

Regularized gradient descent: a nonconvex recipe for fast joint blind deconvolution and demixing

Shuyang Ling and Thomas Strohmer*
Department of Mathematics,
University of California at Davis
Davis CA-95616, USA
{syling,strohmer}@math.ucdavis.edu

March 24, 2017

Abstract

We study the question of extracting a sequence of functions $\{\mathbf{f}_i, \mathbf{g}_i\}_{i=1}^s$ from observing only the sum of their convolutions, i.e., from $\mathbf{y} = \sum_{i=1}^s \mathbf{f}_i * \mathbf{g}_i$. While convex optimization techniques are able to solve this joint blind deconvolution-demixing problem provably and robustly under certain conditions, for medium-size or large-size problems we need computationally faster methods without sacrificing the benefits of mathematical rigor that come with convex methods. In this paper we present a non-convex algorithm which guarantees exact recovery under conditions that are competitive with convex optimization methods, with the additional advantage of being computationally much more efficient. Our two-step algorithm converges to the global minimum linearly and is also robust in the presence of additive noise. While the derived performance bounds are suboptimal in terms of the information-theoretic limit, numerical simulations show remarkable performance even if the number of measurements is close to the number of degrees of freedom. We discuss an application of the proposed framework in wireless communications in connection with the Internet-of-Things.

1 Introduction

Blind deconvolution is the task of estimating two unknown functions from their convolution. While it is a highly ill-posed bilinear inverse problem, blind deconvolution is also an extremely important problem in signal processing [1], communications engineering [32], imaging processing [5], audio processing [21], etc. In this paper, we deal with an even more difficult and more general variation of the blind deconvolution problem, in which we have to extract multiple convolved signals mixed together in one observation signal. This joint blind deconvolution-demixing problem arises in a range of applications such as acoustics [21], dictionary learning [2], and wireless communications [32].

We briefly discuss one such application in more detail. Blind deconvolution/demixing problems are expected to play a vital role in the future Internet-of-Things. The Internet-of-Things will connect billions of wireless devices, which is far more than the current wireless systems can technically and economically accommodate. One of the many challenges in the design of the Internet-of-Things will be its ability to manage the massive number of sporadic traffic generating devices which are most of the time inactive, but regularly access the network for minor updates with no human interaction [36]. This means among others that the overhead caused by the exchange of certain types of information between transmitter and receiver, such as channel estimation, assignment of data slots, etc, has to be avoided as much as possible.

Focusing on the underlying mathematical challenges, we consider a *multi-user communication* scenario where many different users/devices communicate with a common base station, as illustrated in Figure 1. Suppose we have s users and each of them sends a signal \mathbf{g}_i through an

*The authors acknowledge support from the NSF via grants DTRA-DMS 1322393 and DMS 1620455.

unknown channel (which differs from user to user) to a common base station,. We assume that the i -th channel, represented by its impulse response \mathbf{f}_i , does not change during the transmission of the signal \mathbf{g}_i . Therefore \mathbf{f}_i acts as convolution operator, i.e., the signal transmitted by the i -th user arriving at the base station becomes $\mathbf{f}_i * \mathbf{g}_i$, where “*” denotes convolution. The

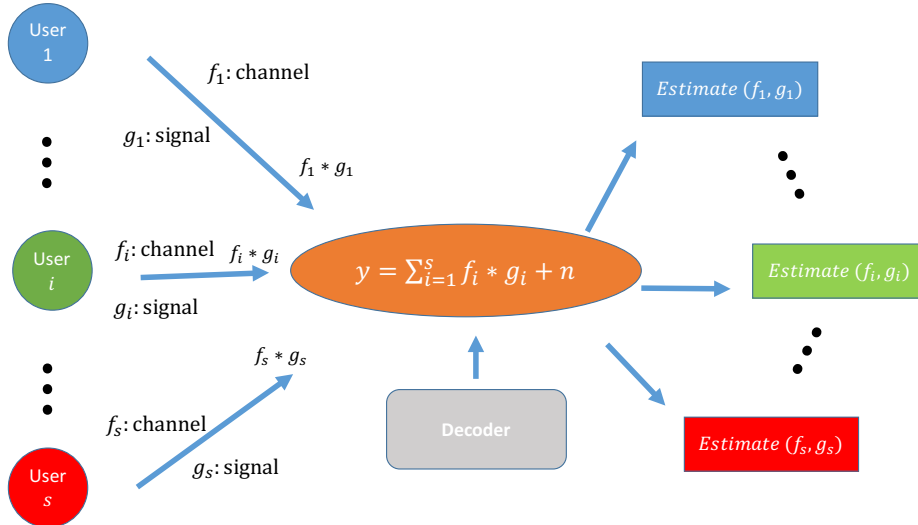


Figure 1: Single-antenna multi-user communication scenario without explicit channel estimation: Each of the s users sends a signal \mathbf{g}_i through an unknown channel \mathbf{f}_i to a common base station. The base station measures the superposition of all those signals, namely, $\mathbf{y} = \sum_{i=1}^s \mathbf{f}_i * \mathbf{g}_i$ (plus noise). The goal is to extract all pairs of $\{(\mathbf{f}_i, \mathbf{g}_i)\}_{i=1}^s$ simultaneously from \mathbf{y} .

antenna at the base station, instead of receiving each individual component $\mathbf{f}_i * \mathbf{g}_i$, is only able to record the superposition of all those signals, namely,

$$\mathbf{y} = \sum_{i=1}^s \mathbf{f}_i * \mathbf{g}_i + \mathbf{n}, \quad (1.1)$$

where \mathbf{n} represents noise. We aim to develop a fast algorithm to simultaneously extract all pairs $\{(\mathbf{f}_i, \mathbf{g}_i)\}_{i=1}^s$ from \mathbf{y} (i.e., estimating the channel/impulse responses \mathbf{f}_i and the signals \mathbf{g}_i jointly) in a numerically efficient and robust way, while keeping the number of required measurements as small as possible.

1.1 State of the art and contributions of this paper

A thorough theoretical analysis concerning the solvability of demixing problems via convex optimization can be found in [23]. There, the authors derive explicit sharp bounds and phase transitions regarding the number of measurements required to successfully demix structured signals (such as sparse signals or low-rank matrices) from a single measurement vector. In principle we could recast the blind deconvolution/demixing problem as the demixing of a sum of rank-one matrices, see (2.3). As such, it seems to fit into the framework analyzed by McCoy and Tropp. However, the setup in [23] differs from ours in a crucial manner. McCoy and Tropp consider as measurement matrices (see the matrices \mathcal{A}_i in (2.3)) full-rank random matrices, while in our setting the measurement matrices are rank-one. This difference fundamentally changes the theoretical analysis. The findings in [23] are therefore not applicable to the problem of joint blind deconvolution/demixing. The compressive principal component analysis in [35] is also a form of demixing problem, but its setting is only vaguely related to ours. There is a large amount of literature on demixing problems, but the vast majority does not have a “blind deconvolution component”, therefore this body of work is only marginally related to the topic of our paper.

Blind deconvolution/demixing problems also appear in convolutional dictionary learning, see e.g. [2]. There, the aim is to factorize an ensemble of input vectors into a linear combination of overcomplete basis elements which are modeled as shift-invariant—the latter property is why the factorization turns into a convolution. The setup is similar to (1.1), but with an additional penalty term to enforce sparsity of the convolving filters. The existing literature on convolutional dictionary learning is mainly focused on empirical results, therefore there is little overlap with our work. But it is an interesting challenge for future research to see whether the approach in this paper can be modified to provide a fast and theoretically sound solver for the sparse convolutional coding problem.

There are numerous papers concerned with blind deconvolution/demixing problems in the area of wireless communications. But the majority of these papers assumes the availability of multiple measurement vectors, which makes the problem significantly easier. Those methods however cannot be applied to the case of a single measurement vector, which is the focus of this paper. Thus there is essentially no overlap of those papers with our work.

Our previous paper [19] solves (1.1) under subspace conditions, i.e., assuming that both \mathbf{f}_i and \mathbf{g}_i belong to known linear subspaces. This contributes to generalizing the pioneering work by Ahmed, Recht, and Romberg [1] from the “single-user” scenario to the “multi-user” scenario. Both [1] and [19] employ a two-step convex approach: first “lifting” [9] is used and then the lifted version of the original bilinear inverse problems is relaxed into a semi-definite program. An improvement of the theoretical bounds in [19] was announced in [25].

While the convex approach is certainly effective and elegant, it can hardly handle large-scale problems. This motivates us to apply a nonconvex optimization approach [8, 18] to this blind-deconvolution-blind-demixing problem. The mathematical challenge, when using non-convex methods, is to derive a rigorous convergence framework with conditions that are competitive with those in a convex framework.

In the last few years several excellent articles have appeared on provably convergent nonconvex optimization applied to various problems in signal processing and machine learning, e.g., matrix completion [15, 14, 29], phase retrieval [8, 11, 28, 3], blind deconvolution [17, 4, 18], dictionary learning [27] and low-rank matrix recovery [30, 34]. In this paper we derive the first nonconvex optimization algorithm to solve (1.1) fast and with rigorous theoretical guarantees concerning exact recovery, convergence rates, as well as robustness for noisy data. Our work can be viewed as a generalization of blind deconvolution [18] ($s = 1$) to the multi-user scenario ($s > 1$).

The idea behind our approach is strongly motivated by the nonconvex optimization algorithm for phase retrieval proposed in [8]. In this foundational paper, the authors use a two-step approach: (i) Construct a good initial guess with a numerically efficient algorithm; (ii) Starting with this initial guess, prove that simple gradient descent will converge to the true solution. Our paper follows a similar two-step scheme. However, the techniques used here are quite different from [8]. Like the matrix completion problem [7], the performance of the algorithm relies heavily and inherently on how much the ground truth signals are aligned with the design matrix. Due to this so-called “incoherence” issue, we need to impose extra constraints, which results in a different construction of the so-called *basin of attraction*. Therefore, influenced by [15, 29, 18], we add penalty terms to control the incoherence and this leads to the regularized gradient descent method, which forms the core of our proposed algorithm.

To the best of our knowledge, our algorithm is the first algorithm for the blind deconvolution/blind demixing problem that is numerically efficient, robust against noise, and comes with rigorous recovery guarantees.

1.2 Notation

For a matrix \mathbf{Z} , $\|\mathbf{Z}\|$ denotes its operator norm and $\|\mathbf{Z}\|_F$ is its the Frobenius norm. For a vector \mathbf{z} , $\|\mathbf{z}\|$ is its Euclidean norm and $\|\mathbf{z}\|_\infty$ is the ℓ_∞ -norm. For both matrices and vectors, \mathbf{Z}^* and \mathbf{z}^* denote their complex conjugate transpose. \bar{z} is the complex conjugate of z . We equip the matrix space $\mathbb{C}^{K \times N}$ with the inner product defined by $\langle \mathbf{U}, \mathbf{V} \rangle := \text{Tr}(\mathbf{U}^* \mathbf{V})$. For a

given vector \mathbf{z} , $\text{diag}(\mathbf{z})$ represents the diagonal matrix whose diagonal entries are \mathbf{z} . For any $z \in \mathbb{R}$, let $z_+ = \frac{z+|z|}{2}$.

2 Preliminaries

Obviously, without any further assumption, it is impossible to solve (1.1). Therefore, we impose the following subspace assumptions throughout our discussion [1, 19].

- **Channel subspace assumption:** Each finite impulse response $\mathbf{f}_i \in \mathbb{C}^L$ is assumed to have *maximum delay spread* K , i.e.,

$$\mathbf{f}_i(n) = 0, \quad \text{for } n > K.$$

- **Signal subspace assumption:** Let $\mathbf{g}_i := \mathbf{C}_i \bar{\mathbf{x}}_i$ be the outcome of the signal $\bar{\mathbf{x}}_i \in \mathbb{C}^N$ encoded by a matrix $\mathbf{C}_i \in \mathbb{C}^{L \times N}$ with $L > N$, where the encoding matrix \mathbf{C}_i is known and assumed to have full rank¹.

Remark 2.1. *Both subspace assumptions are common in various applications. For instance in wireless communications, the channel impulse response can always be modeled to have finite support (or maximum delay spread, as it is called in engineering jargon) due to the physical properties of wave propagation [13]; and the signal subspace assumption is a standard feature found in many current communication systems) [13], including CDMA (where \mathbf{C}_i is known as spreading matrix) and OFDM (where \mathbf{C}_i is known as precoding matrix).*

The specific choice of the encoding matrices \mathbf{C}_i depends on a variety of conditions. In this paper, we derive our theory by assuming that \mathbf{C}_i is a complex Gaussian random matrix, i.e., each entry in \mathbf{C}_i is i.i.d. $\mathcal{CN}(0, 1)$. This assumption, while sometimes imposed in the wireless communications literature, is somewhat unrealistic in practice, due to the lack of a fast algorithm to apply \mathbf{C}_i and due to storage requirements. In practice one would rather choose \mathbf{C}_i to be something like the product of a Hadamard matrix and a diagonal matrix with random binary entries. We hope to address such more structured encoding matrices in our future research. Our numerical simulations (see Section 4) show no difference in the performance of our algorithm for either choice.

Under the two assumptions above, the model actually has a simpler form in the *frequency* domain. We assume throughout the paper that the convolution of finite sequences is circular convolution². By applying the Discrete Fourier Transform (DFT) to (1.1) along with the two assumptions, we have

$$\frac{1}{\sqrt{L}} \mathbf{F} \mathbf{y} = \sum_{i=1}^s \text{diag}(\mathbf{F} \mathbf{h}_i) (\mathbf{F} \mathbf{C}_i \bar{\mathbf{x}}_i) + \frac{1}{\sqrt{L}} \mathbf{F} \mathbf{n}$$

where \mathbf{F} is the $L \times L$ normalized unitary DFT matrix with $\mathbf{F}^* \mathbf{F} = \mathbf{F} \mathbf{F}^* = \mathbf{I}_L$. The noise is assumed to be additive white complex Gaussian noise with $\mathbf{n} \sim \mathcal{CN}(\mathbf{0}, \sigma^2 d_0^2 \mathbf{I}_L)$ where $d_0 = \sqrt{\sum_{i=1}^s \|\mathbf{h}_{i0}\|^2 \|\mathbf{x}_{i0}\|^2}$, and $\{(\mathbf{h}_{i0}, \mathbf{x}_{i0})\}_{i=1}^s$ is the ground truth. We define $d_{i0} = \|\mathbf{h}_{i0} \mathbf{x}_{i0}^*\|_F$ and assume without loss of generality that $\|\mathbf{h}_{i0}\|$ and $\|\mathbf{x}_{i0}\|$ are of the same norm, i.e., $\|\mathbf{h}_{i0}\| = \|\mathbf{x}_{i0}\| = \sqrt{d_{i0}}$. In that way, $\frac{1}{\sigma^2}$ actually is a measure of SNR (signal to noise ratio).

Let $\mathbf{h}_i \in \mathbb{C}^K$ be the first K nonzero entries of \mathbf{f}_i and $\mathbf{B} \in \mathbb{C}^{L \times K}$ be a low-frequency DFT matrix (the first K columns of an $L \times L$ unitary DFT matrix). Then a simple relation holds,

$$\mathbf{F} \mathbf{f}_i = \mathbf{B} \mathbf{h}_i, \quad \mathbf{B}^* \mathbf{B} = \mathbf{I}_K.$$

¹Here we use the conjugate $\bar{\mathbf{x}}_i$ instead of \mathbf{x}_i because it will simplify our notation in later derivations.

²This circular convolution assumption can often be reinforced directly (for example in wireless communications the use of a cyclic prefix in OFDM renders the convolution circular) or indirectly (e.g. via zero-padding). In the first case replacing regular convolution by circular convolution does not introduce any errors at all. In the latter case one introduces an additional approximation error in the inversion which is negligible, since it decays exponentially for impulse responses of finite length [26].

We also denote $\mathbf{A}_i := \overline{\mathbf{F}\mathbf{C}_i}$ and $\mathbf{e} := \frac{1}{\sqrt{L}}\mathbf{F}\mathbf{n}$. Due to the Gaussianity, \mathbf{A}_i also possesses complex Gaussian distribution and so does \mathbf{e} . From now on, instead of focusing on the original model, we consider (with a slight abuse of notation) the following equivalent formulation throughout our discussion:

$$\mathbf{y} = \sum_{i=1}^s \text{diag}(\mathbf{B}\mathbf{h}_i)\overline{\mathbf{A}_i\mathbf{x}_i} + \mathbf{e}, \quad (2.1)$$

where $\mathbf{e} \sim \mathcal{CN}(\mathbf{0}, \frac{\sigma^2 d_0^2}{L}\mathbf{I}_L)$. Our goal here is to estimate all $\{\mathbf{h}_i, \mathbf{x}_i\}_{i=1}^s$ from \mathbf{y}, \mathbf{B} and $\{\mathbf{A}_i\}_{i=1}^s$. Obviously, this is a bilinear inverse problem, i.e., if all $\{\mathbf{h}_i\}_{i=1}^s$ are given, it is a linear inverse problem (the ordinary demixing problem) to recover all $\{\mathbf{x}_i\}_{i=1}^s$, and vice versa. We note that there is a scaling ambiguity in all blind deconvolution problems that cannot be resolved by any reconstruction method without further information. Namely, if the pair $(\mathbf{h}_i, \mathbf{x}_i)$ is a solution then so is $(\alpha\mathbf{h}_i, \alpha^{-1}\mathbf{x}_i)$ for any $\alpha \neq 0$. Therefore, when we talk about exact recovery in the following, then this is understood modulo such a trivial scaling ambiguity.

Before proceeding to our proposed algorithm we introduce some notation to facilitate a more convenient presentation of our approach. Let \mathbf{b}_l be the l -th column of \mathbf{B}^* and \mathbf{a}_{il} be the l -th column of \mathbf{A}_i^* . Based on our assumptions the following properties hold:

$$\sum_{l=1}^L \mathbf{b}_l \mathbf{b}_l^* = \mathbf{I}_K, \quad \|\mathbf{b}_l\|^2 = \frac{K}{L}, \quad \mathbf{a}_{il} \sim \mathcal{CN}(\mathbf{0}, \mathbf{I}_N).$$

Moreover, inspired by the well-known *lifting* idea [9, 1, 6, 20], we define the useful matrix-valued linear operator $\mathcal{A}_i : \mathbb{C}^{K \times N} \rightarrow \mathbb{C}^L$ and its adjoint $\mathcal{A}_i^* : \mathbb{C}^L \rightarrow \mathbb{C}^{K \times N}$ by

$$\mathcal{A}_i(\mathbf{Z}) := \{\mathbf{b}_l^* \mathbf{Z} \mathbf{a}_{il}\}_{l=1}^L, \quad \mathcal{A}_i^*(z) := \sum_{l=1}^L z_l \mathbf{b}_l \mathbf{a}_{il}^* = \mathbf{B}^* \text{diag}(z) \mathbf{A}_i \quad (2.2)$$

for each $1 \leq i \leq s$ under canonical inner product over $\mathbb{C}^{K \times N}$. Therefore, (2.1) can be written in the following equivalent form

$$\mathbf{y} = \sum_{i=1}^s \mathcal{A}_i(\mathbf{h}_i \mathbf{x}_i^*) + \mathbf{e}. \quad (2.3)$$

Hence, we can think of \mathbf{y} as the observation vector obtained from taking *linear* measurements with respect to a set of rank-1 matrices $\{\mathbf{h}_i \mathbf{x}_i^*\}_{i=1}^s$. In fact, with a bit of linear algebra (and ignoring the noise term for the moment), the l -th entry of \mathbf{y} in (2.3) equals the inner product of two block-diagonal matrices:

$$y_l = \left\langle \begin{bmatrix} \mathbf{h}_{1,0} \mathbf{x}_{1,0}^* & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{h}_{2,0} \mathbf{x}_{2,0}^* & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{h}_{s,0} \mathbf{x}_{s,0}^* \end{bmatrix}, \begin{bmatrix} \mathbf{b}_l \mathbf{a}_{1l}^* & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{b}_l \mathbf{a}_{2l}^* & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{b}_l \mathbf{a}_{sl}^* \end{bmatrix} \right\rangle + e_l \quad (2.4)$$

where $y_l = \sum_{i=1}^s \mathbf{b}_l^* \mathbf{h}_{i,0} \mathbf{x}_{i,0}^* \mathbf{a}_{il} + e_l, 1 \leq l \leq L$. In other words, we aim to recover such a block-diagonal matrix (the left-hand side in the inner product (2.4)) from L linear measurements with block structure if $\mathbf{e} = \mathbf{0}$.

By stacking all $\{\mathbf{h}_i\}_{i=1}^s$ (and $\{\mathbf{x}_i\}_{i=1}^s, \{\mathbf{h}_{i,0}\}_{i=1}^s, \{\mathbf{x}_{i,0}\}_{i=1}^s$) into a long column, we let

$$\mathbf{h} := \begin{bmatrix} \mathbf{h}_1 \\ \vdots \\ \mathbf{h}_s \end{bmatrix}, \quad \mathbf{h}_0 := \begin{bmatrix} \mathbf{h}_{1,0} \\ \vdots \\ \mathbf{h}_{s,0} \end{bmatrix} \in \mathbb{C}^{Ks}, \quad \mathbf{x} := \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_s \end{bmatrix}, \quad \mathbf{x}_0 := \begin{bmatrix} \mathbf{x}_{1,0} \\ \vdots \\ \mathbf{x}_{s,0} \end{bmatrix} \in \mathbb{C}^{Ns}. \quad (2.5)$$

We define \mathcal{H} as a bilinear operator which maps a pair $(\mathbf{h}, \mathbf{x}) \in \mathbb{C}^{Ks} \times \mathbb{C}^{Ns}$ into a block diagonal

matrix in $\mathbb{C}^{Ks \times Ns}$, i.e.,

$$\mathcal{H}(\mathbf{h}, \mathbf{x}) := \begin{bmatrix} \mathbf{h}_1 \mathbf{x}_1^* & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{h}_2 \mathbf{x}_2^* & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{h}_s \mathbf{x}_s^* \end{bmatrix} \in \mathbb{C}^{Ks \times Ns}. \quad (2.6)$$

Let $\mathbf{X} := \mathcal{H}(\mathbf{h}, \mathbf{x})$ and $\mathbf{X}_0 := \mathcal{H}(\mathbf{h}_0, \mathbf{x}_0)$ where \mathbf{X}_0 is the ground truth. Define $\mathcal{A}(\mathbf{Z}) : \mathbb{C}^{Ks \times Ns} \rightarrow \mathbb{C}^L$ as

$$\mathcal{A}(\mathbf{Z}) := \sum_{i=1}^s \mathcal{A}_i(\mathbf{Z}_i) \quad (2.7)$$

where $\mathbf{Z} = \text{blkdiag}(\mathbf{Z}_1, \dots, \mathbf{Z}_s)$. Therefore, $\mathcal{A}(\mathcal{H}(\mathbf{h}, \mathbf{x})) = \sum_{i=1}^s \mathcal{A}_i(\mathbf{h}_i \mathbf{x}_i^*)$ and $\mathbf{y} = \mathcal{A}(\mathcal{H}(\mathbf{h}_0, \mathbf{x}_0)) + \mathbf{e}$. The adjoint operator \mathcal{A}^* is defined naturally as

$$\mathcal{A}^*(\mathbf{z}) := \begin{bmatrix} \mathcal{A}_1^*(\mathbf{z}) & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathcal{A}_2^*(\mathbf{z}) & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathcal{A}_s^*(\mathbf{z}) \end{bmatrix} \in \mathbb{C}^{Ks \times Ns}, \quad (2.8)$$

which is a linear map from \mathbb{C}^L to $\mathbb{C}^{Ks \times Ns}$. To measure the approximation error of \mathbf{X}_0 given by \mathbf{X} , we define $\delta(\mathbf{h}, \mathbf{x})$ as the global relative error:

$$\delta(\mathbf{h}, \mathbf{x}) := \frac{\|\mathbf{X} - \mathbf{X}_0\|_F}{\|\mathbf{X}_0\|_F} = \frac{\sqrt{\sum_{i=1}^s \|\mathbf{h}_i \mathbf{x}_i^* - \mathbf{h}_{i0} \mathbf{x}_{i0}^*\|_F^2}}{d_0} = \sqrt{\frac{\sum_{i=1}^s \delta_i^2 d_{i0}^2}{\sum_{i=1}^s d_{i0}^2}} \quad (2.9)$$

where $\delta_i := \delta_i(\mathbf{h}_i, \mathbf{x}_i)$ is the relative error within each component:

$$\delta_i(\mathbf{h}_i, \mathbf{x}_i) := \frac{\|\mathbf{h}_i \mathbf{x}_i^* - \mathbf{h}_{i0} \mathbf{x}_{i0}^*\|_F}{d_{i0}}.$$

Note that δ and δ_i are functions of (\mathbf{h}, \mathbf{x}) and $(\mathbf{h}_i, \mathbf{x}_i)$ respectively and in most cases, we just simply use δ and δ_i if no possibility of confusion exists.

2.1 Convex versus nonconvex approaches

As indicated in (2.4), joint blind deconvolution-demixing can be recast as the task to recover a rank- s block-diagonal matrix from linear measurements. In general, such a low-rank matrix recovery problem is NP-hard. In order to take advantage of the low-rank property of the ground truth, it is natural to adopt convex relaxation by solving a convenient nuclear norm minimization program, i.e.,

$$\min \sum_{i=1}^s \|\mathbf{Z}_i\|_*, \quad \text{s.t.} \quad \sum_{i=1}^s \mathcal{A}_i(\mathbf{Z}_i) = \mathbf{y}. \quad (2.10)$$

The question of when the solution of (2.10) yields exact recovery is answered in our previous work [19], whose main theoretical result is informally summarized in the following theorem.

Theorem 2.2. *Suppose that \mathbf{A}_i are $L \times N$ i.i.d. complex Gaussian matrices and \mathbf{B} is an $L \times K$ partial DFT matrix with $\mathbf{B}^* \mathbf{B} = \mathbf{I}_K$. Then solving (2.10) gives exact recovery if the number of measurements L yields*

$$L \geq C_\gamma s^2 (K + N) \log^3 L$$

with probability at least $1 - L^{-\gamma}$ where C_γ is an absolute scalar only depending on γ linearly.

Numerical simulations in [19] show that the semidefinite program (SDP) in (2.10) is able to estimate all pairs of $\{\mathbf{h}_i, \mathbf{x}_i\}_{i=1}^s$ even when L is very close to $s(K + N)$, i.e., the degree of freedom for the unknowns, although L depends on s quadratically in our theory. However, the computational cost for solving an SDP already become challenging for moderate size problems and too expensive for large scale problems.

Therefore, we try to look for a more efficient nonconvex approach, which hopefully is also reinforced by theory. It seems quite natural to achieve the goal by minimizing the following *non-linear* least squares objective function with respect to (\mathbf{h}, \mathbf{x})

$$F(\mathbf{h}, \mathbf{x}) := \|\mathcal{A}(\mathcal{H}(\mathbf{h}, \mathbf{x}) - \mathbf{y})\|^2 = \left\| \sum_{i=1}^s \mathcal{A}_i(\mathbf{h}_i \mathbf{x}_i^*) - \mathbf{y} \right\|^2. \quad (2.11)$$

In particular, if $\mathbf{e} = \mathbf{0}$, we write

$$F_0(\mathbf{h}, \mathbf{x}) := \left\| \sum_{i=1}^s \mathcal{A}_i(\mathbf{h}_i \mathbf{x}_i^* - \mathbf{h}_{i0} \mathbf{x}_{i0}^*) \right\|^2. \quad (2.12)$$

As also pointed out in [18], this is a highly nonconvex optimization problem. Many of the commonly used algorithms, such as gradient descent or alternating minimization, may not necessarily yield convergence to the global minimum, so that we cannot always hope to obtain the desired solution. Often, those simple algorithms might get stuck in local minima.

2.2 The basin of attraction

Motivated by several excellent recent papers of nonconvex optimization on various signal processing and machine learning problem, we propose our two-step algorithm: (i) Compute an initial guess carefully; (ii) Apply gradient descent to the objective function, starting with the carefully chosen initial guess. One difficulty of understanding nonconvex optimization consists in how to construct the so-called *basin of attraction*, i.e., if the starting point is inside this basin of attraction, the iterates will always stay inside the region and converge to the global minimum. The construction of the basin of attraction varies for different problems [8, 3, 29]. For this problem, similar to [18], the construction follows from the following three observations. Each of these observations suggests the definition of a certain *neighborhood* and the basin of attraction is then defined as the intersection of these three neighborhood sets $\mathcal{N}_d \cap \mathcal{N}_\mu \cap \mathcal{N}_\epsilon$.

1. **Ambiguity of solution:** in fact, we can only recover $(\mathbf{h}_i, \mathbf{x}_i)$ up to a scalar since $(\alpha \mathbf{h}_i, \alpha^{-1} \mathbf{x}_i)$ and $(\mathbf{h}_i, \mathbf{x}_i)$ are both solutions for $\alpha \neq 0$. From a numerical perspective, we want to avoid the scenario when $\|\mathbf{h}_i\| \rightarrow 0$ and $\|\mathbf{x}_i\| \rightarrow \infty$ while $\|\mathbf{h}_i\| \|\mathbf{x}_i\|$ is fixed, which potentially leads to numerical instability. To balance both the norm of $\|\mathbf{h}_i\|$ and $\|\mathbf{x}_i\|$ for all $1 \leq i \leq s$, we define

$$\mathcal{N}_d := \{ \{(\mathbf{h}_i, \mathbf{x}_i)\}_{i=1}^s : \|\mathbf{h}_i\| \leq 2\sqrt{d_{i0}}, \|\mathbf{x}_i\| \leq 2\sqrt{d_{i0}}, 1 \leq i \leq s \},$$

which is a convex set.

2. **Incoherence:** the performance depends on how large/small the incoherence μ_h^2 is, where μ_h^2 is defined by

$$\mu_h^2 := \max_{1 \leq i \leq s} \frac{L \|\mathbf{B} \mathbf{h}_{i0}\|_\infty^2}{\|\mathbf{h}_{i0}\|^2}.$$

The idea is that: *the smaller the μ_h^2 is, the better the performance is.* Let's consider an extreme case: if $\mathbf{B} \mathbf{h}_{i0}$ is highly sparse or spiky, we lose much information on those zero/small entries and cannot hope to get satisfactory recovered signals.

A similar quantity is also introduced in the matrix completion problem [7, 29]. The larger μ_h^2 is, the more \mathbf{h}_{i0} is aligned with one particular row of \mathbf{B} . To control the incoherence between \mathbf{b}_l and \mathbf{h}_i , we define the second neighborhood,

$$\mathcal{N}_\mu := \{ \{ \mathbf{h}_i \}_{i=1}^s : \sqrt{L} \|\mathbf{B} \mathbf{h}_i\|_\infty \leq 4\sqrt{d_{i0}} \mu, 1 \leq i \leq s \}, \quad (2.13)$$

where μ is a parameter and $\mu \geq \mu_h$. Note that \mathcal{N}_μ is also a convex set.

3. **Close to the ground truth:** we also want to construct an initial guess such that it is close to the ground truth, i.e.,

$$\mathcal{N}_\epsilon := \left\{ \{(\mathbf{h}_i, \mathbf{x}_i)\}_{i=1}^s : \delta_i = \frac{\|\mathbf{h}_i \mathbf{x}_i^* - \mathbf{h}_{i0} \mathbf{x}_{i0}^*\|_F}{d_{i0}} \leq \epsilon, 1 \leq i \leq s \right\} \quad (2.14)$$

where ϵ is a predetermined parameter in $(0, \frac{1}{15}]$.

Remark 2.3. To ensure $\delta_i \leq \epsilon$, it suffices to ensure $\delta \leq \frac{\epsilon}{\sqrt{s\kappa}}$ where $\kappa := \frac{\max d_{i0}}{\min d_{i0}} \geq 1$. This is because

$$\frac{1}{s\kappa^2} \sum_{i=1}^s \delta_i^2 \leq \delta^2 \leq \frac{\epsilon^2}{s\kappa^2}$$

which implies $\max_{1 \leq i \leq s} \delta_i \leq \epsilon$.

Remark 2.4. When we say $(\mathbf{h}, \mathbf{x}) \in \mathcal{N}_d, \mathcal{N}_d$ or \mathcal{N}_ϵ , it means for all $i = 1, \dots, s$ we have $(\mathbf{h}_i, \mathbf{x}_i) \in \mathcal{N}_d, \mathcal{N}_\mu$ or \mathcal{N}_ϵ respectively. In particular, $(\mathbf{h}_0, \mathbf{x}_0) \in \mathcal{N}_d \cap \mathcal{N}_\mu \cap \mathcal{N}_\epsilon$.

2.3 Objective function and Wirtinger derivative

To implement the first two observations, we introduce the regularizer $G(\mathbf{h}, \mathbf{x})$, defined as the sum of s components

$$G(\mathbf{h}, \mathbf{x}) := \sum_{i=1}^s G_i(\mathbf{h}_i, \mathbf{x}_i). \quad (2.15)$$

For each component $G_i(\mathbf{h}_i, \mathbf{x}_i)$, we let $\rho \geq d^2 + 2\|e\|^2$, $0.9d_0 \leq d \leq 1.1d_0$, $0.9d_{i0} \leq d_i \leq 1.1d_{i0}$ for all $1 \leq i \leq s$ and

$$G_i := \rho \left[\underbrace{G_0\left(\frac{\|\mathbf{h}_i\|^2}{2d_i}\right) + G_0\left(\frac{\|\mathbf{x}_i\|^2}{2d_i}\right)}_{\mathcal{N}_d} + \underbrace{\sum_{l=1}^L G_0\left(\frac{L|\mathbf{b}_l^* \mathbf{h}_i|^2}{8d_i \mu^2}\right)}_{\mathcal{N}_\mu} \right], \quad (2.16)$$

where $G_0(z) = \max\{z - 1, 0\}^2$. Here both d and $\{d_i\}_{i=1}^s$ are data-driven and well approximated by our spectral initialization procedure; and μ^2 is a tuning parameter which could be estimated if we assume a specific statistical model for the channel (for example, in the widely used Rayleigh fading model, the channel coefficients are assumed to be complex Gaussian). The idea behind G_i is quite straightforward though the formulation is complicated. For each G_i in (2.16), the first two terms try to force the iterates to lie in \mathcal{N}_d and the third term tries to force the iterates to lie in \mathcal{N}_μ . What about the neighborhood \mathcal{N}_ϵ ? A proper choice of the initialization along with gradient descent (keeping the objective function decrease) will ensure that the iterates lie in \mathcal{N}_ϵ .

Finally, we consider the objective function as the sum of nonlinear least squares objective function $F(\mathbf{h}, \mathbf{x})$ in (2.11) and the regularizer $G(\mathbf{h}, \mathbf{x})$,

$$\tilde{F}(\mathbf{h}, \mathbf{x}) := F(\mathbf{h}, \mathbf{x}) + G(\mathbf{h}, \mathbf{x}). \quad (2.17)$$

Note that the input of the function $\tilde{F}(\mathbf{h}, \mathbf{x})$ consists of complex variables, but the output is real-valued (so do $F(\mathbf{h}, \mathbf{x})$ and $G(\mathbf{h}, \mathbf{x})$) and thus simple relations hold

$$\frac{\partial \tilde{F}}{\partial \bar{\mathbf{h}}_i} = \overline{\frac{\partial \tilde{F}}{\partial \mathbf{h}_i}}, \quad \frac{\partial \tilde{F}}{\partial \bar{\mathbf{x}}_i} = \overline{\frac{\partial \tilde{F}}{\partial \mathbf{x}_i}}.$$

Therefore, to minimize this function, it suffices to consider only the gradient of \tilde{F} with respect to $\bar{\mathbf{h}}_i$ and $\bar{\mathbf{x}}_i$, which is also called Wirtinger derivative [8]. The Wirtinger derivatives of

$F(\mathbf{h}, \mathbf{x})$ and $G(\mathbf{h}, \mathbf{x})$ w.r.t. $\bar{\mathbf{h}}_i$ and $\bar{\mathbf{x}}_i$ can be easily computed as follows

$$\nabla F_{\mathbf{h}_i} = \mathcal{A}_i^*(\mathcal{A}(\mathbf{X}) - \mathbf{y}) \mathbf{x}_i = \mathcal{A}_i^*(\mathcal{A}(\mathbf{X} - \mathbf{X}_0) - \mathbf{e}) \mathbf{x}_i, \quad (2.18)$$

$$\nabla F_{\mathbf{x}_i} = (\mathcal{A}_i^*(\mathcal{A}(\mathbf{X}) - \mathbf{y}))^* \mathbf{h}_i = (\mathcal{A}_i^*(\mathcal{A}(\mathbf{X} - \mathbf{X}_0) - \mathbf{e}))^* \mathbf{h}_i, \quad (2.19)$$

$$\nabla G_{\mathbf{h}_i} = \frac{\rho}{2d_i} \left[G'_0 \left(\frac{\|\mathbf{h}_i\|^2}{2d_i} \right) \mathbf{h}_i + \frac{L}{4\mu^2} \sum_{l=1}^L G'_0 \left(\frac{L|\mathbf{b}_l^* \mathbf{h}_i|^2}{8d_i \mu^2} \right) \mathbf{b}_l \mathbf{b}_l^* \mathbf{h}_i \right], \quad (2.20)$$

$$\nabla G_{\mathbf{x}_i} = \frac{\rho}{2d_i} G'_0 \left(\frac{\|\mathbf{x}_i\|^2}{2d_i} \right) \mathbf{x}_i, \quad (2.21)$$

where $\mathcal{A}(\mathbf{X}) = \sum_{i=1}^s \mathcal{A}_i(\mathbf{h}_i \mathbf{x}_i^*)$ and \mathcal{A}^* is defined in (2.8). In short, we denote

$$\nabla \tilde{F}_{\mathbf{h}} := \nabla F_{\mathbf{h}} + \nabla G_{\mathbf{h}}, \quad \nabla F_{\mathbf{h}} := \begin{bmatrix} \nabla F_{\mathbf{h}_1} \\ \vdots \\ \nabla F_{\mathbf{h}_s} \end{bmatrix}, \quad \nabla G_{\mathbf{h}} := \begin{bmatrix} \nabla G_{\mathbf{h}_1} \\ \vdots \\ \nabla G_{\mathbf{h}_s} \end{bmatrix}, \quad (2.22)$$

similar definitions hold for $\nabla \tilde{F}_{\mathbf{x}}$, $\nabla F_{\mathbf{x}}$ and $G_{\mathbf{x}}$. It is easy to see that $\nabla F_{\mathbf{h}} = \mathcal{A}^*(\mathcal{A}(\mathbf{X}) - \mathbf{y})\mathbf{x}$ and $\nabla F_{\mathbf{x}} = (\mathcal{A}^*(\mathcal{A}(\mathbf{X}) - \mathbf{y}))^*\mathbf{h}$.

3 Algorithm and Main Theory

3.1 Two-step algorithm

As mentioned before, the first step is to find a good initial guess $(\mathbf{u}^{(0)}, \mathbf{v}^{(0)}) \in \mathbb{C}^{Ks} \oplus \mathbb{C}^{Ns}$ such that it is inside the basin of attraction. The initialization follows from this key fact:

$$\mathbb{E}(\mathcal{A}_i^*(\mathbf{y})) = \mathbb{E} \left(\mathcal{A}_i^* \left(\sum_{j=1}^s \mathcal{A}_j(\mathbf{h}_{j0} \mathbf{x}_{j0}^*) + \mathbf{e} \right) \right) = \mathbf{h}_{i0} \mathbf{x}_{i0}^*$$

where we use $\mathbf{B}^* \mathbf{B} = \sum_{l=1}^L \mathbf{b}_l \mathbf{b}_l^* = \mathbf{I}_K$, $\mathbb{E}(\mathbf{a}_{il} \mathbf{a}_{il}^*) = \mathbf{I}_N$ and

$$\begin{aligned} \mathbb{E}(\mathcal{A}_i^* \mathcal{A}_i(\mathbf{h}_{i0} \mathbf{x}_{i0}^*)) &= \sum_{l=1}^L \mathbf{b}_l \mathbf{b}_l^* \mathbf{h}_{i0} \mathbf{x}_{i0}^* \mathbb{E}(\mathbf{a}_{il} \mathbf{a}_{il}^*) = \mathbf{h}_{i0} \mathbf{x}_{i0}^*, \\ \mathbb{E}(\mathcal{A}_j^* \mathcal{A}_i(\mathbf{h}_{i0} \mathbf{x}_{i0}^*)) &= \sum_{l=1}^L \mathbf{b}_l \mathbf{b}_l^* \mathbf{h}_{i0} \mathbf{x}_{i0}^* \mathbb{E}(\mathbf{a}_{il} \mathbf{a}_{jl}^*) = \mathbf{0}, \quad \forall j \neq i. \end{aligned}$$

Therefore, it is natural to extract the leading singular value and associated left and right singular vectors from each $\mathcal{A}_i^*(\mathbf{y})$ and use them as (a hopefully good) approximation to $(d_{i0}, \mathbf{h}_{i0}, \mathbf{x}_{i0})$. This idea leads to Algorithm 1, the proof of which is given in Section 6.5. The second step of the algorithm is just to apply gradient descent to \tilde{F} with the initial guess $\{(\mathbf{u}_i^{(0)}, \mathbf{v}_i^{(0)}, d_i)\}_{i=1}^s$ or $(\mathbf{u}^{(0)}, \mathbf{v}^{(0)}, \{d_i\}_{i=1}^s)$, where $\mathbf{u}^{(0)}$ stems from stacking all $\mathbf{u}_i^{(0)}$ into one long vector.

Remark 3.1. For Algorithm 2, we can rewrite each iteration into

$$\mathbf{u}^{(t)} = \mathbf{u}^{(t-1)} - \eta \nabla \tilde{F}_{\mathbf{h}}(\mathbf{u}^{(t-1)}, \mathbf{v}^{(t-1)}), \quad \mathbf{v}^{(t)} = \mathbf{v}^{(t-1)} - \eta \nabla \tilde{F}_{\mathbf{x}}(\mathbf{u}^{(t-1)}, \mathbf{v}^{(t-1)}),$$

where $\nabla \tilde{F}_{\mathbf{h}}$ and $\nabla \tilde{F}_{\mathbf{x}}$ are in (2.22), and

$$\mathbf{u}^{(t)} := \begin{bmatrix} \mathbf{u}_1^{(t)} \\ \vdots \\ \mathbf{u}_s^{(t)} \end{bmatrix}, \quad \mathbf{v}^{(t)} := \begin{bmatrix} \mathbf{v}_1^{(t)} \\ \vdots \\ \mathbf{v}_s^{(t)} \end{bmatrix}.$$

Algorithm 1 Initialization via spectral method and projection

- 1: **for** $i = 1, 2, \dots, s$ **do**
- 2: Compute $\mathcal{A}_i^*(\mathbf{y})$.
- 3: Find the leading singular value, left and right singular vectors of $\mathcal{A}_i^*(\mathbf{y})$, denoted by $(d_i, \hat{\mathbf{h}}_{i0}, \hat{\mathbf{x}}_{i0})$.
- 4: Solve the following optimization problem for $1 \leq i \leq s$:

$$\mathbf{u}_i^{(0)} := \operatorname{argmin}_{\mathbf{z} \in \mathbb{C}^K} \|\mathbf{z} - \sqrt{d_i} \hat{\mathbf{h}}_{i0}\|^2 \text{ s.t. } \sqrt{L} \|\mathbf{B}\mathbf{z}\|_\infty \leq 2\sqrt{d_i} \mu.$$

- 5: Set $\mathbf{v}_i^{(0)} = \sqrt{d_i} \hat{\mathbf{x}}_{i0}$.
 - 6: **end for**
 - 7: Output: $\{(\mathbf{u}_i^{(0)}, \mathbf{v}_i^{(0)}, d_i)\}_{i=1}^s$ or $(\mathbf{u}^{(0)}, \mathbf{v}^{(0)}, \{d_i\}_{i=1}^s)$.
-

Algorithm 2 Wirtinger gradient descent with constant stepsize η

- 1: **Initialization:** obtain $(\mathbf{u}^{(0)}, \mathbf{v}^{(0)}, \{d_i\}_{i=1}^s)$ via Algorithm 1.
 - 2: **for** $t = 1, 2, \dots$, **do**
 - 3: **for** $i = 1, 2, \dots, s$ **do**
 - 4: $\mathbf{u}_i^{(t)} = \mathbf{u}_i^{(t-1)} - \eta \nabla \tilde{F}_{\mathbf{h}_i}(\mathbf{u}_i^{(t-1)}, \mathbf{v}_i^{(t-1)})$,
 - 5: $\mathbf{v}_i^{(t)} = \mathbf{v}_i^{(t-1)} - \eta \nabla \tilde{F}_{\mathbf{x}_i}(\mathbf{u}_i^{(t-1)}, \mathbf{v}_i^{(t-1)})$,
 - 6: **end for**
 - 7: **end for**
-

3.2 Main theorem

Our main findings are summarized as follows: Theorem 3.2 shows that the initial guess given by Algorithm 1 indeed belongs to the basin of attraction. Moreover, d_i also serves as a good approximation of d_{i0} for each i . Theorem 3.3 demonstrates that the regularized Wirtinger gradient descent will guarantee the linear convergence of the iterates and the recovery is exact in the noise-free case and stable in the presence of noise.

Theorem 3.2. *The initialization obtained via Algorithm 1 satisfies*

$$(\mathbf{u}^{(0)}, \mathbf{v}^{(0)}) \in \frac{1}{\sqrt{3}} \mathcal{N}_d \cap \frac{1}{\sqrt{3}} \mathcal{N}_\mu \cap \mathcal{N}_{\frac{2\varepsilon}{5\sqrt{s}\kappa}} \quad (3.1)$$

and

$$0.9d_{i0} \leq d_i \leq 1.1d_{i0}, \quad 0.9d_0 \leq d \leq 1.1d_0, \quad (3.2)$$

holds with probability at least $1 - L^{-\gamma+1}$ if the number of measurements satisfies

$$L \geq C_{\gamma+\log(s)} (\mu_h^2 + \sigma^2) s^2 \kappa^4 \max\{K, N\} \log^2 L / \varepsilon^2. \quad (3.3)$$

Here ε is any predetermined constant in $(0, \frac{1}{15}]$, and C_γ is a constant only linearly depending on γ with $\gamma \geq 1$.

Theorem 3.3. *Starting with the initial value $\mathbf{z}^{(0)} := (\mathbf{u}^{(0)}, \mathbf{v}^{(0)})$ satisfying (3.1), the Algorithm 2 creates a sequence of iterates $(\mathbf{u}^{(t)}, \mathbf{v}^{(t)})$ which converges to the global minimum linearly,*

$$\|\mathcal{H}(\mathbf{u}^{(t)}, \mathbf{v}^{(t)}) - \mathcal{H}(\mathbf{h}_0, \mathbf{x}_0)\|_F \leq \frac{\varepsilon d_0}{\sqrt{2s}\kappa^2} (1 - \eta\omega)^{t/2} + 60\sqrt{s} \|\mathcal{A}^*(\mathbf{e})\| \quad (3.4)$$

with probability at least $1 - L^{-\gamma+1}$ and $\eta\omega = \mathcal{O}((s\kappa d_0(K+N)\log^2 L)^{-1})$ if the number of measurements L satisfies

$$L \geq C_{\gamma+\log(s)} (\mu^2 + \sigma^2) s^2 \kappa^4 \max\{K, N\} \log^2 L / \varepsilon^2. \quad (3.5)$$

In particular, with probability at least $1 - L^{-\gamma+1}$, there holds

$$\|\mathcal{A}^*(\mathbf{e})\| \leq C_0 \sigma d_0 \sqrt{\frac{\gamma s (K+N) (\log^2 L)}{L}}.$$

Remark 3.4. Our previous work [19] shows that the convex approach via semidefinite programming (see (2.10)) requires $L \geq C_0 s^2 (K + \mu_h^2 N) \log^3(L)$ to ensure exact recovery. Later, [25] claimed to improve this result to the near-optimal bound $L \geq C_0 s (K + \mu_h^2 N)$ up to some log-factors. The difference between nonconvex and convex methods lies in the appearance of the condition number κ in (3.5). This is not just an artifact of the proof—empirically we also observe that the value of κ affects the convergence rate of our nonconvex algorithm, see Figure 5.

Remark 3.5. Our theory suggests s^2 -dependence for the number of measurements L , although numerically L in fact depends on s linearly, as shown in Section 4. The reason for s^2 -dependence will be addressed in details in Section 5.2.

Remark 3.6. In the theoretical analysis, we assume that $\mathbf{C}_i/\mathbf{A}_i$ is a Gaussian random matrix. Numerical simulations suggest that this assumption is clearly not necessary. For example, \mathbf{C}_i may be chosen to be a Hadamard-type matrix which is more appropriate and favorable for communications.

Remark 3.7. If $\mathbf{e} = \mathbf{0}$, (3.4) shows that $(\mathbf{u}^{(t)}, \mathbf{v}^{(t)})$ converges to the ground truth at a linear rate. On the other hand, if noise exists, $(\mathbf{u}^{(t)}, \mathbf{v}^{(t)})$ is guaranteed to converge to a point within a small neighborhood of $(\mathbf{h}_0, \mathbf{x}_0)$. More importantly, if the number of measurements L gets larger, $\|\mathcal{A}^*(\mathbf{e})\|$ decays at the rate of $\mathcal{O}(L^{-1/2})$.

4 Numerics

In this section we present a range of numerical simulations to illustrate and complement different aspects of our theoretical framework. We will empirically analyze the number of measurements needed for perfect joint deconvolution/demixing to see how this compares to our theoretical bounds. We will also study the robustness for noisy data. In our simulations we use Gaussian encoding matrices, as in our theorems. But we also more more realistic structured encoding matrices, that are more reminiscent of what one might come across in wireless communications.

While Theorem 3.3 says that the number of measurements L depends *quadratically* on the number of sources s , numerical simulations suggest near-optimal performance. Figure 2 demonstrates that L actually depends linearly on s , i.e., the boundary between success (white) and failure (black) is approximately a linear function of s . In the experiment, $K = N = 50$ are fixed, all \mathbf{A}_i are complex Gaussians and all $(\mathbf{h}_i, \mathbf{x}_i)$ are standard complex Gaussian vectors. For each pair of (L, s) , 25 experiments are performed and we treat the recovery as a success if $\frac{\|\hat{\mathbf{X}} - \mathbf{X}_0\|_F}{\|\mathbf{X}_0\|_F} \leq 10^{-3}$. For our algorithm, we use backtracking to determine the stepsize and the iteration stops either if $\|\mathcal{A}(\mathcal{H}(\mathbf{h}^{(t+1)}, \mathbf{x}^{(t+1)}) - \mathcal{H}(\mathbf{h}^{(t)}, \mathbf{x}^{(t)}))\| < 10^{-6} \|\mathbf{y}\|$ or if the number of iterations reaches 500. The backtracking is based on the Armijo-Goldstein condition [22]. The initial stepsize is chosen to be $\eta = \frac{1}{K+N}$. If $\tilde{F}(\mathbf{z}^{(t)} - \eta \nabla \tilde{F}(\mathbf{z}^{(t)})) > \tilde{F}(\mathbf{z}^{(t)})$, we just divide η by two and use a smaller stepsize.

We see from Figure 2 that the number of measurements for the proposed algorithm to succeed not only seems to depend linearly on the number of sensors, but it is actually rather close to the information-theoretic limit $s(K + N)$. Indeed, the green dashed line in Figure 2, which represents the empirical boundary for the phase transition between success and failure corresponds to a line with slope about $\frac{3}{2}s(K + N)$. It is interesting to compare this empirical performance to the sharp theoretical phase transition bounds one would obtain via convex optimization [10, 23]. Considering the convex approach based on lifting in [19], we can adapt the theoretical framework in [10] to the blind deconvolution/demixing setting, but with one modification. The bounds in [10] rely on Gaussian widths of tangent cones related to the measurement matrices \mathcal{A}_i . Since simply analytic formulas for these expressions seem to be out of reach for the structured rank-one measurement matrices used in our paper, we instead compute the bounds for full-rank Gaussian random matrices, which yields a sharp bound of about $3s(K + N)$ (the corresponding bounds for rank-one sensing matrices will likely have a constant larger than 3). Note that these sharp theoretical bounds predict quite accurately

the empirical behavior of convex methods. Thus our empirical bound for using a non-convex methods compares rather favorably with that of the convex approach.

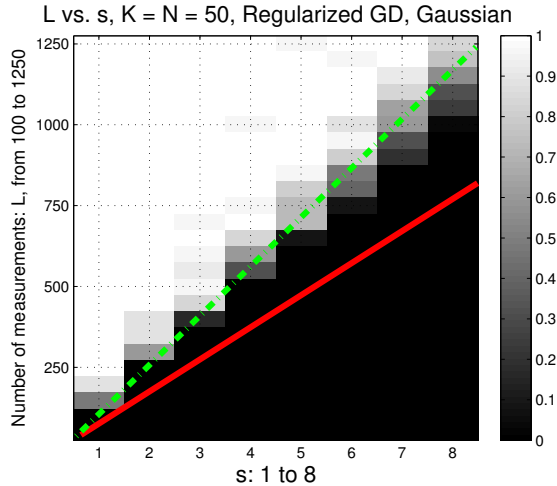


Figure 2: Phase transition plot for empirical recovery performance under different choices of (L, s) where $K = N = 50$ are fixed. Black region: failure; white region: success. The red solid line depicts the number of degrees of freedom and the green dashed line shows the empirical phase transition bound for Algorithm 2.

Similar conclusions can be drawn from Figure 3; there all \mathbf{A}_i are in the form of $\mathbf{A}_i = \mathbf{F}\mathbf{D}_i\mathbf{H}$ where \mathbf{F} is the unitary $L \times L$ DFT matrix, all \mathbf{D}_i are independent diagonal binary ± 1 matrices and \mathbf{H} is an $L \times N$ fixed partial deterministic Hadamard matrix. The purpose of using \mathbf{D}_i is to enhance the incoherence between each channel so that our algorithm is able to tell apart each individual signal and channel. As before we assume Gaussian channels, i.e., $\mathbf{h}_i \sim \mathcal{CN}(\mathbf{0}, \mathbf{I}_K)$. Therefore, our approach does not only work for Gaussian encoding matrices \mathbf{A}_i but also for the matrices that are interesting to real-world applications, although no satisfactory theory has been derived yet for that case. Moreover, due to the structure of \mathbf{A}_i and \mathbf{B} , fast transform algorithms are available, potentially allowing for real-time deployment.

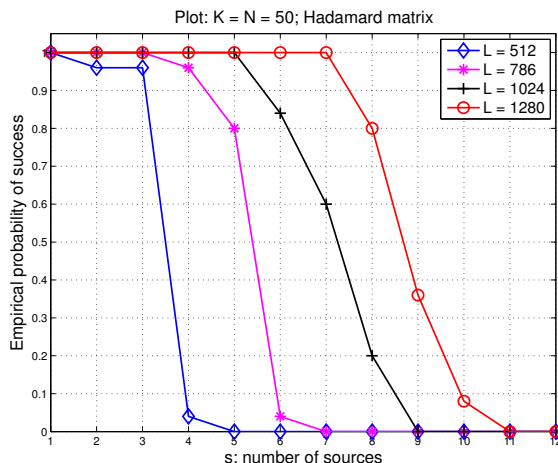


Figure 3: Empirical probability of successful recovery for different pairs of (L, s) when $K = N = 50$ are fixed.

Figure 4 shows the robustness of our algorithm under different levels of noise. We also run 25 samples for each level of SNR and different L and then compute the average relative error. It is easily seen that the relative error scales linearly with the SNR and one unit of increase in SNR (in dB) results in one unit of decrease in the relative error.

Theorem 3.3 suggests that the performance and convergence rate actually depend on the condition number of $\mathbf{X}_0 = \mathcal{H}(\mathbf{h}_0, \mathbf{x}_0)$, i.e., on $\kappa = \frac{\max d_{i0}}{\min d_{i0}}$ where $d_{i0} = \|\mathbf{h}_{i0}\| \|\mathbf{x}_{i0}\|$. Next we

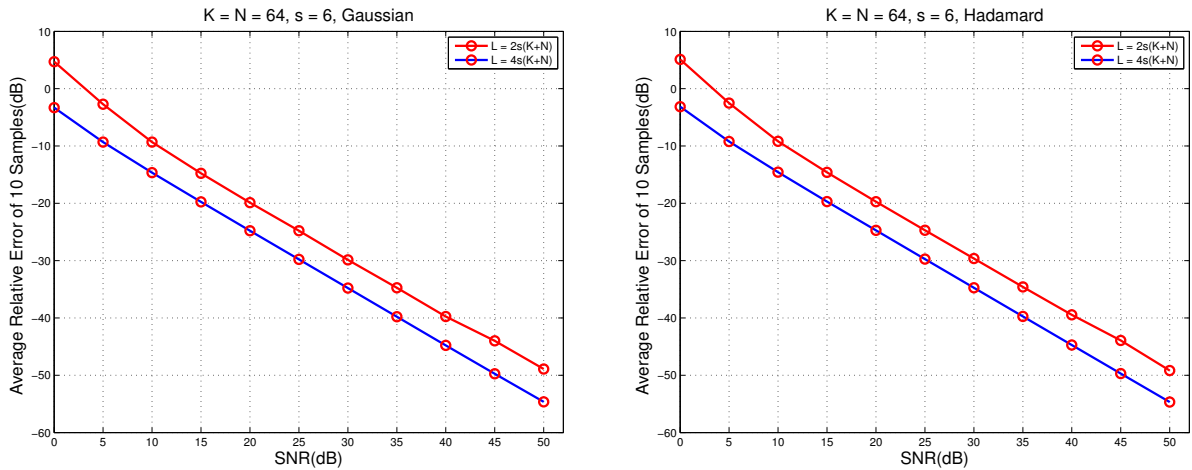


Figure 4: Relative error vs. SNR (dB): $\text{SNR} = 20 \log_{10} \left(\frac{\|y\|}{\|e\|} \right)$.

demonstrate that this dependence on the condition number is not an artifact of the proof, but is indeed also observed empirically. In this experiment, we let $s = 2$ and set for the first component $d_{1,0} = 1$ and for the second one $d_{2,0} = \kappa$ for $\kappa \in \{1, 2, 5\}$. Here, $\kappa = 1$ means that the received signals of both sensors have equal power, whereas $\kappa = 5$ means that the signal received from the second sensor is considerably stronger. The initial stepsize is chosen as $\eta = 1$, followed by the backtracking scheme. Figure 5 shows how the relative error decays with respect to the number of iterations t under different condition number κ and L .

The larger κ is, the slower the convergence rate is, as we see from Figure 5. This may result from two reasons: our spectral initialization may not be able to give a good initial guess for those weak components; moreover, during the gradient descent procedure, the gradient directions for the weak components could be totally dominated/polluted by the strong components. Currently, we still have no effective way of how to deal with this issue of slow convergence when κ is not small. We have to leave this topic for future investigations.

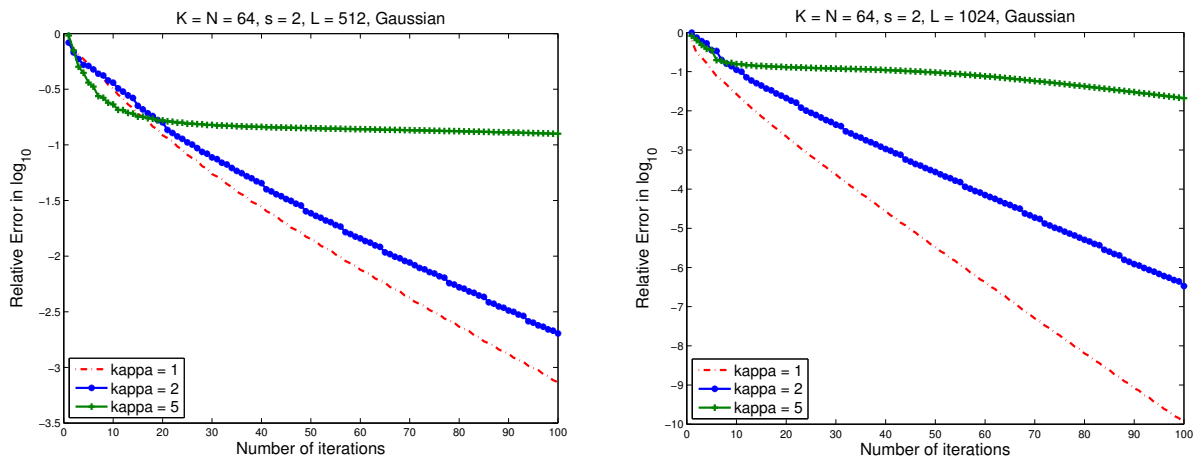


Figure 5: Relative error vs. number of iterations t .

5 Convergence analysis

Our convergence analysis relies on the following four conditions where the first three of them are local properties. We will also briefly discuss how they contribute to the proof of our main theorem. Note that our previous work [18] on blind deconvolution is actually a special case ($s = 1$) of (2.1). Therefore, the proof of Theorem 3.3 follows in part the main ideas in [18]. However, there are still many key differences since we are now dealing with a more complicated

scenario and thus many technical details are much more involved. During the presentation, we will clearly point out both the similarities to and differences from [18].

5.1 Four key conditions

Condition 5.1. Local regularity condition: Let $\mathbf{z} := (\mathbf{h}, \mathbf{x}) \in \mathbb{C}^{s(K+N)}$ and $\nabla \tilde{F}(\mathbf{z}) := \begin{bmatrix} \nabla \tilde{F}_{\mathbf{h}}(\mathbf{z}) \\ \nabla \tilde{F}_{\mathbf{x}}(\mathbf{z}) \end{bmatrix} \in \mathbb{C}^{s(K+N)}$, then

$$\|\nabla \tilde{F}(\mathbf{z})\|^2 \geq \omega[\tilde{F}(\mathbf{z}) - c]_+ \quad (5.1)$$

for $\mathbf{z} \in \mathcal{N}_d \cap \mathcal{N}_\mu \cap \mathcal{N}_\epsilon$ where $\omega = \frac{d_0}{7000}$ and $c = \|\mathbf{e}\|^2 + 2000s\|\mathcal{A}^*(\mathbf{e})\|^2$.

We will prove Condition 5.1 in Section 6.3. Condition 5.1 tells that $\tilde{F}(\mathbf{z}) = 0$ if $\|\nabla \tilde{F}(\mathbf{z})\| = 0$ and $\mathbf{e} = 0$, i.e., all the stationary points inside the basin of attraction are global minima.

Condition 5.2. Local smoothness condition: Let $\mathbf{z} = (\mathbf{h}, \mathbf{x})$ and $\mathbf{w} = (\mathbf{u}, \mathbf{v})$ and there holds

$$\|\tilde{F}(\mathbf{z} + \mathbf{w}) - \tilde{F}(\mathbf{z})\| \leq C_L \|\mathbf{w}\| \quad (5.2)$$

for $\mathbf{z} + \mathbf{w}$ and \mathbf{z} inside $\mathcal{N}_d \cap \mathcal{N}_\mu \cap \mathcal{N}_\epsilon$ where $C_L \approx \mathcal{O}(d_0 s \kappa (1 + \sigma^2)(K + N) \log^2 L)$. The convergence rate is governed by C_L .

The proof of Condition 5.2 can be found in Section 6.4.

Condition 5.3. Local restricted isometry property: Denote $\mathbf{X} = \mathcal{H}(\mathbf{h}, \mathbf{x})$ and $\mathbf{X}_0 = \mathcal{H}(\mathbf{h}_0, \mathbf{x}_0)$. There holds

$$\frac{2}{3} \|\mathbf{X} - \mathbf{X}_0\|_F^2 \leq \|\mathcal{A}(\mathbf{X} - \mathbf{X}_0)\|^2 \leq \frac{3}{2} \|\mathbf{X} - \mathbf{X}_0\|_F^2 \quad (5.3)$$

uniformly all for $(\mathbf{h}, \mathbf{x}) \in \mathcal{N}_d \cap \mathcal{N}_\mu \cap \mathcal{N}_\epsilon$.

Condition 5.3 will be proven in Section 6.2. This condition says that the convergence of the objective function implies the convergence of the iterates.

Condition 5.4. Robustness condition: Let $\varepsilon \leq \frac{1}{15}$ be a predetermined constant. We have

$$\|\mathcal{A}^*(\mathbf{e})\| = \max_{1 \leq i \leq s} \|\mathcal{A}_i^*(\mathbf{e})\| \leq \frac{\varepsilon d_0}{10\sqrt{2s\kappa}}, \quad (5.4)$$

where $\mathbf{e} \sim \mathcal{CN}(0, \frac{\sigma^2 d_0^2}{L})$ if $L \geq C_\gamma \kappa^2 s^2 (K + N) / \varepsilon^2$.

We will prove Condition 5.4 in Section 6.5. We now extract one useful result based on Conditions 5.3 and 5.4. From these two conditions, we are able to produce a good approximation of $F(\mathbf{h}, \mathbf{x})$ for all $(\mathbf{h}, \mathbf{x}) \in \mathcal{N}_d \cap \mathcal{N}_\mu \cap \mathcal{N}_\epsilon$ in terms of δ in (2.9). For $(\mathbf{h}, \mathbf{x}) \in \mathcal{N}_d \cap \mathcal{N}_\mu \cap \mathcal{N}_\epsilon$, the following inequality holds

$$\frac{2}{3} \delta^2 d_0^2 - \frac{\varepsilon \delta d_0^2}{5\sqrt{s\kappa}} + \|\mathbf{e}\|^2 \leq F(\mathbf{h}, \mathbf{x}) \leq \frac{3}{2} \delta^2 d_0^2 + \frac{\varepsilon \delta d_0^2}{5\sqrt{s\kappa}} + \|\mathbf{e}\|^2. \quad (5.5)$$

Note that (5.5) simply follows from

$$F(\mathbf{h}, \mathbf{x}) = \|\mathcal{A}(\mathbf{X} - \mathbf{X}_0)\|_F^2 - 2 \operatorname{Re}(\langle \mathbf{X} - \mathbf{X}_0, \mathcal{A}^*(\mathbf{e}) \rangle) + \|\mathbf{e}\|^2.$$

Note that (5.3) implies $\frac{2}{3} \delta^2 d_0^2 \leq \|\mathcal{A}(\mathbf{X} - \mathbf{X}_0)\|_F^2 \leq \frac{3}{2} \delta^2 d_0^2$. Thus it suffices to estimate the cross-term,

$$\begin{aligned} |\operatorname{Re}(\langle \mathbf{X} - \mathbf{X}_0, \mathcal{A}^*(\mathbf{e}) \rangle)| &\leq \|\mathcal{A}^*(\mathbf{e})\| \|\mathbf{X} - \mathbf{X}_0\|_* = \|\mathcal{A}^*(\mathbf{e})\| \sum_{i=1}^s \|\mathbf{h}_i \mathbf{x}_i^* - \mathbf{h}_{i0} \mathbf{x}_{i0}^*\|_* \\ &\leq \sqrt{2} \|\mathcal{A}^*(\mathbf{e})\| \sum_{i=1}^s \|\mathbf{h}_i \mathbf{x}_i^* - \mathbf{h}_{i0} \mathbf{x}_{i0}^*\|_F \\ &\leq \sqrt{2s} \|\mathcal{A}^*(\mathbf{e})\| \|\mathbf{X} - \mathbf{X}_0\|_F \leq \frac{\varepsilon \delta d_0^2}{10\sqrt{s\kappa}} \end{aligned} \quad (5.6)$$

where $\|\cdot\|_*$ and $\|\cdot\|$ are a pair of dual norms and $\|\mathcal{A}^*(\mathbf{e})\|$ comes from (5.4).

5.2 Outline of the convergence analysis

For the ease of proof, we introduce another neighborhood:

$$\mathcal{N}_{\tilde{F}} = \left\{ (\mathbf{h}, \mathbf{x}) : \tilde{F}(\mathbf{h}, \mathbf{x}) \leq \frac{\varepsilon^2 d_0^2}{3s\kappa^2} + \|\mathbf{e}\|^2 \right\}.$$

Moreover, another reason to consider $\mathcal{N}_{\tilde{F}}$ is based on the fact that gradient descent *only* allows one to make the objective function decrease. In other words, all the iterates $\mathbf{z}^{(t)}$ generated by gradient descent are inside $\mathcal{N}_{\tilde{F}}$ as long as $\mathbf{z}^{(0)} \in \mathcal{N}_{\tilde{F}}$.

On the other hand, it is crucial to note that the decrease of the objective function does not necessarily imply the decrease of the relative error of the iterates. Therefore, we want to construct an initial guess in $\mathcal{N}_\varepsilon \cap \mathcal{N}_{\tilde{F}}$ so that $\mathbf{z}^{(0)}$ is sufficiently close to the ground truth and then analyze the behavior of $\mathbf{z}^{(t)}$.

In the rest of this section, we basically try to prove the following relation:

$$\underbrace{\frac{1}{\sqrt{3}}\mathcal{N}_d \cap \frac{1}{\sqrt{3}}\mathcal{N}_\mu \cap \mathcal{N}_{\frac{2\varepsilon}{5\sqrt{s\kappa}}}}_{\text{Initial guess}} \subset \underbrace{\mathcal{N}_\varepsilon \cap \mathcal{N}_{\tilde{F}}}_{\{\mathbf{z}^{(t)}\}_{t \geq 0} \text{ in } \mathcal{N}_\varepsilon \cap \mathcal{N}_{\tilde{F}}} \subset \underbrace{\mathcal{N}_d \cap \mathcal{N}_\mu \cap \mathcal{N}_\varepsilon}_{\text{Key conditions hold over } \mathcal{N}_d \cap \mathcal{N}_\mu \cap \mathcal{N}_\varepsilon}.$$

Now we give a more detailed explanation of the relation above, which constitutes the main structure of the proof:

1. We will show $\frac{1}{\sqrt{3}}\mathcal{N}_d \cap \frac{1}{\sqrt{3}}\mathcal{N}_\mu \cap \mathcal{N}_{\frac{2\varepsilon}{5\sqrt{s\kappa}}} \subset \mathcal{N}_\varepsilon \cap \mathcal{N}_{\tilde{F}}$ in the proof of Theorem 3.3 in Section 5.3, which is quite straightforward.
2. Lemma 5.5 explains why it holds that $\mathcal{N}_\varepsilon \cap \mathcal{N}_{\tilde{F}} \subset \mathcal{N}_d \cap \mathcal{N}_\mu \cap \mathcal{N}_\varepsilon$ and where the s^2 -bottleneck comes from.
3. Lemma 5.7 implicitly tells us that the iterates $\mathbf{z}^{(t)}$ will remain in $\mathcal{N}_\varepsilon \cap \mathcal{N}_{\tilde{F}}$ if the initial guess $\mathbf{z}^{(0)}$ is inside $\mathcal{N}_\varepsilon \cap \mathcal{N}_{\tilde{F}}$ and $\tilde{F}(\mathbf{z}^{(t)})$ is monotonically decreasing. Lemma 5.8 makes this observation explicit by showing that $\mathbf{z}^{(t)} \in \mathcal{N}_\varepsilon \cap \mathcal{N}_{\tilde{F}}$ implies $\mathbf{z}^{(t+1)} := \mathbf{z}^{(t)} - \eta \nabla \tilde{F}(\mathbf{z}^{(t)}) \in \mathcal{N}_\varepsilon \cap \mathcal{N}_{\tilde{F}}$ if the stepsize η obeys $\eta \leq \frac{1}{C_L}$. Moreover, Lemma 5.8 guarantees sufficient decrease of $\tilde{F}(\mathbf{z}^{(t)})$ in each iteration, which paves the road towards the proof of linear convergence of $\tilde{F}(\mathbf{z}^{(t)})$ and thus $\mathbf{z}^{(t)}$.

Remember that \mathcal{N}_d and \mathcal{N}_μ are both convex sets, and the purpose of introducing regularizers $G_i(\mathbf{h}_i, \mathbf{x}_i)$ is to approximately project the iterates onto $\mathcal{N}_d \cap \mathcal{N}_\mu$. Moreover, we hope that once the iterates are inside \mathcal{N}_ε and inside a sublevel subset $\mathcal{N}_{\tilde{F}}$, they will never escape from $\mathcal{N}_{\tilde{F}} \cap \mathcal{N}_\varepsilon$. Those ideas are fully reflected in the following lemma.

Lemma 5.5. *Assume $0.9d_{i0} \leq d_i \leq 1.1d_{i0}$ and $0.9d_0 \leq d \leq 1.1d_0$, there holds $\mathcal{N}_{\tilde{F}} \subset \mathcal{N}_d \cap \mathcal{N}_\mu$; moreover, under Conditions 5.3 and 5.4, we have $\mathcal{N}_{\tilde{F}} \cap \mathcal{N}_\varepsilon \subset \mathcal{N}_d \cap \mathcal{N}_\mu \cap \mathcal{N}_{\frac{9}{10}\varepsilon}$.*

Proof: If $(\mathbf{h}, \mathbf{x}) \notin \mathcal{N}_d \cap \mathcal{N}_\mu$, by the definition of G in (2.15), at least one component in G exceeds $\rho G_0 \left(\frac{2d_{i0}}{d_i} \right)$. We have

$$\begin{aligned} \tilde{F}(\mathbf{h}, \mathbf{x}) &\geq \rho G_0 \left(\frac{2d_{i0}}{d_i} \right) \geq (d^2 + 2\|\mathbf{e}\|^2) \left(\frac{2d_{i0}}{d_i} - 1 \right)^2 \\ &\geq (2/1.1 - 1)^2 (d^2 + 2\|\mathbf{e}\|^2) \\ &\geq \frac{1}{2}d_0^2 + \|\mathbf{e}\|^2 > \frac{\varepsilon^2 d_0^2}{3s\kappa^2} + \|\mathbf{e}\|^2, \end{aligned}$$

where $\rho \geq d^2 + 2\|\mathbf{e}\|^2$, $0.9d_0 \leq d \leq 1.1d_0$ and $0.9d_{i0} \leq d_i \leq 1.1d_{i0}$. This implies $(\mathbf{h}, \mathbf{x}) \notin \mathcal{N}_{\tilde{F}}$ and hence $\mathcal{N}_{\tilde{F}} \subset \mathcal{N}_d \cap \mathcal{N}_\mu$.

Now we have $(\mathbf{h}, \mathbf{x}) \in \mathcal{N}_d \cap \mathcal{N}_\mu \cap \mathcal{N}_\epsilon$ if $(\mathbf{h}, \mathbf{x}) \in \mathcal{N}_{\tilde{F}} \cap \mathcal{N}_\epsilon$. Applying (5.5) gives

$$\frac{2}{3}\delta^2 d_0^2 - \frac{\epsilon \delta d_0^2}{5\sqrt{s\kappa}} + \|\mathbf{e}\|^2 \leq F(\mathbf{h}, \mathbf{x}) \leq \tilde{F}(\mathbf{h}, \mathbf{x}) \leq \frac{\epsilon^2 d_0^2}{3s\kappa^2} + \|\mathbf{e}\|^2$$

which implies that $\delta \leq \frac{9}{10} \frac{\epsilon}{\sqrt{s\kappa}}$. By definition of δ in (2.9), there holds

$$\frac{81\epsilon^2}{100s\kappa^2} \geq \delta^2 = \frac{\sum_{i=1}^s \delta_i^2 d_{i0}^2}{\sum_{i=1}^s d_{i0}^2} \geq \frac{\sum_{i=1}^s \delta_i^2}{s\kappa^2} \geq \frac{1}{s\kappa^2} \max_{1 \leq i \leq s} \delta_i^2, \quad (5.7)$$

which gives $\delta_i \leq \frac{9}{10}\epsilon$ and $(\mathbf{h}, \mathbf{x}) \in \mathcal{N}_{\frac{9}{10}\epsilon}$. \square

Remark 5.6. *The s^2 -bottleneck comes from (5.7). If $\delta \leq \epsilon$ is small, we cannot guarantee that each δ_i is also smaller than ϵ . Just consider the simplest case when all d_{i0} are the same: then $d_0^2 = \sum_{i=1}^s d_{i0}^2 = s d_{i0}^2$ and there holds*

$$\epsilon^2 \geq \delta^2 = \frac{1}{s} \sum_{i=1}^s \delta_i^2.$$

Obviously, we cannot conclude that $\max \delta_i \leq \epsilon$ but only say that $\delta_i \leq \sqrt{s}\epsilon$. This is why we require $\delta = \mathcal{O}(\frac{\epsilon}{\sqrt{s}})$ to ensure $\delta_i \leq \epsilon$, which gives s^2 -dependence in L .

Lemma 5.7. *Denote $\mathbf{z}_1 = (\mathbf{h}_1, \mathbf{x}_1)$ and $\mathbf{z}_2 = (\mathbf{h}_2, \mathbf{x}_2)$. Let $\mathbf{z}(\lambda) := (1 - \lambda)\mathbf{z}_1 + \lambda\mathbf{z}_2$. If $\mathbf{z}_1 \in \mathcal{N}_\epsilon$ and $\mathbf{z}(\lambda) \in \mathcal{N}_{\tilde{F}}$ for all $\lambda \in [0, 1]$, we have $\mathbf{z}_2 \in \mathcal{N}_\epsilon$.*

Proof: We prove it by contradiction based on $\mathcal{N}_{\tilde{F}} \cap \mathcal{N}_\epsilon \subset \mathcal{N}_{\frac{9}{10}\epsilon}$ in Lemma 5.5. Suppose that $\mathbf{z}_2 \notin \mathcal{N}_\epsilon$ and $\mathbf{z}_1 \in \mathcal{N}_\epsilon$, and there exists $\mathbf{z}(\lambda_0) := (\mathbf{h}(\lambda_0), \mathbf{x}(\lambda_0)) \in \mathcal{N}_\epsilon$ for some $\lambda_0 \in [0, 1]$, such that $\max_{1 \leq i \leq s} \frac{\|\mathbf{h}_i \mathbf{x}_i^* - \mathbf{h}_{i0} \mathbf{x}_{i0}^*\|_F}{d_{i0}} = \epsilon$. Therefore, $\mathbf{z}(\lambda_0) \in \mathcal{N}_{\tilde{F}} \cap \mathcal{N}_\epsilon$ and Lemma 5.5 implies $\max_{1 \leq i \leq s} \frac{\|\mathbf{h}_i \mathbf{x}_i^* - \mathbf{h}_{i0} \mathbf{x}_{i0}^*\|_F}{d_{i0}} \leq \frac{9}{10}\epsilon$, which contradicts $\max_{1 \leq i \leq s} \frac{\|\mathbf{h}_i \mathbf{x}_i^* - \mathbf{h}_{i0} \mathbf{x}_{i0}^*\|_F}{d_{i0}} = \epsilon$. \square

Lemma 5.8. *Let the stepsize $\eta \leq \frac{1}{C_L}$, $\mathbf{z}^{(t)} := (\mathbf{u}^{(t)}, \mathbf{v}^{(t)}) \in \mathbb{C}^{s(K+N)}$ and C_L be the Lipschitz constant of $\nabla \tilde{F}(\mathbf{z})$ over $\mathcal{N}_d \cap \mathcal{N}_\mu \cap \mathcal{N}_\epsilon$ in (5.2). If $\mathbf{z}^{(t)} \in \mathcal{N}_\epsilon \cap \mathcal{N}_{\tilde{F}}$, we have $\mathbf{z}^{(t+1)} \in \mathcal{N}_\epsilon \cap \mathcal{N}_{\tilde{F}}$ and*

$$\tilde{F}(\mathbf{z}^{(t+1)}) \leq \tilde{F}(\mathbf{z}^{(t)}) - \eta \|\nabla \tilde{F}(\mathbf{z}^{(t)})\|^2 \quad (5.8)$$

where $\mathbf{z}^{(t+1)} = \mathbf{z}^{(t)} - \eta \nabla \tilde{F}(\mathbf{z}^{(t)})$.

Remark 5.9. *This lemma tells us that once $\mathbf{z}^{(t)} \in \mathcal{N}_\epsilon \cap \mathcal{N}_{\tilde{F}}$, the next iterate $\mathbf{z}^{(t+1)} = \mathbf{z}^{(t)} - \eta \nabla \tilde{F}(\mathbf{z}^{(t)})$ is also inside $\mathcal{N}_\epsilon \cap \mathcal{N}_{\tilde{F}}$ as long as the stepsize $\eta \leq \frac{1}{C_L}$. In other words, $\mathcal{N}_\epsilon \cap \mathcal{N}_{\tilde{F}}$ is in fact a stronger version of the basin of attraction. Moreover, the objective function will decay sufficiently in each step as long as we can control the lower bound of the $\nabla \tilde{F}$, which is guaranteed by the Local Regularity Condition 5.3.*

Proof: Let $\phi(\tau) := \tilde{F}(\mathbf{z}^{(t)} - \tau \nabla \tilde{F}(\mathbf{z}^{(t)}))$, $\phi(0) = \tilde{F}(\mathbf{z}^{(t)})$ and consider the following quantity:

$$\tau_{\max} := \max\{\mu : \phi(\tau) \leq \tilde{F}(\mathbf{z}^{(t)}), 0 \leq \tau \leq \mu\},$$

where τ_{\max} is the largest stepsize such that the objective function $\tilde{F}(\mathbf{z})$ evaluated at any point over the whole line segment $\{\mathbf{z}^{(t)} - \tau \nabla \tilde{F}(\mathbf{z}^{(t)}), 0 \leq \tau \leq \tau_{\max}\}$ is not greater than $\tilde{F}(\mathbf{z}^{(t)})$. Now we will show $\tau_{\max} \geq \frac{1}{C_L}$. Obviously, if $\|\nabla \tilde{F}(\mathbf{z}^{(t)})\| = 0$, it holds automatically.

Consider $\|\nabla \tilde{F}(\mathbf{z}^{(t)})\| \neq 0$ and assume $\tau_{\max} < \frac{1}{C_L}$. First note that,

$$\frac{d}{d\tau} \phi(\tau) < 0 \implies \tau_{\max} > 0.$$

By the definition of τ_{\max} , there holds $\phi(\tau_{\max}) = \phi(0)$ since $\phi(\tau)$ is a continuous function w.r.t. τ . Lemma 5.7 implies

$$\{\mathbf{z}^{(t)} - \tau \nabla \tilde{F}(\mathbf{z}^{(t)}), 0 \leq \tau \leq \tau_{\max}\} \subseteq \mathcal{N}_\epsilon \cap \mathcal{N}_{\tilde{F}}.$$

Now we apply Lemma 6.19, the modified descent lemma, and obtain

$$\tilde{F}(\mathbf{z}^{(t)} - \tau_{\max} \nabla \tilde{F}(\mathbf{z}^{(t)})) \leq \tilde{F}(\mathbf{z}^{(t)}) - (2\tau_{\max} - C_L \tau_{\max}^2) \|\tilde{F}(\mathbf{z}^{(t)})\|^2 \leq \tilde{F}(\mathbf{z}^{(t)}) - \tau_{\max} \|\tilde{F}(\mathbf{z}^{(t)})\|^2$$

where $C_L \tau_{\max} \leq 1$. In other words, $\phi(\tau_{\max}) \leq \tilde{F}(\mathbf{z}^{(t)} - \tau_{\max} \nabla \tilde{F}(\mathbf{z}^{(t)})) < \tilde{F}(\mathbf{z}^{(t)}) = \phi(0)$ contradicts $\phi(\tau_{\max}) = \phi(0)$.

Therefore, we conclude that $\tau_{\max} \geq \frac{1}{C_L}$. For any $\eta \leq \frac{1}{C_L}$, Lemma 5.7 implies

$$\{\mathbf{z}^{(t)} - \tau \nabla \tilde{F}(\mathbf{z}^{(t)}), 0 \leq \tau \leq \eta\} \subseteq \mathcal{N}_\epsilon \cap \mathcal{N}_{\tilde{F}}$$

and applying Lemma 6.19 gives

$$\tilde{F}(\mathbf{z}^{(t)} - \eta \nabla \tilde{F}(\mathbf{z}^{(t)})) \leq \tilde{F}(\mathbf{z}^{(t)}) - (2\eta - C_L \eta^2) \|\tilde{F}(\mathbf{z}^{(t)})\|^2 \leq \tilde{F}(\mathbf{z}^{(t)}) - \eta \|\tilde{F}(\mathbf{z}^{(t)})\|^2.$$

□

5.3 Proof of Theorem 3.3

Combining all the considerations above, we now prove Theorem 3.3 to conclude this section.

Proof: The proof consists of three parts:

Part I: Proof of $\mathbf{z}^{(0)} := (\mathbf{u}^{(0)}, \mathbf{v}^{(0)}) \in \mathcal{N}_\epsilon \cap \mathcal{N}_{\tilde{F}}$. From the assumption of Theorem 3.3,

$$\mathbf{z}^{(0)} \in \frac{1}{\sqrt{3}} \mathcal{N}_d \cap \frac{1}{\sqrt{3}} \mathcal{N}_\mu \cap \mathcal{N}_{\frac{2\epsilon}{5\sqrt{s\kappa}}}.$$

First we show $G(\mathbf{u}^{(0)}, \mathbf{v}^{(0)}) = 0$: for $0 \leq i \leq s$ and the definition of \mathcal{N}_d and \mathcal{N}_μ ,

$$\frac{\|\mathbf{u}_i^{(0)}\|^2}{2d_i} \leq \frac{2d_{i0}}{3d_i} < 1, \quad \frac{L|\mathbf{b}_i^* \mathbf{u}_i^{(0)}|^2}{8d_i \mu^2} \leq \frac{L}{8d_i \mu^2} \cdot \frac{16d_{i0} \mu^2}{3L} \leq \frac{2d_{i0}}{3d_i} < 1,$$

where $\|\mathbf{u}_i^{(0)}\| \leq \frac{2\sqrt{d_{i0}}}{\sqrt{3}}$, $\sqrt{L} \|\mathbf{B} \mathbf{u}_i^{(0)}\|_\infty \leq \frac{4\sqrt{d_{i0}} \mu}{\sqrt{3}}$ and $\frac{9}{10} d_{i0} \leq d_i \leq \frac{11}{10} d_{i0}$. Therefore

$$G_0 \left(\frac{\|\mathbf{u}_i^{(0)}\|^2}{2d_i} \right) = G_0 \left(\frac{\|\mathbf{v}_i^{(0)}\|^2}{2d_i} \right) = G_0 \left(\frac{L|\mathbf{b}_i^* \mathbf{u}_i^{(0)}|^2}{8d_i \mu^2} \right) = 0$$

for all $1 \leq l \leq L$ and $G(\mathbf{u}^{(0)}, \mathbf{v}^{(0)}) = 0$.

For $\mathbf{z}^{(0)} = (\mathbf{u}^{(0)}, \mathbf{v}^{(0)}) \in \mathcal{N}_{\frac{2\epsilon}{5\sqrt{s\kappa}}}$, we have $\delta(\mathbf{z}^{(0)}) := \frac{\sqrt{\sum_{i=1}^s \delta_i^2 d_{i0}^2}}{d_0} \leq \frac{2\epsilon}{5\sqrt{s\kappa}}$. By (5.5), there holds $\delta(\mathbf{z}^{(0)}) \leq \frac{2\epsilon}{5\sqrt{s\kappa}}$ and $G(\mathbf{u}^{(0)}, \mathbf{v}^{(0)}) = 0$,

$$\tilde{F}(\mathbf{u}^{(0)}, \mathbf{v}^{(0)}) = F(\mathbf{u}^{(0)}, \mathbf{v}^{(0)}) \leq \|\mathbf{e}\|^2 + \frac{3}{2} \delta^2(\mathbf{z}^{(0)}) d_0^2 + \frac{\epsilon \delta(\mathbf{z}^{(0)}) d_0^2}{5\sqrt{s\kappa}} \leq \|\mathbf{e}\|^2 + \frac{\epsilon^2 d_0^2}{3s\kappa^2}$$

and hence $\mathbf{z}^{(0)} = (\mathbf{u}^{(0)}, \mathbf{v}^{(0)}) \in \mathcal{N}_\epsilon \cap \mathcal{N}_{\tilde{F}}$.

Part II: The linear convergence of the objective function $\tilde{F}(\mathbf{z}^{(t)})$. Denote $\mathbf{z}^{(t)} := (\mathbf{u}^{(t)}, \mathbf{v}^{(t)})$. Note that $\mathbf{z}^{(0)} \in \mathcal{N}_\epsilon \cap \mathcal{N}_{\tilde{F}}$, Lemma 5.8 implies $\mathbf{z}^{(t)} \in \mathcal{N}_\epsilon \cap \mathcal{N}_{\tilde{F}}$ for all $t \geq 0$ by induction if $\eta \leq \frac{1}{C_L}$. Moreover, combining Condition 5.1 with Lemma 5.8 leads to

$$\tilde{F}(\mathbf{z}^{(t)}) \leq \tilde{F}(\mathbf{z}^{(t-1)}) - \eta \omega \left[\tilde{F}(\mathbf{z}^{(t-1)}) - c \right]_+, \quad t \geq 1$$

with $c = \|\mathbf{e}\|^2 + a \|\mathcal{A}^*(\mathbf{e})\|^2$ and $a = 2000s$. Therefore, by induction, we have

$$\left[\tilde{F}(\mathbf{z}^{(t)}) - c \right]_+ \leq (1 - \eta \omega) \left[\tilde{F}(\mathbf{z}^{(t-1)}) - c \right]_+ \leq (1 - \eta \omega)^t \left[\tilde{F}(\mathbf{z}^{(0)}) - c \right]_+ \leq \frac{\epsilon^2 d_0^2}{3s\kappa^2} (1 - \eta \omega)^t$$

where $\tilde{F}(\mathbf{z}^{(0)}) \leq \frac{\epsilon^2 d_0^2}{3s\kappa^2} + \|\mathbf{e}\|^2$ and $\left[\tilde{F}(\mathbf{z}^{(0)}) - c \right]_+ \leq \left[\frac{1}{3s\kappa^2} \epsilon^2 d_0^2 - a \|\mathcal{A}^*(\mathbf{e})\|^2 \right]_+ \leq \frac{\epsilon^2 d_0^2}{3s\kappa^2}$. Now we conclude that $\left[\tilde{F}(\mathbf{z}^{(t)}) - c \right]_+$ converges to 0 linearly.

Part III: The linear convergence of the objective function $\tilde{F}(\mathbf{z}^{(t)})$. Denote

$$\delta(\mathbf{z}^{(t)}) := \frac{\|\mathcal{H}(\mathbf{u}^{(t)}, \mathbf{v}^{(t)}) - \mathcal{H}(\mathbf{h}_0, \mathbf{x}_0)\|_F}{d_0}.$$

Note that $\mathbf{z}^{(t)} \in \mathcal{N}_\epsilon \cap \mathcal{N}_{\tilde{F}} \subseteq \mathcal{N}_d \cap \mathcal{N}_\mu \cap \mathcal{N}_\epsilon$ and over $\mathcal{N}_d \cap \mathcal{N}_\mu \cap \mathcal{N}_\epsilon$, there holds $F_0(\mathbf{z}^{(t)}) \geq \frac{2}{3}\delta^2(\mathbf{z}^{(t)})d_0^2$ which follows from Local RIP Condition in (5.3) and $F_0(\mathbf{z}^{(t)})$ defined in (2.12). Moreover

$$\begin{aligned} \tilde{F}(\mathbf{z}^{(t)}) - \|\mathbf{e}\|^2 &\geq F_0(\mathbf{z}^{(t)}) - 2 \operatorname{Re} \left(\langle \mathcal{A}^*(\mathbf{e}), \mathcal{H}(\mathbf{u}^{(0)}, \mathbf{v}^{(0)}) - \mathcal{H}(\mathbf{h}_0, \mathbf{x}_0) \rangle \right) \\ &\geq \frac{2}{3}\delta^2(\mathbf{z}^{(t)})d_0^2 - 2\sqrt{2s}\|\mathcal{A}^*(\mathbf{e})\|\delta(\mathbf{z}^{(t)})d_0 \end{aligned}$$

where $G(\mathbf{z}^{(t)}) \geq 0$ and the second inequality follows from (5.6). There holds

$$\frac{2}{3}\delta^2(\mathbf{z}^{(t)})d_0^2 - 2\sqrt{2s}\|\mathcal{A}^*(\mathbf{e})\|\delta(\mathbf{z}^{(t)})d_0 - a\|\mathcal{A}^*(\mathbf{e})\|^2 \leq \left[\tilde{F}(\mathbf{z}^{(t)}) - c \right]_+ \leq \frac{\varepsilon^2 d_0^2}{3s\kappa^2}(1 - \eta\omega)^t$$

and equivalently,

$$\left| \delta(\mathbf{z}^{(t)})d_0 - \frac{3\sqrt{2}}{2}\|\mathcal{A}^*(\mathbf{e})\| \right|^2 \leq \frac{\varepsilon^2 d_0^2}{2s\kappa^2}(1 - \eta\omega)^t + \left(\frac{3}{2}a + \frac{9}{2} \right) \|\mathcal{A}^*(\mathbf{e})\|^2.$$

Solving the inequality above for $\delta(\mathbf{z}^{(t)})$, we have

$$\begin{aligned} \delta(\mathbf{z}^{(t)})d_0 &\leq \frac{\varepsilon d_0}{\sqrt{2s\kappa^2}}(1 - \eta\omega)^{t/2} + \left(\frac{3\sqrt{2}}{2} + \sqrt{\frac{3}{2}a + \frac{9}{2}} \right) \|\mathcal{A}^*(\mathbf{e})\| \\ &\leq \frac{\varepsilon d_0}{\sqrt{2s\kappa^2}}(1 - \eta\omega)^{t/2} + 60\sqrt{s}\|\mathcal{A}^*(\mathbf{e})\| \end{aligned} \quad (5.9)$$

where $a = 2000s$. Let $d^{(t)} := \sqrt{\sum_{i=1}^s \|\mathbf{u}_i^{(t)}\|^2 \|\mathbf{v}_i^{(t)}\|^2}$ for $t \in \mathbb{Z}_{\geq 0}$. By (5.9) and triangle inequality, we immediately obtain $|d^{(t)} - d_0| \leq \frac{\varepsilon d_0}{\sqrt{2s\kappa^2}}(1 - \eta\omega)^{t/2} + 60\sqrt{s}\|\mathcal{A}^*(\mathbf{e})\|$. \square

6 Proof of the four conditions

This section is devoted to proving the four key conditions introduced in Section 5. The *local smoothness condition* and the *robustness condition* are relatively less challenging to deal with. The more difficult part is to show the *local regularity condition* and the *local isometry property*. The key to solve those problems is to understand how the matrix-valued linear operator \mathcal{A} in (2.7) behaves on block-diagonal matrices, such as $\mathcal{H}(\mathbf{h}, \mathbf{x})$, $\mathcal{H}(\mathbf{h}_0, \mathbf{x}_0)$ and $\mathcal{H}(\mathbf{h}, \mathbf{x}) - \mathcal{H}(\mathbf{h}_0, \mathbf{x}_0)$. In particular, when $s = 1$, all those matrices become rank-1 matrices, which have been well discussed in our previous work [18].

First of all, we define the linear subspace $T_i \subset \mathbb{C}^{K \times N}$ along with its orthogonal complement for $1 \leq i \leq s$ as

$$\begin{aligned} T_i &:= \{ \mathbf{Z}_i \in \mathbb{C}^{K \times N} : \mathbf{Z}_i = \mathbf{h}_{i0}\mathbf{v}_i^* + \mathbf{u}_i\mathbf{x}_{i0}^*, \quad \mathbf{u}_i \in \mathbb{C}^K, \mathbf{v}_i \in \mathbb{C}^N \}, \\ T_i^\perp &:= \left\{ \left(\mathbf{I}_K - \frac{\mathbf{h}_{i0}\mathbf{h}_{i0}^*}{d_{i0}} \right) \mathbf{Z}_i \left(\mathbf{I}_N - \frac{\mathbf{x}_{i0}\mathbf{x}_{i0}^*}{d_{i0}} \right) : \mathbf{Z}_i \in \mathbb{C}^{K \times N} \right\} \end{aligned}$$

where $\|\mathbf{h}_{i0}\| = \|\mathbf{x}_{i0}\| = \sqrt{d_{i0}}$. In particular, $\mathbf{h}_{i0}\mathbf{x}_{i0}^* \in T_i$ for all $1 \leq i \leq s$.

The proof also requires us to consider block-diagonal matrices whose i -th block belongs to T_i (or T_i^\perp). Let $\mathbf{Z} = \operatorname{blkdiag}(\mathbf{Z}_1, \dots, \mathbf{Z}_s) \in \mathbb{C}^{Ks \times Ns}$ be a block-diagonal matrix and say $\mathbf{Z} \in T$ if

$$T := \{ \operatorname{blkdiag}(\{\mathbf{Z}_i\}_{i=1}^s) \mid \mathbf{Z}_i \in T_i \}$$

and $\mathbf{Z} \in T^\perp$ if

$$T^\perp := \{\text{blkdiag}(\{\mathbf{Z}_i\}_{i=1}^s) \mid \mathbf{Z}_i \in T_i^\perp\}$$

where both T and T^\perp are subsets in $\mathbb{C}^{Ks \times Ns}$ and $\mathcal{H}(\mathbf{h}_0, \mathbf{x}_0) \in T$.

Now we take a closer look at a special case of block-diagonal matrices, i.e., $\mathcal{H}(\mathbf{h}, \mathbf{x})$ and calculate its projection onto T and T^\perp respectively and it suffices to consider $\mathcal{P}_{T_i}(\mathbf{h}_i \mathbf{x}_i^*)$ and $\mathcal{P}_{T_i^\perp}(\mathbf{h}_i \mathbf{x}_i^*)$. For each block $\mathbf{h}_i \mathbf{x}_i^*$ and $1 \leq i \leq s$, there are unique orthogonal decompositions

$$\mathbf{h}_i := \alpha_{i1} \mathbf{h}_{i0} + \tilde{\mathbf{h}}_i, \quad \mathbf{x}_i := \alpha_{i2} \mathbf{x}_{i0} + \tilde{\mathbf{x}}_i, \quad (6.1)$$

where $\mathbf{h}_{i0} \perp \tilde{\mathbf{h}}_i$ and $\mathbf{x}_{i0} \perp \tilde{\mathbf{x}}_i$. It is important to note that $\alpha_{i1} = \alpha_{i1}(\mathbf{h}_i) = \frac{\langle \mathbf{h}_{i0}, \mathbf{h}_i \rangle}{d_{i0}}$ and $\alpha_{i2} = \alpha_{i2}(\mathbf{x}_i) = \frac{\langle \mathbf{x}_{i0}, \mathbf{x}_i \rangle}{d_{i0}}$ and thus α_{i1} and α_{i2} are functions of \mathbf{h}_i and \mathbf{x}_i respectively. Immediately, we have the following matrix orthogonal decomposition for $\mathbf{h}_i \mathbf{x}_i^*$ onto T_i and T_i^\perp ,

$$\mathbf{h}_i \mathbf{x}_i^* - \mathbf{h}_{i0} \mathbf{x}_{i0}^* = \underbrace{(\alpha_{i1} \overline{\alpha_{i2}} - 1) \mathbf{h}_{i0} \mathbf{x}_{i0}^* + \overline{\alpha_{i2}} \tilde{\mathbf{h}}_i \mathbf{x}_{i0}^* + \alpha_{i1} \mathbf{h}_{i0} \tilde{\mathbf{x}}_i^*}_{\text{belongs to } T_i} + \underbrace{\tilde{\mathbf{h}}_i \tilde{\mathbf{x}}_i^*}_{\text{belongs to } T_i^\perp} \quad (6.2)$$

where the first three components are in T_i while $\tilde{\mathbf{h}}_i \tilde{\mathbf{x}}_i^* \in T_i^\perp$.

6.1 Key lemmata

From the decomposition in (6.1) and (6.2), we want to analyze how $\|\tilde{\mathbf{h}}_i\|$, $\|\tilde{\mathbf{x}}_i\|$, α_{i1} and α_{i2} depend on $\delta_i = \frac{\|\mathbf{h}_i \mathbf{x}_i^* - \mathbf{h}_{i0} \mathbf{x}_{i0}^*\|_F}{d_{i0}}$ if $\delta_i < 1$. The following lemma answers this question, which can be viewed as an application of singular value/vector perturbation theory [33] applied to rank-1 matrices. From the lemma below, we can see that if $\mathbf{h}_i \mathbf{x}_i^*$ is close to $\mathbf{h}_{i0} \mathbf{x}_{i0}^*$, then $\mathcal{P}_{T_i^\perp}(\mathbf{h}_i \mathbf{x}_i^*)$ is in fact very small (of order $\mathcal{O}(\delta_i^2 d_{i0})$).

Lemma 6.1. (Lemma 5.9 in [18]) *Recall that $\|\mathbf{h}_{i0}\| = \|\mathbf{x}_{i0}\| = \sqrt{d_{i0}}$. If $\delta_i := \frac{\|\mathbf{h}_i \mathbf{x}_i^* - \mathbf{h}_{i0} \mathbf{x}_{i0}^*\|_F}{d_{i0}} < 1$, we have the following useful bounds*

$$|\alpha_{i1}| \leq \frac{\|\mathbf{h}_i\|}{\|\mathbf{h}_{i0}\|}, \quad |\alpha_{i1} \overline{\alpha_{i2}} - 1| \leq \delta_i,$$

and

$$\|\tilde{\mathbf{h}}_i\| \leq \frac{\delta_i}{1 - \delta_i} \|\mathbf{h}_i\|, \quad \|\tilde{\mathbf{x}}_i\| \leq \frac{\delta_i}{1 - \delta_i} \|\mathbf{x}_i\|, \quad \|\tilde{\mathbf{h}}_i\| \|\tilde{\mathbf{x}}_i\| \leq \frac{\delta_i^2}{2(1 - \delta_i)} d_{i0}.$$

Moreover, if $\|\mathbf{h}_i\| \leq 2\sqrt{d_{i0}}$ and $\sqrt{L} \|\mathbf{B} \mathbf{h}_i\|_\infty \leq 4\mu\sqrt{d_{i0}}$, i.e., $\mathbf{h}_i \in \mathcal{N}_d \cap \mathcal{N}_\mu$, we have $\sqrt{L} \|\mathbf{B}_i \tilde{\mathbf{h}}_i\|_\infty \leq 6\mu\sqrt{d_{i0}}$.

Now we start to focus on several results related to the linear operator \mathcal{A} .

Lemma 6.2. (Operator norm of \mathcal{A}). *For \mathcal{A} defined in (2.7), there holds*

$$\|\mathcal{A}\| \leq \sqrt{s(N \log(NL/2) + (\gamma + \log s) \log L)} \quad (6.3)$$

with probability at least $1 - L^{-\gamma}$.

Proof: Note that $\mathcal{A}_i(\mathbf{Z}_i) := \{\mathbf{b}_i^* \mathbf{Z}_i \mathbf{a}_{il}\}_{l=1}^L$ in (2.2). Lemma 1 in [1] implies

$$\|\mathcal{A}_i\| \leq \sqrt{N \log(NL/2) + \gamma' \log L}$$

with probability at least $1 - L^{-\gamma'}$. By taking the union bound over $1 \leq i \leq s$,

$$\max \|\mathcal{A}_i\| \leq \sqrt{N \log(NL/2) + (\gamma + \log s) \log L}$$

with probability at least $1 - sL^{-\gamma - \log s} \geq 1 - L^{-\gamma}$.

For \mathcal{A} defined in (2.7), applying the triangle inequality gives

$$\|\mathcal{A}(\mathbf{Z})\| = \left\| \sum_{i=1}^s \mathcal{A}_i(\mathbf{Z}_i) \right\| \leq \sum_{i=1}^s \|\mathcal{A}_i\| \|\mathbf{Z}_i\|_F \leq \max_{1 \leq i \leq s} \|\mathcal{A}_i\| \sqrt{s \sum_{i=1}^s \|\mathbf{Z}_i\|_F^2} = \sqrt{s} \max_{1 \leq i \leq s} \|\mathcal{A}_i\| \|\mathbf{Z}\|_F$$

where $\mathbf{Z} = \text{blkdiag}(\mathbf{Z}_1, \dots, \mathbf{Z}_s) \in \mathbb{C}^{Ks \times Ns}$. Therefore,

$$\|\mathcal{A}\| \leq \sqrt{s} \max_{1 \leq i \leq s} \|\mathcal{A}_i\| \leq \sqrt{s(N \log(NL/2) + (\gamma + \log s) \log L)}$$

with probability at least $1 - L^{-\gamma}$. \square

Lemma 6.3. (Restricted isometry property for \mathcal{A} on T). *\mathcal{A} restricted on T is well-conditioned, i.e.,*

$$\|\mathcal{P}_T \mathcal{A}^* \mathcal{A} \mathcal{P}_T - \mathcal{P}_T\| \leq \frac{1}{10} \quad (6.4)$$

where \mathcal{P}_T is the projection operator from $\mathbb{C}^{Ks \times Ns}$ onto T , given $L \geq C_\gamma s^2 \max\{K, \mu_h^2 N\} \log^2 L$ with probability at least $1 - L^{-\gamma}$.

Remark 6.4. Here $\mathcal{A} \mathcal{P}_T$ and $\mathcal{P}_T \mathcal{A}^*$ are defined as

$$\mathcal{A} \mathcal{P}_T(\mathbf{Z}) = \sum_{i=1}^s \mathcal{A}_i(\mathcal{P}_{T_i}(\mathbf{Z}_i)), \quad \mathcal{P}_T \mathcal{A}^*(\mathbf{z}) = \text{blkdiag}(\mathcal{P}_{T_1}(\mathcal{A}_1^*(\mathbf{z})), \dots, \mathcal{P}_{T_s}(\mathcal{A}_s^*(\mathbf{z})))$$

respectively where \mathbf{Z} is a block-diagonal matrix and $\mathbf{z} \in \mathbb{C}^L$.

Proof: From Corollary 5.3 and 5.8 in [19], we know that

$$\|\mathcal{P}_{T_i} \mathcal{A}_i^* \mathcal{A}_j \mathcal{P}_{T_j}\| \leq \frac{1}{10s}, \quad \forall i \neq j; \quad \|\mathcal{P}_{T_i} \mathcal{A}_i^* \mathcal{A}_i \mathcal{P}_{T_i} - \mathcal{P}_{T_i}\| \leq \frac{1}{10s}, \quad \forall 1 \leq i \leq s \quad (6.5)$$

with probability at least $1 - L^{-\gamma+1}$ if $L \geq C_\gamma s^2 \max\{K, \mu_h^2 N\} \log^2 L \log(s+1)$.

For any block diagonal matrix $\mathbf{Z} = \text{blkdiag}(\mathbf{Z}_1, \dots, \mathbf{Z}_s) \in \mathbb{C}^{Ks \times Ns}$ and $\mathbf{Z}_i \in \mathbb{C}^{K \times N}$,

$$\begin{aligned} \langle \mathbf{Z}, \mathcal{P}_T \mathcal{A}^* \mathcal{A} \mathcal{P}_T(\mathbf{Z}) - \mathcal{P}_T(\mathbf{Z}) \rangle &= \sum_{1 \leq i, j \leq s} \langle \mathcal{A}_i \mathcal{P}_{T_i}(\mathbf{Z}_i), \mathcal{A}_j \mathcal{P}_{T_j}(\mathbf{Z}_j) \rangle - \|\mathcal{P}_T(\mathbf{Z})\|_F^2 \\ &= \sum_{i=1}^s \langle \mathbf{Z}_i, \mathcal{P}_{T_i} \mathcal{A}_i^* \mathcal{A}_i \mathcal{P}_{T_i}(\mathbf{Z}_i) - \mathcal{P}_{T_i}(\mathbf{Z}_i) \rangle + \sum_{i \neq j} \langle \mathcal{A}_i \mathcal{P}_{T_i}(\mathbf{Z}_i), \mathcal{A}_j \mathcal{P}_{T_j}(\mathbf{Z}_j) \rangle. \end{aligned} \quad (6.6)$$

Using (6.5), the following two inequalities hold,

$$\begin{aligned} |\langle \mathbf{Z}_i, \mathcal{P}_{T_i} \mathcal{A}_i^* \mathcal{A}_i \mathcal{P}_{T_i}(\mathbf{Z}_i) - \mathcal{P}_{T_i}(\mathbf{Z}_i) \rangle| &\leq \|\mathcal{P}_{T_i} \mathcal{A}_i^* \mathcal{A}_i \mathcal{P}_{T_i} - \mathcal{P}_{T_i}\| \|\mathbf{Z}_i\|_F^2 \leq \frac{\|\mathbf{Z}_i\|_F^2}{10s}, \\ |\langle \mathcal{A}_i \mathcal{P}_{T_i}(\mathbf{Z}_i), \mathcal{A}_j \mathcal{P}_{T_j}(\mathbf{Z}_j) \rangle| &\leq \|\mathcal{P}_{T_i} \mathcal{A}_i^* \mathcal{A}_j \mathcal{P}_{T_j}\| \|\mathbf{Z}_i\|_F \|\mathbf{Z}_j\|_F \leq \frac{\|\mathbf{Z}_i\|_F \|\mathbf{Z}_j\|_F}{10s}. \end{aligned}$$

After substituting both estimates into (6.6), we have

$$|\langle \mathbf{Z}, \mathcal{P}_T \mathcal{A}^* \mathcal{A} \mathcal{P}_T(\mathbf{Z}) - \mathcal{P}_T(\mathbf{Z}) \rangle| \leq \sum_{1 \leq i, j \leq s} \frac{\|\mathbf{Z}_i\|_F \|\mathbf{Z}_j\|_F}{10s} \leq \frac{1}{10s} \left(\sum_{i=1}^s \|\mathbf{Z}_i\|_F \right)^2 \leq \frac{\|\mathbf{Z}\|_F^2}{10}.$$

\square

Finally, we show how \mathcal{A} behaves when applied to block-diagonal matrices $\mathbf{X} = \mathcal{H}(\mathbf{h}, \mathbf{x})$.

Lemma 6.5. (*\mathcal{A} restricted on block-diagonal matrices with rank-1 blocks*).

Consider $\mathbf{X} = \mathcal{H}(\mathbf{h}, \mathbf{x})$ and

$$\sigma_{\max}^2(\mathbf{h}, \mathbf{x}) := \max_{1 \leq l \leq L} \sum_{i=1}^s |\mathbf{b}_l^* \mathbf{h}_i|^2 \|\mathbf{x}_i\|^2. \quad (6.7)$$

Conditioned on (6.3), we have

$$\|\mathcal{A}(\mathbf{X})\|^2 \leq \frac{4}{3} \|\mathbf{X}\|_F^2 + 2\sqrt{2s\|\mathbf{X}\|_F^2 \sigma_{\max}^2(\mathbf{h}, \mathbf{x})(K+N) \log L + 8s\sigma_{\max}^2(\mathbf{h}, \mathbf{x})(K+N) \log L}, \quad (6.8)$$

uniformly for any $\mathbf{h} \in \mathbb{C}^{Ks}$ and $\mathbf{x} \in \mathbb{C}^{Ns}$ with probability at least $1 - \frac{1}{\gamma} \exp(-s(K+N))$ if $L \geq C_\gamma s(K+N) \log L$. Here $\|\mathbf{X}\|_F^2 = \|\mathcal{H}(\mathbf{h}, \mathbf{x})\|_F^2 = \sum_{i=1}^s \|\mathbf{h}_i\|^2 \|\mathbf{x}_i\|^2$.

Remark 6.6. Here are a few more explanations and facts about $\sigma_{\max}^2(\mathbf{h}, \mathbf{x})$. Note that $\|\mathcal{A}(\mathbf{X})\|^2$ is the sum of L sub-exponential random variables, i.e.,

$$\|\mathcal{A}(\mathbf{X})\|^2 = \sum_{l=1}^L \left| \sum_{i=1}^s \mathbf{b}_l^* \mathbf{h}_i \mathbf{x}_i^* \mathbf{a}_{il} \right|^2. \quad (6.9)$$

Here $\sigma_{\max}^2(\mathbf{h}, \mathbf{x})$ corresponds to the largest expectation of all those components in $\|\mathcal{A}(\mathbf{X})\|^2$.

For $\sigma_{\max}^2(\mathbf{h}, \mathbf{x})$, without loss of generality, we assume $\|\mathbf{x}_i\| = 1$ for $1 \leq i \leq s$ and let $\mathbf{h} \in \mathbb{C}^{Ks}$ be a unit vector, i.e., $\|\mathbf{h}\|^2 = \sum_{i=1}^s \|\mathbf{h}_i\|^2 = 1$. The bound

$$\frac{1}{L} \leq \sigma_{\max}^2(\mathbf{h}, \mathbf{x}) \leq \frac{K}{L} \quad (6.10)$$

follows from $L\sigma_{\max}^2(\mathbf{h}, \mathbf{x}) \geq \sum_{l=1}^L \sum_{i=1}^s |\mathbf{b}_l^* \mathbf{h}_i|^2 = \|\mathbf{h}\|^2 = 1$.

Moreover, $\sigma_{\max}^2(\mathbf{h}, \mathbf{x})$ and $\sigma_{\max}(\mathbf{h}, \mathbf{x})$ are both Lipschitz functions w.r.t. \mathbf{h} . Now we want to determine their Lipschitz constants. First note that for $\|\mathbf{x}_i\| = 1$, $\sigma_{\max}(\mathbf{h}, \mathbf{x})$ equals

$$\sigma_{\max}(\mathbf{h}, \mathbf{x}) = \max_{1 \leq l \leq L} \|(\mathbf{I}_s \otimes \mathbf{b}_l^*) \mathbf{h}\|$$

where “ \otimes ” denotes Kronecker product. Let $\mathbf{u} \in \mathbb{C}^{Ks}$ be another unit vector and we have

$$\begin{aligned} |\sigma_{\max}(\mathbf{h}, \mathbf{x}) - \sigma_{\max}(\mathbf{u}, \mathbf{x})| &= \left| \max_{1 \leq l \leq L} \|(\mathbf{I}_s \otimes \mathbf{b}_l^*) \mathbf{h}\| - \max_{1 \leq l \leq L} \|(\mathbf{I}_s \otimes \mathbf{b}_l^*) \mathbf{u}\| \right| \\ &= \max_{1 \leq l \leq L} \left| \|(\mathbf{I}_s \otimes \mathbf{b}_l^*) \mathbf{h}\| - \|(\mathbf{I}_s \otimes \mathbf{b}_l^*) \mathbf{u}\| \right| \\ &\leq \max_{1 \leq l \leq L} \|(\mathbf{I}_s \otimes \mathbf{b}_l^*) (\mathbf{h} - \mathbf{u})\| \leq \|\mathbf{h} - \mathbf{u}\| \end{aligned} \quad (6.11)$$

where $\|\mathbf{I}_s \otimes \mathbf{b}_l^*\| = \|\mathbf{b}_l\| \sqrt{\frac{K}{L}} < 1$. For $\sigma_{\max}^2(\mathbf{h}, \mathbf{x})$,

$$\begin{aligned} |\sigma_{\max}^2(\mathbf{h}, \mathbf{x}) - \sigma_{\max}^2(\mathbf{u}, \mathbf{x})| &\leq (\sigma_{\max}(\mathbf{h}, \mathbf{x}) + \sigma_{\max}(\mathbf{u}, \mathbf{x})) \cdot |\sigma_{\max}(\mathbf{h}, \mathbf{x}) - \sigma_{\max}(\mathbf{u}, \mathbf{x})| \\ &\leq \frac{2K}{L} \|\mathbf{h} - \mathbf{u}\| \leq 2\|\mathbf{h} - \mathbf{u}\|. \end{aligned} \quad (6.12)$$

Proof: Without loss of generality, let $\|\mathbf{x}_i\| = 1$ and $\sum_{i=1}^s \|\mathbf{h}_i\|^2 = 1$. It suffices to prove $f(\mathbf{h}, \mathbf{x}) \leq \frac{4}{3}$ for all $(\mathbf{h}, \mathbf{x}) \in \mathbb{C}^{Ks} \oplus \mathbb{C}^{Ns}$ in (2.5) where $f(\mathbf{h}, \mathbf{x})$ is defined as

$$f(\mathbf{h}, \mathbf{x}) := \|\mathcal{A}(\mathbf{X})\|^2 - 2\sqrt{2s\sigma_{\max}^2(\mathbf{h}, \mathbf{x})(K+N) \log L} - 8s\sigma_{\max}^2(\mathbf{h}, \mathbf{x})(K+N) \log L.$$

Part I: Bounds of $\|\mathcal{A}(\mathbf{X})\|^2$ for any fixed (\mathbf{h}, \mathbf{x}) . From (6.9), we already know that $Y = \|\mathcal{A}(\mathbf{X})\|_F^2 = \sum_{i=1}^{2L} c_i \xi_i^2$ where $\{\xi_i\}$ are i.i.d. χ_1^2 random variables and $\mathbf{c} = (c_1, \dots, c_{2L})^T \in \mathbb{R}^{2L}$. More precisely, we can determine $\{c_i\}_{i=1}^{2L}$ as

$$\left| \sum_{i=1}^s \mathbf{b}_l^* \mathbf{h}_i \mathbf{x}^* \mathbf{a}_{il} \right|^2 = c_{2l-1} \xi_{2l-1}^2 + c_{2l} \xi_{2l}^2, \quad c_{2l-1} = c_{2l} = \frac{1}{2} \sum_{i=1}^s |\mathbf{b}_l^* \mathbf{h}_i|^2$$

because $\sum_{i=1}^s \mathbf{b}_l^* \mathbf{h}_i \mathbf{x}^* \mathbf{a}_{il} \sim \mathcal{CN}(0, \sum_{i=1}^s |\mathbf{b}_l^* \mathbf{h}_i|^2)$.

By the Bernstein inequality, there holds

$$\mathbb{P}(Y - \mathbb{E}(Y) \geq t) \leq \exp\left(-\frac{t^2}{8\|\mathbf{c}\|^2}\right) \vee \exp\left(-\frac{t}{8\|\mathbf{c}\|_\infty}\right) \quad (6.13)$$

where $\mathbb{E}(Y) = \|\mathbf{X}\|_F^2 = 1$. In order to apply the Bernstein inequality, we need to estimate $\|\mathbf{c}\|^2$ and $\|\mathbf{c}\|_\infty$ as follows,

$$\begin{aligned} \|\mathbf{c}\|_\infty &= \frac{1}{2} \max_{1 \leq l \leq L} \sum_{i=1}^s |\mathbf{b}_l^* \mathbf{h}_i|^2 = \frac{1}{2} \sigma_{\max}^2(\mathbf{h}, \mathbf{x}), \\ \|\mathbf{c}\|_2^2 &= \frac{1}{2} \sum_{l=1}^L \left| \sum_{i=1}^s |\mathbf{b}_l^* \mathbf{h}_i|^2 \right|^2 \leq \frac{1}{2} \max_{1 \leq l \leq L} \sum_{i=1}^s |\mathbf{b}_l^* \mathbf{h}_i|^2 \leq \frac{1}{2} \sigma_{\max}^2(\mathbf{h}, \mathbf{x}). \end{aligned}$$

Applying (6.13) gives

$$\mathbb{P}(\|\mathcal{A}(\mathbf{X})\|^2 \geq 1 + t) \leq \exp\left(-\frac{t^2}{4\sigma_{\max}^2(\mathbf{h}, \mathbf{x})}\right) \vee \exp\left(-\frac{t}{4\sigma_{\max}^2(\mathbf{h}, \mathbf{x})}\right).$$

In particular, by setting

$$t = g(\mathbf{h}, \mathbf{x}) := 2\sqrt{2s\sigma_{\max}^2(\mathbf{h}, \mathbf{x})(K+N)\log L} + 8s\sigma_{\max}^2(\mathbf{h}, \mathbf{x})(K+N)\log L,$$

we have

$$\mathbb{P}(\|\mathcal{A}(\mathbf{X})\|^2 \geq 1 + g(\mathbf{h}, \mathbf{x})) \leq e^{-2s(K+N)(\log L)}.$$

So far, we have shown that $f(\mathbf{h}, \mathbf{x}) \leq 1$ with probability at least $1 - e^{-2s(K+N)(\log L)}$ for a fixed pair of (\mathbf{h}, \mathbf{x}) .

Part II: Covering argument. Now we will use a covering argument to extend this result for all (\mathbf{h}, \mathbf{x}) and thus prove that $f(\mathbf{h}, \mathbf{x}) \leq \frac{4}{3}$ uniformly for all (\mathbf{h}, \mathbf{x}) .

We start with defining \mathcal{K} and \mathcal{N}_i as ϵ_0 -nets of \mathcal{S}^{Ks-1} and \mathcal{S}^{N-1} for \mathbf{h} and $\mathbf{x}_i, 1 \leq i \leq s$, respectively. The bounds $|\mathcal{K}| \leq (1 + \frac{2}{\epsilon_0})^{2sK}$ and $|\mathcal{N}_i| \leq (1 + \frac{2}{\epsilon_0})^{2N}$ follow from the covering numbers of the sphere (Lemma 5.2 in [31]). Here we let $\mathcal{N} := \mathcal{N}_1 \times \dots \times \mathcal{N}_s$. By taking the union bound over $\mathcal{K} \times \mathcal{N}$, we have that $f(\mathbf{h}, \mathbf{x}) \leq 1$ holds uniformly for all $(\mathbf{h}, \mathbf{x}) \in \mathcal{K} \times \mathcal{N}$ with probability at least

$$1 - (1 + 2/\epsilon_0)^{2s(K+N)} e^{-2s(K+N)\log L} = 1 - e^{-2s(K+N)(\log L - \log(1+2/\epsilon_0))}.$$

For any $(\mathbf{h}, \mathbf{x}) \in \mathcal{S}^{Ks-1} \times \underbrace{\mathcal{S}^{N-1} \times \dots \times \mathcal{S}^{N-1}}_{s \text{ times}}$, we can find a point $(\mathbf{u}, \mathbf{v}) \in \mathcal{K} \times \mathcal{N}$ satisfying

$\|\mathbf{h} - \mathbf{u}\| \leq \epsilon_0$ and $\|\mathbf{x}_i - \mathbf{v}_i\| \leq \epsilon_0$ for all $1 \leq i \leq s$. Conditioned on (6.3), we know that

$$\|\mathcal{A}\|^2 \leq s(N \log(NL/2) + (\gamma + \log s) \log L) \leq s(N + \gamma + \log s) \log L.$$

Now we aim to evaluate $|f(\mathbf{h}, \mathbf{x}) - f(\mathbf{u}, \mathbf{v})|$. First we consider $|f(\mathbf{u}, \mathbf{x}) - f(\mathbf{u}, \mathbf{v})|$. Since $\sigma_{\max}^2(\mathbf{u}, \mathbf{x}) = \sigma_{\max}^2(\mathbf{u}, \mathbf{v})$ if $\|\mathbf{x}_i\| = \|\mathbf{v}_i\| = \|\mathbf{u}\| = 1$ for $1 \leq i \leq s$, we have

$$\begin{aligned} |f(\mathbf{u}, \mathbf{x}) - f(\mathbf{u}, \mathbf{v})| &= \left| \|\mathcal{A}(\mathcal{H}(\mathbf{u}, \mathbf{x}))\|_F^2 - \|\mathcal{A}(\mathcal{H}(\mathbf{u}, \mathbf{v}))\|_F^2 \right| \\ &\leq \|\mathcal{A}(\mathcal{H}(\mathbf{u}, \mathbf{x} - \mathbf{v}))\| \cdot \|\mathcal{A}(\mathcal{H}(\mathbf{u}, \mathbf{x} + \mathbf{v}))\| \\ &\leq \|\mathcal{A}\|^2 \sqrt{\sum_{i=1}^s \|\mathbf{u}_i\|^2 \|\mathbf{x}_i - \mathbf{v}_i\|^2} \sqrt{\sum_{i=1}^s \|\mathbf{u}_i\|^2 \|\mathbf{x}_i + \mathbf{v}_i\|^2} \\ &\leq 2\|\mathcal{A}\|^2 \epsilon_0 \leq 2s(N + \gamma + \log s)(\log L) \epsilon_0 \end{aligned}$$

where the first inequality is due to $\|z_1\|^2 - \|z_2\|^2 \leq \|z_1 - z_2\| \|z_1 + z_2\|$ for any $z_1, z_2 \in \mathbb{C}$.

We proceed to estimate $|f(\mathbf{h}, \mathbf{x}) - f(\mathbf{u}, \mathbf{x})|$ by using (6.12) and (6.11),

$$\begin{aligned} |f(\mathbf{h}, \mathbf{x}) - f(\mathbf{u}, \mathbf{x})| &\leq J_1 + J_2 + J_3 \\ &\leq (2\|\mathcal{A}\|^2 + 2\sqrt{2s(K+N)\log L} + 16s(K+N)\log L)\varepsilon_0 \\ &\leq 25s(K+N+\gamma+\log s)(\log L)\varepsilon_0 \end{aligned}$$

where (6.12) and (6.11) give

$$\begin{aligned} J_1 &= \|\|\mathcal{A}(\mathcal{H}(\mathbf{h}, \mathbf{x}))\|_F^2 - \|\mathcal{A}(\mathcal{H}(\mathbf{u}, \mathbf{x}))\|_F^2\| \leq \|\mathcal{A}(\mathcal{H}(\mathbf{h} - \mathbf{u}, \mathbf{x}))\| \|\mathcal{A}(\mathcal{H}(\mathbf{h} + \mathbf{u}, \mathbf{x}))\| \leq 2\|\mathcal{A}\|^2\varepsilon_0, \\ J_2 &= 2\sqrt{2s(K+N)\log L} \cdot |\sigma_{\max}(\mathbf{h}, \mathbf{x}) - \sigma_{\max}(\mathbf{u}, \mathbf{x})| \leq 2\sqrt{2s(K+N)\log L}\varepsilon_0, \\ J_3 &= 8s(K+N)(\log L) \cdot |\sigma_{\max}^2(\mathbf{h}, \mathbf{x}) - \sigma_{\max}^2(\mathbf{u}, \mathbf{x})| \leq 16s(K+N)(\log L)\varepsilon_0. \end{aligned}$$

Therefore, if $\varepsilon_0 = \frac{1}{81s(N+K+\gamma+\log s)\log L}$, there holds

$$f(\mathbf{h}, \mathbf{x}) \leq f(\mathbf{u}, \mathbf{v}) + \underbrace{|f(\mathbf{u}, \mathbf{x}) - f(\mathbf{u}, \mathbf{v})| + |f(\mathbf{h}, \mathbf{x}) - f(\mathbf{u}, \mathbf{x})|}_{\leq 27s(K+N+\gamma+\log s)(\log L)\varepsilon_0 \leq \frac{1}{3}} \leq \frac{4}{3}$$

for all (\mathbf{h}, \mathbf{x}) uniformly with probability at least $1 - e^{-2s(K+N)(\log L - \log(1+2/\varepsilon_0))}$. By letting $L \geq C_\gamma s(K+N)\log L$ with C_γ reasonably large and $\gamma \geq 1$, we have $\log L - \log(1+2/\varepsilon_0) \geq \frac{1}{2}(1 + \log(\gamma))$ and with probability at least $1 - \frac{1}{\gamma} \exp(-s(K+N))$. \square

6.2 Proof of the local restricted isometry property

Lemma 6.7. *Conditioned on (6.4) and (6.8), the following RIP type of property holds:*

$$\frac{2}{3}\|\mathbf{X} - \mathbf{X}_0\|_F^2 \leq \|\mathcal{A}(\mathbf{X} - \mathbf{X}_0)\|^2 \leq \frac{3}{2}\|\mathbf{X} - \mathbf{X}_0\|_F^2$$

uniformly for all $(\mathbf{h}, \mathbf{x}) \in \mathcal{N}_d \cap \mathcal{N}_\mu \cap \mathcal{N}_\varepsilon$ with $\mu \geq \mu_h$ and $\varepsilon \leq \frac{1}{15}$ if $L \geq C_\gamma \mu^2 s(K+N)\log^2 L$ for some numerical constant C_γ .

Proof: The main idea of the proof follows two steps: decompose $\mathbf{X} - \mathbf{X}_0$ onto T and T^\perp , then apply (6.4) and (6.8) to $\mathcal{P}_T(\mathbf{X} - \mathbf{X}_0)$ and $\mathcal{P}_{T^\perp}(\mathbf{X} - \mathbf{X}_0)$ respectively.

For any $\mathbf{X} = \mathcal{H}(\mathbf{h}, \mathbf{x}) \in \mathcal{N}_\varepsilon$ with $\delta_i \leq \varepsilon \leq \frac{1}{15}$, we can decompose $\mathbf{X} - \mathbf{X}_0$ as the sum of two block diagonal matrices $\mathbf{U} = \text{blkdiag}(\mathbf{U}_i, 1 \leq i \leq s)$ and $\mathbf{V} = \text{blkdiag}(\mathbf{V}_i, 1 \leq i \leq s)$ where each pair of $(\mathbf{U}_i, \mathbf{V}_i)$ corresponds to the orthogonal decomposition of $\mathbf{h}_i \mathbf{x}_i^* - \mathbf{h}_{i0} \mathbf{x}_{i0}^*$,

$$\mathbf{h}_i \mathbf{x}_i^* - \mathbf{h}_{i0} \mathbf{x}_{i0}^* := \underbrace{(\alpha_{i1} \bar{\alpha}_{i2} - 1) \mathbf{h}_{i0} \mathbf{x}_{i0}^* + \bar{\alpha}_{i2} \tilde{\mathbf{h}}_i \mathbf{x}_{i0}^* + \alpha_{i1} \mathbf{h}_{i0} \tilde{\mathbf{x}}_i^*}_{\mathbf{U}_i \in T_i} + \underbrace{\tilde{\mathbf{h}}_i \tilde{\mathbf{x}}_i^*}_{\mathbf{V}_i \in T_i^\perp} \quad (6.14)$$

which has been briefly discussed in (6.1) and (6.2). Note that $\mathcal{A}(\mathbf{X} - \mathbf{X}_0) = \mathcal{A}(\mathbf{U} + \mathbf{V})$ and

$$\|\mathcal{A}(\mathbf{U})\| - \|\mathcal{A}(\mathbf{V})\| \leq \|\mathcal{A}(\mathbf{U} + \mathbf{V})\| \leq \|\mathcal{A}(\mathbf{U})\| + \|\mathcal{A}(\mathbf{V})\|.$$

Therefore, it suffices to have a two-side bound for $\|\mathcal{A}(\mathbf{U})\|$ and an upper bound for $\|\mathcal{A}(\mathbf{V})\|$ where $\mathbf{U} \in T$ and $\mathbf{V} \in T^\perp$ in order to establish the local isometry property.

Estimation of $\|\mathcal{A}(\mathbf{U})\|$: For $\|\mathcal{A}(\mathbf{U})\|$, we know from Lemma 6.3 that

$$\sqrt{\frac{9}{10}}\|\mathbf{U}\|_F \leq \|\mathcal{A}(\mathbf{U})\| \leq \sqrt{\frac{11}{10}}\|\mathbf{U}\|_F \quad (6.15)$$

and hence we only need to compute $\|\mathbf{U}\|_F$. By Lemma 6.1, there also hold $\|\mathbf{V}_i\|_F \leq \frac{\delta_i^2}{2(1-\delta_i)} d_{i0}$ and $\delta_i - \|\mathbf{V}_i\|_F \leq \|\mathbf{U}_i\|_F \leq \delta_i + \|\mathbf{V}_i\|_F$, i.e.,

$$\left(\delta_i - \frac{\delta_i^2}{2(1-\delta_i)} \right) d_{i0} \leq \|\mathbf{U}_i\|_F \leq \left(\delta_i + \frac{\delta_i^2}{2(1-\delta_i)} \right) d_{i0}, \quad 1 \leq i \leq s.$$

With $\|\mathbf{U}\|_F^2 = \sum_{i=1}^s \|\mathbf{U}_i\|_F^2$, it is easy to get $\delta d_0 \left(1 - \frac{\varepsilon}{2(1-\varepsilon)}\right) \leq \|\mathbf{U}\|_F \leq \delta d_0 \left(1 + \frac{\varepsilon}{2(1-\varepsilon)}\right)$. Combined with (6.15), we get

$$\sqrt{\frac{9}{10}} \left(1 - \frac{\varepsilon}{2(1-\varepsilon)}\right) \delta d_0 \leq \|\mathcal{A}(\mathbf{U})\| \leq \sqrt{\frac{11}{10}} \left(1 + \frac{\varepsilon}{2(1-\varepsilon)}\right) \delta d_0. \quad (6.16)$$

Estimation of $\|\mathcal{A}(\mathbf{V})\|$: Note that \mathbf{V} is a block-diagonal matrix with rank-1 block. So applying Lemma 6.5 gives us

$$\|\mathcal{A}(\mathbf{V})\|^2 \leq \frac{4}{3} \|\mathbf{V}\|_F^2 + 2\sqrt{2s\|\mathbf{V}\|_F^2 \sigma_{\max}^2(\tilde{\mathbf{h}}, \tilde{\mathbf{x}})(K+N) \log L} + 8s\sigma_{\max}^2(\tilde{\mathbf{h}}, \tilde{\mathbf{x}})(K+N) \log L \quad (6.17)$$

where $\mathbf{V} = \mathcal{H}(\tilde{\mathbf{h}}, \tilde{\mathbf{x}})$ and $\tilde{\mathbf{h}} = \begin{bmatrix} \tilde{\mathbf{h}}_1 \\ \vdots \\ \tilde{\mathbf{h}}_s \end{bmatrix}$. It suffices to get an estimation of $\|\mathbf{V}\|_F$ and $\sigma_{\max}^2(\tilde{\mathbf{h}}, \tilde{\mathbf{x}})$ to bound $\|\mathcal{A}(\mathbf{V})\|$ in (6.17).

Lemma 6.1 says that $\|\tilde{\mathbf{h}}_i\| \|\tilde{\mathbf{x}}_i\| \leq \frac{\delta_i^2}{2(1-\delta_i)} d_{i0} \leq \frac{\varepsilon}{2(1-\varepsilon)} \delta_i d_{i0}$ if $\varepsilon < 1$. Moreover,

$$\|\tilde{\mathbf{x}}_i\| \leq \frac{\delta_i}{1-\delta_i} \|\mathbf{x}_i\| \leq \frac{2\delta_i}{1-\delta_i} \sqrt{d_{i0}}, \quad \sqrt{L} \|\mathbf{B}\tilde{\mathbf{h}}_i\|_{\infty} \leq 6\mu\sqrt{d_{i0}}, \quad 1 \leq i \leq s \quad (6.18)$$

if (\mathbf{h}, \mathbf{x}) belongs to $\mathcal{N}_d \cap \mathcal{N}_{\mu} \cap \mathcal{N}_{\varepsilon}$. For $\|\mathbf{V}\|_F$,

$$\|\mathbf{V}\|_F = \sqrt{\sum_{i=1}^s \|\mathbf{V}_i\|_F^2} = \sqrt{\sum_{i=1}^s \|\tilde{\mathbf{h}}_i\|^2 \|\tilde{\mathbf{x}}_i\|^2} \leq \frac{\varepsilon \delta d_0}{2(1-\varepsilon)}.$$

Now we aim to get an upper bound for $\sigma_{\max}^2(\tilde{\mathbf{h}}, \tilde{\mathbf{x}})$ by using (6.18),

$$\sigma_{\max}^2(\tilde{\mathbf{h}}, \tilde{\mathbf{x}}) = \max_{1 \leq l \leq L} \sum_{i=1}^s |\mathbf{b}_l^* \tilde{\mathbf{h}}_i|^2 \|\tilde{\mathbf{x}}_i\|^2 \leq C_0 \frac{\mu^2 \sum_{i=1}^s \delta_i^2 d_{i0}^2}{L} = C_0 \frac{\mu^2 \delta^2 d_0^2}{L}.$$

By substituting the estimations of $\|\mathbf{V}\|_F$ and $\sigma_{\max}^2(\tilde{\mathbf{h}}, \tilde{\mathbf{x}})$ into (6.17)

$$\|\mathcal{A}(\mathbf{V})\|^2 \leq \frac{\varepsilon^2 \delta^2 d_0^2}{3(1-\varepsilon)^2} + \frac{\sqrt{2\varepsilon} \delta^2 d_0^2}{1-\varepsilon} \sqrt{\frac{C_0 \mu^2 s (K+N) \log L}{L}} + \frac{8C_0 \mu^2 \delta^2 d_0^2 s (K+N) \log L}{L}. \quad (6.19)$$

By letting $L \geq C_{\gamma} \mu^2 s (K+N) \log^2 L$ with C_{γ} sufficiently large and combining (6.19) and (6.16), we have

$$\sqrt{\frac{2}{3}} \delta d_0 \leq \|\mathcal{A}(\mathbf{U})\| - \|\mathcal{A}(\mathbf{V})\| \leq \|\mathcal{A}(\mathbf{U} + \mathbf{V})\| \leq \|\mathcal{A}(\mathbf{U})\| + \|\mathcal{A}(\mathbf{V})\| \leq \sqrt{\frac{3}{2}} \delta d_0,$$

which gives $\frac{2}{3} \|\mathbf{X} - \mathbf{X}_0\|_F^2 \leq \|\mathcal{A}(\mathbf{X} - \mathbf{X}_0)\|^2 \leq \frac{3}{2} \|\mathbf{X} - \mathbf{X}_0\|_F^2$. \square

6.3 Proof of the local regularity condition

We first introduce a few notations: for all $(\mathbf{h}, \mathbf{x}) \in \mathcal{N}_d \cap \mathcal{N}_{\varepsilon}$, consider $\alpha_{i1}, \alpha_{i2}, \tilde{\mathbf{h}}_i$ and $\tilde{\mathbf{x}}_i$ defined in (6.1) and define

$$\Delta \mathbf{h}_i = \mathbf{h}_i - \alpha_i \mathbf{h}_{i0}, \quad \Delta \mathbf{x}_i = \mathbf{x}_i - \bar{\alpha}_i^{-1} \mathbf{x}_{i0}$$

where

$$\alpha_i(\mathbf{h}_i, \mathbf{x}_i) = \begin{cases} (1 - \delta_0) \alpha_{i1}, & \text{if } \|\mathbf{h}_i\|_2 \geq \|\mathbf{x}_i\|_2 \\ \frac{1}{(1-\delta_0)\bar{\alpha}_{i2}}, & \text{if } \|\mathbf{h}_i\|_2 < \|\mathbf{x}_i\|_2 \end{cases}$$

with

$$\delta_0 := \frac{\delta}{10}. \quad (6.20)$$

The function $\alpha_i(\mathbf{h}_i, \mathbf{x}_i)$ is defined for each block of $\mathbf{X} = \mathcal{H}(\mathbf{h}, \mathbf{x})$. The particular form of $\alpha_i(\mathbf{h}, \mathbf{x})$ serves primarily for proving the Lemma 6.10, i.e., local regularity condition of $G(\mathbf{h}, \mathbf{x})$. We also define

$$\Delta \mathbf{h} := \begin{bmatrix} \mathbf{h}_1 - \alpha_1 \mathbf{h}_{1,0} \\ \vdots \\ \mathbf{h}_s - \alpha_s \mathbf{h}_{s,0} \end{bmatrix} \in \mathbb{C}^{Ks}, \quad \Delta \mathbf{x} := \begin{bmatrix} \mathbf{x}_1 - \alpha_1 \mathbf{x}_{1,0} \\ \vdots \\ \mathbf{x}_s - \alpha_s \mathbf{x}_{s,0} \end{bmatrix} \in \mathbb{C}^{Ns}.$$

The following lemma gives bounds of $\Delta \mathbf{x}_i$ and $\Delta \mathbf{h}_i$.

Lemma 6.8. *For all $(\mathbf{h}, \mathbf{x}) \in \mathcal{N}_d \cap \mathcal{N}_\epsilon$ with $\epsilon \leq \frac{1}{15}$, there hold*

$$\begin{aligned} \max\{\|\Delta \mathbf{h}_i\|_2^2, \|\Delta \mathbf{x}_i\|_2^2\} &\leq (7.5\delta_i^2 + 2.88\delta_0^2)d_{i0}, \\ \|\Delta \mathbf{h}_i\|_2^2 \|\Delta \mathbf{x}_i\|_2^2 &\leq \frac{1}{26}(\delta_i^2 + \delta_0^2)d_{i0}^2. \end{aligned}$$

Moreover, if we assume $(\mathbf{h}_i, \mathbf{x}_i) \in \mathcal{N}_\mu$ additionally, we have $\sqrt{L}\|\mathbf{B}(\Delta \mathbf{h}_i)\|_\infty \leq 6\mu\sqrt{d_{i0}}$.

Proof: We only consider $\|\mathbf{h}_i\|_2 \geq \|\mathbf{x}_i\|_2$ and $\alpha_i = (1 - \delta_0)\alpha_{i1}$, and the other case is exactly the same due to the symmetry. For both $\Delta \mathbf{h}_i$ and $\Delta \mathbf{x}_i$, by definition,

$$\Delta \mathbf{h}_i = \mathbf{h}_i - \alpha_i \mathbf{h}_{i0} = \delta_0 \alpha_{i1} \mathbf{h}_{i0} + \tilde{\mathbf{h}}_i, \quad (6.21)$$

$$\Delta \mathbf{x}_i = \mathbf{x}_i - \frac{1}{(1 - \delta_0)\bar{\alpha}_{i1}} \mathbf{x}_{i0} = \left(\alpha_{i2} - \frac{1}{(1 - \delta_0)\bar{\alpha}_{i1}} \right) \mathbf{x}_{i0} + \tilde{\mathbf{x}}_i, \quad (6.22)$$

where $\mathbf{h}_i = \alpha_{i1} \mathbf{h}_{i0} + \tilde{\mathbf{h}}_i$ and $\mathbf{x}_i = \alpha_{i2} \mathbf{x}_{i0} + \tilde{\mathbf{x}}_i$ come from the orthogonal decomposition in (6.1).

We start with estimating $\|\Delta \mathbf{h}_i\|_2^2$. Note that $\|\mathbf{h}_i\|_2^2 \leq 4d_{i0}$ and $\|\alpha_{i1} \mathbf{h}_{i0}\|_2^2 \leq \|\mathbf{h}_i\|_2^2$ since $(\mathbf{h}, \mathbf{x}) \in \mathcal{N}_d \cap \mathcal{N}_\mu$. By Lemma 6.1, we have

$$\|\Delta \mathbf{h}_i\|_2^2 = \|\tilde{\mathbf{h}}_i\|_2^2 + \delta_0^2 \|\alpha_{i1} \mathbf{h}_{i0}\|_2^2 \leq \left(\left(\frac{\delta_i}{1 - \delta_i} \right)^2 + \delta_0^2 \right) \|\mathbf{h}_i\|_2^2 \leq (4.6\delta_i^2 + 4\delta_0^2)d_{i0}. \quad (6.23)$$

Then we calculate $\|\Delta \mathbf{x}_i\|_2$: from (6.22), we have

$$\|\Delta \mathbf{x}_i\|_2^2 = \left| \alpha_{i2} - \frac{1}{(1 - \delta_0)\bar{\alpha}_{i1}} \right|^2 d_{i0} + \|\tilde{\mathbf{x}}_i\|_2^2 \leq \left| \alpha_{i2} - \frac{1}{(1 - \delta_0)\bar{\alpha}_{i1}} \right|^2 d_{i0} + \frac{4\delta_i^2 d_{i0}}{(1 - \delta_i)^2},$$

where Lemma 6.1 gives $\|\tilde{\mathbf{x}}_i\|_2 \leq \frac{\delta_i}{1 - \delta_i} \|\mathbf{x}_i\|_2 \leq \frac{2\delta_i}{1 - \delta_i} \sqrt{d_{i0}}$ for $(\mathbf{h}, \mathbf{x}) \in \mathcal{N}_d \cap \mathcal{N}_\epsilon$.

So it suffices to estimate $\left| \alpha_{i2} - \frac{1}{(1 - \delta_0)\bar{\alpha}_{i1}} \right|$, which satisfies

$$\left| \alpha_{i2} - \frac{1}{(1 - \delta_0)\bar{\alpha}_{i1}} \right| = \frac{1}{|\alpha_{i1}|} \left| \bar{\alpha}_{i1} \alpha_{i2} - 1 - \frac{\delta_0}{1 - \delta_0} \right| \leq \frac{1}{|\alpha_{i1}|} \left(|(\bar{\alpha}_{i1} \alpha_{i2} - 1)| + \frac{\delta_0}{1 - \delta_0} \right). \quad (6.24)$$

Lemma 6.1 implies that $|\bar{\alpha}_{i1} \alpha_{i2} - 1| \leq \delta_i$, and (6.1) gives

$$|\alpha_{i1}|^2 = \frac{1}{d_{i0}} (\|\mathbf{h}_i\|_2^2 - \|\tilde{\mathbf{h}}_i\|_2^2) \geq \frac{1}{d_{i0}} \left(1 - \frac{\delta_i^2}{(1 - \delta_i)^2} \right) \|\mathbf{h}_i\|_2^2 \geq \left(1 - \frac{\delta_i^2}{(1 - \delta_i)^2} \right) (1 - \epsilon) \quad (6.25)$$

where $\|\tilde{\mathbf{h}}_i\|_2 \leq \frac{\delta_i}{1 - \delta_i} \|\mathbf{h}_i\|_2$ and $\|\mathbf{h}_i\|_2^2 \geq \|\mathbf{h}_i\| \|\mathbf{x}_i\| \geq (1 - \epsilon)d_{i0}$ if $\|\mathbf{h}_i\| \geq \|\mathbf{x}_i\|$. Substituting (6.25) into (6.24) gives

$$\left| \alpha_{i2} - \frac{1}{(1 - \delta_0)\bar{\alpha}_{i1}} \right| \leq \frac{1}{\sqrt{1 - \epsilon}} \left(1 - \frac{\delta_i^2}{(1 - \delta_i)^2} \right)^{-1/2} \left(\delta_i + \frac{\delta_0}{1 - \delta_0} \right) \leq 1.2(\delta_i + \delta_0).$$

Then we have

$$\|\Delta \mathbf{x}_i\|_2^2 \leq \left(1.44(\delta_i + \delta_0)^2 + \frac{4\delta_i^2}{(1 - \delta_i)^2} \right) d_{i0} \leq (7.5\delta_i^2 + 2.88\delta_0^2)d_{i0}. \quad (6.26)$$

Finally, we try to bound $\|\Delta \mathbf{h}_i\|^2 \|\Delta \mathbf{x}_i\|^2$. Lemma 6.1 gives $\|\tilde{\mathbf{h}}_i\|_2 \|\tilde{\mathbf{x}}_i\|_2 \leq \frac{\delta_i^2 d_{i0}}{2(1 - \delta_i)}$ and $|\alpha_{i1}| \leq 2$. Combining them along with (6.21), (6.22), (6.23) and (6.26), we have

$$\begin{aligned} \|\Delta \mathbf{h}_i\|_2^2 \|\Delta \mathbf{x}_i\|_2^2 &\leq \|\tilde{\mathbf{h}}_i\|_2^2 \|\tilde{\mathbf{x}}_i\|_2^2 + \delta_0^2 |\alpha_{i1}|^2 \|\mathbf{h}_{i0}\|_2^2 \|\Delta \mathbf{x}_i\|_2^2 + \left| \alpha_{i2} - \frac{1}{(1 - \delta_0)\bar{\alpha}_{i1}} \right|^2 \|\mathbf{x}_{i0}\|_2^2 \|\Delta \mathbf{h}_i\|_2^2 \\ &\leq \left(\frac{\delta_i^4}{4(1 - \delta_i)^2} + 4\delta_0^2(7.5\delta_i^2 + 2.88\delta_0^2) + 1.44(\delta_i + \delta_0)^2(4.6\delta_i^2 + 4\delta_0^2) \right) d_{i0}^2 \\ &\leq \frac{(\delta_i^2 + \delta_0^2)d_{i0}^2}{26}. \end{aligned}$$

By symmetry, similar results hold for the case $\|\mathbf{h}_i\|_2 < \|\mathbf{x}_i\|_2$ and $\max\{\|\Delta \mathbf{h}_i\|, \|\Delta \mathbf{x}_i\|\} \leq (7.5\delta_i^2 + 2.88\delta_0^2)d_{i0}$.

Next, under the additional assumption $(\mathbf{h}, \mathbf{x}) \in \mathcal{N}_\mu$, we now prove $\sqrt{L}\|\mathbf{B}(\Delta \mathbf{h}_i)\|_\infty \leq 6\mu\sqrt{d_{i0}}$:

Case 1: $\|\mathbf{h}_i\|_2 \geq \|\mathbf{x}_i\|_2$ and $\alpha_i = (1 - \delta_0)\alpha_{i1}$. By Lemma 6.1 gives $|\alpha_{i1}| \leq 2$, which implies

$$\begin{aligned} \sqrt{L}\|\mathbf{B}(\Delta \mathbf{h}_i)\|_\infty &\leq \sqrt{L}\|\mathbf{B}\mathbf{h}_i\|_\infty + (1 - \delta_0)|\alpha_{i1}|\sqrt{L}\|\mathbf{B}\mathbf{h}_{i0}\|_\infty \\ &\leq 4\mu\sqrt{d_{i0}} + 2(1 - \delta_0)\mu_h\sqrt{d_{i0}} \leq 6\mu\sqrt{d_{i0}}. \end{aligned}$$

Case 2: $\|\mathbf{h}_i\|_2 < \|\mathbf{x}_i\|_2$ and $\alpha_i = \frac{1}{(1 - \delta_0)\bar{\alpha}_{i2}}$. Using the same argument as (6.25) gives

$$|\alpha_{i2}|^2 \geq \left(1 - \frac{\delta_i^2}{(1 - \delta_i)^2} \right) (1 - \varepsilon).$$

Therefore,

$$\begin{aligned} \sqrt{L}\|\mathbf{B}(\Delta \mathbf{h}_i)\|_\infty &\leq \sqrt{L}\|\mathbf{B}\mathbf{h}_i\|_\infty + \frac{1}{(1 - \delta_0)|\bar{\alpha}_{i2}|} \sqrt{L}\|\mathbf{B}\mathbf{h}_0\|_\infty \\ &\leq 4\mu\sqrt{d_0} + \left(1 - \frac{\delta_i^2}{(1 - \delta_i)^2} \right)^{-1/2} \frac{\mu_h\sqrt{d_0}}{(1 - \delta_0)\sqrt{1 - \varepsilon}} \leq 6\mu\sqrt{d_0}. \end{aligned}$$

□

Lemma 6.9. (Local Regularity for $F(\mathbf{h}, \mathbf{x})$) *Conditioned on (5.3) and (6.8), the following inequality holds*

$$\operatorname{Re}(\langle \nabla F_{\mathbf{h}}, \Delta \mathbf{h} \rangle + \langle \nabla F_{\mathbf{x}}, \Delta \mathbf{x} \rangle) \geq \frac{\delta^2 d_0^2}{8} - 2\sqrt{s}\delta d_0 \|\mathcal{A}^*(\mathbf{e})\|,$$

uniformly for any $(\mathbf{h}, \mathbf{x}) \in \mathcal{N}_d \cap \mathcal{N}_\mu \cap \mathcal{N}_\varepsilon$ with $\varepsilon \leq \frac{1}{15}$ if $L \geq C\mu^2 s(K + N) \log^2 L$ for some numerical constant C .

Proof: First note that for

$$I_0 = \langle \nabla F_{\mathbf{h}}, \Delta \mathbf{h} \rangle + \overline{\langle \nabla F_{\mathbf{x}}, \Delta \mathbf{x} \rangle} = \sum_{i=1}^s \langle \nabla F_{\mathbf{h}_i}, \Delta \mathbf{h}_i \rangle + \overline{\langle \nabla F_{\mathbf{x}_i}, \Delta \mathbf{x}_i \rangle}.$$

For each component, recall that (2.18) and (2.19), we have

$$\begin{aligned} \langle \nabla F_{\mathbf{h}_i}, \Delta \mathbf{h}_i \rangle + \overline{\langle \nabla F_{\mathbf{x}_i}, \Delta \mathbf{x}_i \rangle} &= \langle \mathcal{A}_i^*(\mathcal{A}(\mathbf{X} - \mathbf{X}_0) - \mathbf{e})\mathbf{x}_i, \Delta \mathbf{h}_i \rangle + \overline{\langle (\mathcal{A}_i^*(\mathcal{A}(\mathbf{X} - \mathbf{X}_0) - \mathbf{e}))^* \mathbf{h}_i, \Delta \mathbf{x}_i \rangle} \\ &= \langle \mathcal{A}(\mathbf{X} - \mathbf{X}_0) - \mathbf{e}, \mathcal{A}_i((\Delta \mathbf{h}_i)\mathbf{x}_i^* + \mathbf{h}_i(\Delta \mathbf{x}_i)^*) \rangle. \end{aligned}$$

Define \mathbf{U}_i and \mathbf{V}_i as

$$\mathbf{U}_i := \alpha_i \mathbf{h}_{i0} (\Delta \mathbf{x}_i)^* + \bar{\alpha}_i^{-1} (\Delta \mathbf{h}_i) \mathbf{x}_{i0}^* \in T_i, \quad \mathbf{V}_i := \Delta \mathbf{h}_i (\Delta \mathbf{x}_i)^*. \quad (6.27)$$

Here \mathbf{V}_i does not necessarily belong to T_i^\perp . From the way of how $\Delta \mathbf{h}_i$, $\Delta \mathbf{x}_i$, \mathbf{U}_i and \mathbf{V}_i are constructed, two simple relations hold:

$$\begin{aligned} \mathbf{h}_i \mathbf{x}_i^* - \mathbf{h}_{i0} \mathbf{x}_{i0}^* &= \mathbf{U}_i + \mathbf{V}_i, \\ (\Delta \mathbf{h}_i) \mathbf{x}_i^* + \mathbf{h}_i (\Delta \mathbf{x}_i)^* &= \mathbf{U}_i + 2\mathbf{V}_i. \end{aligned}$$

Define $\mathbf{U} := \text{blkdiag}(\mathbf{U}_1, \dots, \mathbf{U}_s)$ and $\mathbf{V} := \text{blkdiag}(\mathbf{V}_1, \dots, \mathbf{V}_s)$. I_0 can be simplified to

$$\begin{aligned} I_0 &= \sum_{i=1}^s \langle \nabla F \mathbf{h}_i, \Delta \mathbf{h}_i \rangle + \overline{\langle \nabla F \mathbf{x}_i, \Delta \mathbf{x}_i \rangle} = \sum_{i=1}^s \langle \mathcal{A}(\mathbf{U} + \mathbf{V}) - \mathbf{e}, \mathcal{A}_i(\mathbf{U}_i + 2\mathbf{V}_i) \rangle \\ &= \underbrace{\langle \mathcal{A}(\mathbf{U} + \mathbf{V}), \mathcal{A}(\mathbf{U} + 2\mathbf{V}) \rangle}_{I_{01}} - \underbrace{\langle \mathbf{e}, \mathcal{A}(\mathbf{U} + 2\mathbf{V}) \rangle}_{I_{02}}. \end{aligned}$$

Now we will give a lower bound for $\text{Re}(I_{01})$ and an upper bound for $\text{Re}(I_{02})$ so that the lower bound of $\text{Re}(I_0)$ is obtained. By the Cauchy-Schwarz inequality, $\text{Re}(I_{01})$ has the lower bound

$$\text{Re}(I_{01}) \geq (\|\mathcal{A}(\mathbf{U})\| - \|\mathcal{A}(\mathbf{V})\|)(\|\mathcal{A}(\mathbf{U})\| - 2\|\mathcal{A}(\mathbf{V})\|). \quad (6.28)$$

In the following, we will give an upper bound for $\|\mathcal{A}(\mathbf{V})\|$ and a lower bound for $\|\mathcal{A}(\mathbf{U})\|$.

Upper bound for $\|\mathcal{A}(\mathbf{V})\|$: Note that \mathbf{V} is a block-diagonal matrix with rank-1 blocks, and applying Lemma 6.5 results in

$$\|\mathcal{A}(\mathbf{V})\|^2 \leq \frac{4}{3} \sum_{i=1}^s \|\mathbf{V}\|_F^2 + 2\sigma_{\max}(\Delta \mathbf{h}, \Delta \mathbf{x}) \|\mathbf{V}\|_F \sqrt{2s(K+N) \log L} + 8s\sigma_{\max}^2(\Delta \mathbf{h}, \Delta \mathbf{x})(K+N) \log L.$$

By using Lemma 6.8, we have $\|\Delta \mathbf{h}_i\|^2 \leq (7.5\delta_i^2 + 2.88\delta_0^2)d_{i0}$ and $\sqrt{L}\|\mathbf{B}(\Delta \mathbf{h}_i)\|_\infty \leq 6\mu\sqrt{d_{i0}}$. Substituting them into $\sigma_{\max}^2(\Delta \mathbf{h}, \Delta \mathbf{x})$ gives

$$\sigma_{\max}^2(\Delta \mathbf{h}, \Delta \mathbf{x}) = \max_{1 \leq l \leq L} \left(\sum_{i=1}^s |\mathbf{b}_l^* \Delta \mathbf{h}_i|^2 \|\Delta \mathbf{x}_i\|^2 \right) \leq \frac{36\mu^2}{L} \sum_{i=1}^s (7.5\delta_i^2 + 2.88\delta_0^2) d_{i0}^2 \leq \frac{272\mu^2 \delta^2 d_0^2}{L}.$$

For $\|\mathbf{V}\|_F$, note that $\|\Delta \mathbf{h}_i\|^2 \|\Delta \mathbf{x}_i\|^2 \leq \frac{1}{26}(\delta_i^2 + \delta_0^2)d_{i0}^2$ and thus

$$\|\mathbf{V}\|_F^2 = \sum_{i=1}^s \|\Delta \mathbf{h}_i\|^2 \|\Delta \mathbf{x}_i\|^2 \leq \frac{1}{26} \sum_{i=1}^s (\delta_i^2 + \delta_0^2) d_{i0}^2 \leq \frac{1}{26} \cdot 1.01\delta^2 d_0^2 = \frac{\delta^2 d_0^2}{25}.$$

Then by $\delta \leq \epsilon \leq \frac{1}{15}$ and letting $L \geq C\mu^2 s(K+N) \log^2 L$ for a sufficiently large numerical constant C , there holds

$$\|\mathcal{A}(\mathbf{V})\|^2 \leq \frac{\delta^2 d_0^2}{16} \implies \|\mathcal{A}(\mathbf{V})\| \leq \frac{\delta d_0}{4}. \quad (6.29)$$

Lower bound for $\|\mathcal{A}(\mathbf{U})\|$: By the triangle inequality, $\|\mathbf{U}\|_F \geq \delta d_0 - \frac{1}{5}\delta d_0 \geq \frac{4}{5}\delta d_0$ if $\epsilon \leq \frac{1}{15}$ since $\|\mathbf{V}\|_F \leq 0.2\delta d_0$. Since $\mathbf{U} \in T$, by Lemma 6.3, there holds

$$\|\mathcal{A}(\mathbf{U})\| \geq \sqrt{\frac{9}{10}} \|\mathbf{U}\|_F \geq \frac{3}{4}\delta d_0. \quad (6.30)$$

With the upper bound of $\mathcal{A}(\mathbf{V})$ in (6.29), the lower bound of $\mathcal{A}(\mathbf{U})$ in (6.30), and (6.28), we get $\text{Re}(I_{01}) \geq \frac{\delta^2 d_0^2}{8}$.

Now let us give an upper bound for $\text{Re}(I_{02})$,

$$\begin{aligned} \|I_{02}\| &\leq \|\mathcal{A}^*(\mathbf{e})\| \|\mathbf{U} + 2\mathbf{V}\|_* = \|\mathcal{A}^*(\mathbf{e})\| \sum_{i=1}^s \underbrace{\|\mathbf{U}_i + 2\mathbf{V}_i\|_*}_{\text{rank-2}} \\ &\leq \sqrt{2} \|\mathcal{A}^*(\mathbf{e})\| \sum_{i=1}^s \|\mathbf{U}_i + 2\mathbf{V}_i\|_F \\ &\leq \sqrt{2s} \|\mathcal{A}^*(\mathbf{e})\| \|\mathbf{U} + 2\mathbf{V}\|_F \leq 2\sqrt{s}\delta d_0 \|\mathcal{A}^*(\mathbf{e})\| \end{aligned}$$

where $\|\cdot\|$ and $\|\cdot\|_*$ are a pair of dual norms and

$$\|\mathbf{U} + 2\mathbf{V}\|_F \leq \|\mathbf{U} + \mathbf{V}\|_F + \|\mathbf{V}\|_F \leq \delta d_0 + 0.2\delta d_0 \leq 1.2\delta d_0.$$

Combining the estimation of $\text{Re}(I_{01})$ and $\text{Re}(I_{02})$ above leads to

$$\text{Re}(\langle \nabla F_{\mathbf{h}}, \Delta \mathbf{h} \rangle + \langle \nabla F_{\mathbf{x}}, \Delta \mathbf{x} \rangle) \geq \frac{\delta^2 d_0^2}{8} - 2\sqrt{s}\delta d_0 \|\mathcal{A}^*(\mathbf{e})\|.$$

□

Lemma 6.10. (Local Regularity for $G(\mathbf{h}, \mathbf{x})$) For any $(\mathbf{h}, \mathbf{x}) \in \mathcal{N}_d \cap \mathcal{N}_\epsilon$ with $\epsilon \leq \frac{1}{15}$ and $\frac{9}{10}d_0 \leq d \leq \frac{11}{10}d_0$, $\frac{9}{10}d_{i0} \leq d_i \leq \frac{11}{10}d_{i0}$, the following inequality holds uniformly

$$\text{Re}(\langle \nabla G_{\mathbf{h}_i}, \Delta \mathbf{h}_i \rangle + \langle \nabla G_{\mathbf{x}_i}, \Delta \mathbf{x}_i \rangle) \geq 2\delta_0 \sqrt{\rho G_i(\mathbf{h}_i, \mathbf{x}_i)} = \frac{\delta}{5} \sqrt{\rho G_i(\mathbf{h}_i, \mathbf{x}_i)}, \quad (6.31)$$

where $\rho \geq d^2 + 2\|\mathbf{e}\|^2$. Immediately, we have

$$\text{Re}(\langle \nabla G_{\mathbf{h}}, \Delta \mathbf{h} \rangle + \langle \nabla G_{\mathbf{x}}, \Delta \mathbf{x} \rangle) = \sum_{i=1}^s \text{Re}(\langle \nabla G_{\mathbf{h}_i}, \Delta \mathbf{h}_i \rangle + \langle \nabla G_{\mathbf{x}_i}, \Delta \mathbf{x}_i \rangle) \geq \frac{\delta}{5} \sqrt{\rho G(\mathbf{h}, \mathbf{x})}. \quad (6.32)$$

Remark 6.11. For the local regularity condition for $G(\mathbf{h}, \mathbf{x})$, we use the results from [18] when $s = 1$. This is because each component $G_i(\mathbf{h}, \mathbf{x})$ only depends on $(\mathbf{h}_i, \mathbf{x}_i)$ by definition and thus the lower bound of $\text{Re}(\langle \nabla G_{\mathbf{h}_i}, \Delta \mathbf{h}_i \rangle + \langle \nabla G_{\mathbf{x}_i}, \Delta \mathbf{x}_i \rangle)$ is completely determined by $(\mathbf{h}_i, \mathbf{x}_i)$ and δ_0 , and is independent of s .

Proof: For each $i : 1 \leq i \leq s$, $\nabla G_{\mathbf{h}_i}$ (or $\nabla G_{\mathbf{x}_i}$) only depends on \mathbf{h}_i (or \mathbf{x}_i) and there holds

$$\text{Re}(\langle \nabla G_{\mathbf{h}_i}, \Delta \mathbf{h}_i \rangle + \langle \nabla G_{\mathbf{x}_i}, \Delta \mathbf{x}_i \rangle) \geq 2\delta_0 \sqrt{\rho G_i(\mathbf{h}_i, \mathbf{x}_i)} = \frac{\delta}{5} \sqrt{\rho G_i(\mathbf{h}_i, \mathbf{x}_i)},$$

which follows exactly from Lemma 5.17 in [18]. For (6.32), by definition of $\nabla G_{\mathbf{h}}$ and $\nabla G_{\mathbf{x}}$ in (2.22),

$$\begin{aligned} \text{Re}(\langle \nabla G_{\mathbf{h}}, \Delta \mathbf{h} \rangle + \langle \nabla G_{\mathbf{x}}, \Delta \mathbf{x} \rangle) &= \sum_{i=1}^s \text{Re}(\langle \nabla G_{\mathbf{h}_i}, \Delta \mathbf{h}_i \rangle + \langle \nabla G_{\mathbf{x}_i}, \Delta \mathbf{x}_i \rangle) \\ &\geq \frac{\delta}{5} \sum_{i=1}^s \sqrt{\rho G_i(\mathbf{h}_i, \mathbf{x}_i)} \geq \frac{\delta}{5} \sqrt{\rho G(\mathbf{h}, \mathbf{x})} \end{aligned}$$

where $G(\mathbf{h}, \mathbf{x}) = \sum_{i=1}^s G_i(\mathbf{h}_i, \mathbf{x}_i)$.

□

Lemma 6.12. (Proof of the Local Regularity Condition) Conditioned on (5.3), for the objective function $\tilde{F}(\mathbf{h}, \mathbf{x})$ in (2.17), there exists a positive constant ω such that

$$\|\nabla \tilde{F}(\mathbf{h}, \mathbf{x})\|^2 \geq \omega \left[\tilde{F}(\mathbf{h}, \mathbf{x}) - c \right]_+ \quad (6.33)$$

with $c = \|\mathbf{e}\|^2 + 2000s\|\mathcal{A}^*(\mathbf{e})\|^2$ and $\omega = \frac{d_0}{7000}$ for all $(\mathbf{h}, \mathbf{x}) \in \mathcal{N}_d \cap \mathcal{N}_\mu \cap \mathcal{N}_\epsilon$. Here we set $\rho \geq d^2 + 2\|\mathbf{e}\|^2$.

Proof: Following from Lemma 6.9 and Lemma 6.10, we have

$$\begin{aligned}\operatorname{Re}(\langle \nabla F_{\mathbf{h}}, \Delta \mathbf{h} \rangle + \langle \nabla F_{\mathbf{x}}, \Delta \mathbf{x} \rangle) &\geq \frac{\delta^2 d_0^2}{8} - 2\sqrt{s}\delta d_0 \|\mathcal{A}^*(\mathbf{e})\| \\ \operatorname{Re}(\langle \nabla G_{\mathbf{h}}, \Delta \mathbf{h} \rangle + \langle \nabla G_{\mathbf{x}}, \Delta \mathbf{x} \rangle) &\geq \frac{\delta d}{5} \sqrt{G(\mathbf{h}, \mathbf{x})} \geq \frac{9\delta d_0}{50} \sqrt{G(\mathbf{h}, \mathbf{x})}\end{aligned}$$

for all $(\mathbf{h}, \mathbf{x}) \in \mathcal{N}_d \cap \mathcal{N}_\mu \cap \mathcal{N}_\epsilon$ where $\rho \geq d^2 + 2\|\mathbf{e}\|^2 \geq d^2$ and $\frac{9}{10}d_0 \leq d \leq \frac{11}{10}d_0$. Adding them together gives $\operatorname{Re}(\langle \nabla \tilde{F}_{\mathbf{h}}, \Delta \mathbf{h} \rangle + \langle \nabla \tilde{F}_{\mathbf{x}}, \Delta \mathbf{x} \rangle)$ on the left side. Moreover, Cauchy-Schwarz inequality implies

$$\operatorname{Re}(\langle \nabla \tilde{F}_{\mathbf{h}}, \Delta \mathbf{h} \rangle + \langle \nabla \tilde{F}_{\mathbf{x}}, \Delta \mathbf{x} \rangle) \leq 4\delta\sqrt{d_0} \|\nabla \tilde{F}(\mathbf{h}, \mathbf{x})\|$$

where both $\|\Delta \mathbf{h}\|^2$ and $\|\Delta \mathbf{x}\|^2$ are bounded by $8\delta^2 d_0$ in Lemma 6.8 since

$$\|\Delta \mathbf{h}\|^2 = \sum_{i=1}^s \|\Delta \mathbf{h}_i\|^2 \leq \sum_{i=1}^s (7.5\delta_i^2 + 2.88\delta_0^2) d_{i0} \leq 8\delta^2 d_0.$$

Therefore,

$$\frac{\delta^2 d_0^2}{8} + \frac{9\delta d_0 \sqrt{G(\mathbf{h}, \mathbf{x})}}{50} - 2\sqrt{s}\delta d_0 \|\mathcal{A}^*(\mathbf{e})\| \leq 4\delta\sqrt{d_0} \|\nabla \tilde{F}(\mathbf{h}, \mathbf{x})\|. \quad (6.34)$$

Dividing both sides of (6.34) by δd_0 , we obtain

$$\begin{aligned}\frac{4}{\sqrt{d_0}} \|\nabla \tilde{F}(\mathbf{h}, \mathbf{x})\| &\geq \frac{\delta d_0}{12} + \frac{9}{50} \sqrt{G(\mathbf{h}, \mathbf{x})} + \frac{\delta d_0}{24} - 2\sqrt{s} \|\mathcal{A}^*(\mathbf{e})\| \\ &\geq \frac{1}{6\sqrt{6}} [\sqrt{F_0(\mathbf{h}, \mathbf{x})} + \sqrt{G(\mathbf{h}, \mathbf{x})}] + \frac{\delta d_0}{24} - 2\sqrt{s} \|\mathcal{A}^*(\mathbf{e})\|\end{aligned}$$

where the Local RIP condition (5.3) implies $F_0(\mathbf{h}, \mathbf{x}) \leq \frac{3}{2}\delta^2 d_0^2$ and hence $\frac{\delta d_0}{12} \geq \frac{1}{6\sqrt{6}} \sqrt{F_0(\mathbf{h}, \mathbf{x})}$, where $F_0(\mathbf{h}, \mathbf{x})$ is defined in (2.12).

Note that (5.6) gives

$$\sqrt{2[\operatorname{Re}(\langle \mathcal{A}^*(\mathbf{e}), \mathbf{X} - \mathbf{X}_0 \rangle)]_+} \leq \sqrt{2\sqrt{2s} \|\mathcal{A}^*(\mathbf{e})\| \delta d_0} \leq \frac{\sqrt{6}\delta d_0}{4} + \frac{4\sqrt{s}}{\sqrt{6}} \|\mathcal{A}^*(\mathbf{e})\|. \quad (6.35)$$

By (6.35) and $\tilde{F}(\mathbf{h}, \mathbf{x}) - \|\mathbf{e}\|^2 \leq F_0(\mathbf{h}, \mathbf{x}) + 2[\operatorname{Re}(\langle \mathcal{A}^*(\mathbf{e}), \mathbf{X} - \mathbf{X}_0 \rangle)]_+ + G(\mathbf{h}, \mathbf{x})$, there holds

$$\begin{aligned}\frac{4}{\sqrt{d_0}} \|\nabla \tilde{F}(\mathbf{h}, \mathbf{x})\| &\geq \frac{1}{6\sqrt{6}} \left[\left(\sqrt{F_0(\mathbf{h}, \mathbf{x})} + \sqrt{2[\operatorname{Re}(\langle \mathcal{A}^*(\mathbf{e}), \mathbf{X} - \mathbf{X}_0 \rangle)]_+} + \sqrt{G(\mathbf{h}, \mathbf{x})} \right) \right. \\ &\quad \left. + \frac{\delta d_0}{24} - \frac{1}{6\sqrt{6}} \left(\frac{\sqrt{6}\delta d_0}{4} + \frac{4\sqrt{s}}{\sqrt{6}} \|\mathcal{A}^*(\mathbf{e})\| \right) - 2\sqrt{s} \|\mathcal{A}^*(\mathbf{e})\| \right] \\ &\geq \frac{1}{6\sqrt{6}} \left[\sqrt{[\tilde{F}(\mathbf{h}, \mathbf{x}) - \|\mathbf{e}\|^2]_+} - \sqrt{1000s} \|\mathcal{A}^*(\mathbf{e})\| \right].\end{aligned}$$

For any nonnegative real numbers a and b , we have $[\sqrt{(x-a)_+} - b]_+ + b \geq \sqrt{(x-a)_+}$ and it implies

$$(x-a)_+ \leq 2([\sqrt{(x-a)_+} - b]_+^2 + b^2) \implies [\sqrt{(x-a)_+} - b]_+ \geq \frac{(x-a)_+}{2} - b^2.$$

Therefore, by setting $a = \|\mathbf{e}\|^2$ and $b = \sqrt{1000s} \|\mathcal{A}^*(\mathbf{e})\|$, there holds

$$\begin{aligned}\|\nabla \tilde{F}(\mathbf{h}, \mathbf{x})\|^2 &\geq \frac{d_0}{3500} \left[\frac{\tilde{F}(\mathbf{h}, \mathbf{x}) - \|\mathbf{e}\|^2}{2} - 1000s \|\mathcal{A}^*(\mathbf{e})\|^2 \right]_+ \\ &\geq \frac{d_0}{7000} \left[\tilde{F}(\mathbf{h}, \mathbf{x}) - (\|\mathbf{e}\|^2 + 2000s \|\mathcal{A}^*(\mathbf{e})\|^2) \right]_+.\end{aligned}$$

□

6.4 Local smoothness

Lemma 6.13. *Conditioned on (5.3), (5.4) and (6.3), for any $\mathbf{z} := (\mathbf{h}, \mathbf{x}) \in \mathbb{C}^{(K+N)s}$ and $\mathbf{w} := (\mathbf{u}, \mathbf{v}) \in \mathbb{C}^{(K+N)s}$ such that \mathbf{z} and $\mathbf{z} + \mathbf{w} \in \mathcal{N}_\epsilon \cap \mathcal{N}_{\tilde{F}}$, there holds*

$$\|\nabla \tilde{F}(\mathbf{z} + \mathbf{w}) - \nabla \tilde{F}(\mathbf{z})\| \leq C_L \|\mathbf{w}\|,$$

with

$$C_L \leq \left(10 \|\mathcal{A}\|^2 d_0 + \frac{2\rho}{\min d_{i0}} \left(5 + \frac{2L}{\mu^2} \right) \right)$$

where $\rho \geq d^2 + 2\|\mathbf{e}\|^2$ and $\|\mathcal{A}\| \leq \sqrt{s(N \log(NL/2) + (\gamma + \log s) \log L)}$ holds with probability at least $1 - L^{-\gamma}$ from Lemma 6.2.

In particular, $L = \mathcal{O}((\mu^2 + \sigma^2)s(K + N) \log^2 L)$ and $\|\mathbf{e}\|^2 = \mathcal{O}(\sigma^2 d_0^2)$ follows from $\|\mathbf{e}\|^2 \sim \frac{\sigma^2 d_0^2}{2L} \chi_{2L}^2$ and (6.13). Therefore, C_L can be simplified to

$$C_L = \mathcal{O}(d_0 s \kappa (1 + \sigma^2) (K + N) \log^2 L)$$

by choosing $\rho \approx d^2 + 2\|\mathbf{e}\|^2$.

Proof: By Lemma 5.5, we know that both $\mathbf{z} = (\mathbf{h}, \mathbf{x})$ and $\mathbf{z} + \mathbf{w} = (\mathbf{h} + \mathbf{u}, \mathbf{x} + \mathbf{v}) \in \mathcal{N}_d \cap \mathcal{N}_\mu \cap \mathcal{N}_\epsilon$. Note that

$$\nabla \tilde{F} = (\nabla \tilde{F}_\mathbf{h}, \nabla \tilde{F}_\mathbf{x}) = (\nabla F_\mathbf{h} + \nabla G_\mathbf{h}, \nabla F_\mathbf{x} + \nabla G_\mathbf{x}),$$

where (2.18), (2.19), (2.20) and (2.21) give $\nabla F_\mathbf{h}, \nabla F_\mathbf{x}, \nabla G_\mathbf{h}$ and $\nabla G_\mathbf{x}$. It suffices to find out the Lipschitz constants for all of those four functions.

Step 1: We first estimate the Lipschitz constant for $\nabla F_\mathbf{h}$ and the result can be applied to $\nabla F_\mathbf{x}$ due to symmetry.

$$\begin{aligned} \nabla F_\mathbf{h}(\mathbf{z} + \mathbf{w}) - \nabla F_\mathbf{h}(\mathbf{z}) &= \mathcal{A}^* \mathcal{A}(\mathcal{H}(\mathbf{h} + \mathbf{u}, \mathbf{x} + \mathbf{v}))(\mathbf{x} + \mathbf{v}) - [\mathcal{A}^* \mathcal{A}(\mathcal{H}(\mathbf{h}, \mathbf{x}))\mathbf{x} + \mathcal{A}^*(\mathbf{y})\mathbf{v}] \\ &= \mathcal{A}^*(\mathcal{A}(\mathcal{H}(\mathbf{h} + \mathbf{u}, \mathbf{x} + \mathbf{v}) - \mathcal{H}(\mathbf{h}, \mathbf{x}))) (\mathbf{x} + \mathbf{v}) \\ &\quad + \mathcal{A}^* \mathcal{A}(\mathcal{H}(\mathbf{h}, \mathbf{x}) - \mathcal{H}(\mathbf{h}_0, \mathbf{x}_0)) \mathbf{v} - \mathcal{A}^*(\mathbf{e}) \mathbf{v} \\ &= \mathcal{A}^*(\mathcal{A}(\mathcal{H}(\mathbf{h} + \mathbf{u}, \mathbf{v}) + \mathcal{H}(\mathbf{u}, \mathbf{x}))) (\mathbf{x} + \mathbf{v}) \\ &\quad + \mathcal{A}^* \mathcal{A}(\mathcal{H}(\mathbf{h}, \mathbf{x}) - \mathcal{H}(\mathbf{h}_0, \mathbf{x}_0)) \mathbf{v} - \mathcal{A}^*(\mathbf{e}) \mathbf{v}. \end{aligned}$$

Note that $\|\mathcal{H}(\mathbf{h}, \mathbf{x})\|_F \leq \sqrt{\sum_{i=1}^s \|\mathbf{h}_i\|^2 \|\mathbf{x}_i\|^2} \leq \|\mathbf{h}\| \|\mathbf{x}\|$ and $\mathbf{z}, \mathbf{z} + \mathbf{w} \in \mathcal{N}_d$ directly implies

$$\|\mathcal{H}(\mathbf{u}, \mathbf{x}) + \mathcal{H}(\mathbf{h} + \mathbf{u}, \mathbf{v})\|_F \leq \|\mathbf{u}\| \|\mathbf{x}\| + \|\mathbf{h} + \mathbf{u}\| \|\mathbf{v}\| \leq 2\sqrt{d_0} (\|\mathbf{u}\| + \|\mathbf{v}\|)$$

where $\|\mathbf{h} + \mathbf{u}\| \leq 2\sqrt{d_0}$. Moreover, (5.3) implies

$$\|\mathcal{H}(\mathbf{h}, \mathbf{x}) - \mathcal{H}(\mathbf{h}_0, \mathbf{x}_0)\|_F \leq \epsilon d_0$$

since $\mathbf{z} \in \mathcal{N}_d \cap \mathcal{N}_\mu \cap \mathcal{N}_\epsilon$. Combined with $\|\mathcal{A}^*(\mathbf{e})\| \leq \epsilon d_0$ in (5.4) and $\|\mathbf{x} + \mathbf{v}\| \leq 2\sqrt{d_0}$, we have

$$\begin{aligned} \|\nabla F_\mathbf{h}(\mathbf{z} + \mathbf{w}) - \nabla F_\mathbf{h}(\mathbf{z})\| &\leq 4d_0 \|\mathcal{A}\|^2 (\|\mathbf{u}\| + \|\mathbf{v}\|) + \epsilon d_0 \|\mathcal{A}\|^2 \|\mathbf{v}\| + \epsilon d_0 \|\mathbf{v}\| \\ &\leq 5d_0 \|\mathcal{A}\|^2 (\|\mathbf{u}\| + \|\mathbf{v}\|). \end{aligned} \tag{6.36}$$

Due to the symmetry between $\nabla F_\mathbf{h}$ and $\nabla F_\mathbf{x}$, we have,

$$\|\nabla F_\mathbf{x}(\mathbf{z} + \mathbf{w}) - \nabla F_\mathbf{x}(\mathbf{z})\| \leq 5d_0 \|\mathcal{A}\|^2 (\|\mathbf{u}\| + \|\mathbf{v}\|). \tag{6.37}$$

In other words,

$$\|\nabla F(\mathbf{z} + \mathbf{w}) - \nabla F(\mathbf{z})\| \leq 5\sqrt{2}d_0 \|\mathcal{A}\|^2 (\|\mathbf{u}\| + \|\mathbf{v}\|) \leq 10d_0 \|\mathcal{A}\|^2 \|\mathbf{w}\|$$

where $\|\mathbf{u}\| + \|\mathbf{v}\| \leq \sqrt{2} \|\mathbf{w}\|$.

Step 2: We estimate the upper bound of $\|\nabla G_{\mathbf{x}_i}(\mathbf{z}_i + \mathbf{w}_i) - \nabla G_{\mathbf{x}_i}(\mathbf{z}_i)\|$. Implied by Lemma 5.19 in [18], we have

$$\|\nabla G_{\mathbf{x}_i}(\mathbf{z}_i + \mathbf{w}_i) - \nabla G_{\mathbf{x}_i}(\mathbf{z}_i)\| \leq \frac{5d_{i0}\rho}{d_i^2} \|\mathbf{v}_i\|. \quad (6.38)$$

Step 3: We estimate the upper bound of $\|\nabla G_{\mathbf{h}_i}(\mathbf{z} + \mathbf{w}) - \nabla G_{\mathbf{h}_i}(\mathbf{z})\|$. Denote

$$\begin{aligned} \nabla G_{\mathbf{h}_i}(\mathbf{z} + \mathbf{w}) - \nabla G_{\mathbf{h}_i}(\mathbf{z}) &= \underbrace{\frac{\rho}{2d_i} \left[G'_0 \left(\frac{\|\mathbf{h}_i + \mathbf{u}_i\|^2}{2d_i} \right) (\mathbf{h}_i + \mathbf{u}_i) - G'_0 \left(\frac{\|\mathbf{h}_i\|^2}{2d_i} \right) \mathbf{h}_i \right]}_{\mathbf{j}_1} \\ &\quad + \underbrace{\frac{\rho L}{8d_i\mu^2} \sum_{l=1}^L \left[G'_0 \left(\frac{L|\mathbf{b}_l^*(\mathbf{h}_i + \mathbf{u}_i)|^2}{8d_i\mu^2} \right) \mathbf{b}_l^*(\mathbf{h}_i + \mathbf{u}_i) - G'_0 \left(\frac{L|\mathbf{b}_l^*\mathbf{h}_i|^2}{8d_i\mu^2} \right) \mathbf{b}_l^*\mathbf{h}_i \right]}_{\mathbf{j}_2} \mathbf{b}_l. \end{aligned}$$

Following the same estimation of \mathbf{j}_1 and \mathbf{j}_2 in Lemma 5.19 of [18], we have

$$\|\mathbf{j}_1\| \leq \frac{5d_{i0}\rho}{d_i^2} \|\mathbf{u}_i\|, \quad \|\mathbf{j}_2\| \leq \frac{3\rho L d_{i0}}{2d_i^2 \mu^2} \|\mathbf{u}_i\|. \quad (6.39)$$

Therefore, combining (6.38) and (6.39) gives

$$\begin{aligned} \|\nabla G(\mathbf{z} + \mathbf{w}) - \nabla G(\mathbf{z})\| &= \sqrt{\sum_{i=1}^s (\|\nabla G_{\mathbf{h}_i}(\mathbf{z} + \mathbf{w}) - \nabla G_{\mathbf{h}_i}(\mathbf{z})\|^2 + \|\nabla G_{\mathbf{x}_i}(\mathbf{z} + \mathbf{w}) - \nabla G_{\mathbf{x}_i}(\mathbf{z})\|^2)} \\ &\leq \max \left\{ \frac{5d_{i0}\rho}{d_i^2} + \frac{3\rho L d_{i0}}{2d_i^2 \mu^2} \right\} \sqrt{\sum_{i=1}^s \|\mathbf{u}_i\|^2} + \max \left\{ \frac{5d_{i0}\rho}{d_i^2} \right\} \sqrt{\sum_{i=1}^s \|\mathbf{v}_i\|^2} \\ &\leq \max \left\{ \frac{5d_{i0}\rho}{d_i^2} + \frac{3\rho L d_{i0}}{2d_i^2 \mu^2} \right\} \|\mathbf{u}\| + \max \left\{ \frac{5d_{i0}\rho}{d_i^2} \right\} \|\mathbf{v}\| \\ &\leq \frac{2\rho}{\min d_{i0}} \left(5 + \frac{2L}{\mu^2} \right) \|\mathbf{w}\|. \end{aligned}$$

In summary, the Lipschitz constant C_L of $\tilde{F}(\mathbf{z})$ has an upper bound as follows:

$$\begin{aligned} \|\nabla \tilde{F}(\mathbf{z} + \mathbf{w}) - \nabla \tilde{F}(\mathbf{z})\| &\leq \|\nabla F(\mathbf{z} + \mathbf{w}) - \nabla F(\mathbf{z})\| + \|\nabla G(\mathbf{z} + \mathbf{w}) - \nabla G(\mathbf{z})\| \\ &\leq \left(10\|\mathcal{A}\|^2 d_0 + \frac{2\rho}{\min d_{i0}} \left(5 + \frac{2L}{\mu^2} \right) \right) \|\mathbf{w}\|. \end{aligned}$$

□

6.5 Robustness condition and spectral initialization

In this section, we will prove the robustness condition (5.4) and also Theorem 3.2. To prove (5.4), it suffices to show the following lemma, which is a more general version of (5.4).

Lemma 6.14. *Consider a sequence of Gaussian independent random variable $\mathbf{c} = (c_1, \dots, c_L) \in \mathbb{C}^L$ where $c_l \sim \mathcal{CN}(0, \frac{\lambda_l^2}{L})$ with $\lambda_l \leq \lambda$. Moreover, we assume \mathcal{A}_i in (2.2) is independent of \mathbf{c} . Then there holds*

$$\|\mathcal{A}^*(\mathbf{c})\| = \left\| \max_{1 \leq i \leq s} \|\mathcal{A}_i^*(\mathbf{c})\| \right\| \leq \xi$$

with probability at least $1 - L^{-\gamma}$ if $L \geq C_{\gamma+\log(s)} \left(\frac{\lambda}{\xi} + \frac{\lambda^2}{\xi^2} \right) \max\{K, N\} \log L / \xi^2$.

Proof: It suffices to show that $\max_{1 \leq i \leq s} \|\mathcal{A}_i^*(\mathbf{c})\| \leq \xi$. For each fixed $i : 1 \leq i \leq s$,

$$\mathcal{A}_i^*(\mathbf{c}) = \sum_{l=1}^L c_l \mathbf{b}_l \mathbf{a}_{il}^*$$

The key is to apply the matrix Bernstein inequality (6.52) and we need to estimate $\|\mathcal{Z}_l\|_{\psi_1}$, and the variance of $\sum_{l=1}^L \mathcal{Z}_l$. For each l , $\|c_l \mathbf{b}_l \mathbf{a}_{il}^*\|_{\psi_1} \leq \frac{\lambda \sqrt{KN}}{L}$ follows from (6.57). Moreover, the variance of $\mathcal{A}_i^*(\mathbf{c})$ is bounded by $\frac{\lambda^2 \max\{K, N\}}{L}$ since

$$\begin{aligned} \mathbb{E}[\mathcal{A}_i^*(\mathbf{c})(\mathcal{A}_i^*(\mathbf{c}))^*] &= \sum_{l=1}^L \mathbb{E}(|c_l|^2 \|\mathbf{a}_{il}\|^2) \mathbf{b}_l \mathbf{b}_l^* = \frac{N}{L} \sum_{l=1}^L \lambda_l^2 \mathbf{b}_l \mathbf{b}_l^* \preceq \frac{\lambda^2 N}{L}, \\ \mathbb{E}[(\mathcal{A}_i^*(\mathbf{c}))^* (\mathcal{A}_i^*(\mathbf{c}))] &= \sum_{l=1}^L \|\mathbf{b}_l\|^2 \mathbb{E}(|c_l|^2 \mathbf{a}_{il} \mathbf{a}_{il}^*) = \frac{K}{L^2} \sum_{l=1}^L \lambda_l^2 \mathbf{I}_N \preceq \frac{\lambda^2 K}{L}. \end{aligned}$$

Letting $t = \gamma \log L$ and applying (6.52) leads to

$$\|\mathcal{A}_i^*(\mathbf{c})\| \leq C_0 \max \left\{ \frac{\lambda \sqrt{KN} \log^2 L}{L}, \sqrt{\frac{C_\gamma \lambda^2 \max\{K, N\} \log L}{L}} \right\} \leq \xi.$$

Therefore, by taking the union bound over $1 \leq i \leq s$,

$$\|\mathcal{A}_i^*(\mathbf{c})\| \leq \xi$$

with probability at least $1 - L^{-\gamma}$ if $L \geq C_{\gamma + \log(s)} \left(\frac{\lambda}{\xi} + \frac{\lambda^2}{\xi^2} \right) \max\{K, N\} \log^2 L$. \square

The robustness condition is an immediate result of Lemma 6.14 by setting $\xi = \frac{\varepsilon d_0}{10\sqrt{2s\kappa}}$ and $\lambda = \sigma d_0$.

Corollary 6.15. [Robustness Condition] For $\mathbf{e} \sim \mathcal{CN}(\mathbf{0}, \frac{\sigma^2 d_0^2}{L} \mathbf{I}_L)$

$$\|\mathcal{A}_i^*(\mathbf{e})\| \leq \frac{\varepsilon d_0}{10\sqrt{2s\kappa}}, \quad \forall 1 \leq i \leq s$$

with probability at least $1 - L^{-\gamma}$ if $L \geq C_\gamma \left(\frac{s^2 \kappa^2 \sigma^2}{\varepsilon^2} + \frac{s\kappa\sigma}{\varepsilon} \right) \max\{K, N\} \log L$.

Lemma 6.16. For $\mathbf{e} \sim \mathcal{CN}(\mathbf{0}, \frac{\sigma^2 d_0^2}{L} \mathbf{I}_L)$, there holds

$$\|\mathcal{A}_i^*(\mathbf{y}) - \mathbf{h}_{i0} \mathbf{x}_{i0}^*\| \leq \xi d_{i0}, \quad \forall 1 \leq i \leq s \quad (6.40)$$

with probability at least $1 - L^{-\gamma}$ if $L \geq C_{\gamma + \log(s)} s \kappa^2 (\mu_h^2 + \sigma^2) \max\{K, N\} \log L / \xi^2$.

Remark 6.17. The success of the initialization algorithm completely relies on the lemma above. As mentioned in Section 3, $\mathbb{E}(\mathcal{A}_i^*(\mathbf{y})) = \mathbf{h}_{i0} \mathbf{x}_{i0}^*$ and Lemma 6.40 confirms that $\mathcal{A}_i^*(\mathbf{y})$ is close to $\mathbf{h}_{i0} \mathbf{x}_{i0}^*$ in operator norm and hence the spectral method is able to give us a reliable initialization.

Proof: Note that

$$\mathcal{A}_i^*(\mathbf{y}) = \mathcal{A}_i^* \mathcal{A}_i(\mathbf{h}_{i0} \mathbf{x}_{i0}^*) + \mathcal{A}_i^*(\mathbf{w}_i)$$

where

$$\mathbf{w}_i = \mathbf{y} - \mathcal{A}_i(\mathbf{h}_{i0} \mathbf{x}_{i0}^*) = \sum_{j \neq i} \mathcal{A}_j(\mathbf{h}_{j0} \mathbf{x}_{j0}^*) + \mathbf{e} \quad (6.41)$$

is independent of \mathcal{A}_i . The proof consists of two parts: 1. show that $\|\mathcal{A}_i^* \mathcal{A}_i(\mathbf{h}_{i0} \mathbf{x}_{i0}^*) - \mathbf{h}_{i0} \mathbf{x}_{i0}^*\| \leq \frac{\xi d_{i0}}{2}$; 2. prove that $\|\mathcal{A}_i^*(\mathbf{w}_i)\| \leq \frac{\xi d_{i0}}{2}$.

Part I: Following from the definition of \mathcal{A}_i and \mathcal{A}_i^* in (2.2),

$$\mathcal{A}_i^* \mathcal{A}_i(\mathbf{h}_{i0} \mathbf{x}_{i0}^*) - \mathbf{h}_{i0} \mathbf{x}_{i0}^* = \sum_{l=1}^L \underbrace{\mathbf{b}_l \mathbf{b}_l^* \mathbf{h}_{i0} \mathbf{x}_{i0}^* (\mathbf{a}_{il} \mathbf{a}_{il}^* - \mathbf{I}_N)}_{\text{defined as } \mathcal{Z}_l}.$$

where $\mathbf{B}^* \mathbf{B} = \mathbf{I}_K$. The sub-exponential norm of \mathcal{Z}_l is bounded by

$$\|\mathcal{Z}_l\|_{\psi_1} \leq \max_{1 \leq l \leq L} \|\mathbf{b}_l\| \|\mathbf{b}_l^* \mathbf{h}_{i0}\| \|(\mathbf{a}_{il} \mathbf{a}_{il}^* - \mathbf{I}_N) \mathbf{x}_{i0}\|_{\psi_1} \leq \frac{\mu \sqrt{KN} d_{i0}}{L}$$

where $\|\mathbf{b}_l\| = \sqrt{\frac{K}{L}}$, $\max_l |\mathbf{b}_l^* \mathbf{h}_{i0}|^2 \leq \frac{\mu^2 d_{i0}}{L}$ and $\|(\mathbf{a}_{il} \mathbf{a}_{il}^* - \mathbf{I}_N) \mathbf{x}_{i0}\|_{\psi_1} \leq \sqrt{N d_{i0}}$ follows from (6.55).

We proceed to estimate the variance of $\sum_{l=1}^L \mathcal{Z}_l$ by using (6.54) and (6.56):

$$\begin{aligned} \left\| \sum_{l=1}^L \mathbb{E}(\mathcal{Z}_l \mathcal{Z}_l^*) \right\| &= \left\| \sum |\mathbf{b}_l^* \mathbf{h}_{i0}|^2 \mathbf{x}_{i0}^* \mathbb{E}(\mathbf{a}_{il} \mathbf{a}_{il}^* - \mathbf{I}_N)^2 \mathbf{x}_{i0} \mathbf{b}_l \mathbf{b}_l^* \right\| \leq \frac{\mu^2 N d_{i0}^2}{L}, \\ \left\| \sum_{l=1}^L \mathbb{E}(\mathcal{Z}_l^* \mathcal{Z}_l) \right\| &= \frac{K}{L} \left\| \sum_{l=1}^L |\mathbf{b}_l^* \mathbf{h}_{i0}|^2 \mathbb{E}[(\mathbf{a}_{il} \mathbf{a}_{il}^* - \mathbf{I}_N) \mathbf{x}_{i0} \mathbf{x}_{i0}^* (\mathbf{a}_{il} \mathbf{a}_{il}^* - \mathbf{I}_N)] \right\| \leq \frac{K d_{i0}^2}{L}. \end{aligned}$$

Therefore, the variance of $\sum_{l=1}^L \mathcal{Z}_l$ is bounded by $\frac{\max\{K, \mu_h^2 N\} d_{i0}^2}{L}$. By applying matrix Bernstein inequality (6.52) and taking the union bound over all i , we prove that

$$\|\mathcal{A}_i^* \mathcal{A}_i(\mathbf{h}_{i0} \mathbf{x}_{i0}^*) - \mathbf{h}_{i0} \mathbf{x}_{i0}^*\| \leq \frac{\xi d_{i0}}{2}, \quad \forall 1 \leq i \leq s$$

holds with probability at least $1 - L^{-\gamma+1}$ if $L \geq C_{\gamma+\log(s)} \max\{K, \mu_h^2 N\} \log L / \xi^2$.

Part II: For each $1 \leq l \leq L$, the l -th entry of \mathbf{w}_i in (6.41), i.e., $(\mathbf{w}_i)_l = \sum_{j \neq i} \mathbf{b}_l^* \mathbf{h}_{j0} \mathbf{x}_{j0}^* \mathbf{a}_{jl} + e_l$, is independent of $\mathbf{b}_l^* \mathbf{h}_{i0} \mathbf{x}_{i0}^* \mathbf{a}_{il}$ and obeys $\mathcal{CN}(0, \frac{\sigma^2}{L})$. Here

$$\begin{aligned} \sigma_{il}^2 &= L \mathbb{E} |(\mathbf{w}_i)_l|^2 = L \sum_{j \neq i} |\mathbf{b}_l^* \mathbf{h}_{j0}|^2 \|\mathbf{x}_{j0}\|^2 + \sigma^2 \|\mathbf{X}_0\|_F^2 \\ &\leq \mu_h^2 \sum_{j \neq i} \|\mathbf{h}_{j0}\|^2 \|\mathbf{x}_{j0}\|^2 + \sigma^2 \|\mathbf{X}_0\|_F^2 \leq (\mu_h^2 + \sigma^2) \|\mathbf{X}_0\|_F^2. \end{aligned}$$

This gives $\max_{i,l} \sigma_{il}^2 \leq (\mu_h^2 + \sigma^2) \|\mathbf{X}_0\|_F^2$. Thanks to the independence between \mathbf{w}_i and \mathcal{A}_i , applying Lemma 6.14 results in

$$\|\mathcal{A}_i^*(\mathbf{w}_i)\| \leq \frac{\xi d_{i0}}{2} \tag{6.42}$$

with probability $1 - L^{-\gamma+1}$ if $L \geq C \max\left(\frac{(\mu_h^2 + \sigma^2) \|\mathbf{X}_0\|_F^2}{\xi^2 d_{i0}^2}, \frac{\sqrt{\mu_h^2 + \sigma^2} \|\mathbf{X}_0\|_F}{\xi d_{i0}}\right) \max\{K, N\} \log L$.

Therefore, combining (6.41) with (6.42), we get

$$\|\mathcal{A}_i^*(\mathbf{y}) - \mathbf{h}_{i0} \mathbf{x}_{i0}^*\| \leq \|\mathcal{A}_i^* \mathcal{A}_i(\mathbf{h}_{i0} \mathbf{x}_{i0}^*) - \mathbf{h}_{i0} \mathbf{x}_{i0}^*\| + \|\mathcal{A}_i^*(\mathbf{w}_i)\| \leq \xi d_{i0}$$

for all $1 \leq i \leq s$ with probability at least $1 - L^{-\gamma+1}$ if

$$L \geq C_{\gamma+\log(s)} (\mu_h^2 + \sigma^2) s \kappa^2 \max\{K, N\} \log L / \xi^2$$

where $\|\mathbf{X}_0\|_F / d_{i0} \leq \sqrt{s} \kappa$. □

Before moving to the proof of Theorem 3.2, we introduce a property about the projection onto a closed convex set.

Lemma 6.18 (Theorem 2.8 in [12]). *Let $Q := \{\mathbf{w} \in \mathbb{C}^K | \sqrt{L} \|\mathbf{B}\mathbf{w}\|_\infty \leq 2\sqrt{d}\mu\}$ be a closed nonempty convex set. There holds*

$$\operatorname{Re}(\langle \mathbf{z} - \mathcal{P}_Q(\mathbf{z}), \mathbf{w} - \mathcal{P}_Q(\mathbf{z}) \rangle) \leq 0, \quad \forall \mathbf{w} \in Q, \mathbf{z} \in \mathbb{C}^K$$

where $\mathcal{P}_Q(\mathbf{z})$ is the projection of \mathbf{z} onto Q .

With this lemma, we can easily see

$$\|\mathbf{z} - \mathbf{w}\|^2 = \|\mathbf{z} - \mathcal{P}_Q(\mathbf{z})\|^2 + \|\mathcal{P}_Q(\mathbf{z}) - \mathbf{w}\|^2 + 2\operatorname{Re}(\langle \mathbf{z} - \mathcal{P}_Q(\mathbf{z}), \mathcal{P}_Q(\mathbf{z}) - \mathbf{w} \rangle) \geq \|\mathcal{P}_Q(\mathbf{z}) - \mathbf{w}\|^2 \quad (6.43)$$

for all $\mathbf{z} \in \mathbb{C}^K$ and $\mathbf{w} \in Q$. It means that projection onto nonempty closed convex set is non-expansive. Now we present the proof of Theorem 3.2.

Proof of Theorem 3.2. By choosing $L \geq C_{\gamma+\log(s)}(\mu_h^2 + \sigma^2)s^2\kappa^4 \max\{K, N\} \log L/\varepsilon^2$, we have

$$\|\mathcal{A}_i^*(\mathbf{y}) - \mathbf{h}_{i0}\mathbf{x}_{i0}^*\| \leq \xi d_{i0}, \quad \forall 1 \leq i \leq s \quad (6.44)$$

where $\xi = \frac{\varepsilon}{10\sqrt{2s\kappa}}$.

By applying the triangle inequality to (6.44), it is easy to see that

$$(1 - \xi)d_{i0} \leq d_i \leq (1 + \xi)d_{i0}, \quad |d_i - d_{i0}| \leq \xi d_{i0} \leq \frac{\varepsilon d_{i0}}{10\sqrt{2s\kappa}} < \frac{d_{i0}}{10}, \quad (6.45)$$

which gives $\frac{9}{10}d_{i0} \leq d_i \leq \frac{11}{10}d_{i0}$ where $d_i = \|\mathcal{A}_i^*(\mathbf{y})\|$ according to Algorithm 1.

Part I: Proof of $(\mathbf{u}^{(0)}, \mathbf{v}^{(0)}) \in \frac{1}{\sqrt{3}}\mathcal{N}_d \cap \frac{1}{\sqrt{3}}\mathcal{N}_\mu$ Note that $\mathbf{v}_i^{(0)} = \sqrt{d_i}\|\hat{\mathbf{x}}_{i0}\| = \sqrt{d_i}$ where $\hat{\mathbf{x}}_{i0}$ is the leading right singular vector of $\mathcal{A}_i^*(\mathbf{y})$. Therefore,

$$\|\mathbf{v}_i^{(0)}\| = \sqrt{d_i}\|\hat{\mathbf{x}}_{i0}\| = \sqrt{d_i} \leq \sqrt{(1 + \xi)d_{i0}} \leq \frac{2}{\sqrt{3}}\sqrt{d_{i0}}, \quad \forall 1 \leq i \leq s$$

which implies $\{\mathbf{v}_i^{(0)}\}_{i=1}^s \in \frac{1}{\sqrt{3}}\mathcal{N}_d$.

Now we will prove that $\mathbf{u}_i^{(0)} \in \frac{1}{\sqrt{3}}\mathcal{N}_d \cap \frac{1}{\sqrt{3}}\mathcal{N}_\mu$ by Lemma 6.18. By Algorithm 1, $\mathbf{u}_i^{(0)}$ is the minimizer to the function $f(\mathbf{z}) = \frac{1}{2}\|\mathbf{z} - \sqrt{d_i}\hat{\mathbf{h}}_{i0}\|^2$ over $Q_i := \{\mathbf{z} | \sqrt{L}\|\mathbf{B}\mathbf{z}\|_\infty \leq 2\sqrt{d_i}\mu\}$. Obviously, by definition, $\mathbf{u}_i^{(0)}$ is the projection of $\sqrt{d_i}\hat{\mathbf{h}}_{i0}$ onto Q_i . Note that $\mathbf{u}_i^{(0)} \in Q_i$ implies $\sqrt{L}\|\mathbf{B}\mathbf{u}_i^{(0)}\|_\infty \leq 2\sqrt{d_i}\mu \leq 2\sqrt{(1 + \xi)d_{i0}}\mu \leq \frac{4\sqrt{d_{i0}}\mu}{\sqrt{3}}$ and hence $\mathbf{u}_i^{(0)} \in \frac{1}{\sqrt{3}}\mathcal{N}_\mu$.

Moreover, due to (6.43), there holds

$$\|\sqrt{d_i}\hat{\mathbf{h}}_{i0} - \mathbf{w}\|^2 \geq \|\mathbf{u}_i^{(0)} - \mathbf{w}\|^2, \quad \forall \mathbf{w} \in Q_i \quad (6.46)$$

In particular, let $\mathbf{w} = \mathbf{0} \in Q_i$ and immediately we have

$$\|\mathbf{u}_i^{(0)}\|^2 \leq d_i \leq \frac{4}{3} \implies \mathbf{u}_i^{(0)} \in \frac{1}{\sqrt{3}}\mathcal{N}_\mu.$$

In other words, $\{(\mathbf{u}_i^{(0)}, \mathbf{v}_i^{(0)})\}_{i=1}^s \in \frac{1}{\sqrt{3}}\mathcal{N}_d \cap \frac{1}{\sqrt{3}}\mathcal{N}_\mu$.

Part II: Proof of $(\mathbf{u}^{(0)}, \mathbf{v}^{(0)}) \in \mathcal{N}_{\frac{2\varepsilon}{5\sqrt{s\kappa}}}$ We will show $\|\mathbf{u}_i^{(0)}(\mathbf{v}_i^{(0)})^* - \mathbf{h}_{i0}\mathbf{x}_{i0}^*\|_F \leq 4\xi d_{i0}$ for

$1 \leq i \leq s$ so that $\frac{\|\mathbf{u}_i^{(0)}(\mathbf{v}_i^{(0)})^* - \mathbf{h}_{i0}\mathbf{x}_{i0}^*\|_F}{d_{i0}} \leq \frac{2\varepsilon}{5\sqrt{s\kappa}}$.

First note that $\sigma_j(\mathcal{A}_i^*(\mathbf{y})) \leq \xi d_{i0}$ for all $j \geq 2$, which follows from Weyl's inequality [24] for singular values where $\sigma_j(\mathcal{A}_i^*(\mathbf{y}))$ denotes the j -th largest singular value of $\mathcal{A}_i^*(\mathbf{y})$. Hence there holds

$$\|d_i\hat{\mathbf{h}}_{i0}\hat{\mathbf{x}}_{i0}^* - \mathbf{h}_{i0}\mathbf{x}_{i0}^*\| \leq \|\mathcal{A}_i^*(\mathbf{y}) - d_i\hat{\mathbf{h}}_{i0}\hat{\mathbf{x}}_{i0}^*\| + \|\mathcal{A}_i^*(\mathbf{y}) - \mathbf{h}_{i0}\mathbf{x}_{i0}^*\| \leq 2\xi d_{i0}. \quad (6.47)$$

On the other hand, for any i ,

$$\begin{aligned} \left\| \left(\mathbf{I}_K - \frac{\mathbf{h}_{i0}\mathbf{h}_{i0}^*}{d_{i0}} \right) \hat{\mathbf{h}}_{i0} \right\| &= \left\| \left(\mathbf{I}_K - \frac{\mathbf{h}_{i0}\mathbf{h}_{i0}^*}{d_{i0}} \right) \hat{\mathbf{h}}_{i0} \hat{\mathbf{x}}_{i0}^* \hat{\mathbf{x}}_{i0} \hat{\mathbf{h}}_{i0}^* \right\| \\ &= \left\| \left(\mathbf{I}_K - \frac{\mathbf{h}_{i0}\mathbf{h}_{i0}^*}{d_{i0}} \right) \left[\frac{1}{d_{i0}} ((\mathcal{A}_i^*(\mathbf{y}) - d_i \hat{\mathbf{h}}_{i0} \hat{\mathbf{x}}_{i0}^*) + \hat{\mathbf{h}}_{i0} \hat{\mathbf{x}}_{i0}^* - \frac{\mathbf{h}_{i0}\mathbf{x}_{i0}^*}{d_{i0}}) \right] \hat{\mathbf{x}}_{i0} \hat{\mathbf{h}}_{i0}^* \right\| \\ &= \frac{1}{d_{i0}} \|\mathcal{A}_i^*(\mathbf{y}) - \mathbf{h}_{i0}\mathbf{x}_{i0}^*\| + \left| \frac{d_i}{d_{i0}} - 1 \right| \leq 2\xi \end{aligned}$$

where $(\mathbf{I}_K - \frac{\mathbf{h}_{i0}\mathbf{h}_{i0}^*}{d_{i0}})\mathbf{h}_{i0}\mathbf{x}_{i0}^* = \mathbf{0}$ and $(\mathcal{A}_i^*(\mathbf{y}) - d_i \hat{\mathbf{h}}_{i0} \hat{\mathbf{x}}_{i0}^*)\hat{\mathbf{x}}_{i0} \hat{\mathbf{h}}_{i0}^* = \mathbf{0}$. Therefore, we have

$$\left\| \hat{\mathbf{h}}_{i0} - \frac{\mathbf{h}_{i0}^* \hat{\mathbf{h}}_{i0}}{d_{i0}} \mathbf{h}_{i0} \right\| \leq 2\xi, \quad \|\sqrt{d_i} \hat{\mathbf{h}}_{i0} - t_{i0} \mathbf{h}_{i0}\| \leq 2\sqrt{d_i} \xi, \quad (6.48)$$

where $t_{i0} = \frac{\sqrt{d_i} \mathbf{h}_{i0}^* \hat{\mathbf{h}}_{i0}}{d_{i0}}$ and $|t_{i0}| \leq \sqrt{d_i/d_{i0}} < \sqrt{2}$. If we substitute \mathbf{w} by $t_{i0}\mathbf{h}_{i0} \in Q_i$ into (6.46),

$$\|\sqrt{d_i} \hat{\mathbf{h}}_{i0} - t_{i0} \mathbf{h}_{i0}\| \geq \|\mathbf{u}_i^{(0)} - t_{i0} \mathbf{h}_{i0}\|. \quad (6.49)$$

where $t_{i0}\mathbf{h}_{i0} \in Q_i$ follows from $\sqrt{L}|t_{i0}|\|\mathbf{B}\mathbf{h}_{i0}\|_\infty \leq |t_{i0}|\sqrt{d_{i0}}\mu_h \leq \sqrt{2d_{i0}}\mu$.

Now we are ready to estimate $\|\mathbf{u}_i^{(0)}(\mathbf{v}_i^{(0)})^* - \mathbf{h}_{i0}\mathbf{x}_{i0}^*\|_F$ as follows,

$$\begin{aligned} \|\mathbf{u}_i^{(0)}(\mathbf{v}_i^{(0)})^* - \mathbf{h}_{i0}\mathbf{x}_{i0}^*\|_F &\leq \|\mathbf{u}_i^{(0)}(\mathbf{v}_i^{(0)})^* - t_{i0}\mathbf{h}_{i0}(\mathbf{v}_i^{(0)})^*\|_F + \|t_{i0}\mathbf{h}_{i0}(\mathbf{v}_i^{(0)})^* - \mathbf{h}_{i0}\mathbf{x}_{i0}^*\|_F \\ &\leq \underbrace{\|\mathbf{u}_i^{(0)} - t_{i0}\mathbf{h}_{i0}\|}_{I_1} \|\mathbf{v}_i^{(0)}\| + \underbrace{\left\| \frac{d_i}{d_{i0}} \mathbf{h}_{i0} \mathbf{h}_{i0}^* \hat{\mathbf{h}}_{i0} \hat{\mathbf{x}}_{i0}^* - \mathbf{h}_{i0} \mathbf{x}_{i0}^* \right\|_F}_{I_2}. \end{aligned}$$

Here $I_1 \leq 2\xi d_i$ because $\|\mathbf{v}_i^{(0)}\| = \sqrt{d_i}$ and $\|\mathbf{u}_i^{(0)} - t_{i0}\mathbf{h}_{i0}\| \leq 2\sqrt{d_i}\xi$ follows from (6.48) and (6.49). For I_2 , there holds

$$I_2 = \left\| \frac{\mathbf{h}_{i0}\mathbf{h}_{i0}^*}{d_{i0}} \left(d_i \hat{\mathbf{h}}_{i0} \hat{\mathbf{x}}_{i0}^* - \mathbf{h}_{i0}\mathbf{x}_{i0}^* \right) \right\|_F \leq \|d_i \hat{\mathbf{h}}_{i0} \hat{\mathbf{x}}_{i0}^* - \mathbf{h}_{i0}\mathbf{x}_{i0}^*\|_F \leq 2\sqrt{2}\xi d_{i0},$$

which follows from (6.47). Therefore,

$$\|\mathbf{u}_i^{(0)}(\mathbf{v}_i^{(0)})^* - \mathbf{h}_{i0}\mathbf{x}_{i0}^*\|_F \leq 2\xi d_i + 2\sqrt{2}\xi d_{i0} \leq 5\xi d_{i0} \leq \frac{2\varepsilon d_{i0}}{5\sqrt{s\kappa}}.$$

□

Appendix

Descent Lemma

Lemma 6.19 (Lemma 6.1 in [18]). *If $f(\mathbf{z}, \bar{\mathbf{z}})$ is a continuously differentiable real-valued function with two complex variables \mathbf{z} and $\bar{\mathbf{z}}$, (for simplicity, we just denote $f(\mathbf{z}, \bar{\mathbf{z}})$ by $f(\mathbf{z})$ and keep in the mind that $f(\mathbf{z})$ only assumes real values) for $\mathbf{z} := (\mathbf{h}, \mathbf{x}) \in \mathcal{N}_\varepsilon \cap \mathcal{N}_{\bar{\mathbf{F}}}$. Suppose that there exists a constant C_L such that*

$$\|\nabla f(\mathbf{z} + t\Delta\mathbf{z}) - \nabla f(\mathbf{z})\| \leq C_L t \|\Delta\mathbf{z}\|, \quad \forall 0 \leq t \leq 1,$$

for all $\mathbf{z} \in \mathcal{N}_\varepsilon \cap \mathcal{N}_{\bar{\mathbf{F}}}$ and $\Delta\mathbf{z}$ such that $\mathbf{z} + t\Delta\mathbf{z} \in \mathcal{N}_\varepsilon \cap \mathcal{N}_{\bar{\mathbf{F}}}$ and $0 \leq t \leq 1$. Then

$$f(\mathbf{z} + \Delta\mathbf{z}) \leq f(\mathbf{z}) + 2 \operatorname{Re}((\Delta\mathbf{z})^T \bar{\nabla} f(\mathbf{z})) + C_L \|\Delta\mathbf{z}\|^2$$

where $\bar{\nabla} f(\mathbf{z}) := \frac{\partial f(\mathbf{z}, \bar{\mathbf{z}})}{\partial \bar{\mathbf{z}}}$ is the complex conjugate of $\nabla f(\mathbf{z}) = \frac{\partial f(\mathbf{z}, \bar{\mathbf{z}})}{\partial \mathbf{z}}$.

Concentration inequality

We define the matrix ψ_1 -norm via

$$\|\mathbf{Z}\|_{\psi_1} := \inf_{u \geq 0} \{\mathbb{E}[\exp(\|\mathbf{Z}\|/u)] \leq 2\}.$$

Theorem 6.20. [16] *Consider a finite sequence of \mathcal{Z}_l of independent centered random matrices with dimension $M_1 \times M_2$. Assume that $R := \max_{1 \leq l \leq L} \|\mathcal{Z}_l\|_{\psi_1}$ and introduce the random matrix*

$$\mathcal{S} = \sum_{l=1}^L \mathcal{Z}_l. \quad (6.50)$$

Compute the variance parameter

$$\sigma_0^2 = \max \left\{ \left\| \sum_{l=1}^L \mathbb{E}(\mathcal{Z}_l \mathcal{Z}_l^*) \right\|, \left\| \sum_{l=1}^L \mathbb{E}(\mathcal{Z}_l^* \mathcal{Z}_l) \right\| \right\}, \quad (6.51)$$

then for all $t \geq 0$

$$\|\mathcal{S}\| \leq C_0 \max \left\{ \sigma_0 \sqrt{t + \log(M_1 + M_2)}, R \log \left(\frac{\sqrt{LR}}{\sigma_0} \right) (t + \log(M_1 + M_2)) \right\} \quad (6.52)$$

with probability at least $1 - e^{-t}$ where C_0 is an absolute constant.

Lemma 6.21 (Lemma 10-13 in [1], Lemma 12.4 in [19]). *Let $\mathbf{u} \in \mathbb{C}^n \sim \mathcal{CN}(\mathbf{0}, \mathbf{I}_n)$, then $\|\mathbf{u}\|^2 \sim \frac{1}{2} \chi_{2n}^2$ and*

$$\|\|\mathbf{u}\|^2\|_{\psi_1} = \|\langle \mathbf{u}, \mathbf{u} \rangle\|_{\psi_1} \leq Cn \quad (6.53)$$

and

$$\mathbb{E}(\mathbf{u}\mathbf{u}^* - \mathbf{I}_n)^2 = n\mathbf{I}_n. \quad (6.54)$$

Let $\mathbf{q} \in \mathbb{C}^n$ be any deterministic vector, then the following properties hold

$$\|(\mathbf{u}\mathbf{u}^* - \mathbf{I})\mathbf{q}\|_{\psi_1} \leq C\sqrt{n}\|\mathbf{q}\|, \quad (6.55)$$

$$\mathbb{E}[(\mathbf{u}\mathbf{u}^* - \mathbf{I})\mathbf{q}\mathbf{q}^*(\mathbf{u}\mathbf{u}^* - \mathbf{I})] = \|\mathbf{q}\|^2 \mathbf{I}_n. \quad (6.56)$$

Let $\mathbf{v} \sim \mathcal{CN}(\mathbf{0}, \mathbf{I}_m)$ be a complex Gaussian random vector in \mathbb{C}^m , independent of \mathbf{u} , then

$$\|\|\mathbf{u}\| \cdot \|\mathbf{v}\|\|_{\psi_1} \leq C\sqrt{mn}. \quad (6.57)$$

Acknowledgement

S.Ling would like to thank Felix Krahmer and Dominik Stöger for the discussion about [25], and also thank Ju Sun for pointing out the connection between convolutional dictionary learning and this work.

References

- [1] A. Ahmed, B. Recht, and J. Romberg. Blind deconvolution using convex programming. *Information Theory, IEEE Transactions on*, 60(3):1711–1732, 2014.
- [2] H. Bristow, A. Eriksson, and S. Lucey. Fast convolutional sparse coding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 391–398, 2013.
- [3] T. T. Cai, X. Li, Z. Ma, et al. Optimal rates of convergence for noisy sparse phase retrieval via thresholded wirtinger flow. *The Annals of Statistics*, 44(5):2221–2251, 2016.

- [4] V. Cambbareri and L. Jacques. Through the haze: A non-convex approach to blind calibration for linear random sensing models. *arXiv preprint arXiv:1610.09028*, 2016.
- [5] P. Campisi and K. Egiazarian. *Blind image deconvolution: theory and applications*. CRC press, 2007.
- [6] E. Candès, Y. Eldar, T. Strohmer, and V. Voroninski. Phase retrieval via matrix completion. *SIAM J. Imag. Sci.*, 6(1):199–225, 2013.
- [7] E. Candès and B. Recht. Exact matrix completion via convex optimization. *Foundations of Computational Mathematics*, 9(6):717–772, 2009.
- [8] E. J. Candès, X. Li, and M. Soltanolkotabi. Phase retrieval via Wirtinger flow: Theory and algorithms. *Information Theory, IEEE Transactions on*, 61(4):1985–2007, 2015.
- [9] E. J. Candès, T. Strohmer, and V. Voroninski. Phaselift: Exact and stable signal recovery from magnitude measurements via convex programming. *Communications on Pure and Applied Mathematics*, 66(8):1241–1274, 2013.
- [10] V. Chandrasekaran, B. Recht, P. Parrilo, and A. Willsky. The convex geometry of linear inverse problems. *Foundations of Computational Mathematics*, 12(6):805–849, 2012.
- [11] Y. Chen and E. Candès. Solving random quadratic systems of equations is nearly as easy as solving linear systems. In *Advances in Neural Information Processing Systems*, pages 739–747, 2015.
- [12] R. Escalante and M. Raydan. *Alternating projection methods*, volume 8. SIAM, 2011.
- [13] A. Goldsmith. *Wireless communications*. Cambridge University Press, 2005.
- [14] R. Keshavan, A. Montanari, and S. Oh. Matrix completion from noisy entries. In *Advances in Neural Information Processing Systems*, pages 952–960, 2009.
- [15] R. H. Keshavan, A. Montanari, and S. Oh. Matrix completion from a few entries. *Information Theory, IEEE Transactions on*, 56(6):2980–2998, 2010.
- [16] V. Koltchinskii et al. Von Neumann entropy penalization and low-rank matrix estimation. *The Annals of Statistics*, 39(6):2936–2973, 2011.
- [17] K. Lee, Y. Li, M. Junge, and Y. Bresler. Blind recovery of sparse signals from subsampled convolution. *IEEE Transactions on Information Theory*, 2016.
- [18] X. Li, S. Ling, T. Strohmer, and K. Wei. Rapid, robust, and reliable blind deconvolution via nonconvex optimization. *arXiv preprint arXiv:1606.04933*, 2016.
- [19] S. Ling and T. Strohmer. Blind deconvolution meets blind demixing: Algorithms and performance bounds. *arXiv preprint arXiv:1512.07730*, 2015.
- [20] S. Ling and T. Strohmer. Self-calibration and biconvex compressive sensing. *Inverse Problems*, 31(11):115002, 2015.
- [21] J. Liu, J. Xin, Y. Qi, F.-G. Zheng, et al. A time domain algorithm for blind separation of convolutive sound mixtures and L_1 constrained minimization of cross correlations. *Communications in Mathematical Sciences*, 7(1):109–128, 2009.
- [22] D. G. Luenberger and Y. Ye. *Linear and nonlinear programming*, volume 228. Springer, 2015.
- [23] M. B. McCoy and J. A. Tropp. Achievable performance of convex demixing. Technical report, Caltech, 2017, Paper dated Feb. 2013. ACM Technical Report 2017-02.

- [24] G. W. Stewart. Perturbation theory for the singular value decomposition. *Technical report CS-TR 2539, University of Maryland*, September 1990.
- [25] D. Stöger, P. Jung, and F. Kraemer. Blind deconvolution and compressed sensing. In *Compressed Sensing Theory and its Applications to Radar, Sonar and Remote Sensing (CoSeRa), 2016 4th International Workshop on*, pages 24–27. IEEE, 2016.
- [26] T. Strohmer. Four short stories about Toeplitz matrix calculations. *Linear Algebra Appl.*, 343/344:321–344, 2002. Special issue on structured and infinite systems of linear equations.
- [27] J. Sun, Q. Qu, and J. Wright. Complete dictionary recovery over the sphere i: Overview and the geometric picture. *IEEE Transactions on Information Theory*, 2016.
- [28] J. Sun, Q. Qu, and J. Wright. A geometric analysis of phase retrieval. *arXiv preprint arXiv:1602.06664*, 2016.
- [29] R. Sun and Z.-Q. Luo. Guaranteed matrix completion via non-convex factorization. *IEEE Transactions on Information Theory*, 62(11):6535–6579, 2016.
- [30] S. Tu, R. Boczar, M. Simchowitz, M. Soltanolkotabi, and B. Recht. Low-rank solutions of linear matrix equations via procrustes flow. In *Proceedings of The 33rd International Conference on Machine Learning*, pages 964–973, 2016.
- [31] R. Vershynin. Introduction to the non-asymptotic analysis of random matrices. In Y. C. Eldar and G. Kutyniok, editors, *Compressed Sensing: Theory and Applications*, chapter 5. Cambridge University Press, 2012.
- [32] X. Wang and H. V. Poor. Blind equalization and multiuser detection in dispersive CDMA channels. *Communications, IEEE Transactions on*, 46(1):91–103, 1998.
- [33] P.-Å. Wedin. Perturbation bounds in connection with singular value decomposition. *BIT Numerical Mathematics*, 12(1):99–111, 1972.
- [34] K. Wei, J.-F. Cai, T. F. Chan, and S. Leung. Guarantees of riemannian optimization for low rank matrix recovery. *SIAM Journal on Matrix Analysis and Applications*, 37(3):1198–1222, 2016.
- [35] J. Wright, A. Ganesh, K. Min, and Y. Ma. Compressive principal component pursuit. *Information and Inference*, 2(1):32–68, 2013.
- [36] G. Wunder, H. Boche, T. Strohmer, and P. Jung. Sparse signal processing concepts for efficient 5G system design. *IEEE Access*, 3:195–208, 2015.