l)}\|_2 + \|\alpha \check{x}^0 - \check{x}^{0,(l)}\|_2 \right\} \tag{233}$$

allows us to obtain

$$\begin{aligned} & \|\beta^{0,(l)} \check{h}^0 - \check{h}^{0,(l)}\|_2 + \|\beta^{0,(l)} \check{x}^0 - \check{x}^{0,(l)}\|_2 \\ & \leq 2c_1 \left\{ \left\| (M - M^{(l)}) \check{x}^{0,(l)} \right\|_2 + \left\| \check{h}^{0,(l)H} (M - M^{(l)}) \right\|_2 \right\}. \end{aligned} \tag{234}$$

It then boils down to controlling the two terms on the right-hand side of (234). By construction,

$$M - M^{(l)} = b_l b_l^H h^* x^{*H} a_l a_l^H.$$

- To bound the first term, observe that

$$\begin{aligned} \left\| (M - M^{(l)}) \check{x}^{0,(l)} \right\|_2 &= \left\| \mathbf{b}_l \mathbf{b}_l^H \mathbf{h}^* \mathbf{x}^{*H} \mathbf{a}_l \mathbf{a}_l^H \check{x}^{0,(l)} \right\|_2 = \|\mathbf{b}_l\|_2 \left| \mathbf{b}_l^H \mathbf{h}^* \right| \left| \mathbf{a}_l^H \mathbf{x}^* \right| \cdot \left| \mathbf{a}_l^H \check{x}^{0,(l)} \right| \\ &\leq 30 \frac{\mu}{\sqrt{m}} \cdot \sqrt{\frac{K \log^2 m}{m}}, \end{aligned} \tag{235}$$

where we use the fact that  $\|\mathbf{b}_l\|_2 = \sqrt{K/m}$ , the incoherence condition (36), bound (189), and the fact that with probability exceeding  $1 - O(m^{-10})$ ,

$$\max_{1 \leq l \leq m} \left| \mathbf{a}_l^H \check{x}^{0,(l)} \right| \leq 5\sqrt{\log m},$$

due to the independence between  $\check{x}^{0,(l)}$  and  $\mathbf{a}_l$ .

- To bound the second term, for any  $\tilde{\alpha}$  obeying  $|\tilde{\alpha}| = 1$ , one has

$$\begin{aligned} &\left\| \check{\mathbf{h}}^{0,(l)H} (M - M^{(l)}) \right\|_2 \\ &= \left\| \check{\mathbf{h}}^{0,(l)H} \mathbf{b}_l \mathbf{b}_l^H \mathbf{h}^* \mathbf{x}^{*H} \mathbf{a}_l \mathbf{a}_l^H \right\|_2 = \|\mathbf{a}_l\|_2 \left| \mathbf{b}_l^H \mathbf{h}^* \right| \left| \mathbf{a}_l^H \mathbf{x}^* \right| \cdot \left| \mathbf{b}_l^H \check{\mathbf{h}}^{0,(l)} \right| \\ &\stackrel{(i)}{\leq} 3\sqrt{K} \cdot \frac{\mu}{\sqrt{m}} \cdot 5\sqrt{\log m} \cdot \left| \mathbf{b}_l^H \check{\mathbf{h}}^{0,(l)} \right| \\ &\stackrel{(ii)}{\leq} 15\sqrt{\frac{\mu^2 K \log m}{m}} |\tilde{\alpha} \mathbf{b}_l^H \check{\mathbf{h}}^0| + 15\sqrt{\frac{\mu^2 K \log m}{m}} \left| \mathbf{b}_l^H (\tilde{\alpha} \check{\mathbf{h}}^0 - \check{\mathbf{h}}^{0,(l)}) \right| \\ &\stackrel{(iii)}{\leq} 15\sqrt{\frac{\mu^2 K \log m}{m}} \left| \mathbf{b}_l^H \check{\mathbf{h}}^0 \right| + 15\sqrt{\frac{\mu^2 K \log m}{m}} \cdot \sqrt{\frac{K}{m}} \left\| \tilde{\alpha} \check{\mathbf{h}}^0 - \check{\mathbf{h}}^{0,(l)} \right\|_2. \end{aligned}$$

Here, (i) arises from the incoherence condition (36) together with bounds (189) and (190), inequality (ii) comes from the triangle inequality, and the last line (iii) holds since  $\|\mathbf{b}_l\|_2 = \sqrt{K/m}$  and  $|\tilde{\alpha}| = 1$ .

Substitution of the above bounds into (234) yields

$$\begin{aligned} &\left\| \beta^{0,(l)} \check{\mathbf{h}}^0 - \check{\mathbf{h}}^{0,(l)} \right\|_2 + \left\| \beta^{0,(l)} \check{\mathbf{x}}^0 - \check{\mathbf{x}}^{0,(l)} \right\|_2 \\ &\leq 2c_1 \left\{ 30 \frac{\mu}{\sqrt{m}} \cdot \sqrt{\frac{K \log^2 m}{m}} + 15\sqrt{\frac{\mu^2 K \log m}{m}} \left| \mathbf{b}_l^H \check{\mathbf{h}}^0 \right| \right. \\ &\quad \left. + 15\sqrt{\frac{\mu^2 K \log m}{m}} \cdot \sqrt{\frac{K}{m}} \left\| \tilde{\alpha} \check{\mathbf{h}}^0 - \check{\mathbf{h}}^{0,(l)} \right\|_2 \right\}. \end{aligned}$$

Since the previous inequality holds for all  $|\tilde{\alpha}| = 1$ , we can choose  $\tilde{\alpha} = \beta^{0,(l)}$  and rearrange terms to get

$$\begin{aligned} & \left( 1 - 30c_1 \sqrt{\frac{\mu^2 K \log m}{m}} \sqrt{\frac{K}{m}} \right) \left( \|\beta^{0,(l)} \check{\mathbf{h}}^0 - \check{\mathbf{h}}^{0,(l)}\|_2 + \|\beta^{0,(l)} \check{\mathbf{x}}^0 - \check{\mathbf{x}}^{0,(l)}\|_2 \right) \\ & \leq 60c_1 \frac{\mu}{\sqrt{m}} \cdot \sqrt{\frac{K \log^2 m}{m}} + 30c_1 \sqrt{\frac{\mu^2 K \log m}{m}} |b_l^H \check{\mathbf{h}}^0|. \end{aligned}$$

Under the condition that  $m \gg \mu K \log^{1/2} m$ , one has  $1 - 30c_1 \sqrt{\mu^2 K \log m/m} \cdot \sqrt{K/m} \geq \frac{1}{2}$ , and therefore,

$$\begin{aligned} & \|\beta^{0,(l)} \check{\mathbf{h}}^0 - \check{\mathbf{h}}^{0,(l)}\|_2 + \|\beta^{0,(l)} \check{\mathbf{x}}^0 - \check{\mathbf{x}}^{0,(l)}\|_2 \\ & \leq 120c_1 \frac{\mu}{\sqrt{m}} \cdot \sqrt{\frac{K \log^2 m}{m}} + 60c_1 \sqrt{\frac{\mu^2 K \log m}{m}} |b_l^H \check{\mathbf{h}}^0|, \end{aligned}$$

which immediately implies that

$$\begin{aligned} & \max_{1 \leq l \leq m} \left\{ \|\beta^{0,(l)} \check{\mathbf{h}}^0 - \check{\mathbf{h}}^{0,(l)}\|_2 + \|\beta^{0,(l)} \check{\mathbf{x}}^0 - \check{\mathbf{x}}^{0,(l)}\|_2 \right\} \\ & \leq 120c_1 \frac{\mu}{\sqrt{m}} \cdot \sqrt{\frac{K \log^2 m}{m}} + 60c_1 \sqrt{\frac{\mu^2 K \log m}{m}} \max_{1 \leq l \leq m} |b_l^H \check{\mathbf{h}}^0|. \end{aligned} \tag{236}$$

We then move on to  $|b_l^H \check{\mathbf{h}}^0|$ . The aim is to show that  $\max_{1 \leq l \leq m} |b_l^H \check{\mathbf{h}}^0|$  can also be upper bounded by the left-hand side of (236). By construction, we have  $M \check{\mathbf{x}}^0 = \sigma_1(M) \check{\mathbf{h}}^0$ , which further leads to

$$\begin{aligned} |b_l^H \check{\mathbf{h}}^0| &= \frac{1}{\sigma_1(M)} |b_l^H M \check{\mathbf{x}}^0| \\ &\stackrel{(i)}{\leq} 2 \left| \sum_{j=1}^m (b_l^H b_j) b_j^H \mathbf{h}^* \mathbf{x}^{*H} a_j a_j^H \check{\mathbf{x}}^0 \right| \\ &\leq 2 \left( \sum_{j=1}^m |b_l^H b_j| \right) \max_{1 \leq j \leq m} \left\{ |b_j^H \mathbf{h}^*| |a_j^H \mathbf{x}^*| |a_j^H \check{\mathbf{x}}^0| \right\} \\ &\stackrel{(ii)}{\leq} 8 \log m \cdot \frac{\mu}{\sqrt{m}} \cdot (5\sqrt{\log m}) \max_{1 \leq j \leq m} \left\{ |a_j^H \check{\mathbf{x}}^{0,(j)}| + \|a_j\|_2 \|\beta^{0,(j)} \check{\mathbf{x}}^0 - \check{\mathbf{x}}^{0,(j)}\|_2 \right\} \\ &\leq 200 \frac{\mu \log^2 m}{\sqrt{m}} + 120 \sqrt{\frac{\mu^2 K \log^3 m}{m}} \max_{1 \leq j \leq m} \|\beta^{0,(j)} \check{\mathbf{x}}^0 - \check{\mathbf{x}}^{0,(j)}\|_2, \end{aligned} \tag{237}$$



where  $\beta^{0,(j)}$  is as defined in (233). Here, (i) comes from the lower bound  $\sigma_1(\mathbf{M}) \geq 1/2$ . Bound (ii) follows by combining the incoherence condition (36), bound (189), the triangle inequality, as well as the estimate  $\sum_{j=1}^m \|\mathbf{b}_l^H \mathbf{b}_j\| \leq 4 \log m$  from Lemma 48. The last line uses the upper estimate  $\max_{1 \leq j \leq m} \|\mathbf{a}_j^H \check{\mathbf{x}}^{0,(j)}\| \leq 5\sqrt{\log m}$  and (190). Our bound (237) further implies

$$\max_{1 \leq l \leq m} \|\mathbf{b}_l^H \check{\mathbf{h}}^0\| \leq 200 \frac{\mu \log^2 m}{\sqrt{m}} + 120 \sqrt{\frac{\mu^2 K \log^3 m}{m}} \max_{1 \leq j \leq m} \|\beta^{0,(j)} \check{\mathbf{x}}^0 - \check{\mathbf{x}}^{0,(j)}\|_2. \tag{238}$$

The above bound (238) taken together with (236) gives

$$\begin{aligned} \max_{1 \leq l \leq m} \left\{ \|\beta^{0,(l)} \check{\mathbf{h}}^0 - \check{\mathbf{h}}^{0,(l)}\|_2 + \|\beta^{0,(l)} \check{\mathbf{x}}^0 - \check{\mathbf{x}}^{0,(l)}\|_2 \right\} &\leq 120c_1 \frac{\mu}{\sqrt{m}} \cdot \sqrt{\frac{K \log^2 m}{m}} \\ &+ 60c_1 \sqrt{\frac{\mu^2 K \log m}{m}} \left( 200 \frac{\mu \log^2 m}{\sqrt{m}} + 120 \sqrt{\frac{\mu^2 K \log^3 m}{m}} \max_{1 \leq j \leq m} \|\beta^{0,(j)} \check{\mathbf{x}}^0 - \check{\mathbf{x}}^{0,(j)}\|_2 \right). \end{aligned} \tag{239}$$

As long as  $m \gg \mu^2 K \log^2 m$ , we have  $60c_1 \sqrt{\mu^2 K \log m/m} \cdot 120 \sqrt{\mu^2 K \log^3 m/m} \leq 1/2$ . Rearranging terms, we are left with

$$\max_{1 \leq l \leq m} \left\{ \|\beta^{0,(l)} \check{\mathbf{h}}^0 - \check{\mathbf{h}}^{0,(l)}\|_2 + \|\beta^{0,(l)} \check{\mathbf{x}}^0 - \check{\mathbf{x}}^{0,(l)}\|_2 \right\} \leq c_3 \frac{\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 K \log^5 m}{m}} \tag{240}$$

for some constant  $c_3 > 0$ . Further, this bound combined with (238) yields

$$\max_{1 \leq l \leq m} \|\mathbf{b}_l^H \check{\mathbf{h}}^0\| \leq 200 \frac{\mu \log^2 m}{\sqrt{m}} + 120 \sqrt{\frac{\mu^2 K \log^3 m}{m}} \cdot c_3 \frac{\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 K \log^5 m}{m}} \leq c_2 \frac{\mu \log^2 m}{\sqrt{m}} \tag{241}$$

for some constant  $c_2 > 0$ , with the proviso that  $m \gg \mu^2 K \log^2 m$ .

We now translate the preceding bounds to the scaled version. Recall from bound (231) that

$$1/2 \leq 1 - \xi \leq \|\mathbf{M}\| = \sigma_1(\mathbf{M}) \leq 1 + \xi \leq 2, \tag{242}$$

as long as  $\xi \leq 1/2$ . For any  $\alpha \in \mathbb{C}$  with  $|\alpha| = 1$ ,  $\alpha \check{\mathbf{h}}^0, \alpha \check{\mathbf{x}}^0$  are still the leading left and right singular vectors of  $\mathbf{M}$ . Hence, we can use Lemma 60 to derive that

$$\begin{aligned} & \left| \sigma_1(\mathbf{M}) - \sigma_1(\mathbf{M}^{(l)}) \right| \\ & \leq \left\| (\mathbf{M} - \mathbf{M}^{(l)}) \check{\mathbf{x}}^{0,(l)} \right\|_2 + \left\{ \left\| \alpha \check{\mathbf{h}}^0 - \check{\mathbf{h}}^{0,(l)} \right\|_2 + \left\| \alpha \check{\mathbf{x}}^0 - \check{\mathbf{x}}^{0,(l)} \right\|_2 \right\} \|\mathbf{M}\| \\ & \leq \left\| (\mathbf{M} - \mathbf{M}^{(l)}) \check{\mathbf{x}}^{0,(l)} \right\|_2 + 2 \left\{ \left\| \alpha \check{\mathbf{h}}^0 - \check{\mathbf{h}}^{0,(l)} \right\|_2 + \left\| \alpha \check{\mathbf{x}}^0 - \check{\mathbf{x}}^{0,(l)} \right\|_2 \right\} \end{aligned}$$

and

$$\begin{aligned} & \left\| \alpha \mathbf{h}^0 - \mathbf{h}^{0,(l)} \right\|_2 + \left\| \alpha \mathbf{x}^0 - \mathbf{x}^{0,(l)} \right\|_2 \\ &= \left\| \sqrt{\sigma_1(\mathbf{M})} \left( \alpha \check{\mathbf{h}}^0 \right) - \sqrt{\sigma_1(\mathbf{M}^{(l)})} \check{\mathbf{h}}^{0,(l)} \right\|_2 + \left\| \sqrt{\sigma_1(\mathbf{M})} \alpha \check{\mathbf{x}}^0 - \sqrt{\sigma_1(\mathbf{M}^{(l)})} \check{\mathbf{x}}^{0,(l)} \right\|_2 \\ &\leq \sqrt{\sigma_1(\mathbf{M})} \left\{ \left\| \alpha \check{\mathbf{h}}^0 - \check{\mathbf{h}}^{0,(l)} \right\|_2 + \left\| \alpha \check{\mathbf{x}}^0 - \check{\mathbf{x}}^{0,(l)} \right\|_2 \right\} + \frac{2 \left| \sigma_1(\mathbf{M}) - \sigma_1(\mathbf{M}^{(l)}) \right|}{\sqrt{\sigma_1(\mathbf{M})} + \sqrt{\sigma_1(\mathbf{M}^{(l)})}} \\ &\leq \sqrt{2} \left\{ \left\| \alpha \check{\mathbf{h}}^0 - \check{\mathbf{h}}^{0,(l)} \right\|_2 + \left\| \alpha \check{\mathbf{x}}^0 - \check{\mathbf{x}}^{0,(l)} \right\|_2 \right\} + \sqrt{2} \left| \sigma_1(\mathbf{M}) - \sigma_1(\mathbf{M}^{(l)}) \right|. \end{aligned}$$

Taking the previous two bounds collectively yields

$$\begin{aligned} & \left\| \alpha \mathbf{h}^0 - \mathbf{h}^{0,(l)} \right\|_2 + \left\| \alpha \mathbf{x}^0 - \mathbf{x}^{0,(l)} \right\|_2 \leq \sqrt{2} \left\| (\mathbf{M} - \mathbf{M}^{(l)}) \check{\mathbf{x}}^{0,(l)} \right\|_2 \\ & \quad + 6 \left\{ \left\| \alpha \check{\mathbf{h}}^0 - \check{\mathbf{h}}^{0,(l)} \right\|_2 + \left\| \alpha \check{\mathbf{x}}^0 - \check{\mathbf{x}}^{0,(l)} \right\|_2 \right\}, \end{aligned}$$

which together with (235) and (240) implies

$$\min_{\alpha \in \mathbb{C}, |\alpha|=1} \left\{ \left\| \alpha \mathbf{h}^0 - \mathbf{h}^{0,(l)} \right\|_2 + \left\| \alpha \mathbf{x}^0 - \mathbf{x}^{0,(l)} \right\|_2 \right\} \leq c_5 \frac{\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 K \log^5 m}{m}} \tag{243}$$

for some constant  $c_5 > 0$ , as long as  $\xi$  is sufficiently small. Moreover, we have

$$\left\| \frac{1}{\alpha^0} \mathbf{h}^0 - \frac{\alpha}{\alpha^0} \mathbf{h}^{0,(l)} \right\|_2 + \left\| \alpha^0 \mathbf{x}^0 - \alpha \alpha^0 \mathbf{x}^{0,(l)} \right\|_2 \leq 2 \left\{ \left\| \mathbf{h}^0 - \alpha \mathbf{h}^{0,(l)} \right\|_2 + \left\| \mathbf{x}^0 - \alpha \mathbf{x}^{0,(l)} \right\|_2 \right\}$$

for any  $|\alpha| = 1$ , where  $\alpha^0$  is defined in (38) and, according to Lemma 19, satisfies

$$1/2 \leq |\alpha^0| \leq 2. \tag{244}$$

Therefore,

$$\begin{aligned} & \min_{\alpha \in \mathbb{C}, |\alpha|=1} \sqrt{\left\| \frac{1}{\alpha^0} \mathbf{h}^0 - \frac{\alpha}{\alpha^0} \mathbf{h}^{0,(l)} \right\|_2^2 + \left\| \alpha^0 \mathbf{x}^0 - \alpha \alpha^0 \mathbf{x}^{0,(l)} \right\|_2^2} \\ & \leq \min_{\alpha \in \mathbb{C}, |\alpha|=1} \left\{ \left\| \frac{1}{\alpha^0} \mathbf{h}^0 - \frac{\alpha}{\alpha^0} \mathbf{h}^{0,(l)} \right\|_2 + \left\| \alpha^0 \mathbf{x}^0 - \alpha \alpha^0 \mathbf{x}^{0,(l)} \right\|_2 \right\} \\ & \leq 2 \min_{\alpha \in \mathbb{C}, |\alpha|=1} \left\{ \left\| \mathbf{h}^0 - \alpha \mathbf{h}^{0,(l)} \right\|_2 + \left\| \mathbf{x}^0 - \alpha \mathbf{x}^{0,(l)} \right\|_2 \right\} \\ & \leq 2c_5 \frac{\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 K \log^5 m}{m}}. \end{aligned}$$

Furthermore, we have

$$\begin{aligned} \text{dist}(z^{0,(l)}, \tilde{z}^0) &= \min_{\alpha \in \mathbb{C}} \sqrt{\left\| \frac{1}{\alpha} \mathbf{h}^{0,(l)} - \frac{1}{\alpha^0} \mathbf{h}^0 \right\|_2^2 + \left\| \alpha \mathbf{x}^{0,(l)} - \alpha^0 \mathbf{x}^0 \right\|_2^2} \\ &\leq \min_{\alpha \in \mathbb{C}, |\alpha|=1} \sqrt{\left\| \frac{1}{\alpha^0} \mathbf{h}^0 - \frac{\alpha}{\alpha^0} \mathbf{h}^{0,(l)} \right\|_2^2 + \left\| \alpha^0 \mathbf{x}^0 - \alpha \alpha^0 \mathbf{x}^{0,(l)} \right\|_2^2} \\ &\leq 2c_5 \frac{\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 K \log^5 m}{m}}, \end{aligned} \tag{245}$$

where the second line follows since the latter is minimizing over a smaller feasible set. This completes the proof for claim (96).

Regarding  $|\mathbf{b}_l^H \tilde{\mathbf{h}}^0|$ , one first sees that

$$|\mathbf{b}_l^H \mathbf{h}^0| = \left| \sqrt{\sigma_1(\mathbf{M})} \mathbf{b}_l^H \check{\mathbf{h}}^0 \right| \leq \sqrt{2} c_2 \frac{\mu \log^2 m}{\sqrt{m}},$$

where the last relation holds due to (241) and (242). Hence, using the property (244), we have

$$|\mathbf{b}_l^H \tilde{\mathbf{h}}^0| = \left| \mathbf{b}_l^H \frac{1}{\alpha^0} \mathbf{h}^0 \right| \leq \left| \frac{1}{\alpha^0} \right| |\mathbf{b}_l^H \mathbf{h}^0| \leq 2\sqrt{2} c_2 \frac{\mu \log^2 m}{\sqrt{m}},$$

which finishes the proof of claim (97).

Before concluding this section, we note a by-product of the proof. Specifically, we can establish the claim required in (226) using many results derived in this section. This is formally stated in the following lemma.

**Lemma 30** Fix any small constant  $c > 0$ . Suppose the number of samples obeys  $m \gg \tau K \log^4 m$ . Then with probability at least  $1 - O(m^{-10})$ , we have

$$\max_{1 \leq j \leq \tau} \left| (\mathbf{b}_j - \mathbf{b}_1)^H \tilde{\mathbf{h}}^0 \right| \leq c \frac{\mu}{\sqrt{m}} \log m.$$

**Proof** Instate the notation and hypotheses in Appendix C.6. Recognize that

$$\begin{aligned} \left| (\mathbf{b}_j - \mathbf{b}_1)^H \tilde{\mathbf{h}}^0 \right| &= \left| (\mathbf{b}_j - \mathbf{b}_1)^H \frac{1}{\alpha^0} \mathbf{h}^0 \right| = \left| (\mathbf{b}_j - \mathbf{b}_1)^H \frac{1}{\alpha^0} \sqrt{\sigma_1(\mathbf{M})} \check{\mathbf{h}}^0 \right| \\ &\leq \left| \frac{1}{\alpha^0} \right| \sqrt{\sigma_1(\mathbf{M})} \left| (\mathbf{b}_j - \mathbf{b}_1)^H \check{\mathbf{h}}^0 \right| \\ &\leq 4 \left| (\mathbf{b}_j - \mathbf{b}_1)^H \check{\mathbf{h}}^0 \right|, \end{aligned}$$

where the last inequality comes from (242) and (244). It thus suffices to prove that  $\left| (\mathbf{b}_j - \mathbf{b}_1)^H \check{\mathbf{h}}^0 \right| \leq c \mu \log m / \sqrt{m}$  for some  $c > 0$  small enough. To this end, it can be

seen that

$$\begin{aligned}
 |(\mathbf{b}_j - \mathbf{b}_1)^H \check{\mathbf{h}}^0| &= \frac{1}{\sigma_1(\mathbf{M})} |(\mathbf{b}_j - \mathbf{b}_1)^H \mathbf{M} \check{\mathbf{x}}^0| \\
 &\leq 2 \left| \sum_{k=1}^m (\mathbf{b}_j - \mathbf{b}_1)^H \mathbf{b}_k \mathbf{b}_k^H \mathbf{h}^* \mathbf{x}^{*H} \mathbf{a}_k \mathbf{a}_k^H \check{\mathbf{x}}^0 \right| \\
 &\leq 2 \left( \sum_{k=1}^m |(\mathbf{b}_j - \mathbf{b}_1)^H \mathbf{b}_k| \right) \max_{1 \leq k \leq m} \left\{ |\mathbf{b}_k^H \mathbf{h}^*| |\mathbf{a}_k^H \mathbf{x}^*| |\mathbf{a}_k^H \check{\mathbf{x}}^0| \right\} \\
 &\stackrel{(i)}{\leq} c \frac{1}{\log^2 m} \cdot \frac{\mu}{\sqrt{m}} \cdot (5\sqrt{\log m}) \max_{1 \leq j \leq m} \left\{ |\mathbf{a}_j^H \check{\mathbf{x}}^{0,(j)}| \right. \\
 &\quad \left. + \|\mathbf{a}_j\|_2 \|\alpha^{0,(j)} \check{\mathbf{x}}^0 - \check{\mathbf{x}}^{0,(j)}\|_2 \right\} \\
 &\stackrel{(ii)}{\lesssim} c \frac{\mu}{\sqrt{m}} \frac{1}{\log m} \leq c \frac{\mu}{\sqrt{m}} \log m, \tag{246}
 \end{aligned}$$

where (i) comes from Lemma 50, the incoherence condition (36), and estimate (189). The last line (ii) holds since we have already established (see (237) and (240))

$$\max_{1 \leq j \leq m} \left\{ |\mathbf{a}_j^H \check{\mathbf{x}}^{0,(j)}| + \|\mathbf{a}_j\|_2 \|\alpha^{0,(j)} \check{\mathbf{x}}^0 - \check{\mathbf{x}}^{0,(j)}\|_2 \right\} \lesssim \sqrt{\log m}.$$

The proof is then complete. □

### C.7 Proof of Lemma 21

Recall that  $\alpha^0$  and  $\alpha^{0,(l)}$  are the alignment parameters between  $\mathbf{z}^0$  and  $\mathbf{z}^*$ , and between  $\mathbf{z}^{0,(l)}$  and  $\mathbf{z}^*$ , respectively, that is,

$$\begin{aligned}
 \alpha^0 &:= \operatorname{argmin}_{\alpha \in \mathbb{C}} \left\{ \left\| \frac{1}{\alpha} \mathbf{h}^0 - \mathbf{h}^* \right\|_2^2 + \|\alpha \mathbf{x}^0 - \mathbf{x}^*\|_2^2 \right\}, \\
 \alpha^{0,(l)} &:= \operatorname{argmin}_{\alpha \in \mathbb{C}} \left\{ \left\| \frac{1}{\alpha} \mathbf{h}^{0,(l)} - \mathbf{h}^* \right\|_2^2 + \|\alpha \mathbf{x}^{0,(l)} - \mathbf{x}^*\|_2^2 \right\}.
 \end{aligned}$$

Also, we let

$$\alpha_{\text{mutual}}^{0,(l)} := \operatorname{argmin}_{\alpha \in \mathbb{C}} \left\{ \left\| \frac{1}{\alpha} \mathbf{h}^{0,(l)} - \frac{1}{\alpha^0} \mathbf{h}^0 \right\|_2^2 + \|\alpha \mathbf{x}^{0,(l)} - \alpha^0 \mathbf{x}^0\|_2^2 \right\}.$$

The triangle inequality together with (94) and (245) then tells us that

$$\sqrt{\left\| \frac{1}{\alpha_{\text{mutual}}^{0,(l)}} \mathbf{h}^{0,(l)} - \mathbf{h}^* \right\|_2^2 + \|\alpha_{\text{mutual}}^{0,(l)} \mathbf{x}^{0,(l)} - \mathbf{x}^*\|_2^2}$$

$$\begin{aligned}
 &\leq \sqrt{\left\| \frac{1}{\alpha^0} \mathbf{h}^0 - \frac{1}{\alpha_{\text{mutual}}^{0,(l)}} \mathbf{h}^{0,(l)} \right\|_2^2 + \left\| \alpha^0 \mathbf{x}^0 - \alpha_{\text{mutual}}^{0,(l)} \mathbf{x}^{0,(l)} \right\|_2^2} \\
 &\quad + \sqrt{\left\| \frac{1}{\alpha^0} \mathbf{h}^0 - \mathbf{h}^* \right\|_2^2 + \left\| \alpha^0 \mathbf{x}^0 - \mathbf{x}^* \right\|_2^2} \\
 &\leq 2c_5 \frac{\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 K \log^5 m}{m}} + C_1 \frac{1}{\log^2 m} \\
 &\leq 2C_1 \frac{1}{\log^2 m},
 \end{aligned}$$

where the last relation holds as long as  $m \gg \mu^2 \sqrt{K} \log^{9/2} m$ .

Let

$$\mathbf{x}_1 = \alpha^0 \mathbf{x}^0, \quad \mathbf{h}_1 = \frac{1}{\alpha^0} \mathbf{h}^0 \quad \text{and} \quad \mathbf{x}_2 = \alpha_{\text{mutual}}^{0,(l)} \mathbf{x}^{0,(l)}, \quad \mathbf{h}_2 = \frac{1}{\alpha_{\text{mutual}}^{0,(l)}} \mathbf{h}^{0,(l)}.$$

It is easy to see that  $\mathbf{x}_1, \mathbf{h}_1, \mathbf{x}_2, \mathbf{h}_2$  satisfy the assumptions in Lemma 55, which implies

$$\begin{aligned}
 &\sqrt{\left\| \frac{1}{\alpha_{0,(l)}} \mathbf{h}^{0,(l)} - \frac{1}{\alpha^0} \mathbf{h}^0 \right\|_2^2 + \left\| \alpha_{0,(l)} \mathbf{x}^{0,(l)} - \alpha^0 \mathbf{x}^0 \right\|_2^2} \\
 &\lesssim \sqrt{\left\| \frac{1}{\alpha^0} \mathbf{h}^0 - \frac{1}{\alpha_{\text{mutual}}^{0,(l)}} \mathbf{h}^{0,(l)} \right\|_2^2 + \left\| \alpha^0 \mathbf{x}^0 - \alpha_{\text{mutual}}^{0,(l)} \mathbf{x}^{0,(l)} \right\|_2^2} \\
 &\lesssim \frac{\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 K \log^5 m}{m}}, \tag{247}
 \end{aligned}$$

where the last line comes from (245). With this upper estimate at hand, we are now ready to show that with high probability,

$$\begin{aligned}
 &\left| \mathbf{a}_l^H (\alpha^0 \mathbf{x}^0 - \mathbf{x}^*) \right| \stackrel{(i)}{\leq} \left| \mathbf{a}_l^H (\alpha^{0,(l)} \mathbf{x}^{0,(l)} - \mathbf{x}^*) \right| + \left| \mathbf{a}_l^H (\alpha^0 \mathbf{x}^0 - \alpha^{0,(l)} \mathbf{x}^{0,(l)}) \right| \\
 &\stackrel{(ii)}{\leq} 5\sqrt{\log m} \left\| \alpha^{0,(l)} \mathbf{x}^{0,(l)} - \mathbf{x}^* \right\|_2 + \|\mathbf{a}_l\|_2 \left\| \alpha^0 \mathbf{x}^0 - \alpha^{0,(l)} \mathbf{x}^{0,(l)} \right\|_2 \\
 &\stackrel{(iii)}{\lesssim} \sqrt{\log m} \cdot \frac{1}{\log^2 m} + \sqrt{K} \frac{\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 K \log^5 m}{m}} \\
 &\stackrel{(iv)}{\lesssim} \frac{1}{\log^{3/2} m},
 \end{aligned}$$

where (i) follows from the triangle inequality, (ii) uses Cauchy–Schwarz and the independence between  $\mathbf{x}^{0,(l)}$  and  $\mathbf{a}_l$ , (iii) holds because of (95) and (247) under the condition  $m \gg \mu^2 K \log^6 m$ , and (iv) holds true as long as  $m \gg \mu^2 K \log^4 m$ .

## D Technical Lemmas

### D.1 Technical Lemmas for Phase Retrieval

#### D.1.1 Matrix Concentration Inequalities

**Lemma 31** *Suppose that  $\mathbf{a}_j \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$  for every  $1 \leq j \leq m$ . Fix any small constant  $\delta > 0$ . With probability at least  $1 - C_2 e^{-c_2 m}$ , one has*

$$\left\| \frac{1}{m} \sum_{j=1}^m \mathbf{a}_j \mathbf{a}_j^\top - \mathbf{I}_n \right\| \leq \delta,$$

as long as  $m \geq c_0 n$  for some sufficiently large constant  $c_0 > 0$ . Here,  $C_2, c_2 > 0$  are some universal constants.

**Proof** This is an immediate consequence of [116, Corollary 5.35]. □

**Lemma 32** *Suppose that  $\mathbf{a}_j \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ , for every  $1 \leq j \leq m$ . Fix any small constant  $\delta > 0$ . With probability at least  $1 - O(n^{-10})$ , we have*

$$\left\| \frac{1}{m} \sum_{j=1}^m (\mathbf{a}_j^\top \mathbf{x}^*)^2 \mathbf{a}_j \mathbf{a}_j^\top - (\|\mathbf{x}^*\|_2^2 \mathbf{I}_n + 2\mathbf{x}^* \mathbf{x}^{*\top}) \right\| \leq \delta \|\mathbf{x}^*\|_2^2,$$

provided that  $m \geq c_0 n \log n$  for some sufficiently large constant  $c_0 > 0$ .

**Proof** This is adapted from [18, Lemma 7.4]. □

**Lemma 33** *Suppose that  $\mathbf{a}_j \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ , for every  $1 \leq j \leq m$ . Fix any small constant  $\delta > 0$  and any constant  $C > 0$ . Suppose  $m \geq c_0 n$  for some sufficiently large constant  $c_0 > 0$ . Then with probability at least  $1 - C_2 e^{-c_2 m}$ ,*

$$\left\| \frac{1}{m} \sum_{j=1}^m (\mathbf{a}_j^\top \mathbf{x})^2 \mathbb{1}_{\{|\mathbf{a}_j^\top \mathbf{x}| \leq C\}} \mathbf{a}_j \mathbf{a}_j^\top - (\beta_1 \mathbf{x} \mathbf{x}^\top + \beta_2 \|\mathbf{x}\|_2^2 \mathbf{I}_n) \right\| \leq \delta \|\mathbf{x}\|_2^2, \quad \forall \mathbf{x} \in \mathbb{R}^n,$$

holds for some absolute constants  $c_2, C_2 > 0$ , where

$$\beta_1 := \mathbb{E} \left[ \xi^4 \mathbb{1}_{\{|\xi| \leq C\}} \right] - \mathbb{E} \left[ \xi^2 \mathbb{1}_{\{|\xi| \leq C\}} \right] \quad \text{and} \quad \beta_2 = \mathbb{E} \left[ \xi^2 \mathbb{1}_{\{|\xi| \leq C\}} \right]$$

with  $\xi$  being a standard Gaussian random variable.

**Proof** This is supplied in [25, supplementary material]. □

### D.1.2 Matrix Perturbation Bounds

**Lemma 34** *Let  $\lambda_1(A)$ ,  $\mathbf{u}$  be the leading eigenvalue and eigenvector of a symmetric matrix  $A$ , respectively, and  $\lambda_1(\tilde{A})$ ,  $\tilde{\mathbf{u}}$  be the leading eigenvalue and eigenvector of a symmetric matrix  $\tilde{A}$ , respectively. Suppose that  $\lambda_1(A)$ ,  $\lambda_1(\tilde{A})$ ,  $\|A\|$ ,  $\|\tilde{A}\| \in [C_1, C_2]$  for some  $C_1, C_2 > 0$ . Then,*

$$\left\| \sqrt{\lambda_1(A)} \mathbf{u} - \sqrt{\lambda_1(\tilde{A})} \tilde{\mathbf{u}} \right\|_2 \leq \frac{\|(A - \tilde{A})\mathbf{u}\|_2}{2\sqrt{C_1}} + \left( \sqrt{C_2} + \frac{C_2}{\sqrt{C_1}} \right) \|\mathbf{u} - \tilde{\mathbf{u}}\|_2.$$

**Proof** Observe that

$$\begin{aligned} & \left\| \sqrt{\lambda_1(A)} \mathbf{u} - \sqrt{\lambda_1(\tilde{A})} \tilde{\mathbf{u}} \right\|_2 \\ & \leq \left\| \sqrt{\lambda_1(A)} \mathbf{u} - \sqrt{\lambda_1(\tilde{A})} \mathbf{u} \right\|_2 + \left\| \sqrt{\lambda_1(\tilde{A})} \mathbf{u} - \sqrt{\lambda_1(\tilde{A})} \tilde{\mathbf{u}} \right\|_2 \\ & \leq \left| \sqrt{\lambda_1(A)} - \sqrt{\lambda_1(\tilde{A})} \right| + \sqrt{\lambda_1(\tilde{A})} \|\mathbf{u} - \tilde{\mathbf{u}}\|_2, \end{aligned} \tag{248}$$

where the last inequality follows since  $\|\mathbf{u}\|_2 = 1$ . Using the identity  $\sqrt{a} - \sqrt{b} = (a - b)/(\sqrt{a} + \sqrt{b})$ , we have

$$\left| \sqrt{\lambda_1(A)} - \sqrt{\lambda_1(\tilde{A})} \right| = \frac{|\lambda_1(A) - \lambda_1(\tilde{A})|}{\left| \sqrt{\lambda_1(A)} + \sqrt{\lambda_1(\tilde{A})} \right|} \leq \frac{|\lambda_1(A) - \lambda_1(\tilde{A})|}{2\sqrt{C_1}},$$

where the last inequality comes from our assumptions on  $\lambda_1(A)$  and  $\lambda_1(\tilde{A})$ . This combined with (248) yields

$$\left\| \sqrt{\lambda_1(A)} \mathbf{u} - \sqrt{\lambda_1(\tilde{A})} \tilde{\mathbf{u}} \right\|_2 \leq \frac{|\lambda_1(A) - \lambda_1(\tilde{A})|}{2\sqrt{C_1}} + \sqrt{C_2} \|\mathbf{u} - \tilde{\mathbf{u}}\|_2. \tag{249}$$

To control  $|\lambda_1(A) - \lambda_1(\tilde{A})|$ , use the relationship between the eigenvalue and the eigenvector to obtain

$$\begin{aligned} |\lambda_1(A) - \lambda_1(\tilde{A})| &= \left| \mathbf{u}^\top A \mathbf{u} - \tilde{\mathbf{u}}^\top \tilde{A} \tilde{\mathbf{u}} \right| \\ &\leq \left| \mathbf{u}^\top (A - \tilde{A}) \mathbf{u} \right| + \left| \mathbf{u}^\top \tilde{A} \mathbf{u} - \tilde{\mathbf{u}}^\top \tilde{A} \mathbf{u} \right| + \left| \tilde{\mathbf{u}}^\top \tilde{A} \mathbf{u} - \tilde{\mathbf{u}}^\top \tilde{A} \tilde{\mathbf{u}} \right| \\ &\leq \|(A - \tilde{A})\mathbf{u}\|_2 + 2 \|\mathbf{u} - \tilde{\mathbf{u}}\|_2 \|\tilde{A}\|, \end{aligned}$$

which together with (249) gives

$$\begin{aligned} \left\| \sqrt{\lambda_1(\mathbf{A})} \mathbf{u} - \sqrt{\lambda_1(\tilde{\mathbf{A}})} \tilde{\mathbf{u}} \right\|_2 &\leq \frac{\|(\mathbf{A} - \tilde{\mathbf{A}})\mathbf{u}\|_2 + 2\|\mathbf{u} - \tilde{\mathbf{u}}\|_2 \|\tilde{\mathbf{A}}\|}{2\sqrt{C_1}} + \sqrt{C_2} \|\mathbf{u} - \tilde{\mathbf{u}}\|_2 \\ &\leq \frac{\|(\mathbf{A} - \tilde{\mathbf{A}})\mathbf{u}\|_2}{2\sqrt{C_1}} + \left( \frac{C_2}{\sqrt{C_1}} + \sqrt{C_2} \right) \|\mathbf{u} - \tilde{\mathbf{u}}\|_2 \end{aligned}$$

as claimed. □

## D.2 Technical Lemmas for Matrix Completion

### D.2.1 Orthogonal Procrustes Problem

The orthogonal Procrustes problem is a matrix approximation problem which seeks an orthogonal matrix  $\mathbf{R}$  to best “align” two matrices  $\mathbf{A}$  and  $\mathbf{B}$ . Specifically, for  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times r}$ , define  $\hat{\mathbf{R}}$  to be the minimizer of

$$\text{minimize}_{\mathbf{R} \in \mathcal{O}^{r \times r}} \|\mathbf{A}\mathbf{R} - \mathbf{B}\|_F. \tag{250}$$

The first lemma is concerned with the characterization of the minimizer  $\hat{\mathbf{R}}$  of (250).

**Lemma 35** *For  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times r}$ ,  $\hat{\mathbf{R}}$  is the minimizer of (250) if and only if  $\hat{\mathbf{R}}^\top \mathbf{A}^\top \mathbf{B}$  is symmetric and positive semidefinite.*

**Proof** This is an immediate consequence of [112, Theorem 2]. □

Let  $\mathbf{A}^\top \mathbf{B} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$  be the singular value decomposition of  $\mathbf{A}^\top \mathbf{B} \in \mathbb{R}^{r \times r}$ . It is easy to check that  $\hat{\mathbf{R}} := \mathbf{U}\mathbf{V}^\top$  satisfies the conditions that  $\hat{\mathbf{R}}^\top \mathbf{A}^\top \mathbf{B}$  is both symmetric and positive semidefinite. In view of Lemma 35,  $\hat{\mathbf{R}} = \mathbf{U}\mathbf{V}^\top$  is the minimizer of (250). In the special case when  $\mathbf{C} := \mathbf{A}^\top \mathbf{B}$  is invertible,  $\hat{\mathbf{R}}$  enjoys the following equivalent form:

$$\hat{\mathbf{R}} = \hat{\mathbf{H}}(\mathbf{C}) := \mathbf{C} \left( \mathbf{C}^\top \mathbf{C} \right)^{-1/2}, \tag{251}$$

where  $\hat{\mathbf{H}}(\cdot)$  is an  $\mathbb{R}^{r \times r}$ -valued function on  $\mathbb{R}^{r \times r}$ . This motivates us to look at the perturbation bounds for the matrix-valued function  $\hat{\mathbf{H}}(\cdot)$ , which is formulated in the following lemma.

**Lemma 36** *Let  $\mathbf{C} \in \mathbb{R}^{r \times r}$  be a nonsingular matrix. Then for any matrix  $\mathbf{E} \in \mathbb{R}^{r \times r}$  with  $\|\mathbf{E}\| \leq \sigma_{\min}(\mathbf{C})$  and any unitarily invariant norm  $\|\cdot\|$ , one has*

$$\|\hat{\mathbf{H}}(\mathbf{C} + \mathbf{E}) - \hat{\mathbf{H}}(\mathbf{C})\| \leq \frac{2}{\sigma_{r-1}(\mathbf{C}) + \sigma_r(\mathbf{C})} \|\mathbf{E}\|,$$

where  $\hat{\mathbf{H}}(\cdot)$  is defined above.

**Proof** This is an immediate consequence of [85, Theorem 2.3]. □



With Lemma 36 in place, we are ready to present the following bounds on two matrices after “aligning” them with  $X^*$ .

**Lemma 37** *Instate the notation in Sect. 3.2. Suppose  $X_1, X_2 \in \mathbb{R}^{n \times r}$  are two matrices such that*

$$\|X_1 - X^*\| \|X^*\| \leq \sigma_{\min}/2, \tag{252a}$$

$$\|X_1 - X_2\| \|X^*\| \leq \sigma_{\min}/4. \tag{252b}$$

Denote

$$R_1 := \operatorname{argmin}_{R \in \mathcal{O}^{r \times r}} \|X_1 R - X^*\|_F \quad \text{and} \quad R_2 := \operatorname{argmin}_{R \in \mathcal{O}^{r \times r}} \|X_2 R - X^*\|_F.$$

Then the following two inequalities hold true:

$$\|X_1 R_1 - X_2 R_2\| \leq 5\kappa \|X_1 - X_2\| \quad \text{and} \quad \|X_1 R_1 - X_2 R_2\|_F \leq 5\kappa \|X_1 - X_2\|_F.$$

**Proof** Before proving the claims, we first gather some immediate consequences of assumptions (252). Denote  $C = X_1^\top X^*$  and  $E = (X_2 - X_1)^\top X^*$ . It is easily seen that  $C$  is invertible since

$$\|C - X_1^\top X^*\| \leq \|X_1 - X^*\| \|X^*\| \stackrel{(i)}{\leq} \sigma_{\min}/2 \quad \stackrel{(ii)}{\implies} \quad \sigma_r(C) \geq \sigma_{\min}/2, \tag{253}$$

where (i) follows from assumption (252a) and (ii) is a direct application of Weyl’s inequality. In addition,  $C + E = X_2^\top X^*$  is also invertible since

$$\|E\| \leq \|X_1 - X_2\| \|X^*\| \stackrel{(i)}{\leq} \sigma_{\min}/4 \stackrel{(ii)}{<} \sigma_r(C),$$

where (i) arises from assumption (252b) and (ii) holds because of (253). When both  $C$  and  $C + E$  are invertible, the orthonormal matrices  $R_1$  and  $R_2$  admit closed-form expressions as follows

$$R_1 = C \left( C^\top C \right)^{-1/2} \quad \text{and} \quad R_2 = (C + E) \left[ (C + E)^\top (C + E) \right]^{-1/2}.$$

Moreover, we have the following bound on  $\|X_1\|$ :

$$\|X_1\| \stackrel{(i)}{\leq} \|X_1 - X^*\| + \|X^*\| \stackrel{(ii)}{\leq} \frac{\sigma_{\min}}{2 \|X^*\|} + \|X^*\| \leq \frac{\sigma_{\max}}{2 \|X^*\|} + \|X^*\| \stackrel{(iii)}{\leq} 2 \|X^*\|, \tag{254}$$

where (i) is the triangle inequality, (ii) uses assumption (252a), and (iii) arises from the fact that  $\|X^*\| = \sqrt{\sigma_{\max}}$ .

With these in place, we turn to establishing the claimed bounds. We will focus on the upper bound on  $\|X_1 R_1 - X_2 R_2\|_F$ , as the bound on  $\|X_1 R_1 - X_2 R_2\|$  can be

easily obtained using the same argument. Simple algebra reveals that

$$\begin{aligned} \|X_1 R_1 - X_2 R_2\|_F &= \|(X_1 - X_2) R_2 + X_1 (R_1 - R_2)\|_F \\ &\leq \|X_1 - X_2\|_F + \|X_1\| \|R_1 - R_2\|_F \\ &\leq \|X_1 - X_2\|_F + 2 \|X^*\| \|R_1 - R_2\|_F, \end{aligned} \tag{255}$$

where the first inequality uses the fact that  $\|R_2\| = 1$  and the last inequality comes from (254). An application of Lemma 36 leads us to conclude that

$$\begin{aligned} \|R_1 - R_2\|_F &\leq \frac{2}{\sigma_r(C) + \sigma_{r-1}(C)} \|E\|_F \\ &\leq \frac{2}{\sigma_{\min}} \|(X_2 - X_1)^\top X^*\|_F \end{aligned} \tag{256}$$

$$\leq \frac{2}{\sigma_{\min}} \|X_2 - X_1\|_F \|X^*\|, \tag{257}$$

where (256) utilizes (253). Combine (255) and (257) to reach

$$\begin{aligned} \|X_1 R_1 - X_2 R_2\|_F &\leq \|X_1 - X_2\|_F + \frac{4}{\sigma_{\min}} \|X_2 - X_1\|_F \|X^*\|^2 \\ &\leq (1 + 4\kappa) \|X_1 - X_2\|_F, \end{aligned}$$

which finishes the proof by noting that  $\kappa \geq 1$ . □

### D.2.2 Matrix Concentration Inequalities

This section collects various measure concentration results regarding the Bernoulli random variables  $\{\delta_{j,k}\}_{1 \leq j,k \leq n}$ , which is ubiquitous in the analysis for matrix completion.

**Lemma 38** *Fix any small constant  $\delta > 0$ , and suppose that  $m \gg \delta^{-2} \mu nr \log n$ . Then with probability exceeding  $1 - O(n^{-10})$ , one has*

$$(1 - \delta) \|B\|_F \leq \frac{1}{\sqrt{p}} \|\mathcal{P}_\Omega(B)\|_F \leq (1 + \delta) \|B\|_F$$

which holds simultaneously for all  $B \in \mathbb{R}^{n \times n}$  lying within the tangent space of  $M^*$ .

**Proof** This result has been established in [19, Section 4.2] for asymmetric sampling patterns (where each  $(i, j), i \neq j$ , is included in  $\Omega$  independently). It is straightforward to extend the proof and the result to symmetric sampling patterns (where each  $(i, j), i \geq j$ , is included in  $\Omega$  independently). We omit the proof for conciseness. □

**Lemma 39** Fix a matrix  $\mathbf{M} \in \mathbb{R}^{n \times n}$ . Suppose  $n^2 p \geq c_0 n \log n$  for some sufficiently large constant  $c_0 > 0$ . With probability at least  $1 - O(n^{-10})$ , one has

$$\left\| \frac{1}{p} \mathcal{P}_\Omega(\mathbf{M}) - \mathbf{M} \right\| \leq C \sqrt{\frac{n}{p}} \|\mathbf{M}\|_\infty,$$

where  $C > 0$  is some absolute constant.

**Proof** See [64, Lemma 3.2]. Similar to Lemma 38, the result therein was provided for the asymmetric sampling patterns but can be easily extended to the symmetric case.  $\square$

**Lemma 40** Recall from Sect. 3.2 that  $\mathbf{E} \in \mathbb{R}^{n \times n}$  is the symmetric noise matrix. Suppose the sample size obeys  $n^2 p \geq c_0 n \log^2 n$  for some sufficiently large constant  $c_0 > 0$ . With probability at least  $1 - O(n^{-10})$ , one has

$$\left\| \frac{1}{p} \mathcal{P}_\Omega(\mathbf{E}) \right\| \leq C \sigma \sqrt{\frac{n}{p}},$$

where  $C > 0$  is some universal constant.

**Proof** See [32, Lemma 11].  $\square$

**Lemma 41** Fix some matrix  $\mathbf{A} \in \mathbb{R}^{n \times r}$  with  $n \geq 2r$  and some  $1 \leq l \leq n$ . Suppose  $\{\delta_{l,j}\}_{1 \leq j \leq n}$  are independent Bernoulli random variables with means  $\{p_j\}_{1 \leq j \leq n}$  no more than  $p$ . Define

$$\mathbf{G}_l(\mathbf{A}) := \left[ \delta_{l,1} \mathbf{A}_{1,\cdot}^\top, \delta_{l,2} \mathbf{A}_{2,\cdot}^\top, \dots, \delta_{l,n} \mathbf{A}_{n,\cdot}^\top \right] \in \mathbb{R}^{r \times n}.$$

Then one has

$$\text{Median} [\|\mathbf{G}_l(\mathbf{A})\|] \leq \sqrt{p \|\mathbf{A}\|^2 + \sqrt{2p \|\mathbf{A}\|_{2,\infty}^2 \|\mathbf{A}\|^2 \log(4r)} + \frac{2 \|\mathbf{A}\|_{2,\infty}^2}{3} \log(4r)}$$

and for any constant  $C \geq 3$ , with probability exceeding  $1 - n^{-(1.5C-1)}$

$$\left\| \sum_{j=1}^n (\delta_{l,j} - p) \mathbf{A}_{j,\cdot}^\top \cdot \mathbf{A}_{j,\cdot} \right\| \leq C \left( \sqrt{p \|\mathbf{A}\|_{2,\infty}^2 \|\mathbf{A}\|^2 \log n + \|\mathbf{A}\|_{2,\infty}^2 \log n} \right),$$

and

$$\|\mathbf{G}_l(\mathbf{A})\| \leq \sqrt{p \|\mathbf{A}\|^2 + C \left( \sqrt{p \|\mathbf{A}\|_{2,\infty}^2 \|\mathbf{A}\|^2 \log n + \|\mathbf{A}\|_{2,\infty}^2 \log n} \right)}.$$

**Proof** By the definition of  $G_l(\mathbf{A})$  and the triangle inequality, one has

$$\begin{aligned} \|G_l(\mathbf{A})\|^2 &= \|G_l(\mathbf{A}) G_l(\mathbf{A})^\top\| = \left\| \sum_{j=1}^n \delta_{l,j} \mathbf{A}_{j,\cdot}^\top \mathbf{A}_{j,\cdot} \right\| \\ &\leq \left\| \sum_{j=1}^n (\delta_{l,j} - p_j) \mathbf{A}_{j,\cdot}^\top \mathbf{A}_{j,\cdot} \right\| + p \|\mathbf{A}\|^2. \end{aligned}$$

Therefore, it suffices to control the first term. It can be seen that  $\{(\delta_{l,j} - p_j) \mathbf{A}_{j,\cdot}^\top \mathbf{A}_{j,\cdot}\}_{1 \leq j \leq n}$  are i.i.d. zero-mean random matrices. Letting

$$\begin{aligned} L &:= \max_{1 \leq j \leq n} \|(\delta_{l,j} - p_j) \mathbf{A}_{j,\cdot}^\top \mathbf{A}_{j,\cdot}\| \leq \|\mathbf{A}\|_{2,\infty}^2 \\ \text{and } V &:= \left\| \sum_{j=1}^n \mathbb{E} \left[ (\delta_{l,j} - p_j)^2 \mathbf{A}_{j,\cdot}^\top \mathbf{A}_{j,\cdot} \mathbf{A}_{j,\cdot}^\top \mathbf{A}_{j,\cdot} \right] \right\| \\ &\leq \mathbb{E} \left[ (\delta_{l,j} - p_j)^2 \right] \|\mathbf{A}\|_{2,\infty}^2 \left\| \sum_{j=1}^n \mathbf{A}_{j,\cdot}^\top \mathbf{A}_{j,\cdot} \right\| \leq p \|\mathbf{A}\|_{2,\infty}^2 \|\mathbf{A}\|^2 \end{aligned}$$

and invoking matrix Bernstein’s inequality [114, Theorem 6.1.1], one has for all  $t \geq 0$ ,

$$\mathbb{P} \left\{ \left\| \sum_{j=1}^n (\delta_{l,j} - p_j) \mathbf{A}_{j,\cdot}^\top \mathbf{A}_{j,\cdot} \right\| \geq t \right\} \leq 2r \cdot \exp \left( \frac{-t^2/2}{p \|\mathbf{A}\|_{2,\infty}^2 \|\mathbf{A}\|^2 + \|\mathbf{A}\|_{2,\infty}^2 \cdot t/3} \right). \tag{258}$$

We can thus find an upper bound on  $\text{Median} \left[ \left\| \sum_{j=1}^n (\delta_{l,j} - p_j) \mathbf{A}_{j,\cdot}^\top \mathbf{A}_{j,\cdot} \right\| \right]$  by finding a value  $t$  that ensures the right-hand side of (258) is smaller than  $1/2$ . Using this strategy and some simple calculations, we get

$$\text{Median} \left[ \left\| \sum_{j=1}^n (\delta_{l,j} - p_j) \mathbf{A}_{j,\cdot}^\top \mathbf{A}_{j,\cdot} \right\| \right] \leq \sqrt{2p \|\mathbf{A}\|_{2,\infty}^2 \|\mathbf{A}\|^2 \log(4r)} + \frac{2 \|\mathbf{A}\|_{2,\infty}^2}{3} \log(4r)$$

and for any  $C \geq 3$ ,

$$\left\| \sum_{j=1}^n (\delta_{l,j} - p_j) \mathbf{A}_{j,\cdot}^\top \mathbf{A}_{j,\cdot} \right\| \leq C \left( \sqrt{p \|\mathbf{A}\|_{2,\infty}^2 \|\mathbf{A}\|^2 \log n} + \|\mathbf{A}\|_{2,\infty}^2 \log n \right)$$

holds with probability at least  $1 - n^{-(1.5C-1)}$ . As a consequence, we have

$$\text{Median} [\|G_l(\mathbf{A})\|] \leq \sqrt{p \|\mathbf{A}\|^2 + \sqrt{2p \|\mathbf{A}\|_{2,\infty}^2 \|\mathbf{A}\|^2 \log(4r)} + \frac{2 \|\mathbf{A}\|_{2,\infty}^2 \log(4r)}{3},$$

and with probability exceeding  $1 - n^{-(1.5C-1)}$ ,

$$\|G_l(\mathbf{A})\|^2 \leq p \|\mathbf{A}\|^2 + C \left( \sqrt{p \|\mathbf{A}\|_{2,\infty}^2 \|\mathbf{A}\|^2 \log n} + \|\mathbf{A}\|_{2,\infty}^2 \log n \right).$$

This completes the proof. □

**Lemma 42** *Let  $\{\delta_{l,j}\}_{1 \leq l \leq j \leq n}$  be i.i.d. Bernoulli random variables with mean  $p$  and  $\delta_{l,j} = \delta_{j,l}$ . For any  $\mathbf{\Delta} \in \mathbb{R}^{n \times r}$ , define*

$$G_l(\mathbf{\Delta}) := \left[ \delta_{l,1} \mathbf{\Delta}_{1,\cdot}^\top, \delta_{l,2} \mathbf{\Delta}_{2,\cdot}^\top, \dots, \delta_{l,n} \mathbf{\Delta}_{n,\cdot}^\top \right] \in \mathbb{R}^{r \times n}.$$

*Suppose the sample size obeys  $n^2 p \gg \kappa \mu r n \log^2 n$ . Then for any  $k > 0$  and  $\alpha > 0$  large enough, with probability at least  $1 - c_1 e^{-\alpha C n r \log n / 2}$ ,*

$$\sum_{l=1}^n \mathbb{1}_{\{\|G_l(\mathbf{\Delta})\| \geq 4\sqrt{p}\psi + 2\sqrt{kr}\xi\}} \leq \frac{2\alpha n \log n}{k}$$

*holds simultaneously for all  $\mathbf{\Delta} \in \mathbb{R}^{n \times r}$  obeying*

$$\begin{aligned} \|\mathbf{\Delta}\|_{2,\infty} &\leq C_5 \rho^t \mu r \sqrt{\frac{\log n}{np}} \|\mathbf{X}^*\|_{2,\infty} + C_8 \sigma \sqrt{\frac{n \log n}{p}} \|\mathbf{X}^*\|_{2,\infty} := \xi \\ \text{and} \quad \|\mathbf{\Delta}\| &\leq C_9 \rho^t \mu r \frac{1}{\sqrt{np}} \|\mathbf{X}^*\| + C_{10} \sigma \sqrt{\frac{n}{p}} \|\mathbf{X}^*\| := \psi, \end{aligned}$$

*where  $c_1, C_5, C_8, C_9, C_{10} > 0$  are some absolute constants.*

**Proof** For simplicity of presentation, we will prove the claim for the asymmetric case where  $\{\delta_{l,j}\}_{1 \leq l, j \leq n}$  are independent. The results immediately carry over to the symmetric case as claimed in this lemma. To see this, note that we can always divide  $G_l(\mathbf{\Delta})$  into

$$G_l(\mathbf{\Delta}) = G_l^{\text{upper}}(\mathbf{\Delta}) + G_l^{\text{lower}}(\mathbf{\Delta}),$$

where all nonzero components of  $G_l^{\text{upper}}(\mathbf{\Delta})$  come from the upper triangular part (those blocks with  $l \leq j$ ), while all nonzero components of  $G_l^{\text{lower}}(\mathbf{\Delta})$  are from the lower triangular part (those blocks with  $l > j$ ). We can then look at  $\{G_l^{\text{upper}}(\mathbf{\Delta}) \mid 1 \leq l \leq n\}$  and  $\{G_l^{\text{lower}}(\mathbf{\Delta}) \mid 1 \leq l \leq n\}$  separately using the argument we develop for the asymmetric case. From now on, we assume that  $\{\delta_{l,j}\}_{1 \leq l, j \leq n}$  are independent.

Suppose for the moment that  $\mathbf{\Delta}$  is statistically independent of  $\{\delta_{l,j}\}$ . Clearly, for any  $\mathbf{\Delta}, \tilde{\mathbf{\Delta}} \in \mathbb{R}^{n \times r}$ ,

$$\begin{aligned} \left| \|\mathbf{G}_l(\mathbf{\Delta})\| - \|\mathbf{G}_l(\tilde{\mathbf{\Delta}})\| \right| &\leq \|\mathbf{G}_l(\mathbf{\Delta}) - \mathbf{G}_l(\tilde{\mathbf{\Delta}})\| \leq \|\mathbf{G}_l(\mathbf{\Delta}) - \mathbf{G}_l(\tilde{\mathbf{\Delta}})\|_F \\ &\leq \sqrt{\sum_{j=1}^n \|\mathbf{\Delta}_{j,\cdot} - \tilde{\mathbf{\Delta}}_{j,\cdot}\|_2^2} \\ &:= d(\mathbf{\Delta}, \tilde{\mathbf{\Delta}}), \end{aligned}$$

which implies that  $\|\mathbf{G}_l(\mathbf{\Delta})\|$  is 1-Lipschitz with respect to the metric  $d(\cdot, \cdot)$ . Moreover,

$$\max_{1 \leq j \leq n} \|\delta_{l,j} \mathbf{\Delta}_{j,\cdot}\|_2 \leq \|\mathbf{\Delta}\|_{2,\infty} \leq \xi$$

according to our assumption. Hence, Talagrand’s inequality [24, Proposition 1] reveals the existence of some absolute constants  $C, c > 0$  such that for all  $\lambda > 0$

$$\mathbb{P}\{\|\mathbf{G}_l(\mathbf{\Delta})\| - \text{Median}[\|\mathbf{G}_l(\mathbf{\Delta})\|] \geq \lambda \xi\} \leq C \exp(-c\lambda^2). \tag{259}$$

We then proceed to control  $\text{Median}[\|\mathbf{G}_l(\mathbf{\Delta})\|]$ . A direct application of Lemma 41 yields

$$\text{Median}[\|\mathbf{G}_l(\mathbf{\Delta})\|] \leq \sqrt{2p\psi^2 + \sqrt{p \log(4r)}\xi\psi + \frac{2\xi^2}{3} \log(4r)} \leq 2\sqrt{p}\psi,$$

where the last relation holds since  $p\psi^2 \gg \xi^2 \log r$ , which follows by combining the definitions of  $\psi$  and  $\xi$ , the sample size condition  $np \gg \kappa\mu r \log^2 n$ , and the incoherence condition (114). Thus, substitution into (259) and taking  $\lambda = \sqrt{kr}$  give

$$\mathbb{P}\left\{\|\mathbf{G}_l(\mathbf{\Delta})\| \geq 2\sqrt{p}\psi + \sqrt{kr}\xi\right\} \leq C \exp(-ckr) \tag{260}$$

for any  $k \geq 0$ . Furthermore, invoking [4, Corollary A.1.14] and using bound (260), one has

$$\mathbb{P}\left(\sum_{l=1}^n \mathbb{1}_{\{\|\mathbf{G}_l(\mathbf{\Delta})\| \geq 2\sqrt{p}\psi + \sqrt{kr}\xi\}} \geq tnC \exp(-ckr)\right) \leq 2 \exp\left(-\frac{t \log t}{2} nC \exp(-ckr)\right)$$

for any  $t \geq 6$ . Choose  $t = \alpha \log n / [kC \exp(-ckr)] \geq 6$  to obtain

$$\mathbb{P}\left(\sum_{l=1}^n \mathbb{1}_{\{\|\mathbf{G}_l(\mathbf{\Delta})\| \geq 2\sqrt{p}\psi + \sqrt{kr}\xi\}} \geq \frac{\alpha n \log n}{k}\right) \leq 2 \exp\left(-\frac{\alpha C}{2} nr \log n\right). \tag{261}$$

So far, we have demonstrated that for any fixed  $\Delta$  obeying our assumptions,  $\sum_{l=1}^n \mathbb{1}_{\{\|G_l(\Delta)\| \geq 2\sqrt{p}\psi + \sqrt{kr}\xi\}}$  is well controlled with exponentially high probability. In order to extend the results to all feasible  $\Delta$ , we resort to the standard  $\epsilon$ -net argument. Clearly, due to the homogeneity property of  $\|G_l(\Delta)\|$ , it suffices to restrict attention to the following set:

$$\mathcal{S} = \{\Delta \mid \min\{\xi, \psi\} \leq \|\Delta\| \leq \psi\}, \tag{262}$$

where  $\psi/\xi \lesssim \|X^*\|/\|X^*\|_{2,\infty} \lesssim \sqrt{n}$ . We then proceed with the following steps.

1. Introduce the auxiliary function

$$\chi_l(\Delta) = \begin{cases} 1, & \text{if } \|G_l(\Delta)\| \geq 4\sqrt{p}\psi + 2\sqrt{kr}\xi, \\ \frac{\|G_l(\Delta)\| - 2\sqrt{p}\psi - \sqrt{kr}\xi}{2\sqrt{p}\psi + \sqrt{kr}\xi}, & \text{if } \|G_l(\Delta)\| \in [2\sqrt{p}\psi + \sqrt{kr}\xi, 4\sqrt{p}\psi + 2\sqrt{kr}\xi], \\ 0, & \text{else.} \end{cases}$$

Clearly, this function is sandwiched between two indicator functions

$$\mathbb{1}_{\{\|G_l(\Delta)\| \geq 4\sqrt{p}\psi + 2\sqrt{kr}\xi\}} \leq \chi_l(\Delta) \leq \mathbb{1}_{\{\|G_l(\Delta)\| \geq 2\sqrt{p}\psi + \sqrt{kr}\xi\}}.$$

Note that  $\chi_l$  is more convenient to work with due to continuity.

2. Consider an  $\epsilon$ -net  $\mathcal{N}_\epsilon$  [111, Section 2.3.1] of the set  $\mathcal{S}$  as defined in (262). For any  $\epsilon = 1/n^{O(1)}$ , one can find such a net with cardinality  $\log |\mathcal{N}_\epsilon| \lesssim nr \log n$ . Apply the union bound and (261) to yield

$$\begin{aligned} & \mathbb{P}\left(\sum_{l=1}^n \chi_l(\Delta) \geq \frac{\alpha n \log n}{k}, \forall \Delta \in \mathcal{N}_\epsilon\right) \\ & \leq \mathbb{P}\left(\sum_{l=1}^n \mathbb{1}_{\{\|G_l(\Delta)\| \geq 2\sqrt{p}\psi + \sqrt{kr}\xi\}} \geq \frac{\alpha n \log n}{k}, \forall \Delta \in \mathcal{N}_\epsilon\right) \\ & \leq 2|\mathcal{N}_\epsilon| \exp\left(-\frac{\alpha C}{2} nr \log n\right) \leq 2 \exp\left(-\frac{\alpha C}{4} nr \log n\right), \end{aligned}$$

as long as  $\alpha$  is chosen to be sufficiently large.

3. One can then use the continuity argument to extend the bound to all  $\Delta$  outside the  $\epsilon$ -net, i.e., with exponentially high probability,

$$\begin{aligned} & \sum_{l=1}^n \chi_l(\Delta) \leq \frac{2\alpha n \log n}{k}, \quad \forall \Delta \in \mathcal{S} \\ \implies & \sum_{l=1}^n \mathbb{1}_{\{\|G_l(\Delta)\| \geq 4\sqrt{p}\psi + 2\sqrt{kr}\xi\}} \leq \sum_{l=1}^n \chi_l(\Delta) \leq \frac{2\alpha n \log n}{k}, \quad \forall \Delta \in \mathcal{S}. \end{aligned}$$

This is fairly standard (see, e.g., [111, Section 2.3.1]) and is thus omitted here.

We have thus concluded the proof. □

**Lemma 43** *Suppose the sample size obeys  $n^2 p \geq C\kappa\mu rn \log n$  for some sufficiently large constant  $C > 0$ . Then with probability at least  $1 - O(n^{-10})$ ,*

$$\left\| \frac{1}{p} \mathcal{P}_\Omega \left( \mathbf{X} \mathbf{X}^\top - \mathbf{X}^* \mathbf{X}^{*\top} \right) \right\| \leq 2n\epsilon^2 \|\mathbf{X}^*\|_{2,\infty}^2 + 4\epsilon\sqrt{n} \log n \|\mathbf{X}^*\|_{2,\infty} \|\mathbf{X}^*\|$$

holds simultaneously for all  $\mathbf{X} \in \mathbb{R}^{n \times r}$  satisfying

$$\|\mathbf{X} - \mathbf{X}^*\|_{2,\infty} \leq \epsilon \|\mathbf{X}^*\|_{2,\infty}, \tag{263}$$

where  $\epsilon > 0$  is any fixed constant.

**Proof** To simplify the notations hereafter, we denote  $\Delta := \mathbf{X} - \mathbf{X}^*$ . With this notation in place, one can decompose

$$\mathbf{X} \mathbf{X}^\top - \mathbf{X}^* \mathbf{X}^{*\top} = \Delta \mathbf{X}^{*\top} + \mathbf{X}^* \Delta^\top + \Delta \Delta^\top,$$

which together with the triangle inequality implies that

$$\begin{aligned} \left\| \frac{1}{p} \mathcal{P}_\Omega \left( \mathbf{X} \mathbf{X}^\top - \mathbf{X}^* \mathbf{X}^{*\top} \right) \right\| &\leq \left\| \frac{1}{p} \mathcal{P}_\Omega \left( \Delta \mathbf{X}^{*\top} \right) \right\| + \left\| \frac{1}{p} \mathcal{P}_\Omega \left( \mathbf{X}^* \Delta^\top \right) \right\| + \left\| \frac{1}{p} \mathcal{P}_\Omega \left( \Delta \Delta^\top \right) \right\| \\ &= \underbrace{\left\| \frac{1}{p} \mathcal{P}_\Omega \left( \Delta \mathbf{X}^{*\top} \right) \right\|}_{:=\alpha_1} + 2 \underbrace{\left\| \frac{1}{p} \mathcal{P}_\Omega \left( \Delta \mathbf{X}^{*\top} \right) \right\|}_{:=\alpha_2}. \end{aligned} \tag{264}$$

In the sequel, we bound  $\alpha_1$  and  $\alpha_2$  separately.

1. Recall from [84, Theorem 2.5] the elementary inequality that

$$\|\mathbf{C}\| \leq \|\mathbf{C}\|, \tag{265}$$

where  $|\mathbf{C}| := [c_{i,j}]_{1 \leq i,j \leq n}$  for any matrix  $\mathbf{C} = [c_{i,j}]_{1 \leq i,j \leq n}$ . In addition, for any matrix  $\mathbf{D} := [d_{i,j}]_{1 \leq i,j \leq n}$  such that  $|d_{i,j}| \geq |c_{i,j}|$  for all  $i$  and  $j$ , one has  $\|\mathbf{C}\| \leq \|\mathbf{D}\|$ . Therefore,

$$\alpha_1 \leq \left\| \frac{1}{p} \mathcal{P}_\Omega \left( \mathbf{C} \right) \right\| \leq \|\Delta\|_{2,\infty}^2 \left\| \frac{1}{p} \mathcal{P}_\Omega \left( \mathbf{1} \mathbf{1}^\top \right) \right\|.$$

Lemma 39 then tells us that with probability at least  $1 - O(n^{-10})$ ,

$$\left\| \frac{1}{p} \mathcal{P}_\Omega \left( \mathbf{1} \mathbf{1}^\top \right) - \mathbf{1} \mathbf{1}^\top \right\| \leq C \sqrt{\frac{n}{p}} \tag{266}$$

for some universal constant  $C > 0$ , as long as  $p \gg \log n/n$ . This together with the triangle inequality yields

$$\left\| \frac{1}{p} \mathcal{P}_\Omega \left( \mathbf{1} \mathbf{1}^\top \right) \right\| \leq \left\| \frac{1}{p} \mathcal{P}_\Omega \left( \mathbf{1} \mathbf{1}^\top \right) - \mathbf{1} \mathbf{1}^\top \right\| + \|\mathbf{1} \mathbf{1}^\top\| \leq C \sqrt{\frac{n}{p}} + n \leq 2n, \tag{267}$$



provided that  $p \gg 1/n$ . Putting together the previous bounds, we arrive at

$$\alpha_1 \leq 2n \|\Delta\|_{2,\infty}^2. \tag{268}$$

2. Regarding the second term  $\alpha_2$ , apply the elementary inequality (265) once again to get

$$\|\mathcal{P}_\Omega(\Delta X^{*\top})\| \leq \|\mathcal{P}_\Omega(|\Delta X^{*\top}|)\|,$$

which motivates us to look at  $\|\mathcal{P}_\Omega(|\Delta X^{*\top}|)\|$  instead. A key step of this part is to take advantage of the  $\ell_{2,\infty}$  norm constraint of  $\mathcal{P}_\Omega(|\Delta X^{*\top}|)$ . Specifically, we claim for the moment that with probability exceeding  $1 - O(n^{-10})$ ,

$$\|\mathcal{P}_\Omega(|\Delta X^{*\top}|)\|_{2,\infty}^2 \leq 2p\sigma_{\max} \|\Delta\|_{2,\infty}^2 := \theta \tag{269}$$

holds under our sample size condition. In addition, we also have the following trivial  $\ell_\infty$  norm bound

$$\|\mathcal{P}_\Omega(|\Delta X^{*\top}|)\|_\infty \leq \|\Delta\|_{2,\infty} \|X^*\|_{2,\infty} := \gamma. \tag{270}$$

In what follows, for simplicity of presentation, we will denote

$$A := \mathcal{P}_\Omega(|\Delta X^{*\top}|). \tag{271}$$

(a) To facilitate the analysis of  $\|A\|$ , we first introduce  $k_0 + 1 = \frac{1}{2} \log(\kappa\mu r)$  auxiliary matrices<sup>9</sup>  $B_s \in \mathbb{R}^{n \times n}$  that satisfy

$$\|A\| \leq \|B_{k_0}\| + \sum_{s=0}^{k_0-1} \|B_s\|. \tag{272}$$

To be precise, each  $B_s$  is defined such that

$$[B_s]_{j,k} = \begin{cases} \frac{1}{2^s} \gamma, & \text{if } A_{j,k} \in (\frac{1}{2^{s+1}} \gamma, \frac{1}{2^s} \gamma], \\ 0, & \text{else,} \end{cases} \quad \text{for } 0 \leq s \leq k_0 - 1 \quad \text{and}$$

$$[B_{k_0}]_{j,k} = \begin{cases} \frac{1}{2^{k_0}} \gamma, & \text{if } A_{j,k} \leq \frac{1}{2^{k_0}} \gamma, \\ 0, & \text{else,} \end{cases}$$

which clearly satisfy (272); in words,  $B_s$  is constructed by rounding up those entries of  $A$  within a prescribed magnitude interval. Thus, it suffices to bound

<sup>9</sup> For simplicity, we assume  $\frac{1}{2} \log(\kappa\mu r)$  is an integer. The argument here can be easily adapted to the case when  $\frac{1}{2} \log(\kappa\mu r)$  is not an integer.

$\|B_s\|$  for every  $s$ . To this end, we start with  $s = k_0$  and use the definition of  $B_{k_0}$  to get

$$\begin{aligned} \|B_{k_0}\| &\stackrel{(i)}{\leq} \|B_{k_0}\|_\infty \sqrt{(2np)^2} \stackrel{(ii)}{\leq} 4np \frac{1}{\sqrt{\kappa\mu r}} \|\Delta\|_{2,\infty} \|X^*\|_{2,\infty} \\ &\stackrel{(iii)}{\leq} 4\sqrt{np} \|\Delta\|_{2,\infty} \|X^*\|, \end{aligned}$$

where (i) arises from Lemma 44, with  $2np$  being a crude upper bound on the number of nonzero entries in each row and each column. This can be derived by applying the standard Chernoff bound on  $\Omega$ . The second inequality (ii) relies on the definitions of  $\gamma$  and  $k_0$ . The last one (iii) follows from the incoherence condition (114). Besides, for any  $0 \leq s \leq k_0 - 1$ , by construction one has

$$\|B_s\|_{2,\infty}^2 \leq 4\theta = 8p\sigma_{\max} \|\Delta\|_{2,\infty}^2 \quad \text{and} \quad \|B_s\|_\infty = \frac{1}{2^s} \gamma,$$

where  $\theta$  is as defined in (269). Here, we have used the fact that the magnitude of each entry of  $B_s$  is at most two times that of  $A$ . An immediate implication is that there are at most

$$\frac{\|B_s\|_{2,\infty}^2}{\|B_s\|_\infty^2} \leq \frac{8p\sigma_{\max} \|\Delta\|_{2,\infty}^2}{\left(\frac{1}{2^s} \gamma\right)^2} := k_r$$

nonzero entries in each row of  $B_s$  and at most

$$k_c = 2np$$

nonzero entries in each column of  $B_s$ , where  $k_c$  is derived from the standard Chernoff bound on  $\Omega$ . Utilizing Lemma 44 once more, we discover that

$$\begin{aligned} \|B_s\| &\leq \|B_s\|_\infty \sqrt{k_r k_c} = \frac{1}{2^s} \gamma \sqrt{k_r k_c} = \sqrt{16np^2 \sigma_{\max} \|\Delta\|_{2,\infty}^2} \\ &= 4\sqrt{np} \|\Delta\|_{2,\infty} \|X^*\| \end{aligned}$$

for each  $0 \leq s \leq k_0 - 1$ . Combining all, we arrive at

$$\begin{aligned} \|A\| &\leq \sum_{s=0}^{k_0-1} \|B_s\| + \|B_{k_0}\| \leq (k_0 + 1) 4\sqrt{np} \|\Delta\|_{2,\infty} \|X^*\| \\ &\leq 2\sqrt{np} \log(\kappa\mu r) \|\Delta\|_{2,\infty} \|X^*\| \\ &\leq 2\sqrt{np} \log n \|\Delta\|_{2,\infty} \|X^*\|, \end{aligned}$$

where the last relation holds under the condition  $n \geq \kappa\mu r$ . This further gives

$$\alpha_2 \leq \frac{1}{p} \|A\| \leq 2\sqrt{n} \log n \|\Delta\|_{2,\infty} \|X^*\|. \tag{273}$$

(b) In order to finish the proof of this part, we need to justify claim (269). Observe that

$$\begin{aligned} \left\| \left[ \mathcal{P}_\Omega \left( \left| \Delta \mathbf{X}^{\star\top} \right| \right) \right]_{l,\cdot} \right\|_2^2 &= \sum_{j=1}^n \left( \Delta_{l,\cdot} \mathbf{X}_{j,\cdot}^{\star\top} \delta_{l,j} \right)^2 \\ &= \Delta_{l,\cdot} \left( \sum_{j=1}^n \delta_{l,j} \mathbf{X}_{j,\cdot}^{\star\top} \mathbf{X}_{j,\cdot}^{\star} \right) \Delta_{l,\cdot}^\top \\ &\leq \|\Delta\|_{2,\infty}^2 \left\| \sum_{j=1}^n \delta_{l,j} \mathbf{X}_{j,\cdot}^{\star\top} \mathbf{X}_{j,\cdot}^{\star} \right\| \end{aligned} \tag{274}$$

for every  $1 \leq l \leq n$ , where  $\delta_{l,j}$  indicates whether the entry with the index  $(l, j)$  is observed or not. Invoke Lemma 41 to yield

$$\begin{aligned} \left\| \sum_{j=1}^n \delta_{l,j} \mathbf{X}_{j,\cdot}^{\star\top} \mathbf{X}_{j,\cdot}^{\star} \right\| &= \left\| \left[ \delta_{l,1} \mathbf{X}_{1,\cdot}^{\star\top}, \delta_{l,2} \mathbf{X}_{2,\cdot}^{\star\top}, \dots, \delta_{l,n} \mathbf{X}_{n,\cdot}^{\star\top} \right] \right\|^2 \\ &\leq p\sigma_{\max} + C \left( \sqrt{p \|\mathbf{X}^{\star}\|_{2,\infty}^2 \|\mathbf{X}^{\star}\|^2 \log n} + \|\mathbf{X}^{\star}\|_{2,\infty}^2 \log n \right) \\ &\leq \left( p + C \sqrt{\frac{p\kappa\mu r \log n}{n}} + C \frac{\kappa\mu r \log n}{n} \right) \sigma_{\max} \\ &\leq 2p\sigma_{\max}, \end{aligned} \tag{275}$$

with high probability, as soon as  $np \gg \kappa\mu r \log n$ . Combining (274) and (275) yields

$$\left\| \left[ \mathcal{P}_\Omega \left( \left| \Delta \mathbf{X}^{\star\top} \right| \right) \right]_{l,\cdot} \right\|_2^2 \leq 2p\sigma_{\max} \|\Delta\|_{2,\infty}^2, \quad 1 \leq l \leq n,$$

as claimed in (269).

3. Taken together, the preceding bounds (264), (268), and (273) yield

$$\left\| \frac{1}{p} \mathcal{P}_\Omega \left( \mathbf{X} \mathbf{X}^\top - \mathbf{X}^{\star} \mathbf{X}^{\star\top} \right) \right\| \leq \alpha_1 + 2\alpha_2 \leq 2n \|\Delta\|_{2,\infty}^2 + 4\sqrt{n} \log n \|\Delta\|_{2,\infty} \|\mathbf{X}^{\star}\|.$$

The proof is completed by substituting the assumption  $\|\Delta\|_{2,\infty} \leq \epsilon \|\mathbf{X}^{\star}\|_{2,\infty}$ . □

In the end of this subsection, we record a useful lemma to bound the spectral norm of a sparse Bernoulli matrix.

**Lemma 44** *Let  $\mathbf{A} \in \{0, 1\}^{n_1 \times n_2}$  be a binary matrix, and suppose that there are at most  $k_r$  and  $k_c$  nonzero entries in each row and column of  $\mathbf{A}$ , respectively. Then one has  $\|\mathbf{A}\| \leq \sqrt{k_c k_r}$ .*

**Proof** This immediately follows from the elementary inequality  $\|\mathbf{A}\|^2 \leq \|\mathbf{A}\|_{1 \rightarrow 1} \|\mathbf{A}\|_{\infty \rightarrow \infty}$  (see [56, equation (1.11)]), where  $\|\mathbf{A}\|_{1 \rightarrow 1}$  and  $\|\mathbf{A}\|_{\infty \rightarrow \infty}$  are the induced 1-norm (or maximum absolute column sum norm) and the induced  $\infty$ -norm (or maximum absolute row sum norm), respectively. □

### D.2.3 Matrix Perturbation Bounds

**Lemma 45** *Let  $M \in \mathbb{R}^{n \times n}$  be a symmetric matrix with the top- $r$  eigendecomposition  $U \Sigma U^\top$ . Assume  $\|M - M^*\| \leq \sigma_{\min}/2$ , and denote*

$$\widehat{Q} := \operatorname{argmin}_{R \in \mathcal{O}^{r \times r}} \|UR - U^*\|_F.$$

*Then there is some numerical constant  $c_3 > 0$  such that*

$$\|U \widehat{Q} - U^*\| \leq \frac{c_3}{\sigma_{\min}} \|M - M^*\|.$$

**Proof** Define  $Q = U^\top U^*$ . The triangle inequality gives

$$\|U \widehat{Q} - U^*\| \leq \|U(\widehat{Q} - Q)\| + \|UQ - U^*\| \leq \|\widehat{Q} - Q\| + \|UU^\top U^* - U^*\|. \tag{276}$$

[1, Lemma 3] asserts that

$$\|\widehat{Q} - Q\| \leq 4(\|M - M^*\| / \sigma_{\min})^2$$

as long as  $\|M - M^*\| \leq \sigma_{\min}/2$ . For the remaining term in (276), one can use  $U^{*\top} U^* = I_r$  to obtain

$$\|UU^\top U^* - U^*\| = \|UU^\top U^* - U^*U^{*\top} U^*\| \leq \|UU^\top - U^*U^{*\top}\|,$$

which together with the Davis–Kahan  $\sin \Theta$  theorem [39] reveals that

$$\|UU^\top U^* - U^*\| \leq \frac{c_2}{\sigma_{\min}} \|M - M^*\|$$

for some constant  $c_2 > 0$ . Combine the estimates on  $\|\widehat{Q} - Q\|$ ,  $\|UU^\top U^* - U^*\|$  and (276) to reach

$$\|U \widehat{Q} - U^*\| \leq \left( \frac{4}{\sigma_{\min}} \|M - M^*\| \right)^2 + \frac{c_2}{\sigma_{\min}} \|M - M^*\| \leq \frac{c_3}{\sigma_{\min}} \|M - M^*\|$$

for some numerical constant  $c_3 > 0$ , where we have utilized the fact that  $\|M - M^*\| / \sigma_{\min} \leq 1/2$ . □

**Lemma 46** *Let  $M, \widetilde{M} \in \mathbb{R}^{n \times n}$  be two symmetric matrices with top- $r$  eigendecompositions  $U \Sigma U^\top$  and  $\widetilde{U} \widetilde{\Sigma} \widetilde{U}^\top$ , respectively. Assume  $\|M - M^*\| \leq \sigma_{\min}/4$  and  $\|\widetilde{M} - M^*\| \leq \sigma_{\min}/4$ , and suppose  $\sigma_{\max}/\sigma_{\min}$  is bounded by some constant  $c_1 > 0$ ,*

with  $\sigma_{\max}$  and  $\sigma_{\min}$  the largest and the smallest singular values of  $\mathbf{M}^*$ , respectively. If we denote

$$\mathbf{Q} := \operatorname{argmin}_{\mathbf{R} \in \mathcal{O}^{r \times r}} \|\mathbf{U}\mathbf{R} - \tilde{\mathbf{U}}\|_{\text{F}},$$

then there exists some numerical constant  $c_3 > 0$  such that

$$\begin{aligned} \|\Sigma^{1/2} \mathbf{Q} - \mathbf{Q} \tilde{\Sigma}^{1/2}\| &\leq \frac{c_3}{\sqrt{\sigma_{\min}}} \|\tilde{\mathbf{M}} - \mathbf{M}\| \quad \text{and} \\ \|\Sigma^{1/2} \mathbf{Q} - \mathbf{Q} \tilde{\Sigma}^{1/2}\|_{\text{F}} &\leq \frac{c_3}{\sqrt{\sigma_{\min}}} \|(\tilde{\mathbf{M}} - \mathbf{M})\mathbf{U}\|_{\text{F}}. \end{aligned}$$

**Proof** Here, we focus on the Frobenius norm; the bound on the operator norm follows from the same argument, and hence we omit the proof. Since  $\|\cdot\|_{\text{F}}$  is unitarily invariant, we have

$$\|\Sigma^{1/2} \mathbf{Q} - \mathbf{Q} \tilde{\Sigma}^{1/2}\|_{\text{F}} = \|\mathbf{Q}^{\text{T}} \Sigma^{1/2} \mathbf{Q} - \tilde{\Sigma}^{1/2}\|_{\text{F}},$$

where  $\mathbf{Q}^{\text{T}} \Sigma^{1/2} \mathbf{Q}$  and  $\tilde{\Sigma}^{1/2}$  are the matrix square roots of  $\mathbf{Q}^{\text{T}} \Sigma \mathbf{Q}$  and  $\tilde{\Sigma}$ , respectively. In view of the matrix square root perturbation bound [97, Lemma 2.1],

$$\begin{aligned} \|\Sigma^{1/2} \mathbf{Q} - \mathbf{Q} \tilde{\Sigma}^{1/2}\|_{\text{F}} &\leq \frac{1}{\sigma_{\min}[(\Sigma)^{1/2}] + \sigma_{\min}[(\tilde{\Sigma})^{1/2}]} \|\mathbf{Q}^{\text{T}} \Sigma \mathbf{Q} - \tilde{\Sigma}\|_{\text{F}} \\ &\leq \frac{1}{\sqrt{\sigma_{\min}}} \|\mathbf{Q}^{\text{T}} \Sigma \mathbf{Q} - \tilde{\Sigma}\|_{\text{F}}, \end{aligned} \tag{277}$$

where the last inequality follows from the lower estimates

$$\sigma_{\min}(\Sigma) \geq \sigma_{\min}(\Sigma^*) - \|\mathbf{M} - \mathbf{M}^*\| \geq \sigma_{\min}/4$$

and, similarly,  $\sigma_{\min}(\tilde{\Sigma}) \geq \sigma_{\min}/4$ . Recognizing that  $\Sigma = \mathbf{U}^{\text{T}} \mathbf{M} \mathbf{U}$  and  $\tilde{\Sigma} = \tilde{\mathbf{U}}^{\text{T}} \tilde{\mathbf{M}} \tilde{\mathbf{U}}$ , one gets

$$\begin{aligned} \|\mathbf{Q}^{\text{T}} \Sigma \mathbf{Q} - \tilde{\Sigma}\|_{\text{F}} &= \|(\mathbf{U} \mathbf{Q})^{\text{T}} \mathbf{M} (\mathbf{U} \mathbf{Q}) - \tilde{\mathbf{U}}^{\text{T}} \tilde{\mathbf{M}} \tilde{\mathbf{U}}\|_{\text{F}} \\ &\leq \|(\mathbf{U} \mathbf{Q})^{\text{T}} \mathbf{M} (\mathbf{U} \mathbf{Q}) - (\mathbf{U} \mathbf{Q})^{\text{T}} \tilde{\mathbf{M}} (\mathbf{U} \mathbf{Q})\|_{\text{F}} + \|(\mathbf{U} \mathbf{Q})^{\text{T}} \tilde{\mathbf{M}} (\mathbf{U} \mathbf{Q}) - \tilde{\mathbf{U}}^{\text{T}} \tilde{\mathbf{M}} (\mathbf{U} \mathbf{Q})\|_{\text{F}} \\ &\quad + \|\tilde{\mathbf{U}}^{\text{T}} \tilde{\mathbf{M}} (\mathbf{U} \mathbf{Q}) - \tilde{\mathbf{U}}^{\text{T}} \tilde{\mathbf{M}} \tilde{\mathbf{U}}\|_{\text{F}} \\ &\leq \|(\tilde{\mathbf{M}} - \mathbf{M})\mathbf{U}\|_{\text{F}} + 2\|\mathbf{U} \mathbf{Q} - \tilde{\mathbf{U}}\|_{\text{F}} \|\tilde{\mathbf{M}}\| \leq \|(\tilde{\mathbf{M}} - \mathbf{M})\mathbf{U}\|_{\text{F}} + 4\sigma_{\max} \|\mathbf{U} \mathbf{Q} - \tilde{\mathbf{U}}\|_{\text{F}}, \end{aligned} \tag{278}$$

where the last relation holds due to the upper estimate

$$\|\tilde{\mathbf{M}}\| \leq \|\mathbf{M}^*\| + \|\tilde{\mathbf{M}} - \mathbf{M}^*\| \leq \sigma_{\max} + \sigma_{\min}/4 \leq 2\sigma_{\max}.$$

Invoke the Davis–Kahan  $\sin\Theta$  theorem [39] to obtain

$$\|UQ - \tilde{U}\|_F \leq \frac{c_2}{\sigma_r(M) - \sigma_{r+1}(\tilde{M})} \|(\tilde{M} - M)U\|_F \leq \frac{2c_2}{\sigma_{\min}} \|(\tilde{M} - M)U\|_F, \tag{279}$$

for some constant  $c_2 > 0$ , where the last inequality follows from the bounds

$$\begin{aligned} \sigma_r(M) &\geq \sigma_r(M^*) - \|M - M^*\| \geq 3\sigma_{\min}/4, \\ \sigma_{r+1}(\tilde{M}) &\leq \sigma_{r+1}(M^*) + \|\tilde{M} - M^*\| \leq \sigma_{\min}/4. \end{aligned}$$

Combine (277), (278), (279), and the fact  $\sigma_{\max}/\sigma_{\min} \leq c_1$  to reach

$$\|\Sigma^{1/2}Q - Q\tilde{\Sigma}^{1/2}\|_F \leq \frac{c_3}{\sqrt{\sigma_{\min}}} \|(\tilde{M} - M)U\|_F$$

for some constant  $c_3 > 0$ . □

**Lemma 47** *Let  $M \in \mathbb{R}^{n \times n}$  be a symmetric matrix with the top- $r$  eigendecomposition  $U\Sigma U^T$ . Denote  $X = U\Sigma^{1/2}$  and  $X^* = U^*(\Sigma^*)^{1/2}$ , and define*

$$\hat{Q} := \operatorname{argmin}_{R \in \mathcal{O}^{r \times r}} \|UR - U^*\|_F \quad \text{and} \quad \hat{H} := \operatorname{argmin}_{R \in \mathcal{O}^{r \times r}} \|XR - X^*\|_F.$$

Assume  $\|M - M^*\| \leq \sigma_{\min}/2$ , and suppose  $\sigma_{\max}/\sigma_{\min}$  is bounded by some constant  $c_1 > 0$ . Then there exists a numerical constant  $c_3 > 0$  such that

$$\|\hat{Q} - \hat{H}\| \leq \frac{c_3}{\sigma_{\min}} \|M - M^*\|.$$

**Proof** We first collect several useful facts about the spectrum of  $\Sigma$ . Weyl’s inequality tells us that  $\|\Sigma - \Sigma^*\| \leq \|M - M^*\| \leq \sigma_{\min}/2$ , which further implies that

$$\sigma_r(\Sigma) \geq \sigma_r(\Sigma^*) - \|\Sigma - \Sigma^*\| \geq \sigma_{\min}/2 \quad \text{and} \quad \|\Sigma\| \leq \|\Sigma^*\| + \|\Sigma - \Sigma^*\| \leq 2\sigma_{\max}.$$

Denote

$$Q = U^T U^* \quad \text{and} \quad H = X^T X^*.$$

Simple algebra yields

$$\begin{aligned} H &= \Sigma^{1/2}Q(\Sigma^*)^{1/2} = \underbrace{\Sigma^{1/2}(Q - \hat{Q})(\Sigma^*)^{1/2} + (\Sigma^{1/2}\hat{Q} - \hat{Q}\Sigma^{1/2})(\Sigma^*)^{1/2}}_{:=E} \\ &\quad + \underbrace{\hat{Q}(\Sigma\Sigma^*)^{1/2}}_{:=A}. \end{aligned}$$

It can be easily seen that  $\sigma_{r-1}(\mathbf{A}) \geq \sigma_r(\mathbf{A}) \geq \sigma_{\min}/2$ , and

$$\begin{aligned} \|\mathbf{E}\| &\leq \|\boldsymbol{\Sigma}^{1/2}\| \cdot \|\mathbf{Q} - \widehat{\mathbf{Q}}\| \cdot \|(\boldsymbol{\Sigma}^*)^{1/2}\| + \|\boldsymbol{\Sigma}^{1/2}\widehat{\mathbf{Q}} - \widehat{\mathbf{Q}}\boldsymbol{\Sigma}^{1/2}\| \cdot \|(\boldsymbol{\Sigma}^*)^{1/2}\| \\ &\leq 2\sigma_{\max} \underbrace{\|\mathbf{Q} - \widehat{\mathbf{Q}}\|}_{:=\alpha} + \sqrt{\sigma_{\max}} \underbrace{\|\boldsymbol{\Sigma}^{1/2}\widehat{\mathbf{Q}} - \widehat{\mathbf{Q}}\boldsymbol{\Sigma}^{1/2}\|}_{:=\beta}, \end{aligned}$$

which can be controlled as follows.

- Regarding  $\alpha$ , use [1, Lemma 3] to reach

$$\alpha = \|\mathbf{Q} - \widehat{\mathbf{Q}}\| \leq 4 \|\mathbf{M} - \mathbf{M}^*\|^2 / \sigma_{\min}^2.$$

- For  $\beta$ , one has

$$\beta \stackrel{(i)}{=} \|\widehat{\mathbf{Q}}^\top \boldsymbol{\Sigma}^{1/2} \widehat{\mathbf{Q}} - \boldsymbol{\Sigma}^{1/2}\| \stackrel{(ii)}{\leq} \frac{1}{2\sigma_r(\boldsymbol{\Sigma}^{1/2})} \|\widehat{\mathbf{Q}}^\top \boldsymbol{\Sigma} \widehat{\mathbf{Q}} - \boldsymbol{\Sigma}\| \stackrel{(iii)}{=} \frac{1}{2\sigma_r(\boldsymbol{\Sigma}^{1/2})} \|\boldsymbol{\Sigma} \widehat{\mathbf{Q}} - \widehat{\mathbf{Q}} \boldsymbol{\Sigma}\|,$$

where (i) and (iii) come from the unitary invariance of  $\|\cdot\|$  and (ii) follows from the matrix square root perturbation bound [97, Lemma 2.1]. We can further take the triangle inequality to obtain

$$\begin{aligned} \|\boldsymbol{\Sigma} \widehat{\mathbf{Q}} - \widehat{\mathbf{Q}} \boldsymbol{\Sigma}\| &= \|\boldsymbol{\Sigma} \mathbf{Q} - \mathbf{Q} \boldsymbol{\Sigma} + \boldsymbol{\Sigma}(\widehat{\mathbf{Q}} - \mathbf{Q}) - (\widehat{\mathbf{Q}} - \mathbf{Q}) \boldsymbol{\Sigma}\| \\ &\leq \|\boldsymbol{\Sigma} \mathbf{Q} - \mathbf{Q} \boldsymbol{\Sigma}\| + 2 \|\boldsymbol{\Sigma}\| \|\mathbf{Q} - \widehat{\mathbf{Q}}\| \\ &= \left\| \mathbf{U}(\mathbf{M} - \mathbf{M}^*) \mathbf{U}^{*\top} + \mathbf{Q}(\boldsymbol{\Sigma}^* - \boldsymbol{\Sigma}) \right\| + 2 \|\boldsymbol{\Sigma}\| \|\mathbf{Q} - \widehat{\mathbf{Q}}\| \\ &\leq \left\| \mathbf{U}(\mathbf{M} - \mathbf{M}^*) \mathbf{U}^{*\top} \right\| + \|\mathbf{Q}(\boldsymbol{\Sigma}^* - \boldsymbol{\Sigma})\| + 2 \|\boldsymbol{\Sigma}\| \|\mathbf{Q} - \widehat{\mathbf{Q}}\| \\ &\leq 2 \|\mathbf{M} - \mathbf{M}^*\| + 4\sigma_{\max} \alpha, \end{aligned}$$

where the last inequality uses the Weyl’s inequality  $\|\boldsymbol{\Sigma}^* - \boldsymbol{\Sigma}\| \leq \|\mathbf{M} - \mathbf{M}^*\|$  and the fact that  $\|\boldsymbol{\Sigma}\| \leq 2\sigma_{\max}$ .

- Rearrange the previous bounds to arrive at

$$\|\mathbf{E}\| \leq 2\sigma_{\max} \alpha + \sqrt{\sigma_{\max}} \frac{1}{\sqrt{\sigma_{\min}}} (2 \|\mathbf{M} - \mathbf{M}^*\| + 4\sigma_{\max} \alpha) \leq c_2 \|\mathbf{M} - \mathbf{M}^*\|$$

for some numerical constant  $c_2 > 0$ , where we have used the assumption that  $\sigma_{\max}/\sigma_{\min}$  is bounded.

Recognizing that  $\widehat{\mathbf{Q}} = \text{sgn}(\mathbf{A})$  (see definition in (177)), we are ready to invoke Lemma 36 to deduce that

$$\|\widehat{\mathbf{Q}} - \widehat{\mathbf{H}}\| \leq \frac{2}{\sigma_{r-1}(\mathbf{A}) + \sigma_r(\mathbf{A})} \|\mathbf{E}\| \leq \frac{c_3}{\sigma_{\min}} \|\mathbf{M} - \mathbf{M}^*\|$$

for some constant  $c_3 > 0$ . □

### D.3 Technical Lemmas for Blind Deconvolution

#### D.3.1 Wirtinger Calculus

In this section, we formally prove the fundamental theorem of calculus and the mean value form of Taylor’s theorem under the Wirtinger calculus; see (283) and (284), respectively.

Let  $f : \mathbb{C}^n \rightarrow \mathbb{R}$  be a real-valued function. Denote  $\mathbf{z} = \mathbf{x} + i\mathbf{y} \in \mathbb{C}^n$ , and then  $f(\cdot)$  can alternatively be viewed as a function  $\mathbb{R}^{2n} \rightarrow \mathbb{R}$ . There is a one-to-one mapping connecting the Wirtinger derivatives and the conventional derivatives [69]:

$$\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \mathbf{J}^{-1} \begin{bmatrix} \mathbf{z} \\ \bar{\mathbf{z}} \end{bmatrix}, \tag{280a}$$

$$\nabla_{\mathbb{R}} f \left( \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \right) = \mathbf{J}^H \nabla_{\mathbb{C}} f \left( \begin{bmatrix} \mathbf{z} \\ \bar{\mathbf{z}} \end{bmatrix} \right), \tag{280b}$$

$$\nabla_{\mathbb{R}}^2 f \left( \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \right) = \mathbf{J}^H \nabla_{\mathbb{C}}^2 f \left( \begin{bmatrix} \mathbf{z} \\ \bar{\mathbf{z}} \end{bmatrix} \right) \mathbf{J}, \tag{280c}$$

where the subscripts  $\mathbb{R}$  and  $\mathbb{C}$  represent calculus in the real (conventional) sense and in the complex (Wirtinger) sense, respectively, and

$$\mathbf{J} = \begin{bmatrix} \mathbf{I}_n & i\mathbf{I}_n \\ \mathbf{I}_n & -i\mathbf{I}_n \end{bmatrix}.$$

With these relationships in place, we are ready to verify the fundamental theorem of calculus using the Wirtinger derivatives. Recall from [70, Chapter XIII, Theorem 4.2] that

$$\nabla_{\mathbb{R}} f \left( \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{y}_1 \end{bmatrix} \right) - \nabla_{\mathbb{R}} f \left( \begin{bmatrix} \mathbf{x}_2 \\ \mathbf{y}_2 \end{bmatrix} \right) = \left[ \int_0^1 \nabla_{\mathbb{R}}^2 f \left( \begin{bmatrix} \mathbf{x}(\tau) \\ \mathbf{y}(\tau) \end{bmatrix} \right) d\tau \right] \left( \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{y}_1 \end{bmatrix} - \begin{bmatrix} \mathbf{x}_2 \\ \mathbf{y}_2 \end{bmatrix} \right), \tag{281}$$

where

$$\begin{bmatrix} \mathbf{x}(\tau) \\ \mathbf{y}(\tau) \end{bmatrix} := \begin{bmatrix} \mathbf{x}_2 \\ \mathbf{y}_2 \end{bmatrix} + \tau \left( \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{y}_1 \end{bmatrix} - \begin{bmatrix} \mathbf{x}_2 \\ \mathbf{y}_2 \end{bmatrix} \right).$$

Substitute identities (280) into (281) to arrive at

$$\begin{aligned} & \mathbf{J}^H \nabla_{\mathbb{C}} f \left( \begin{bmatrix} \mathbf{z}_1 \\ \bar{\mathbf{z}}_1 \end{bmatrix} \right) - \mathbf{J}^H \nabla_{\mathbb{C}} f \left( \begin{bmatrix} \mathbf{z}_2 \\ \bar{\mathbf{z}}_2 \end{bmatrix} \right) \\ &= \mathbf{J}^H \left[ \int_0^1 \nabla_{\mathbb{C}}^2 f \left( \begin{bmatrix} \mathbf{z}(\tau) \\ \bar{\mathbf{z}}(\tau) \end{bmatrix} \right) d\tau \right] \mathbf{J} \mathbf{J}^{-1} \left( \begin{bmatrix} \mathbf{z}_1 \\ \bar{\mathbf{z}}_1 \end{bmatrix} - \begin{bmatrix} \mathbf{z}_2 \\ \bar{\mathbf{z}}_2 \end{bmatrix} \right) \\ &= \mathbf{J}^H \left[ \int_0^1 \nabla_{\mathbb{C}}^2 f \left( \begin{bmatrix} \mathbf{z}(\tau) \\ \bar{\mathbf{z}}(\tau) \end{bmatrix} \right) d\tau \right] \left( \begin{bmatrix} \mathbf{z}_1 \\ \bar{\mathbf{z}}_1 \end{bmatrix} - \begin{bmatrix} \mathbf{z}_2 \\ \bar{\mathbf{z}}_2 \end{bmatrix} \right), \end{aligned} \tag{282}$$



where  $z_1 = x_1 + i y_1, z_2 = x_2 + i y_2$  and

$$\begin{bmatrix} z(\tau) \\ z(\tau) \end{bmatrix} := \begin{bmatrix} z_2 \\ z_2 \end{bmatrix} + \tau \left( \begin{bmatrix} z_1 \\ z_1 \end{bmatrix} - \begin{bmatrix} z_2 \\ z_2 \end{bmatrix} \right).$$

Simplification of (282) gives

$$\nabla_{\mathbb{C}} f \left( \begin{bmatrix} z_1 \\ z_1 \end{bmatrix} \right) - \nabla_{\mathbb{C}} f \left( \begin{bmatrix} z_2 \\ z_2 \end{bmatrix} \right) = \left[ \int_0^1 \nabla_{\mathbb{C}}^2 f \left( \begin{bmatrix} z(\tau) \\ z(\tau) \end{bmatrix} \right) d\tau \right] \left( \begin{bmatrix} z_1 \\ z_1 \end{bmatrix} - \begin{bmatrix} z_2 \\ z_2 \end{bmatrix} \right). \tag{283}$$

Repeating the above arguments, one can also show that

$$f(z_1) - f(z_2) = \nabla_{\mathbb{C}} f(z_2)^H \begin{bmatrix} z_1 - z_2 \\ z_1 - z_2 \end{bmatrix} + \frac{1}{2} \begin{bmatrix} z_1 - z_2 \\ z_1 - z_2 \end{bmatrix}^H \nabla_{\mathbb{C}}^2 f(\tilde{z}) \begin{bmatrix} z_1 - z_2 \\ z_1 - z_2 \end{bmatrix}, \tag{284}$$

where  $\tilde{z}$  is some point lying on the vector connecting  $z_1$  and  $z_2$ . This is the mean value form of Taylor’s theorem under the Wirtinger calculus.

### D.3.2 Discrete Fourier Transform Matrices

Let  $\mathbf{B} \in \mathbb{C}^{m \times K}$  be the first  $K$  columns of a discrete Fourier transform (DFT) matrix  $\mathbf{F} \in \mathbb{C}^{m \times m}$ , and denote by  $\mathbf{b}_l$  the  $l$ th column of the matrix  $\mathbf{B}^H$ . By definition,

$$\mathbf{b}_l = \frac{1}{\sqrt{m}} \left( 1, \omega^{(l-1)}, \omega^{2(l-1)}, \dots, \omega^{(K-1)(l-1)} \right)^H,$$

where  $\omega := e^{-i \frac{2\pi}{m}}$  with  $i$  representing the imaginary unit. It is seen that for any  $j \neq l$ ,

$$\begin{aligned} \mathbf{b}_l^H \mathbf{b}_j &= \frac{1}{m} \sum_{k=0}^{K-1} \omega^{k(l-1)} \cdot \overline{\omega^{k(j-1)}} \stackrel{(i)}{=} \frac{1}{m} \sum_{k=0}^{K-1} \omega^{k(l-1)} \cdot \omega^{k(1-j)} = \frac{1}{m} \sum_{k=0}^{K-1} \left( \omega^{l-j} \right)^k \\ &\stackrel{(ii)}{=} \frac{1}{m} \frac{1 - \omega^{K(l-j)}}{1 - \omega^{l-j}}. \end{aligned} \tag{285}$$

Here, (i) uses  $\overline{\omega^\alpha} = \omega^{-\alpha}$  for all  $\alpha \in \mathbb{R}$ , while the last identity (ii) follows from the formula for the sum of a finite geometric series when  $\omega^{l-j} \neq 1$ . This leads to the following lemma.

**Lemma 48** For any  $m \geq 3$  and any  $1 \leq l \leq m$ , we have

$$\sum_{j=1}^m \left| \mathbf{b}_l^H \mathbf{b}_j \right| \leq 4 \log m.$$

**Proof** We first make use of identity (285) to obtain

$$\sum_{j=1}^m \left| \mathbf{b}_l^H \mathbf{b}_j \right| = \|\mathbf{b}_l\|_2^2 + \frac{1}{m} \sum_{j:j \neq l}^m \left| \frac{1 - \omega^{K(l-j)}}{1 - \omega^{l-j}} \right| = \frac{K}{m} + \frac{1}{m} \sum_{j:j \neq l}^m \left| \frac{\sin \left[ K(l-j) \frac{\pi}{m} \right]}{\sin \left[ (l-j) \frac{\pi}{m} \right]} \right|,$$

where the last identity follows since  $\|\mathbf{b}_l\|_2^2 = K/m$  and, for all  $\alpha \in \mathbb{R}$ ,

$$\left| 1 - \omega^\alpha \right| = \left| 1 - e^{-i \frac{2\pi}{m} \alpha} \right| = \left| e^{-i \frac{\pi}{m} \alpha} \left( e^{i \frac{\pi}{m} \alpha} - e^{-i \frac{\pi}{m} \alpha} \right) \right| = 2 \left| \sin \left( \alpha \frac{\pi}{m} \right) \right|. \tag{286}$$

Without loss of generality, we focus on the case when  $l = 1$  in the sequel. Recall that for  $c > 0$ , we denote by  $\lfloor c \rfloor$  the largest integer that does not exceed  $c$ . We can continue the derivation to get

$$\begin{aligned} \sum_{j=1}^m \left| \mathbf{b}_1^H \mathbf{b}_j \right| &= \frac{K}{m} + \frac{1}{m} \sum_{j=2}^m \left| \frac{\sin \left[ K(1-j) \frac{\pi}{m} \right]}{\sin \left[ (1-j) \frac{\pi}{m} \right]} \right| \stackrel{(i)}{\leq} \frac{1}{m} \sum_{j=2}^m \left| \frac{1}{\sin \left[ (j-1) \frac{\pi}{m} \right]} \right| + \frac{K}{m} \\ &= \frac{1}{m} \left( \sum_{j=2}^{\lfloor \frac{m}{2} \rfloor + 1} \left| \frac{1}{\sin \left[ (j-1) \frac{\pi}{m} \right]} \right| + \sum_{j=\lfloor \frac{m}{2} \rfloor + 2}^m \left| \frac{1}{\sin \left[ (j-1) \frac{\pi}{m} \right]} \right| \right) + \frac{K}{m} \\ &\stackrel{(ii)}{=} \frac{1}{m} \left( \sum_{j=2}^{\lfloor \frac{m}{2} \rfloor + 1} \left| \frac{1}{\sin \left[ (j-1) \frac{\pi}{m} \right]} \right| + \sum_{j=\lfloor \frac{m}{2} \rfloor + 2}^m \left| \frac{1}{\sin \left[ (m+1-j) \frac{\pi}{m} \right]} \right| \right) + \frac{K}{m}, \end{aligned}$$

where (i) follows from  $\left| \sin \left( K(1-j) \frac{\pi}{m} \right) \right| \leq 1$  and  $|\sin(x)| = |\sin(-x)|$ , and (ii) relies on the fact that  $\sin(x) = \sin(\pi - x)$ . The property that  $\sin(x) \geq x/2$  for any  $x \in [0, \pi/2]$  allows one to further derive

$$\begin{aligned} \sum_{j=1}^m \left| \mathbf{b}_1^H \mathbf{b}_j \right| &\leq \frac{1}{m} \left( \sum_{j=2}^{\lfloor \frac{m}{2} \rfloor + 1} \frac{2m}{(j-1)\pi} + \sum_{j=\lfloor \frac{m}{2} \rfloor + 2}^m \frac{2m}{(m+1-j)\pi} \right) \\ &\quad + \frac{K}{m} = \frac{2}{\pi} \left( \sum_{k=1}^{\lfloor \frac{m}{2} \rfloor} \frac{1}{k} + \sum_{k=1}^{\lfloor \frac{m+1}{2} \rfloor - 1} \frac{1}{k} \right) + \frac{K}{m} \\ &\stackrel{(i)}{\leq} \frac{4}{\pi} \sum_{k=1}^m \frac{1}{k} + \frac{K}{m} \stackrel{(ii)}{\leq} \frac{4}{\pi} (1 + \log m) + 1 \stackrel{(iii)}{\leq} 4 \log m, \end{aligned}$$

where in (i) we extend the range of the summation, (ii) uses the elementary inequality  $\sum_{k=1}^m k^{-1} \leq 1 + \log m$ , and (iii) holds true as long as  $m \geq 3$ .  $\square$

The next lemma considers the difference of two inner products, namely  $(\mathbf{b}_l - \mathbf{b}_1)^H \mathbf{b}_j$ .

**Lemma 49** For all  $0 \leq l - 1 \leq \tau \leq \lfloor \frac{m}{10} \rfloor$ , we have

$$|(\mathbf{b}_l - \mathbf{b}_1)^H \mathbf{b}_j| \leq \begin{cases} \frac{4\tau}{(j-l)} \frac{K}{m} + \frac{8\tau/\pi}{(j-l)^2} & \text{for } l + \tau \leq j \leq \lfloor \frac{m}{2} \rfloor + 1, \\ \frac{4\tau}{m-(j-l)} \frac{K}{m} + \frac{8\tau/\pi}{[m-(j-1)]^2} & \text{for } \lfloor \frac{m}{2} \rfloor + l \leq j \leq m - \tau. \end{cases}$$

In addition, for any  $j$  and  $l$ , the following uniform upper bound holds

$$|(\mathbf{b}_l - \mathbf{b}_1)^H \mathbf{b}_j| \leq 2 \frac{K}{m}.$$

**Proof** Given (285), we can obtain for  $j \neq l$  and  $j \neq 1$ ,

$$\begin{aligned} |(\mathbf{b}_l - \mathbf{b}_1)^H \mathbf{b}_j| &= \frac{1}{m} \left| \frac{1 - \omega^{K(l-j)}}{1 - \omega^{l-j}} - \frac{1 - \omega^{K(1-j)}}{1 - \omega^{1-j}} \right| \\ &= \frac{1}{m} \left| \frac{1 - \omega^{K(l-j)}}{1 - \omega^{l-j}} - \frac{1 - \omega^{K(1-j)}}{1 - \omega^{l-j}} + \frac{1 - \omega^{K(1-j)}}{1 - \omega^{l-j}} - \frac{1 - \omega^{K(1-j)}}{1 - \omega^{1-j}} \right| \\ &= \frac{1}{m} \left| \frac{\omega^{K(1-j)} - \omega^{K(l-j)}}{1 - \omega^{l-j}} + (\omega^{l-j} - \omega^{1-j}) \frac{1 - \omega^{K(1-j)}}{(1 - \omega^{l-j})(1 - \omega^{1-j})} \right| \\ &\leq \frac{1}{m} \left| \frac{1 - \omega^{K(l-1)}}{1 - \omega^{l-j}} \right| + \frac{2}{m} \left| (1 - \omega^{1-l}) \frac{1}{(1 - \omega^{l-j})(1 - \omega^{1-j})} \right|, \end{aligned}$$

where the last line is due to the triangle inequality and  $|\omega^\alpha| = 1$  for all  $\alpha \in \mathbb{R}$ . Identity (286) allows us to rewrite this bound as

$$|(\mathbf{b}_l - \mathbf{b}_1)^H \mathbf{b}_j| \leq \frac{1}{m} \left| \frac{1}{\sin \left[ (l-j) \frac{\pi}{m} \right]} \right| \left\{ \left| \sin \left[ K(l-1) \frac{\pi}{m} \right] \right| + \left| \frac{\sin \left[ (1-l) \frac{\pi}{m} \right]}{\sin \left[ (1-j) \frac{\pi}{m} \right]} \right| \right\}. \tag{287}$$

Combined with the fact that  $|\sin x| \leq |x|$  for all  $x \in \mathbb{R}$ , we can upper bound (287) as

$$|(\mathbf{b}_l - \mathbf{b}_1)^H \mathbf{b}_j| \leq \frac{1}{m} \left| \frac{1}{\sin \left[ (l-j) \frac{\pi}{m} \right]} \right| \left\{ 2K\tau \frac{\pi}{m} + \left| \frac{2\tau \frac{\pi}{m}}{\sin \left[ (1-j) \frac{\pi}{m} \right]} \right| \right\},$$

where we also utilize the assumption  $0 \leq l - 1 \leq \tau$ . Then for  $l + \tau \leq j \leq \lfloor m/2 \rfloor + 1$ , one has

$$\left| (l-j) \frac{\pi}{m} \right| \leq \frac{\pi}{2} \quad \text{and} \quad \left| (1-j) \frac{\pi}{m} \right| \leq \frac{\pi}{2}.$$

Therefore, utilizing the property  $\sin(x) \geq x/2$  for any  $x \in [0, \pi/2]$ , we arrive at

$$|(\mathbf{b}_l - \mathbf{b}_1)^H \mathbf{b}_j| \leq \frac{2}{(j-l)\pi} \left( 2K\tau \frac{\pi}{m} + \frac{4\tau}{j-1} \right) \leq \frac{4\tau}{(j-l)} \frac{K}{m} + \frac{8\tau/\pi}{(j-l)^2},$$

where the last inequality holds since  $j - 1 > j - l$ . Similarly, we can obtain the upper bound for  $\lfloor m/2 \rfloor + l \leq j \leq m - \tau$  using nearly identical argument (which is omitted for brevity).

The uniform upper bound can be justified as follows

$$\left| (\mathbf{b}_l - \mathbf{b}_1)^H \mathbf{b}_j \right| \leq (\|\mathbf{b}_l\|_2 + \|\mathbf{b}_1\|_2) \|\mathbf{b}_j\|_2 \leq 2K/m.$$

The last relation holds since  $\|\mathbf{b}_l\|_2^2 = K/m$  for all  $1 \leq l \leq m$ . □

Next, we list two consequences of the above estimates in Lemmas 50 and 51.

**Lemma 50** Fix any constant  $c > 0$  that is independent of  $m$  and  $K$ . Suppose  $m \geq C\tau K \log^4 m$  for some sufficiently large constant  $C > 0$ , which solely depends on  $c$ . If  $0 \leq l - 1 \leq \tau$ , then one has

$$\sum_{j=1}^m \left| (\mathbf{b}_l - \mathbf{b}_1)^H \mathbf{b}_j \right| \leq \frac{c}{\log^2 m}.$$

**Proof** For some constant  $c_0 > 0$ , we can split the index set  $[m]$  into the following three disjoint sets

$$\begin{aligned} \mathcal{A}_1 &= \left\{ j : l + c_0\tau \log^2 m \leq j \leq \left\lfloor \frac{m}{2} \right\rfloor \right\}, \\ \mathcal{A}_2 &= \left\{ j : \left\lfloor \frac{m}{2} \right\rfloor + l \leq j \leq m - c_0\tau \log^2 m \right\}, \\ \text{and } \mathcal{A}_3 &= [m] \setminus (\mathcal{A}_1 \cup \mathcal{A}_2). \end{aligned}$$

With this decomposition in place, we can write

$$\sum_{j=1}^m \left| (\mathbf{b}_l - \mathbf{b}_1)^H \mathbf{b}_j \right| = \sum_{j \in \mathcal{A}_1} \left| (\mathbf{b}_l - \mathbf{b}_1)^H \mathbf{b}_j \right| + \sum_{j \in \mathcal{A}_2} \left| (\mathbf{b}_l - \mathbf{b}_1)^H \mathbf{b}_j \right| + \sum_{j \in \mathcal{A}_3} \left| (\mathbf{b}_l - \mathbf{b}_1)^H \mathbf{b}_j \right|.$$

We first look at  $\mathcal{A}_1$ . By Lemma 49, one has for any  $j \in \mathcal{A}_1$ ,

$$\left| (\mathbf{b}_l - \mathbf{b}_1)^H \mathbf{b}_j \right| \leq \frac{4\tau}{j-l} \frac{K}{m} + \frac{8\tau/\pi}{(j-l)^2},$$

and hence

$$\begin{aligned} \sum_{j \in \mathcal{A}_1} \left| (\mathbf{b}_l - \mathbf{b}_1)^H \mathbf{b}_j \right| &\leq \sum_{j=l+c_0\tau \log^2 m}^{\lfloor \frac{m}{2} \rfloor + 1} \left( \frac{4\tau}{j-l} \frac{K}{m} + \frac{8\tau/\pi}{(j-l)^2} \right) \\ &\leq \frac{4\tau K}{m} \sum_{k=1}^m \frac{1}{k} + \frac{8\tau}{\pi} \sum_{k=c_0\tau \log^2 m}^m \frac{1}{k^2} \end{aligned}$$

$$\leq 8\tau \frac{K}{m} \log m + \frac{16\tau}{\pi} \frac{1}{c_0\tau \log^2 m},$$

where the last inequality arises from  $\sum_{k=1}^m k^{-1} \leq 1 + \log m \leq 2 \log m$  and  $\sum_{k=c}^m k^{-2} \leq 2/c$ .

Similarly, for  $j \in \mathcal{A}_2$ , we have

$$\left| (\mathbf{b}_l - \mathbf{b}_1)^H \mathbf{b}_j \right| \leq \frac{4\tau}{m - (j - l)} \frac{K}{m} + \frac{8\tau/\pi}{[m - (j - 1)]^2},$$

which in turn implies

$$\sum_{j \in \mathcal{A}_2} \left| (\mathbf{b}_l - \mathbf{b}_1)^H \mathbf{b}_j \right| \leq 8\tau \frac{K}{m} \log m + \frac{16\tau}{\pi} \frac{1}{c_0\tau \log^2 m}.$$

Regarding  $j \in \mathcal{A}_3$ , we observe that

$$|\mathcal{A}_3| \leq 2 \left( c_0\tau \log^2 m + l \right) \leq 2 \left( c_0\tau \log^2 m + \tau + 1 \right) \leq 4c_0\tau \log^2 m.$$

This together with the simple bound  $\left| (\mathbf{b}_l - \mathbf{b}_1)^H \mathbf{b}_j \right| \leq 2K/m$  gives

$$\sum_{j \in \mathcal{A}_3} \left| (\mathbf{b}_l - \mathbf{b}_1)^H \mathbf{b}_j \right| \leq 2 \frac{K}{m} |\mathcal{A}_3| \leq \frac{8c_0\tau K \log^2 m}{m}.$$

The previous three estimates taken collectively yield

$$\sum_{j=1}^m \left| (\mathbf{b}_l - \mathbf{b}_1)^H \mathbf{b}_j \right| \leq \frac{16\tau K \log m}{m} + \frac{32\tau}{\pi} \frac{1}{c_0\tau \log^2 m} + \frac{8c_0\tau K \log^2 m}{m} \leq c \frac{1}{\log^2 m}$$

as long as  $c_0 \geq (32/\pi) \cdot (1/c)$  and  $m \geq 8c_0\tau K \log^4 m/c$ . □

**Lemma 51** Fix any constant  $c > 0$  that is independent of  $m$  and  $K$ . Consider an integer  $\tau > 0$ , and suppose that  $m \geq C\tau K \log m$  for some large constant  $C > 0$ , which depends solely on  $c$ . Then we have

$$\sum_{k=0}^{\lfloor m/\tau \rfloor} \sqrt{\sum_{j=1}^{\tau} \left| \mathbf{b}_1^H (\mathbf{b}_{k\tau+j} - \mathbf{b}_{k\tau+1}) \right|^2} \leq \frac{c}{\sqrt{\tau}}.$$

**Proof** The proof strategy is similar to the one used in Lemma 50. First, notice that

$$\left| \mathbf{b}_1^H (\mathbf{b}_{k\tau+j} - \mathbf{b}_{k\tau+1}) \right| = \left| (\mathbf{b}_m - \mathbf{b}_{m+1-j})^H \mathbf{b}_{k\tau} \right|.$$

As before, for some  $c_1 > 0$ , we can split the index set  $\{1, \dots, \lfloor m/\tau \rfloor\}$  into three disjoint sets

$$\begin{aligned} \mathcal{B}_1 &= \left\{ k : c_1 \leq k \leq \left\lfloor \left( \left\lfloor \frac{m}{2} \right\rfloor + 1 - j \right) / \tau \right\rfloor \right\}, \\ \mathcal{B}_2 &= \left\{ k : \left\lfloor \left( \left\lfloor \frac{m}{2} \right\rfloor + 1 - j \right) / \tau \right\rfloor + 1 \leq k \leq \lfloor (m + 1 - j) / \tau \rfloor - c_1 \right\}, \\ \text{and } \mathcal{B}_3 &= \left\{ 1, \dots, \left\lfloor \frac{m}{\tau} \right\rfloor \right\} \setminus (\mathcal{B}_1 \cup \mathcal{B}_2), \end{aligned}$$

where  $1 \leq j \leq \tau$ .

By Lemma 49, one has

$$\left| (\mathbf{b}_m - \mathbf{b}_{m+1-j})^H \mathbf{b}_{k\tau} \right| \leq \frac{4\tau}{k\tau} \frac{K}{m} + \frac{8\tau/\pi}{(k\tau)^2}, \quad k \in \mathcal{B}_1.$$

Hence, for any  $k \in \mathcal{B}_1$ ,

$$\sqrt{\sum_{j=1}^{\tau} \left| \mathbf{b}_1^H (\mathbf{b}_{k\tau+j} - \mathbf{b}_{k\tau+1}) \right|^2} \leq \sqrt{\tau} \left( \frac{4\tau}{k\tau} \frac{K}{m} + \frac{8\tau/\pi}{(k\tau)^2} \right) = \sqrt{\tau} \left( \frac{4}{k} \frac{K}{m} + \frac{8/\pi}{k^2\tau} \right),$$

which further implies that

$$\begin{aligned} \sum_{k \in \mathcal{B}_1} \sqrt{\sum_{j=1}^{\tau} \left| \mathbf{b}_1^H (\mathbf{b}_{k\tau+j} - \mathbf{b}_{k\tau+1}) \right|^2} &\leq \sqrt{\tau} \sum_{k=c_1}^m \left( \frac{4}{k} \frac{K}{m} + \frac{8/\pi}{k^2\tau} \right) \\ &\leq 8\sqrt{\tau} \frac{K \log m}{m} + \frac{16}{\pi} \frac{1}{\sqrt{\tau} c_1}, \end{aligned}$$

where the last inequality follows since  $\sum_{k=1}^m k^{-1} \leq 2 \log m$  and  $\sum_{k=c_1}^m k^{-2} \leq 2/c_1$ . A similar bound can be obtained for  $k \in \mathcal{B}_2$ .

For the remaining set  $\mathcal{B}_3$ , observe that

$$|\mathcal{B}_3| \leq 2c_1.$$

This together with the crude upper bound  $\left| (\mathbf{b}_l - \mathbf{b}_1)^H \mathbf{b}_j \right| \leq 2K/m$  gives

$$\begin{aligned} \sum_{k \in \mathcal{B}_3} \sqrt{\sum_{j=1}^{\tau} \left| \mathbf{b}_1^H (\mathbf{b}_{k\tau+j} - \mathbf{b}_{k\tau+1}) \right|^2} &\leq |\mathcal{B}_3| \sqrt{\tau \max_j \left| \mathbf{b}_1^H (\mathbf{b}_{k\tau+j} - \mathbf{b}_{k\tau+1}) \right|^2} \\ &\leq |\mathcal{B}_3| \sqrt{\tau} \cdot \frac{2K}{m} \leq \frac{4c_1 \sqrt{\tau} K}{m}. \end{aligned}$$

The previous estimates taken collectively yield

$$\sum_{k=0}^{\lfloor m/\tau \rfloor} \sqrt{\sum_{j=1}^{\tau} \left| \mathbf{b}_1^H (\mathbf{b}_{k\tau+j} - \mathbf{b}_{k\tau+1}) \right|^2} \leq 2 \left( 8\sqrt{\tau} \frac{K \log m}{m} + \frac{16}{\pi} \frac{1}{\sqrt{\tau}} \frac{1}{c_1} \right) + \frac{4c_1\sqrt{\tau}K}{m} \leq c \frac{1}{\sqrt{\tau}},$$

as long as  $c_1 \gg 1/c$  and  $m/(c_1 \tau K \log m) \gg 1/c$ . □

### D.3.3 Complex-Valued Alignment

Let  $g_{\mathbf{h},\mathbf{x}}(\cdot) : \mathbb{C} \rightarrow \mathbb{R}$  be a real-valued function defined as

$$g_{\mathbf{h},\mathbf{x}}(\alpha) := \left\| \frac{1}{\alpha} \mathbf{h} - \mathbf{h}^* \right\|_2^2 + \left\| \alpha \mathbf{x} - \mathbf{x}^* \right\|_2^2,$$

which is the key function in definition (34). Therefore, the alignment parameter of  $(\mathbf{h}, \mathbf{x})$  to  $(\mathbf{h}^*, \mathbf{x}^*)$  is the minimizer of  $g_{\mathbf{h},\mathbf{x}}(\alpha)$ . This section is devoted to studying various properties of  $g_{\mathbf{h},\mathbf{x}}(\cdot)$ . To begin with, the Wirtinger gradient and Hessian of  $g_{\mathbf{h},\mathbf{x}}(\cdot)$  can be calculated as

$$\nabla g_{\mathbf{h},\mathbf{x}}(\alpha) = \begin{bmatrix} \frac{\partial g_{\mathbf{h},\mathbf{x}}(\alpha, \bar{\alpha})}{\partial \alpha} \\ \frac{\partial g_{\mathbf{h},\mathbf{x}}(\alpha, \bar{\alpha})}{\partial \bar{\alpha}} \end{bmatrix} = \begin{bmatrix} \alpha \|\mathbf{x}\|_2^2 - \mathbf{x}^H \mathbf{x}^* - \alpha^{-1} (\bar{\alpha})^{-2} \|\mathbf{h}\|_2^2 + (\bar{\alpha})^{-2} \mathbf{h}^H \mathbf{h} \\ \bar{\alpha} \|\mathbf{x}\|_2^2 - \mathbf{x}^{*H} \mathbf{x} - (\bar{\alpha})^{-1} \alpha^{-2} \|\mathbf{h}\|_2^2 + \alpha^{-2} \mathbf{h}^H \mathbf{h}^* \end{bmatrix}; \tag{288}$$

$$\nabla^2 g_{\mathbf{h},\mathbf{x}}(\alpha) = \begin{bmatrix} \|\mathbf{x}\|_2^2 + |\alpha|^{-4} \|\mathbf{h}\|_2^2 & 2\alpha^{-1} (\bar{\alpha})^{-3} \|\mathbf{h}\|_2^2 - 2(\bar{\alpha})^{-3} \mathbf{h}^H \mathbf{h} \\ 2(\bar{\alpha})^{-1} \alpha^{-3} \|\mathbf{h}\|_2^2 - 2\alpha^{-3} \mathbf{h}^H \mathbf{h}^* & \|\mathbf{x}\|_2^2 + |\alpha|^{-4} \|\mathbf{h}\|_2^2 \end{bmatrix}. \tag{289}$$

The first lemma reveals that, as long as  $(\frac{1}{\beta} \mathbf{h}, \beta \mathbf{x})$  is sufficiently close to  $(\mathbf{h}^*, \mathbf{x}^*)$ , the minimizer of  $g_{\mathbf{h},\mathbf{x}}(\alpha)$  cannot be far away from  $\beta$ .

**Lemma 52** *Assume there exists  $\beta \in \mathbb{C}$  with  $1/2 \leq |\beta| \leq 3/2$  such that  $\max \left\{ \left\| \frac{1}{\beta} \mathbf{h} - \mathbf{h}^* \right\|_2, \|\beta \mathbf{x} - \mathbf{x}^*\|_2 \right\} \leq \delta \leq 1/4$ . Denote by  $\hat{\alpha}$  the minimizer of  $g_{\mathbf{h},\mathbf{x}}(\alpha)$ , and then we necessarily have*

$$|\hat{\alpha}| - |\beta| \leq |\hat{\alpha} - \beta| \leq 18\delta.$$

**Proof** The first inequality is a direct consequence of the triangle inequality. Hence, we concentrate on the second one. Notice that by assumption,

$$g_{\mathbf{h},\mathbf{x}}(\beta) = \left\| \frac{1}{\beta} \mathbf{h} - \mathbf{h}^* \right\|_2^2 + \|\beta \mathbf{x} - \mathbf{x}^*\|_2^2 \leq 2\delta^2, \tag{290}$$

which immediately implies that  $g_{\mathbf{h},\mathbf{x}}(\hat{\alpha}) \leq 2\delta^2$ . It thus suffices to show that for any  $\alpha$  obeying  $|\alpha - \beta| > 18\delta$ , one has  $g_{\mathbf{h},\mathbf{x}}(\alpha) > 2\delta^2$ , and hence, it cannot be the minimizer.

To this end, we lower bound  $g_{h,x}(\alpha)$  as follows:

$$\begin{aligned} g_{h,x}(\alpha) &\geq \|\alpha x - x^*\|_2^2 = \|(\alpha - \beta)x + (\beta x - x^*)\|_2^2 \\ &= |\alpha - \beta|^2 \|x\|_2^2 + \|\beta x - x^*\|_2^2 + 2\operatorname{Re}\left[(\alpha - \beta)(\beta x - x^*)^H x\right] \\ &\geq |\alpha - \beta|^2 \|x\|_2^2 - 2|\alpha - \beta| \left|(\beta x - x^*)^H x\right|. \end{aligned}$$

Given that  $\|\beta x - x^*\|_2 \leq \delta \leq 1/4$  and  $\|x^*\|_2 = 1$ , we have

$$\|\beta x\|_2 \geq \|x^*\|_2 - \|\beta x - x^*\|_2 \geq 1 - \delta \geq 3/4,$$

which together with the fact that  $1/2 \leq |\beta| \leq 3/2$  implies

$$\|x\|_2 \geq 1/2 \quad \text{and} \quad \|x\|_2 \leq 2$$

and

$$\left|(\beta x - x^*)^H x\right| \leq \|\beta x - x^*\|_2 \|x\|_2 \leq 2\delta.$$

Taking the previous estimates collectively yields

$$g_{h,x}(\alpha) \geq \frac{1}{4} |\alpha - \beta|^2 - 4\delta |\alpha - \beta|.$$

It is self-evident that once  $|\alpha - \beta| > 18\delta$ , one gets  $g_{h,x}(\alpha) > 2\delta^2$ , and hence,  $\alpha$  cannot be the minimizer as  $g_{h,x}(\alpha) > g_{h,x}(\beta)$  according to (290). This concludes the proof.  $\square$

The next lemma reveals the local strong convexity of  $g_{h,x}(\alpha)$  when  $\alpha$  is close to one.

**Lemma 53** *Assume that  $\max\{\|h - h^*\|_2, \|x - x^*\|_2\} \leq \delta$  for some sufficiently small constant  $\delta > 0$ . Then, for any  $\alpha$  satisfying  $|\alpha - 1| \leq 18\delta$  and any  $u, v \in \mathbb{C}$ , one has*

$$\begin{bmatrix} u^H & v^H \end{bmatrix} \nabla^2 g_{h,x}(\alpha) \begin{bmatrix} u \\ v \end{bmatrix} \geq \frac{1}{2} (|u|^2 + |v|^2),$$

where  $\nabla^2 g_{h,x}(\cdot)$  stands for the Wirtinger Hessian of  $g_{h,x}(\cdot)$ .

**Proof** For simplicity of presentation, we use  $g_{h,x}(\alpha, \bar{\alpha})$  and  $g_{h,x}(\alpha)$  interchangeably. By (289), for any  $u, v \in \mathbb{C}$ , one has

$$\begin{bmatrix} u^H & v^H \end{bmatrix} \nabla^2 g_{h,x}(\alpha) \begin{bmatrix} u \\ v \end{bmatrix} = \underbrace{\left(\|x\|_2^2 + |\alpha|^{-4} \|h\|_2^2\right)}_{:=\beta_1} (|u|^2 + |v|^2)$$



$$+ 2 \operatorname{Re} \underbrace{\left[ u^H v \left( 2\alpha^{-1} (\bar{\alpha})^{-3} \|\mathbf{h}\|_2^2 - 2 (\bar{\alpha})^{-3} \mathbf{h}^{*H} \mathbf{h} \right) \right]}_{:=\beta_2}.$$

We would like to demonstrate that this is at least on the order of  $|u|^2 + |v|^2$ . We first develop a lower bound on  $\beta_1$ . Given the assumption that  $\max \{ \|\mathbf{h} - \mathbf{h}^*\|_2, \|\mathbf{x} - \mathbf{x}^*\|_2 \} \leq \delta$ , one necessarily has

$$1 - \delta \leq \|\mathbf{x}\|_2 \leq 1 + \delta \quad \text{and} \quad 1 - \delta \leq \|\mathbf{h}\|_2 \leq 1 + \delta.$$

Thus, for any  $\alpha$  obeying  $|\alpha - 1| \leq 18\delta$ , one has

$$\beta_1 \geq \left( 1 + |\alpha|^{-4} \right) (1 - \delta)^2 \geq \left( 1 + (1 + 18\delta)^{-4} \right) (1 - \delta)^2 \geq 1$$

as long as  $\delta > 0$  is sufficiently small. Regarding the second term  $\beta_2$ , we utilize the conditions  $|\alpha - 1| \leq 18\delta$ ,  $\|\mathbf{x}\|_2 \leq 1 + \delta$  and  $\|\mathbf{h}\|_2 \leq 1 + \delta$  to get

$$\begin{aligned} |\beta_2| &\leq 2 |u| |v| |\alpha|^{-3} \left| \alpha^{-1} \|\mathbf{h}\|_2^2 - \mathbf{h}^{*H} \mathbf{h} \right| \\ &= 2 |u| |v| |\alpha|^{-3} \left| (\alpha^{-1} - 1) \|\mathbf{h}\|_2^2 - (\mathbf{h}^* - \mathbf{h})^H \mathbf{h} \right| \\ &\leq 2 |u| |v| |\alpha|^{-3} \left( |\alpha^{-1} - 1| \|\mathbf{h}\|_2^2 + \|\mathbf{h} - \mathbf{h}^*\|_2 \|\mathbf{h}\|_2 \right) \\ &\leq 2 |u| |v| (1 - 18\delta)^{-3} \left( \frac{18\delta}{1 - 18\delta} (1 + \delta)^2 + \delta (1 + \delta) \right) \\ &\lesssim \delta (|u|^2 + |v|^2), \end{aligned}$$

where the last relation holds since  $2 |u| |v| \leq |u|^2 + |v|^2$  and  $\delta > 0$  is sufficiently small. Combining the previous bounds on  $\beta_1$  and  $\beta_2$ , we arrive at

$$\left[ u^H, v^H \right] \nabla^2 g_{\mathbf{h}, \mathbf{x}}(\alpha) \begin{bmatrix} u \\ v \end{bmatrix} \geq (1 - O(\delta)) (|u|^2 + |v|^2) \geq \frac{1}{2} (|u|^2 + |v|^2)$$

as long as  $\delta$  is sufficiently small. This completes the proof. □

Additionally, in a local region surrounding the optimizer, the alignment parameter is Lipschitz continuous; namely, the difference of the alignment parameters associated with two distinct vector pairs is at most proportional to the  $\ell_2$  distance between the two vector pairs involved, as demonstrated below.

**Lemma 54** *Suppose that the vectors  $\mathbf{x}_1, \mathbf{x}_2, \mathbf{h}_1, \mathbf{h}_2 \in \mathbb{C}^K$  satisfy*

$$\max \{ \|\mathbf{x}_1 - \mathbf{x}^*\|_2, \|\mathbf{h}_1 - \mathbf{h}^*\|_2, \|\mathbf{x}_2 - \mathbf{x}^*\|_2, \|\mathbf{h}_2 - \mathbf{h}^*\|_2 \} \leq \delta \leq 1/4 \quad (291)$$

for some sufficiently small constant  $\delta > 0$ . Denote by  $\alpha_1$  and  $\alpha_2$  the minimizers of  $g_{\mathbf{h}_1, \mathbf{x}_1}(\alpha)$  and  $g_{\mathbf{h}_2, \mathbf{x}_2}(\alpha)$ , respectively. Then we have

$$|\alpha_1 - \alpha_2| \lesssim \|\mathbf{x}_1 - \mathbf{x}_2\|_2 + \|\mathbf{h}_1 - \mathbf{h}_2\|_2.$$

**Proof** Since  $\alpha_1$  minimizes  $g_{h_1, x_1}(\alpha)$ , the mean value form of Taylor's theorem (see Appendix D.3.1) gives

$$\begin{aligned} g_{h_1, x_1}(\alpha_2) &\geq g_{h_1, x_1}(\alpha_1) \\ &= g_{h_1, x_1}(\alpha_2) + \nabla g_{h_1, x_1}(\alpha_2)^H \left[ \frac{\alpha_1 - \alpha_2}{\alpha_1 - \alpha_2} \right] \\ &\quad + \frac{1}{2} (\overline{\alpha_1 - \alpha_2}, \alpha_1 - \alpha_2) \nabla^2 g_{h_1, x_1}(\tilde{\alpha}) \left[ \frac{\alpha_1 - \alpha_2}{\alpha_1 - \alpha_2} \right], \end{aligned}$$

where  $\tilde{\alpha}$  is some complex number lying between  $\alpha_1$  and  $\alpha_2$ , and  $\nabla g_{h_1, x_1}$  and  $\nabla^2 g_{h_1, x_1}$  are the Wirtinger gradient and Hessian of  $g_{h_1, x_1}(\cdot)$ , respectively. Rearrange the previous inequality to obtain

$$|\alpha_1 - \alpha_2| \lesssim \frac{\|\nabla g_{h_1, x_1}(\alpha_2)\|_2}{\lambda_{\min}(\nabla^2 g_{h_1, x_1}(\tilde{\alpha}))} \quad (292)$$

as long as  $\lambda_{\min}(\nabla^2 g_{h_1, x_1}(\tilde{\alpha})) > 0$ . This calls for evaluation of the Wirtinger gradient and Hessian of  $g_{h_1, x_1}(\cdot)$ .

Regarding the Wirtinger Hessian, by assumption (291), we can invoke Lemma 52 with  $\beta = 1$  to reach  $\max\{|\alpha_1 - 1|, |\alpha_2 - 1|\} \leq 18\delta$ . This together with Lemma 53 implies

$$\lambda_{\min}(\nabla^2 g_{h_1, x_1}(\tilde{\alpha})) \geq 1/2,$$

since  $\tilde{\alpha}$  lies between  $\alpha_1$  and  $\alpha_2$ .

For the Wirtinger gradient, since  $\alpha_2$  is the minimizer of  $g_{h_2, x_2}(\alpha)$ , the first-order optimality condition [69, equation (38)] requires  $\nabla g_{h_2, x_2}(\alpha_2) = \mathbf{0}$ , which gives

$$\|\nabla g_{h_1, x_1}(\alpha_2)\|_2 = \|\nabla g_{h_1, x_1}(\alpha_2) - \nabla g_{h_2, x_2}(\alpha_2)\|_2.$$

Plug in the gradient expression (288) to reach

$$\begin{aligned} &\|\nabla g_{h_1, x_1}(\alpha_2) - \nabla g_{h_2, x_2}(\alpha_2)\|_2 \\ &= \sqrt{2} \left| \left[ \alpha_2 \|\mathbf{x}_1\|_2^2 - \mathbf{x}_1^H \mathbf{x}^* - \alpha_2^{-1} (\overline{\alpha_2})^{-2} \|\mathbf{h}_1\|_2^2 + (\overline{\alpha_2})^{-2} \mathbf{h}^{*H} \mathbf{h}_1 \right] \right. \\ &\quad \left. - \left[ \alpha_2 \|\mathbf{x}_2\|_2^2 - \mathbf{x}_2^H \mathbf{x}^* - \alpha_2^{-1} (\overline{\alpha_2})^{-2} \|\mathbf{h}_2\|_2^2 + (\overline{\alpha_2})^{-2} \mathbf{h}^{*H} \mathbf{h}_2 \right] \right| \\ &\lesssim |\alpha_2| \left| \|\mathbf{x}_1\|_2^2 - \|\mathbf{x}_2\|_2^2 \right| + \left| \mathbf{x}_1^H \mathbf{x}^* - \mathbf{x}_2^H \mathbf{x}^* \right| + \frac{1}{|\alpha_2|^3} \left| \|\mathbf{h}_1\|_2^2 - \|\mathbf{h}_2\|_2^2 \right| + \frac{1}{|\alpha_2|^2} \left| \mathbf{h}^{*H} \mathbf{h}_1 - \mathbf{h}^{*H} \mathbf{h}_2 \right| \\ &\lesssim |\alpha_2| \left| \|\mathbf{x}_1\|_2^2 - \|\mathbf{x}_2\|_2^2 \right| + \|\mathbf{x}_1 - \mathbf{x}_2\|_2 + \frac{1}{|\alpha_2|^3} \left| \|\mathbf{h}_1\|_2^2 - \|\mathbf{h}_2\|_2^2 \right| + \frac{1}{|\alpha_2|^2} \|\mathbf{h}_1 - \mathbf{h}_2\|_2, \end{aligned}$$

where the last line follows from the triangle inequality. It is straightforward to see that

$$1/2 \leq |\alpha_2| \leq 2, \quad \left| \|\mathbf{x}_1\|_2^2 - \|\mathbf{x}_2\|_2^2 \right| \lesssim \|\mathbf{x}_1 - \mathbf{x}_2\|_2, \quad \left| \|\mathbf{h}_1\|_2^2 - \|\mathbf{h}_2\|_2^2 \right| \lesssim \|\mathbf{h}_1 - \mathbf{h}_2\|_2$$

under condition (291) and assumption  $\|\mathbf{x}^*\|_2 = \|\mathbf{h}^*\|_2 = 1$ , where the first inequality follows from Lemma 52. Taking these estimates together reveals that

$$\|\nabla g_{h_1, x_1}(\alpha_2) - \nabla g_{h_2, x_2}(\alpha_2)\|_2 \lesssim \|\mathbf{x}_1 - \mathbf{x}_2\|_2 + \|\mathbf{h}_1 - \mathbf{h}_2\|_2.$$

The proof is accomplished by substituting the two bounds on the gradient and the Hessian into (292). □

Further, if two vector pairs are both close to the optimizer, then their distance after alignment (w.r.t. the optimizer) cannot be much larger than their distance without alignment, as revealed by the following lemma.

**Lemma 55** *Suppose that the vectors  $\mathbf{x}_1, \mathbf{x}_2, \mathbf{h}_1, \mathbf{h}_2 \in \mathbb{C}^K$  satisfy*

$$\max \{ \|\mathbf{x}_1 - \mathbf{x}^*\|_2, \|\mathbf{h}_1 - \mathbf{h}^*\|_2, \|\mathbf{x}_2 - \mathbf{x}^*\|_2, \|\mathbf{h}_2 - \mathbf{h}^*\|_2 \} \leq \delta \leq 1/4 \quad (293)$$

for some sufficiently small constant  $\delta > 0$ . Denote by  $\alpha_1$  and  $\alpha_2$  the minimizers of  $g_{h_1, x_1}(\alpha)$  and  $g_{h_2, x_2}(\alpha)$ , respectively. Then we have

$$\|\alpha_1 \mathbf{x}_1 - \alpha_2 \mathbf{x}_2\|_2^2 + \left\| \frac{1}{\alpha_1} \mathbf{h}_1 - \frac{1}{\alpha_2} \mathbf{h}_2 \right\|_2^2 \lesssim \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 + \|\mathbf{h}_1 - \mathbf{h}_2\|_2^2.$$

**Proof** To start with, we control the magnitudes of  $\alpha_1$  and  $\alpha_2$ . Lemma 52 together with assumption (293) guarantees that

$$1/2 \leq |\alpha_1| \leq 2 \quad \text{and} \quad 1/2 \leq |\alpha_2| \leq 2.$$

Now we can prove the lemma. The triangle inequality gives

$$\begin{aligned} \|\alpha_1 \mathbf{x}_1 - \alpha_2 \mathbf{x}_2\|_2 &= \|\alpha_1 (\mathbf{x}_1 - \mathbf{x}_2) + (\alpha_1 - \alpha_2) \mathbf{x}_2\|_2 \\ &\leq |\alpha_1| \|\mathbf{x}_1 - \mathbf{x}_2\|_2 + |\alpha_1 - \alpha_2| \|\mathbf{x}_2\|_2 \\ &\stackrel{(i)}{\leq} 2 \|\mathbf{x}_1 - \mathbf{x}_2\|_2 + 2 |\alpha_1 - \alpha_2| \\ &\stackrel{(ii)}{\lesssim} \|\mathbf{x}_1 - \mathbf{x}_2\|_2 + \|\mathbf{h}_1 - \mathbf{h}_2\|_2, \end{aligned}$$

where (i) holds since  $|\alpha_1| \leq 2$  and  $\|\mathbf{x}_2\|_2 \leq 1 + \delta \leq 2$ , and (ii) arises from Lemma 54 that  $|\alpha_1 - \alpha_2| \lesssim \|\mathbf{x}_1 - \mathbf{x}_2\|_2 + \|\mathbf{h}_1 - \mathbf{h}_2\|_2$ . Similarly,

$$\begin{aligned} \left\| \frac{1}{\alpha_1} \mathbf{h}_1 - \frac{1}{\alpha_2} \mathbf{h}_2 \right\|_2 &= \left\| \frac{1}{\alpha_1} (\mathbf{h}_1 - \mathbf{h}_2) + \left( \frac{1}{\alpha_1} - \frac{1}{\alpha_2} \right) \mathbf{h}_2 \right\|_2 \\ &\leq \left| \frac{1}{\alpha_1} \right| \|\mathbf{h}_1 - \mathbf{h}_2\|_2 + \left| \frac{1}{\alpha_1} - \frac{1}{\alpha_2} \right| \|\mathbf{h}_2\|_2 \\ &\leq 2 \|\mathbf{h}_1 - \mathbf{h}_2\|_2 + 2 \frac{|\alpha_1 - \alpha_2|}{|\alpha_1 \alpha_2|} \\ &\lesssim \|\mathbf{x}_1 - \mathbf{x}_2\|_2 + \|\mathbf{h}_1 - \mathbf{h}_2\|_2, \end{aligned}$$

where the last inequality comes from Lemma 54 as well as the facts that  $|\alpha_1| \geq 1/2$  and  $|\alpha_2| \geq 1/2$  as shown above. Combining all of the above bounds and recognizing that  $\|\mathbf{x}_1 - \mathbf{x}_2\|_2 + \|\mathbf{h}_1 - \mathbf{h}_2\|_2 \leq \sqrt{2\|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 + 2\|\mathbf{h}_1 - \mathbf{h}_2\|_2^2}$ , we conclude the proof.  $\square$

Finally, there is a useful identity associated with the minimizer of  $\tilde{g}(\alpha)$  as defined below.

**Lemma 56** For any  $\mathbf{h}_1, \mathbf{h}_2, \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{C}^K$ , denote

$$\alpha^\# := \arg \min_{\alpha} \tilde{g}(\alpha), \quad \text{where} \quad \tilde{g}(\alpha) := \left\| \frac{1}{\alpha} \mathbf{h}_1 - \mathbf{h}_2 \right\|_2^2 + \|\alpha \mathbf{x}_1 - \mathbf{x}_2\|_2^2.$$

Let  $\tilde{\mathbf{x}}_1 = \alpha^\# \mathbf{x}_1$  and  $\tilde{\mathbf{h}}_1 = \frac{1}{\alpha^\#} \mathbf{h}_1$ , then we have

$$\|\tilde{\mathbf{x}}_1 - \mathbf{x}_2\|_2^2 + \mathbf{x}_2^H (\tilde{\mathbf{x}}_1 - \mathbf{x}_2) = \|\tilde{\mathbf{h}}_1 - \mathbf{h}_2\|_2^2 + (\tilde{\mathbf{h}}_1 - \mathbf{h}_2)^H \mathbf{h}_2.$$

**Proof** We can rewrite the function  $\tilde{g}(\alpha)$  as

$$\begin{aligned} \tilde{g}(\alpha) &= |\alpha|^2 \|\mathbf{x}_1\|_2^2 + \|\mathbf{x}_2\|_2^2 - (\alpha \mathbf{x}_1)^H \mathbf{x}_2 - \mathbf{x}_2^H (\alpha \mathbf{x}_1) + \left| \frac{1}{\alpha} \right|^2 \|\mathbf{h}_1\|_2^2 + \|\mathbf{h}_2\|_2^2 \\ &\quad - \left( \frac{1}{\alpha} \mathbf{h}_1 \right)^H \mathbf{h}_2 - \mathbf{h}_2^H \left( \frac{1}{\alpha} \mathbf{h}_1 \right) \\ &= \bar{\alpha} \alpha \|\mathbf{x}_1\|_2^2 + \|\mathbf{x}_2\|_2^2 - \bar{\alpha} \mathbf{x}_1^H \mathbf{x}_2 - \alpha \mathbf{x}_2^H \mathbf{x}_1 + \frac{1}{\bar{\alpha} \alpha} \|\mathbf{h}_1\|_2^2 + \|\mathbf{h}_2\|_2^2 \\ &\quad - \frac{1}{\alpha} \mathbf{h}_1^H \mathbf{h}_2 - \frac{1}{\bar{\alpha}} \mathbf{h}_2^H \mathbf{h}_1. \end{aligned}$$

The first-order optimality condition [69, equation (38)] requires

$$\left. \frac{\partial \tilde{g}}{\partial \alpha} \right|_{\alpha=\alpha^\#} = \alpha^\# \|\mathbf{x}_1\|_2^2 - \mathbf{x}_1^H \mathbf{x}_2 + \frac{1}{\alpha^\#} \left( -\frac{1}{\alpha^{\#2}} \right) \|\mathbf{h}_1\|_2^2 - \left( -\frac{1}{\alpha^{\#2}} \right) \mathbf{h}_2^H \mathbf{h}_1 = 0,$$

which further simplifies to

$$\|\tilde{\mathbf{x}}_1\|_2^2 - \tilde{\mathbf{x}}_1^H \mathbf{x}_2 = \|\tilde{\mathbf{h}}_1\|_2^2 - \mathbf{h}_2^H \tilde{\mathbf{h}}_1$$

since  $\tilde{\mathbf{x}}_1 = \alpha^\# \mathbf{x}_1$ ,  $\tilde{\mathbf{h}}_1 = \frac{1}{\alpha^\#} \mathbf{h}_1$ , and  $\alpha^\# \neq 0$  (otherwise  $\tilde{g}(\alpha^\#) = \infty$  and cannot be the minimizer). Furthermore, this condition is equivalent to

$$\tilde{\mathbf{x}}_1^H (\tilde{\mathbf{x}}_1 - \mathbf{x}_2) = (\tilde{\mathbf{h}}_1 - \mathbf{h}_2)^H \tilde{\mathbf{h}}_1.$$

Recognizing that

$$\tilde{\mathbf{x}}_1^H (\tilde{\mathbf{x}}_1 - \mathbf{x}_2) = \mathbf{x}_2^H (\tilde{\mathbf{x}}_1 - \mathbf{x}_2) + (\tilde{\mathbf{x}}_1 - \mathbf{x}_2)^H (\tilde{\mathbf{x}}_1 - \mathbf{x}_2) = \mathbf{x}_2^H (\tilde{\mathbf{x}}_1 - \mathbf{x}_2) + \|\tilde{\mathbf{x}}_1 - \mathbf{x}_2\|_2^2,$$

$$\tilde{\mathbf{h}}_1^H(\tilde{\mathbf{h}}_1 - \mathbf{h}_2) = \mathbf{h}_2^H(\tilde{\mathbf{h}}_1 - \mathbf{h}_2) + (\tilde{\mathbf{h}}_1 - \mathbf{h}_2)^H(\tilde{\mathbf{h}}_1 - \mathbf{h}_2) = \mathbf{h}_2^H(\tilde{\mathbf{h}}_1 - \mathbf{h}_2) + \|\tilde{\mathbf{h}}_1 - \mathbf{h}_2\|_2^2,$$

we arrive at the desired identity. □

### D.3.4 Matrix Concentration Inequalities

The proof for blind deconvolution is largely built upon the concentration of random matrices that are functions of  $\{\mathbf{a}_j \mathbf{a}_j^H\}$ . In this subsection, we collect the measure concentration results for various forms of random matrices that we encounter in the analysis.

**Lemma 57** *Suppose  $\mathbf{a}_j \stackrel{i.i.d.}{\sim} \mathcal{N}(\mathbf{0}, \frac{1}{2} \mathbf{I}_K) + i\mathcal{N}(\mathbf{0}, \frac{1}{2} \mathbf{I}_K)$  for every  $1 \leq j \leq m$ , and  $\{c_j\}_{1 \leq j \leq m}$  are a set of fixed numbers. Then there exist some universal constants  $\tilde{C}_1, \tilde{C}_2 > 0$  such that for all  $t \geq 0$*

$$\mathbb{P}\left(\left\|\sum_{j=1}^m c_j \mathbf{a}_j \mathbf{a}_j^H - \mathbf{I}_K\right\| \geq t\right) \leq 2 \exp\left(\tilde{C}_1 K - \tilde{C}_2 \min\left\{\frac{t}{\max_j |c_j|}, \frac{t^2}{\sum_{j=1}^m c_j^2}\right\}\right).$$

**Proof** This is a simple variant of [116, Theorem 5.39], which uses the Bernstein inequality and the standard covering argument. Hence, we omit its proof. □

**Lemma 58** *Suppose  $\mathbf{a}_j \stackrel{i.i.d.}{\sim} \mathcal{N}(\mathbf{0}, \frac{1}{2} \mathbf{I}_K) + i\mathcal{N}(\mathbf{0}, \frac{1}{2} \mathbf{I}_K)$  for every  $1 \leq j \leq m$ . Then there exist some absolute constants  $\tilde{C}_1, \tilde{C}_2, \tilde{C}_3 > 0$  such that for all  $\max\{1, 3\tilde{C}_1 K / \tilde{C}_2\} / m \leq \varepsilon \leq 1$ , one has*

$$\mathbb{P}\left(\sup_{|J| \leq \varepsilon m} \left\|\sum_{j \in J} \mathbf{a}_j \mathbf{a}_j^H\right\| \geq 4\tilde{C}_3 \varepsilon m \log \frac{e}{\varepsilon}\right) \leq 2 \exp\left(-\frac{\tilde{C}_2 \tilde{C}_3}{3} \varepsilon m \log \frac{e}{\varepsilon}\right),$$

where  $J \subseteq [m]$  and  $|J|$  denotes its cardinality.

**Proof** The proof relies on Lemma 57 and the union bound. First, invoke Lemma 57 to see that for any fixed  $J \subseteq [m]$  and for all  $t \geq 0$ , we have

$$\mathbb{P}\left(\left\|\sum_{j \in J} (\mathbf{a}_j \mathbf{a}_j^H - \mathbf{I}_K)\right\| \geq |J| t\right) \leq 2 \exp\left(\tilde{C}_1 K - \tilde{C}_2 |J| \min\{t, t^2\}\right), \tag{294}$$

for some constants  $\tilde{C}_1, \tilde{C}_2 > 0$ , and as a result,

$$\begin{aligned} &\mathbb{P}\left(\sup_{|J| \leq \varepsilon m} \left\|\sum_{j \in J} \mathbf{a}_j \mathbf{a}_j^H\right\| \geq \lceil \varepsilon m \rceil (1 + t)\right) \\ &\stackrel{(i)}{\leq} \mathbb{P}\left(\sup_{|J| = \lceil \varepsilon m \rceil} \left\|\sum_{j \in J} \mathbf{a}_j \mathbf{a}_j^H\right\| \geq \lceil \varepsilon m \rceil (1 + t)\right) \end{aligned}$$

$$\begin{aligned} &\leq \mathbb{P} \left( \sup_{|J|=\lceil \varepsilon m \rceil} \left\| \sum_{j \in J} (\mathbf{a}_j \mathbf{a}_j^H - \mathbf{I}_K) \right\| \geq \lceil \varepsilon m \rceil t \right) \\ &\stackrel{\text{(ii)}}{\leq} \binom{m}{\lceil \varepsilon m \rceil} \cdot 2 \exp \left( \tilde{C}_1 K - \tilde{C}_2 \lceil \varepsilon m \rceil \min \{t, t^2\} \right), \end{aligned}$$

where  $\lceil c \rceil$  denotes the smallest integer that is no smaller than  $c$ . Here, (i) holds since we take the supremum over a larger set and (ii) results from (294) and the union bound. Apply the elementary inequality  $\binom{n}{k} \leq (en/k)^k$  for any  $0 \leq k \leq n$  to obtain

$$\begin{aligned} &\mathbb{P} \left( \sup_{|J| \leq \varepsilon m} \left\| \sum_{j \in J} \mathbf{a}_j \mathbf{a}_j^H \right\| \geq \lceil \varepsilon m \rceil (1+t) \right) \\ &\leq 2 \left( \frac{\varepsilon m}{\lceil \varepsilon m \rceil} \right)^{\lceil \varepsilon m \rceil} \exp \left( \tilde{C}_1 K - \tilde{C}_2 \lceil \varepsilon m \rceil \min \{t, t^2\} \right) \\ &\leq 2 \left( \frac{e}{\varepsilon} \right)^{2\varepsilon m} \exp \left( \tilde{C}_1 K - \tilde{C}_2 \varepsilon m \min \{t, t^2\} \right) \\ &= 2 \exp \left[ \tilde{C}_1 K - \varepsilon m \left( \tilde{C}_2 \min \{t, t^2\} - 2 \log(e/\varepsilon) \right) \right], \end{aligned} \tag{295}$$

where the second inequality uses  $\varepsilon m \leq \lceil \varepsilon m \rceil \leq 2\varepsilon m$  whenever  $1/m \leq \varepsilon \leq 1$ .

The proof is then completed by taking  $\tilde{C}_3 \geq \max\{1, 6/\tilde{C}_2\}$  and  $t = \tilde{C}_3 \log(e/\varepsilon)$ . To see this, it is easy to check that  $\min\{t, t^2\} = t$  since  $t \geq 1$ . In addition, one has  $\tilde{C}_1 K \leq \tilde{C}_2 \varepsilon m / 3 \leq \tilde{C}_2 \varepsilon m t / 3$ , and  $2 \log(e/\varepsilon) \leq \tilde{C}_2 t / 3$ . Combine the estimates above with (295) to arrive at

$$\begin{aligned} &\mathbb{P} \left( \sup_{|J| \leq \varepsilon m} \left\| \sum_{j \in J} \mathbf{a}_j \mathbf{a}_j^H \right\| \geq 4\tilde{C}_3 \varepsilon m \log(e/\varepsilon) \right) \\ &\stackrel{\text{(i)}}{\leq} \mathbb{P} \left( \sup_{|J| \leq \varepsilon m} \left\| \sum_{j \in J} \mathbf{a}_j \mathbf{a}_j^H \right\| \geq \lceil \varepsilon m \rceil (1+t) \right) \\ &\leq 2 \exp \left[ \tilde{C}_1 K - \varepsilon m \left( \tilde{C}_2 \min \{t, t^2\} - 2 \log(e/\varepsilon) \right) \right] \\ &\stackrel{\text{(ii)}}{\leq} 2 \exp \left( -\varepsilon m \tilde{C}_2 t / 3 \right) = 2 \exp \left( -\frac{\tilde{C}_2 \tilde{C}_3}{3} \varepsilon m \log(e/\varepsilon) \right) \end{aligned}$$

as claimed. Here, (i) holds due to the facts that  $\lceil \varepsilon m \rceil \leq 2\varepsilon m$  and  $1+t \leq 2t \leq 2\tilde{C}_3 \log(e/\varepsilon)$ . Inequality (ii) arises from the estimates listed above.  $\square$

**Lemma 59** *Suppose  $m \gg K \log^3 m$ . With probability exceeding  $1 - O(m^{-10})$ , we have*

$$\left\| \sum_{j=1}^m \left| \mathbf{a}_j^H \mathbf{x}^* \right|^2 \mathbf{b}_j \mathbf{b}_j^H - \mathbf{I}_K \right\| \lesssim \sqrt{\frac{K}{m} \log m}.$$

**Proof** The identity  $\sum_{j=1}^m \mathbf{b}_j \mathbf{b}_j^H = \mathbf{I}_K$  allows us to rewrite the quantity on the left-hand side as

$$\left\| \sum_{j=1}^m \left| \mathbf{a}_j^H \mathbf{x}^\star \right|^2 \mathbf{b}_j \mathbf{b}_j^H - \mathbf{I}_K \right\| = \left\| \sum_{j=1}^m \underbrace{\left( \left| \mathbf{a}_j^H \mathbf{x}^\star \right|^2 - 1 \right)}_{:= \mathbf{Z}_j} \mathbf{b}_j \mathbf{b}_j^H \right\|,$$

where the  $\mathbf{Z}_j$ 's are independent zero-mean random matrices. To control the above spectral norm, we resort to the matrix Bernstein inequality [66, Theorem 2.7]. To this end, we first need to upper bound the sub-exponential norm  $\|\cdot\|_{\psi_1}$  (see definition in [116]) of each summand  $\mathbf{Z}_j$ , i.e.,

$$\|\|\mathbf{Z}_j\|\|_{\psi_1} = \|\mathbf{b}_j\|_2^2 \left\| \left| \mathbf{a}_j^H \mathbf{x}^\star \right|^2 - 1 \right\|_{\psi_1} \lesssim \|\mathbf{b}_j\|_2^2 \left\| \left| \mathbf{a}_j^H \mathbf{x}^\star \right|^2 \right\|_{\psi_1} \lesssim \frac{K}{m},$$

where we make use of the facts that

$$\|\mathbf{b}_j\|_2^2 = K/m \quad \text{and} \quad \left\| \left| \mathbf{a}_j^H \mathbf{x}^\star \right|^2 \right\|_{\psi_1} \lesssim 1.$$

We further need to bound the variance parameter, that is,

$$\begin{aligned} \sigma_0^2 &:= \left\| \mathbb{E} \left[ \sum_{j=1}^m \mathbf{Z}_j \mathbf{Z}_j^H \right] \right\| = \left\| \mathbb{E} \left[ \sum_{j=1}^m \left( \left| \mathbf{a}_j^H \mathbf{x}^\star \right|^2 - 1 \right)^2 \mathbf{b}_j \mathbf{b}_j^H \mathbf{b}_j \mathbf{b}_j^H \right] \right\| \\ &\lesssim \left\| \sum_{j=1}^m \mathbf{b}_j \mathbf{b}_j^H \mathbf{b}_j \mathbf{b}_j^H \right\| = \frac{K}{m} \left\| \sum_{j=1}^m \mathbf{b}_j \mathbf{b}_j^H \right\| = \frac{K}{m}, \end{aligned}$$

where the second line arises since  $\mathbb{E}[(|\mathbf{a}_j^H \mathbf{x}^\star|^2 - 1)^2] \asymp 1$ ,  $\|\mathbf{b}_j\|_2^2 = K/m$ , and  $\sum_{j=1}^m \mathbf{b}_j \mathbf{b}_j^H = \mathbf{I}_K$ . A direct application of the matrix Bernstein inequality [66, Theorem 2.7] leads us to conclude that with probability exceeding  $1 - O(m^{-10})$ ,

$$\left\| \sum_{j=1}^m \mathbf{Z}_j \right\| \lesssim \max \left\{ \sqrt{\frac{K}{m} \log m}, \frac{K}{m} \log^2 m \right\} \asymp \sqrt{\frac{K}{m} \log m},$$

where the last relation holds under the assumption that  $m \gg K \log^3 m$ . □

### D.3.5 Matrix Perturbation Bounds

We also need the following perturbation bound on the top singular vectors of a given matrix. The following lemma is parallel to Lemma 34.

**Lemma 60** Let  $\sigma_1(\mathbf{A})$ ,  $\mathbf{u}$ , and  $\mathbf{v}$  be the leading singular value, left and right singular vectors of  $\mathbf{A}$ , respectively, and let  $\sigma_1(\tilde{\mathbf{A}})$ ,  $\tilde{\mathbf{u}}$ , and  $\tilde{\mathbf{v}}$  be the leading singular value, left and right singular vectors of  $\tilde{\mathbf{A}}$ , respectively. Suppose  $\sigma_1(\mathbf{A})$  and  $\sigma_1(\tilde{\mathbf{A}})$  are not identically zero, and then one has

$$\begin{aligned} |\sigma_1(\mathbf{A}) - \sigma_1(\tilde{\mathbf{A}})| &\leq \|(\mathbf{A} - \tilde{\mathbf{A}})\mathbf{v}\|_2 + (\|\mathbf{u} - \tilde{\mathbf{u}}\|_2 + \|\mathbf{v} - \tilde{\mathbf{v}}\|_2) \|\tilde{\mathbf{A}}\|; \\ \left\| \sqrt{\sigma_1(\mathbf{A})} \mathbf{u} - \sqrt{\sigma_1(\tilde{\mathbf{A}})} \tilde{\mathbf{u}} \right\|_2 &+ \left\| \sqrt{\sigma_1(\mathbf{A})} \mathbf{v} - \sqrt{\sigma_1(\tilde{\mathbf{A}})} \tilde{\mathbf{v}} \right\|_2 \\ &\leq \sqrt{\sigma_1(\mathbf{A})} (\|\mathbf{u} - \tilde{\mathbf{u}}\|_2 + \|\mathbf{v} - \tilde{\mathbf{v}}\|_2) + \frac{2|\sigma_1(\mathbf{A}) - \sigma_1(\tilde{\mathbf{A}})|}{\sqrt{\sigma_1(\mathbf{A})} + \sqrt{\sigma_1(\tilde{\mathbf{A}})}}. \end{aligned}$$

**Proof** The first claim follows since

$$\begin{aligned} |\sigma_1(\mathbf{A}) - \sigma_1(\tilde{\mathbf{A}})| &= \left| \mathbf{u}^H \mathbf{A} \mathbf{v} - \tilde{\mathbf{u}}^H \tilde{\mathbf{A}} \tilde{\mathbf{v}} \right| \\ &\leq \left| \mathbf{u}^H (\mathbf{A} - \tilde{\mathbf{A}}) \mathbf{v} \right| + \left| \mathbf{u}^H \tilde{\mathbf{A}} \mathbf{v} - \tilde{\mathbf{u}}^H \tilde{\mathbf{A}} \tilde{\mathbf{v}} \right| + \left| \tilde{\mathbf{u}}^H \tilde{\mathbf{A}} \mathbf{v} - \tilde{\mathbf{u}}^H \tilde{\mathbf{A}} \tilde{\mathbf{v}} \right| \\ &\leq \|(\mathbf{A} - \tilde{\mathbf{A}})\mathbf{v}\|_2 + \|\mathbf{u} - \tilde{\mathbf{u}}\|_2 \|\tilde{\mathbf{A}}\| + \|\tilde{\mathbf{A}}\| \|\mathbf{v} - \tilde{\mathbf{v}}\|_2. \end{aligned}$$

With regard to the second claim, we see that

$$\begin{aligned} \left\| \sqrt{\sigma_1(\mathbf{A})} \mathbf{u} - \sqrt{\sigma_1(\tilde{\mathbf{A}})} \tilde{\mathbf{u}} \right\|_2 &\leq \left\| \sqrt{\sigma_1(\mathbf{A})} \mathbf{u} - \sqrt{\sigma_1(\mathbf{A})} \tilde{\mathbf{u}} \right\|_2 + \left\| \sqrt{\sigma_1(\mathbf{A})} \tilde{\mathbf{u}} - \sqrt{\sigma_1(\tilde{\mathbf{A}})} \tilde{\mathbf{u}} \right\|_2 \\ &= \sqrt{\sigma_1(\mathbf{A})} \|\mathbf{u} - \tilde{\mathbf{u}}\|_2 + \left| \sqrt{\sigma_1(\mathbf{A})} - \sqrt{\sigma_1(\tilde{\mathbf{A}})} \right| \|\tilde{\mathbf{u}}\|_2 \\ &= \sqrt{\sigma_1(\mathbf{A})} \|\mathbf{u} - \tilde{\mathbf{u}}\|_2 + \frac{|\sigma_1(\mathbf{A}) - \sigma_1(\tilde{\mathbf{A}})|}{\sqrt{\sigma_1(\mathbf{A})} + \sqrt{\sigma_1(\tilde{\mathbf{A}})}}. \end{aligned}$$

Similarly, one can obtain

$$\left\| \sqrt{\sigma_1(\mathbf{A})} \mathbf{v} - \sqrt{\sigma_1(\tilde{\mathbf{A}})} \tilde{\mathbf{v}} \right\|_2 \leq \sqrt{\sigma_1(\mathbf{A})} \|\mathbf{v} - \tilde{\mathbf{v}}\|_2 + \frac{|\sigma_1(\mathbf{A}) - \sigma_1(\tilde{\mathbf{A}})|}{\sqrt{\sigma_1(\mathbf{A})} + \sqrt{\sigma_1(\tilde{\mathbf{A}})}}.$$

Add these two inequalities to complete the proof.  $\square$

## References

1. Abbe, E., Fan, J., Wang, K., Zhong, Y.: Entrywise eigenvector analysis of random matrices with low expected rank. arXiv preprint [arXiv:1709.09565](https://arxiv.org/abs/1709.09565) (2017)
2. Aghasi, A., Ahmed, A., Hand, P., Joshi, B.: Branchhull: Convex bilinear inversion from the entrywise product of signals with known signs. *Applied and Computational Harmonic Analysis* (2019)
3. Ahmed, A., Recht, B., Romberg, J.: Blind deconvolution using convex programming. *IEEE Transactions on Information Theory* **60**(3), 1711–1732 (2014)
4. Alon, N., Spencer, J.H.: *The Probabilistic Method* (3rd Edition). Wiley (2008)
5. Bahmani, S., Romberg, J.: Phase retrieval meets statistical learning theory: A flexible convex relaxation. In: *Artificial Intelligence and Statistics*, pp. 252–260 (2017)
6. Bendory, T., Eldar, Y.C., Boumal, N.: Non-convex phase retrieval from STFT measurements. *IEEE Transactions on Information Theory* (2017)



7. Bhojanapalli, S., Neyshabur, B., Srebro, N.: Global optimality of local search for low rank matrix recovery. In: *Advances in Neural Information Processing Systems*, pp. 3873–3881 (2016)
8. Bousquet, O., Elisseeff, A.: Stability and generalization. *Journal of Machine Learning Research* **2**(Mar), 499–526 (2002)
9. Bubeck, S.: Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning* **8**(3–4), 231–357 (2015)
10. Cai, J.F., Liu, H., Wang, Y.: Fast rank-one alternating minimization algorithm for phase retrieval. *Journal of Scientific Computing* **79**(1), 128–147 (2019)
11. Cai, T., Zhang, A.: ROP: Matrix recovery via rank-one projections. *The Annals of Statistics* **43**(1), 102–138 (2015)
12. Cai, T.T., Li, X., Ma, Z.: Optimal rates of convergence for noisy sparse phase retrieval via thresholded Wirtinger flow. *The Annals of Statistics* **44**(5), 2221–2251 (2016)
13. Candès, E., Plan, Y.: A probabilistic and RIPless theory of compressed sensing. *IEEE Transactions on Information Theory* **57**(11), 7235–7254 (2011). <https://doi.org/10.1109/TIT.2011.2161794>
14. Candès, E., Tao, T.: The power of convex relaxation: Near-optimal matrix completion. *IEEE Transactions on Information Theory* **56**(5), 2053–2080 (2010)
15. Candès, E.J., Eldar, Y.C., Strohmer, T., Voroninski, V.: Phase retrieval via matrix completion. *SIAM Journal on Imaging Sciences* **6**(1), 199–225 (2013)
16. Candès, E.J., Li, X.: Solving quadratic equations via PhaseLift when there are about as many equations as unknowns. *Foundations of Computational Mathematics* **14**(5), 1017–1026 (2014)
17. Candès, E.J., Li, X., Ma, Y., Wright, J.: Robust principal component analysis? *Journal of ACM* **58**(3), 11:1–11:37 (2011)
18. Candès, E.J., Li, X., Soltanolkotabi, M.: Phase retrieval via Wirtinger flow: Theory and algorithms. *IEEE Transactions on Information Theory* **61**(4), 1985–2007 (2015)
19. Candès, E.J., Recht, B.: Exact matrix completion via convex optimization. *Foundations of Computational Mathematics* **9**(6), 717–772 (2009)
20. Candès, E.J., Strohmer, T., Voroninski, V.: Phaselift: Exact and stable signal recovery from magnitude measurements via convex programming. *Communications on Pure and Applied Mathematics* **66**(8), 1017–1026 (2013)
21. Chandrasekaran, V., Sanghavi, S., Parrilo, P.A., Willsky, A.S.: Rank-sparsity incoherence for matrix decomposition. *SIAM Journal on Optimization* **21**(2), 572–596 (2011)
22. Chen, P., Fannjiang, A., Liu, G.R.: Phase retrieval with one or two diffraction patterns by alternating projections with the null initialization. *Journal of Fourier Analysis and Applications*, pp. 1–40 (2015)
23. Chen, Y.: Incoherence-optimal matrix completion. *IEEE Transactions on Information Theory* **61**(5), 2909–2923 (2015)
24. Chen, Y., Candès, E.: The projected power method: An efficient algorithm for joint alignment from pairwise differences. *Communications on Pure and Applied Mathematics* **71**(8), 1648–1714 (2018)
25. Chen, Y., Candès, E.J.: Solving random quadratic systems of equations is nearly as easy as solving linear systems. *Communications on Pure and Applied Mathematics* **70**(5), 822–883 (2017). <https://doi.org/10.1002/cpa.21638>
26. Chen, Y., Cheng, C., Fan, J.: Asymmetry helps: Eigenvalue and eigenvector analyses of asymmetrically perturbed low-rank matrices. *arXiv preprint arXiv:1811.12804* (2018)
27. Chen, Y., Chi, Y., Fan, J., Ma, C.: Gradient descent with random initialization: Fast global convergence for nonconvex phase retrieval. *Mathematical Programming* **176**(1–2), 5–37 (2019)
28. Chen, Y., Chi, Y., Fan, J., Ma, C., Yan, Y.: Noisy matrix completion: Understanding statistical guarantees for convex relaxation via nonconvex optimization. *arXiv preprint arXiv:1902.07698* (2019)
29. Chen, Y., Chi, Y., Goldsmith, A.J.: Exact and stable covariance estimation from quadratic sampling via convex programming. *IEEE Transactions on Information Theory* **61**(7), 4034–4059 (2015)
30. Chen, Y., Fan, J., Ma, C., Wang, K.: Spectral method and regularized MLE are both optimal for top- $K$  ranking. *Annals of Statistics* **47**(4), 2204–2235 (2019)
31. Chen, Y., Fan, J., Ma, C., Yan, Y.: Inference and uncertainty quantification for noisy matrix completion. *arXiv preprint arXiv:1906.04159* (2019)
32. Chen, Y., Wainwright, M.J.: Fast low-rank estimation by projected gradient descent: General statistical and algorithmic guarantees. *arXiv preprint arXiv:1509.03025* (2015)
33. Chen, Y., Yi, X., Caramanis, C.: A convex formulation for mixed regression with two components: Minimax optimal rates. In: *Conference on Learning Theory*, pp. 560–604 (2014)

34. Cherapanamjeri, Y., Jain, P., Netrapalli, P.: Thresholding based outlier robust PCA. In: Conference on Learning Theory, pp. 593–628 (2017)
35. Chi, Y.: Guaranteed blind sparse spikes deconvolution via lifting and convex optimization. *IEEE Journal of Selected Topics in Signal Processing* **10**(4), 782–794 (2016)
36. Chi, Y., Lu, Y.M.: Kaczmarz method for solving quadratic equations. *IEEE Signal Processing Letters* **23**(9), 1183–1187 (2016)
37. Chi, Y., Lu, Y.M., Chen, Y.: Nonconvex optimization meets low-rank matrix factorization: An overview. arXiv preprint [arXiv:1809.09573](https://arxiv.org/abs/1809.09573) (2018)
38. Davenport, M.A., Romberg, J.: An overview of low-rank matrix recovery from incomplete observations. *IEEE Journal of Selected Topics in Signal Processing* **10**(4), 608–622 (2016)
39. Davis, C., Kahan, W.M.: The rotation of eigenvectors by a perturbation. iii. *SIAM Journal on Numerical Analysis* **7**(1), 1–46 (1970)
40. Davis, D., Drusvyatskiy, D., Paquette, C.: The nonsmooth landscape of phase retrieval. arXiv preprint [arXiv:1711.03247](https://arxiv.org/abs/1711.03247) (2017)
41. Dhifallah, O., Thrampoulidis, C., Lu, Y.M.: Phase retrieval via linear programming: Fundamental limits and algorithmic improvements. arXiv preprint [arXiv:1710.05234](https://arxiv.org/abs/1710.05234) (2017)
42. Dopico, F.M.: A note on  $\sin \Theta$  theorems for singular subspace variations. *BIT* **40**(2), 395–403 (2000). <https://doi.org/10.1023/A:1022303426500>.
43. Duchi, J.C., Ruan, F.: Solving (most) of a set of quadratic equalities: Composite optimization for robust phase retrieval. *Information and Inference* (2018)
44. El Karoui, N.: On the impact of predictor geometry on the performance on high-dimensional ridge-regularized generalized robust regression estimators. *Probability Theory and Related Fields* pp. 1–81 (2015)
45. El Karoui, N., Bean, D., Bickel, P.J., Lim, C., Yu, B.: On robust regression with high-dimensional predictors. *Proceedings of the National Academy of Sciences* **110**(36), 14557–14562 (2013)
46. Fan, J., Ma, C., Zhong, Y.: A selective overview of deep learning. arXiv preprint [arXiv:1904.05526](https://arxiv.org/abs/1904.05526) (2019)
47. Gao, B., Xu, Z.: Phase retrieval using Gauss-Newton method. arXiv preprint [arXiv:1606.08135](https://arxiv.org/abs/1606.08135) (2016)
48. Ge, R., Lee, J.D., Ma, T.: Matrix completion has no spurious local minimum. In: *Advances in Neural Information Processing Systems*, pp. 2973–2981 (2016)
49. Ge, R., Ma, T.: On the optimization landscape of tensor decompositions. In: *Advances in Neural Information Processing Systems*, pp. 3653–3663 (2017)
50. Goldstein, T., Studer, C.: Phasemax: Convex phase retrieval via basis pursuit. *IEEE Transactions on Information Theory* **64**(4), 2675–2689 (2018)
51. Gross, D.: Recovering low-rank matrices from few coefficients in any basis. *IEEE Transactions on Information Theory* **57**(3), 1548–1566 (2011)
52. Gunasekar, S., Woodworth, B.E., Bhojanapalli, S., Neyshabur, B., Srebro, N.: Implicit regularization in matrix factorization. In: *Advances in Neural Information Processing Systems*, pp. 6151–6159 (2017)
53. Hand, P., Voroninski, V.: An elementary proof of convex phase retrieval in the natural parameter space via the linear program phasemax. *Communications in Mathematical Sciences* **16**(7), 2047–2051 (2018)
54. Hardt, M., Wootters, M.: Fast matrix completion without the condition number. *Conference on Learning Theory*, pp. 638–678 (2014)
55. Hastie, T., Mazumder, R., Lee, J.D., Zadeh, R.: Matrix completion and low-rank SVD via fast alternating least squares. *Journal of Machine Learning Research* **16**, 3367–3402 (2015)
56. Higham, N.J.: Estimating the matrix  $p$ -norm. *Numerische Mathematik* **62**(1), 539–555 (1992)
57. Hsu, D., Kakade, S.M., Zhang, T.: A tail inequality for quadratic forms of subgaussian random vectors. *Electron. Commun. Probab.* **17**, no. 52, 6 (2012). <https://doi.org/10.1214/ECP.v17-2079>.
58. Huang, W., Hand, P.: Blind deconvolution by a steepest descent algorithm on a quotient manifold. *SIAM Journal on Imaging Sciences* **11**(4), 2757–2785 (2018)
59. Jaganathan, K., Eldar, Y.C., Hassibi, B.: Phase retrieval: An overview of recent developments. arXiv preprint [arXiv:1510.07713](https://arxiv.org/abs/1510.07713) (2015)
60. Jain, P., Netrapalli, P.: Fast exact matrix completion with finite samples. In: *Conference on Learning Theory*, pp. 1007–1034 (2015)
61. Jain, P., Netrapalli, P., Sanghavi, S.: Low-rank matrix completion using alternating minimization. In: *ACM symposium on Theory of computing*, pp. 665–674 (2013)

62. Javanmard, A., Montanari, A., et al.: Debiasing the lasso: Optimal sample size for gaussian designs. *The Annals of Statistics* **46**(6A), 2593–2622 (2018)
63. Jin, C., Kakade, S.M., Netrapalli, P.: Provable efficient online matrix completion via non-convex stochastic gradient descent. In: *Advances in Neural Information Processing Systems*, pp. 4520–4528 (2016)
64. Keshavan, R.H., Montanari, A., Oh, S.: Matrix completion from a few entries. *IEEE Transactions on Information Theory* **56**(6), 2980–2998 (2010)
65. Keshavan, R.H., Montanari, A., Oh, S.: Matrix completion from noisy entries. *J. Mach. Learn. Res.* **11**, 2057–2078 (2010)
66. Koltchinskii, V.: Oracle inequalities in empirical risk minimization and sparse recovery problems, *Lecture Notes in Mathematics*, vol. 2033. Springer, Heidelberg (2011). <https://doi.org/10.1007/978-3-642-22147-7>
67. Koltchinskii, V., Lounici, K., Tsybakov, A.B.: Nuclear-norm penalization and optimal rates for noisy low-rank matrix completion. *Ann. Statist.* **39**(5), 2302–2329 (2011). <https://doi.org/10.1214/11-AOS894>.
68. Kolte, R., Özgür, A.: Phase retrieval via incremental truncated Wirtinger flow. arXiv preprint [arXiv:1606.03196](https://arxiv.org/abs/1606.03196) (2016)
69. Kreuz-Delgado, K.: The complex gradient operator and the CR-calculus. arXiv preprint [arXiv:0906.4835](https://arxiv.org/abs/0906.4835) (2009)
70. Lang, S.: Real and functional analysis. Springer-Verlag, New York, **10**, 11–13 (1993)
71. Lee, K., Bresler, Y.: Admira: Atomic decomposition for minimum rank approximation. *IEEE Transactions on Information Theory* **56**(9), 4402–4416 (2010)
72. Lee, K., Li, Y., Junge, M., Bresler, Y.: Blind recovery of sparse signals from subsampled convolution. *IEEE Transactions on Information Theory* **63**(2), 802–821 (2017)
73. Lee, K., Tian, N., Romberg, J.: Fast and guaranteed blind multichannel deconvolution under a bilinear system model. *IEEE Transactions on Information Theory* **64**(7), 4792–4818 (2018)
74. Lerman, G., Maunu, T.: Fast, robust and non-convex subspace recovery. *Information and Inference: A Journal of the IMA* **7**(2), 277–336 (2017)
75. Li, Q., Tang, G.: The nonconvex geometry of low-rank matrix optimizations with general objective functions. In: *2017 IEEE Global Conference on Signal and Information Processing (GlobalSIP)*, pp. 1235–1239. IEEE (2017)
76. Li, X., Ling, S., Strohmer, T., Wei, K.: Rapid, robust, and reliable blind deconvolution via nonconvex optimization. *Applied and computational harmonic analysis* (2018)
77. Li, X., Wang, Z., Lu, J., Arora, R., Haupt, J., Liu, H., Zhao, T.: Symmetry, saddle points, and global geometry of nonconvex matrix factorization. arXiv preprint [arXiv:1612.09296](https://arxiv.org/abs/1612.09296) (2016)
78. Li, Y., Lee, K., Bresler, Y.: Blind gain and phase calibration for low-dimensional or sparse signal sensing via power iteration. In: *Sampling Theory and Applications (SampTA), 2017 International Conference on*, pp. 119–123. IEEE (2017)
79. Li, Y., Ma, C., Chen, Y., Chi, Y.: Nonconvex matrix factorization from rank-one measurements. arXiv preprint [arXiv:1802.06286](https://arxiv.org/abs/1802.06286) (2018)
80. Lin, J., Camoriano, R., Rosasco, L.: Generalization properties and implicit regularization for multiple passes SGM. In: *International Conference on Machine Learning*, pp. 2340–2348 (2016)
81. Ling, S., Strohmer, T.: Self-calibration and biconvex compressive sensing. *Inverse Problems* **31**(11), 115002 (2015)
82. Ling, S., Strohmer, T.: Regularized gradient descent: a non-convex recipe for fast joint blind deconvolution and demixing. *Information and Inference: A Journal of the IMA* **8**(1), 1–49 (2018)
83. Lu, Y.M., Li, G.: Phase transitions of spectral initialization for high-dimensional nonconvex estimation. arXiv preprint [arXiv:1702.06435](https://arxiv.org/abs/1702.06435) (2017)
84. Mathias, R.: The spectral norm of a nonnegative matrix. *Linear Algebra Appl.* **139**, 269–284 (1990). [https://doi.org/10.1016/0024-3795\(90\)90403-Y](https://doi.org/10.1016/0024-3795(90)90403-Y).
85. Mathias, R.: Perturbation bounds for the polar decomposition. *SIAM Journal on Matrix Analysis and Applications* **14**(2), 588–597 (1993)
86. Maunu, T., Zhang, T., Lerman, G.: A well-tempered landscape for non-convex robust subspace recovery. *Journal of Machine Learning Research* **20**(37), 1–59 (2019)
87. Mei, S., Bai, Y., Montanari, A.: The landscape of empirical risk for nonconvex losses. *The Annals of Statistics* **46**(6A), 2747–2774 (2018)

88. Mondelli, M., Montanari, A.: Fundamental limits of weak recovery with applications to phase retrieval. *Foundations of Computational Mathematics*, pp. 1–71 (2017)
89. Negahban, S., Wainwright, M.J.: Restricted strong convexity and weighted matrix completion: optimal bounds with noise. *J. Mach. Learn. Res.* **13**, 1665–1697 (2012)
90. Netrapalli, P., Jain, P., Sanghavi, S.: Phase retrieval using alternating minimization. *Advances in Neural Information Processing Systems (NIPS)* (2013)
91. Netrapalli, P., Niranjan, U., Sanghavi, S., Anandkumar, A., Jain, P.: Non-convex robust PCA. In: *Advances in Neural Information Processing Systems*, pp. 1107–1115 (2014)
92. Qing, Q., Zhang, Y., Eldar, Y., Wright, J.: Convolutional phase retrieval via gradient descent. *Neural Information Processing Systems* (2017)
93. Recht, B.: A simpler approach to matrix completion. *Journal of Machine Learning Research* **12**(Dec), 3413–3430 (2011)
94. Recht, B., Fazel, M., Parrilo, P.A.: Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM Review* **52**(3), 471–501 (2010)
95. Rudelson, M., Vershynin, R., et al.: Hanson-Wright inequality and sub-Gaussian concentration. *Electronic Communications in Probability* **18** (2013)
96. Sanghavi, S., Ward, R., White, C.D.: The local convexity of solving systems of quadratic equations. *Results in Mathematics* **71**(3–4), 569–608 (2017)
97. Schmitt, B.A.: Perturbation bounds for matrix square roots and Pythagorean sums. *Linear Algebra Appl.* **174**, 215–227 (1992). [https://doi.org/10.1016/0024-3795\(92\)90052-C](https://doi.org/10.1016/0024-3795(92)90052-C).
98. Schniter, P., Rangan, S.: Compressive phase retrieval via generalized approximate message passing. *IEEE Transactions on Signal Processing* **63**(4), 1043–1055 (2015)
99. Schudy, W., Sviridenko, M.: Concentration and moment inequalities for polynomials of independent random variables. In: *Symposium on Discrete Algorithms*, pp. 437–446. ACM, New York (2012)
100. Shechtman, Y., Beck, A., Eldar, Y.C.: GESPAR: Efficient phase retrieval of sparse signals. *IEEE Transactions on Signal Processing* **62**(4), 928–938 (2014)
101. Soltanolkotabi, M.: Algorithms and theory for clustering and nonconvex quadratic programming. Ph.D. thesis, Stanford University (2014)
102. Soltanolkotabi, M.: Structured signal recovery from quadratic measurements: Breaking sample complexity barriers via nonconvex optimization. *IEEE Transactions on Information Theory* **65**(4), 2374–2400 (2019)
103. Soltanolkotabi, M., Javanmard, A., Lee, J.D.: Theoretical insights into the optimization landscape of over-parameterized shallow neural networks. *IEEE Transactions on Information Theory* **65**(2), 742–769 (2019)
104. Soudry, D., Hoffer, E., Nacson, M.S., Gunasekar, S., Srebro, N.: The implicit bias of gradient descent on separable data. *The Journal of Machine Learning Research* **19**(1), 2822–2878 (2018)
105. Sun, J., Qu, Q., Wright, J.: A geometric analysis of phase retrieval. In: *Information Theory (ISIT), 2016 IEEE International Symposium on*, pp. 2379–2383. IEEE (2016)
106. Sun, J., Qu, Q., Wright, J.: Complete dictionary recovery over the sphere i: Overview and the geometric picture. *IEEE Transactions on Information Theory* **63**(2), 853–884 (2017)
107. Sun, R., Luo, Z.Q.: Guaranteed matrix completion via non-convex factorization. *IEEE Transactions on Information Theory* **62**(11), 6535–6579 (2016)
108. Sur, P., Chen, Y., Candès, E.J.: The likelihood ratio test in high-dimensional logistic regression is asymptotically a rescaled chi-square. *arXiv preprint arXiv:1706.01191*, accepted to *Probability Theory and Related Fields* (2017)
109. Tan, Y.S., Vershynin, R.: Phase retrieval via randomized kaczmarz: Theoretical guarantees. *Information and Inference: A Journal of the IMA* **8**(1), 97–123 (2018)
110. Tanner, J., Wei, K.: Low rank matrix completion by alternating steepest descent methods. *Applied and Computational Harmonic Analysis* **40**(2), 417–429 (2016)
111. Tao, T.: *Topics in Random Matrix Theory*. Graduate Studies in Mathematics. American Mathematical Society, Providence, Rhode Island (2012)
112. Ten Berge, J.M.: Orthogonal procrustes rotation for two or more matrices. *Psychometrika* **42**(2), 267–276 (1977)
113. Tropp, J.A.: Convex recovery of a structured signal from independent random linear measurements. In: *Sampling Theory, a Renaissance*, pp. 67–101. Springer (2015)
114. Tropp, J.A.: An introduction to matrix concentration inequalities. *Found. Trends Mach. Learn.* **8**(1–2), 1–230 (2015). <https://doi.org/10.1561/22000000048>.

115. Tu, S., Boczar, R., Simchowitz, M., Soltanolkotabi, M., Recht, B.: Low-rank solutions of linear matrix equations via procrustes flow. In: International Conference on Machine Learning, pp. 964–973. JMLR. org (2016)
116. Vershynin, R.: Introduction to the non-asymptotic analysis of random matrices. Compressed Sensing, Theory and Applications, pp. 210–268 (2012)
117. Wang, G., Giannakis, G., Saad, Y., Chen, J.: Solving most systems of random quadratic equations. In: Advances in Neural Information Processing Systems, pp. 1867–1877 (2017)
118. Wang, G., Giannakis, G.B., Eldar, Y.C.: Solving systems of random quadratic equations via truncated amplitude flow. IEEE Transactions on Information Theory (2017)
119. Wang, G., Zhang, L., Giannakis, G.B., Akçakaya, M., Chen, J.: Sparse phase retrieval via truncated amplitude flow. IEEE Transactions on Signal Processing **66**(2), 479–491 (2018)
120. Wang, L., Chi, Y.: Blind deconvolution from multiple sparse inputs. IEEE Signal Processing Letters **23**(10), 1384–1388 (2016)
121. Wedin, P.Å.: Perturbation bounds in connection with singular value decomposition. BIT Numerical Mathematics **12**(1), 99–111 (1972)
122. Wei, K.: Solving systems of phaseless equations via Kaczmarz methods: A proof of concept study. Inverse Problems **31**(12), 125008 (2015)
123. Wei, K., Cai, J.F., Chan, T.F., Leung, S.: Guarantees of riemannian optimization for low rank matrix recovery. SIAM Journal on Matrix Analysis and Applications **37**(3), 1198–1222 (2016)
124. Yu, Y., Wang, T., Samworth, R.J.: A useful variant of the Davis-Kahan theorem for statisticians. Biometrika **102**(2), 315–323 (2015). <https://doi.org/10.1093/biomet/asv008>
125. Zhang, C., Bengio, S., Hardt, M., Recht, B., Vinyals, O.: Understanding deep learning requires rethinking generalization. International Conference on Learning Representations (2017)
126. Zhang, H., Chi, Y., Liang, Y.: Provable non-convex phase retrieval with outliers: Median truncated Wirtinger flow. In: International conference on machine learning, pp. 1022–1031 (2016)
127. Zhang, H., Zhou, Y., Liang, Y., Chi, Y.: A nonconvex approach for phase retrieval: Reshaped wirtinger flow and incremental algorithms. Journal of Machine Learning Research (2017)
128. Zhang, Y., Lau, Y., Kuo, H.w., Cheung, S., Pasupathy, A., Wright, J.: On the global geometry of sphere-constrained sparse blind deconvolution. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 4894–4902 (2017)
129. Zhao, T., Wang, Z., Liu, H.: A nonconvex optimization framework for low rank matrix estimation. In: Advances in Neural Information Processing Systems, pp. 559–567 (2015)
130. Zheng, Q., Lafferty, J.: A convergent gradient descent algorithm for rank minimization and semidefinite programming from random linear measurements. In: Advances in Neural Information Processing Systems, pp. 109–117 (2015)
131. Zheng, Q., Lafferty, J.: Convergence analysis for rectangular matrix completion using Burer-Monteiro factorization and gradient descent. arXiv preprint [arXiv:1605.07051](https://arxiv.org/abs/1605.07051) (2016)
132. Zhong, K., Song, Z., Jain, P., Bartlett, P.L., Dhillon, I.S.: Recovery guarantees for one-hidden-layer neural networks. In: International Conference on Machine Learning, pp. 4140–4149. JMLR. org (2017)
133. Zhong, Y., Boumal, N.: Near-optimal bounds for phase synchronization. SIAM Journal on Optimization **28**(2), 989–1016 (2018)

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

## Affiliations

Cong Ma<sup>1</sup> · Kaizheng Wang<sup>1</sup> · Yuejie Chi<sup>2</sup> · Yuxin Chen<sup>3</sup>

✉ Yuxin Chen  
yuxin.chen@princeton.edu

Cong Ma  
congma@princeton.edu

Kaizheng Wang  
kaizheng@princeton.edu

Yuejie Chi  
yuejiechi@cmu.edu

- <sup>1</sup> Department of Operations Research and Financial Engineering, Princeton University, Princeton, NJ 08544, USA
- <sup>2</sup> Department of Electrical and Computer Engineering, Carnegie Mellon University, Pittsburgh, PA 15213, USA
- <sup>3</sup> Department of Electrical Engineering, Princeton University, Princeton, NJ 08544, USA