
Blind Demixing via Wirtinger Flow with Random Initialization

Jialin Dong
dongjl@shanghaitech.edu.cn
ShanghaiTech University

Yuanming Shi
shiyym@shanghaitech.edu.cn
ShanghaiTech University

Abstract

This paper concerns the problem of demixing a series of source signals from the sum of bilinear measurements. This problem spans diverse areas such as communication, imaging processing, machine learning, etc. However, semidefinite programming for blind demixing is prohibitive to large-scale problems due to high computational complexity and storage cost. Although several efficient algorithms have been developed recently that enjoy the benefits of fast convergence rates and even regularization free, they still call for spectral initialization. To find simple initialization approach that works equally well as spectral initialization, we propose to solve blind demixing problem via Wirtinger flow with *random initialization*, which yields a natural implementation. To reveal the efficiency of this algorithm, we provide the global convergence guarantee concerning randomly initialized Wirtinger flow for blind demixing. Specifically, it shows that with sufficient samples, the iterates of randomly initialized Wirtinger flow can enter a local region that enjoys strong convexity and strong smoothness within a few iterations at the first stage. At the second stage, iterates of randomly initialized Wirtinger flow further converge linearly to the ground truth.

1 INTRODUCTION

Suppose we are given an observation vector $\mathbf{y} \in \mathbb{C}^m$ in frequency domain generated from the sum of bilinear measurements of unknown vectors $\mathbf{x}_i^{\natural} \in \mathbb{C}^N$, $\mathbf{h}_i^{\natural} \in \mathbb{C}^K$,

$i = 1, \dots, s$, i.e.,

$$y_j = \sum_{i=1}^s \mathbf{b}_j^* \mathbf{h}_i^{\natural} \mathbf{x}_i^{\natural} \mathbf{a}_{ij}, \quad 1 \leq j \leq m, \quad (1)$$

where $\{\mathbf{a}_{ij}\} \in \mathbb{C}^N$, $\{\mathbf{b}_j\} \in \mathbb{C}^K$ are design vectors. Here, the first K columns of the matrix \mathbf{F} form the matrix $\mathbf{B} := [\mathbf{b}_1, \dots, \mathbf{b}_m]^* \in \mathbb{C}^{m \times K}$, where $\mathbf{F} \in \mathbb{C}^{m \times m}$ is the unitary discrete Fourier transform (DFT) matrix with $\mathbf{F}\mathbf{F}^* = \mathbf{I}_m$. Our goal is to recover $\{\mathbf{x}_i^{\natural}\}$ and $\{\mathbf{h}_i^{\natural}\}$ from the sum of bilinear measurements, which is known as *blind demixing* [1, 2].

This problem has spanned a wide scope of applications ranging from imaging processing [3, 4] and machine learning [5, 6] to communication [7, 8]. Specifically, by solving the blind demixing problem, both original images and corresponding convolutional kernels can be recovered from a blurred image [3]. This problem has also been exploited to demix and deconvolve calcium imaging recordings of neuronal ensembles [4]. Blind demixing also finds applications in machine learning [5, 6] where feature maps convolve with filters individually and added overall users to approximate the input vector. In the context of communication, this problem ensures to recover the source signal from distinguishing users without knowing channel vectors, thereby enabling low-latency communication [7].

Although blind demixing plays vital role in various areas, solving it is generally highly intractable. Moreover, some applications of blind demixing problem involves large-scale data, hence efficient algorithms with optimal guarantees are needed urgently to recover signal vectors from a single observation vector. Instead of computationally expensive convex method [1], a non-convex algorithm, regularized gradient descent [9], has been recently proposed to efficiently solve blind demixing problem. This method, however, requires extra regularization and still yields conservative computational optimality guarantees. To elude the regularization and develop an algorithm with progressive computational guarantees, the least square estimation ap-

Proceedings of the 22nd International Conference on Artificial Intelligence and Statistics (AISTATS) 2019, Naha, Okinawa, Japan. PMLR: Volume 89. Copyright 2019 by the author(s).

proach was proposed in [2]:

$$\mathcal{P} : \underset{\{\mathbf{h}_i\}, \{\mathbf{x}_i\}}{\text{minimize}} \quad f(\mathbf{h}, \mathbf{x}) := \sum_{j=1}^m \left| \sum_{i=1}^s \mathbf{b}_j^* \mathbf{h}_i \mathbf{x}_i^* \mathbf{a}_{ij} - y_j \right|^2, \quad (2)$$

which is solved via Wirtinger flow with spectral initialization by harnessing the benefits of low computational complexity, regularization-free, fast convergence rate with aggressive step size and computational optimality guarantees [2]. Nevertheless, this algorithm with theoretical guarantees depends on the spectral initialization by computing the largest singular vector of the data matrix [2]. To find a natural implementation for the practitioners, we propose to solve the blind demixing problem via randomly initialized Wirtinger flow. Our goal, in this paper, is to confirm the efficiency of this algorithm with theoretical guarantees. Specifically, we shall verify that the iterates of randomly initialized Wirtinger flow are able to achieve a local region that enjoys strong convexity and strong smoothness within a few iteration at the first stage. At the next stage, we further show that the iterates of the randomly initialized Wirtinger flow linearly converges to the ground truth.

1.1 State-of-the-Art Algorithms

Despite the general intractability of blind demixing, it can be effectively solved by several algorithms under the proper statistical models. In particular, semidefinite programming was developed in [1] to solve the blind demixing problem by lifting the bilinear model into the matrix space. However, it is computationally prohibitive for solving large-scale problem due to the high computation and storage cost. To address this issue, the nonconvex algorithm, e.g., regularized gradient descent with spectral initialization [9], was further developed to optimize the variables in the natural space. Nevertheless, the theoretical guarantees for the regularized gradient [9] provide pessimistic convergence rate and require carefully-designed initialization. The Riemannian trust-region optimization algorithm without regularization was further proposed in [7] to improve the convergence rate. However, the second-order algorithm brings unique challenges in providing statistical guarantees. Recently, theoretical guarantees concerning regularization-freed Wirtinger flow with spectral initialization for blind demixing was provided in [2]. However, this regularization-freed method still calls for spectral initialization. In this paper, we aim to explore the random initialization strategy for natural implementation and theoretically verify the efficiency of the randomly initialized Wirtinger flow.

Based on the random initialization strategy, a line of research studies the benign global landscapes and aims

to design generic saddle-point escaping algorithms, e.g., noisy stochastic gradient descent [10], trust-region method [11], perturbed gradient descent [12]. With sufficient sample size, these algorithms are guaranteed to converge globally for phase retrieval [11], matrix recovery [13], matrix sensing [13], robust PCA [14] and shallow neural networks [15] where all local minima are as good as global and all the saddle points are strict. However, the theoretical results developed in [10, 11, 12, 13, 14, 15] are fairly general and may yield pessimistic convergence rate guarantees. Moreover, these saddle-point escaping algorithms are more complicated for implementation than the natural vanilla gradient descent or Wirtinger flow. To advance the theoretical analysis for gradient descent with random initialization, the fast global convergence guarantee concerning randomly initialized gradient descent for phase retrieval has been recently provided in [16].

1.2 Wirtinger Flow with Random Initialization

In this paper, we solve the blind demixing problem via randomly initialized Wirtinger flow by harnessing the benefits of computational efficiency, regularization-free and careful initialization free. In this paper, our main contribution is to provide the global convergence guarantee for the randomly initialized Wirtinger flow.

Wirtinger flow with random initialization is an iterative algorithm with vanilla gradient descent update procedure, i.e., without regularization. Specifically, the gradient step of Wirtinger flow is represented by the notion of Wirtinger derivatives [17], i.e., the derivatives of real-valued functions over complex variables. To simplify the presentation, we denote $f(\mathbf{z}) := f(\mathbf{h}, \mathbf{x})$, where

$$\mathbf{z} = \begin{bmatrix} \mathbf{z}_1 \\ \vdots \\ \mathbf{z}_s \end{bmatrix} \in \mathbb{C}^{s(N+K)} \quad \text{with} \quad \mathbf{z}_i = \begin{bmatrix} \mathbf{h}_i \\ \mathbf{x}_i \end{bmatrix} \in \mathbb{C}^{N+K}. \quad (3)$$

Furthermore, we define the discrepancy between the estimate \mathbf{z} and the ground truth \mathbf{z}^\natural as the distance function, given as

$$\text{dist}(\mathbf{z}, \mathbf{z}^\natural) = \left(\sum_{i=1}^s \text{dist}^2(\mathbf{z}_i, \mathbf{z}_i^\natural) \right)^{1/2}, \quad (4)$$

where

$$\text{dist}^2(\mathbf{z}_i, \mathbf{z}_i^\natural) = \min_{\alpha_i \in \mathbb{C}} (\| \frac{1}{\alpha_i} \mathbf{h}_i - \mathbf{h}_i^\natural \|_2^2 + \| \alpha_i \mathbf{x}_i - \mathbf{x}_i^\natural \|_2^2) / d_i$$

for $i = 1, \dots, s$. Here, $d_i = \| \mathbf{h}_i^\natural \|_2^2 + \| \mathbf{x}_i^\natural \|_2^2$ and each α_i is the alignment parameter.

Let \mathbf{A}^* and \mathbf{a}^* denote the conjugate transpose of matrix \mathbf{A} vector \mathbf{a} respectively. For each $i = 1, \dots, s$,

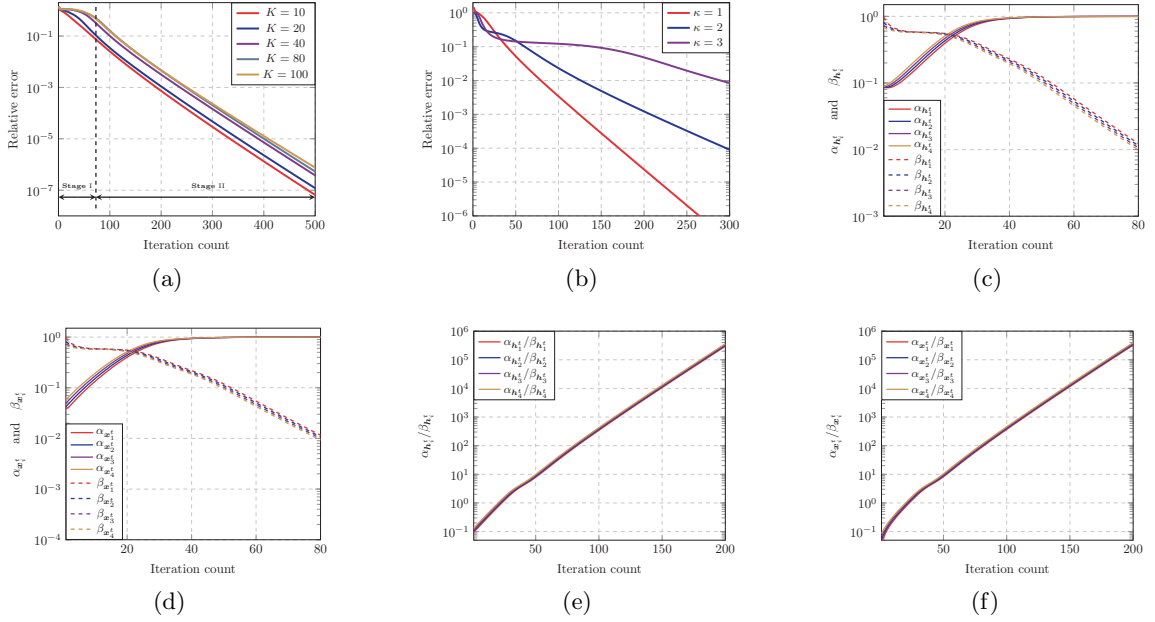


Figure 1: Numerical results.

$\nabla_{\mathbf{h}_i} f(\mathbf{z})$ and $\nabla_{\mathbf{x}_i} f(\mathbf{z})$ denote the Wirtinger gradient of $f(\mathbf{z})$ with respect to \mathbf{h}_i and \mathbf{x}_i respectively:

$$\nabla_{\mathbf{h}_i} f(\mathbf{z}) = \sum_{j=1}^m \left(\sum_{k=1}^s \mathbf{b}_j^* \mathbf{h}_k \mathbf{x}_k^* \mathbf{a}_{kj} - y_j \right) \mathbf{b}_j \mathbf{a}_{ij}^* \mathbf{x}_i, \quad (5a)$$

$$\nabla_{\mathbf{x}_i} f(\mathbf{z}) = \sum_{j=1}^m \left(\sum_{k=1}^s \mathbf{h}_k^* \mathbf{b}_j \mathbf{a}_{kj}^* \mathbf{x}_k - \bar{y}_j \right) \mathbf{a}_{ij} \mathbf{b}_j^* \mathbf{h}_i. \quad (5b)$$

In light of the Wirtinger gradient (5), the update rule of Wirtinger flow is given by

$$\begin{bmatrix} \mathbf{h}_i^{t+1} \\ \mathbf{x}_i^{t+1} \end{bmatrix} = \begin{bmatrix} \mathbf{h}_i^t \\ \mathbf{x}_i^t \end{bmatrix} - \eta \begin{bmatrix} \frac{1}{\|\mathbf{x}_i^t\|_2} \nabla_{\mathbf{h}_i} f(\mathbf{z}^t) \\ \frac{1}{\|\mathbf{h}_i^t\|_2} \nabla_{\mathbf{x}_i} f(\mathbf{z}^t) \end{bmatrix}, \quad i = 1, \dots, s, \quad (6)$$

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

1.3 Numerical Results

The power of randomly initialized WF for solving problem \mathcal{P} (2) will be illustrated by numerical example. The ground truth signals and initial points are randomly generated as

$$\mathbf{h}_i^{\natural} \sim \mathcal{N}(\mathbf{0}, K^{-1} \mathbf{I}_K), \quad \mathbf{x}_i^{\natural} \sim \mathcal{N}(\mathbf{0}, N^{-1} \mathbf{I}_N), \quad (7)$$

$$\mathbf{h}_i^0 \sim \mathcal{N}(\mathbf{0}, K^{-1} \mathbf{I}_K), \quad \mathbf{x}_i^0 \sim \mathcal{N}(\mathbf{0}, N^{-1} \mathbf{I}_N), \quad (8)$$

for $i = 1, \dots, s$. In all simulations, we set $K = N$ and normalize $\|\mathbf{h}_i^{\natural}\|_2 = \|\mathbf{x}_i^{\natural}\|_2 = 1$ for $i = 1, \dots, s$. Specifically, for each $K \in \{10, 20, 40, 80, 100\}$, $s = 10$ and $m = 50K$, the design vectors \mathbf{a}_{ij} 's follow $\mathbf{a}_{ij} \sim$

$\mathcal{N}(\mathbf{0}, \frac{1}{2} \mathbf{I}_N) + i\mathcal{N}(\mathbf{0}, \frac{1}{2} \mathbf{I}_N)$. With the chosen step size $\eta = 0.1$ in all settings, Fig. 1(a) shows the relative error, i.e., $\sum_{i=1}^s \|\mathbf{h}_i^t \mathbf{x}_i^{t*} - \mathbf{h}_i^{\natural} \mathbf{x}_i^{\natural*}\|_F / \sum_{i=1}^s \|\mathbf{h}_i^{\natural} \mathbf{x}_i^{\natural*}\|_F$, versus the iteration count, where $\|\mathbf{X}\|_F$ denotes the Frobenius norm of the matrix \mathbf{X} . We observe that, the iterations of randomly initialized Wirtinger flow can be separately into two stages: Stage I: within first tens of iterations, the relative error maintains nearly flat, Stage II: the relative error enjoys linear convergence rate which rarely changes as the problem size varies.

We further illustrate that the performance and convergence rate of the Wirtinger flow with random initialization depend on the condition number, i.e., $\kappa := \frac{\max_i \|\mathbf{x}_i^{\natural}\|_2}{\min_i \|\mathbf{x}_i^{\natural}\|_2}$. In this experiment, let $K = 50$, $m = 800$, $s = 2$, the step size be $\eta = 0.2$. We then set for the first component $\|\mathbf{h}_1^{\natural}\|_2 = \|\mathbf{x}_1^{\natural}\|_2 = 1$ and for the second one $\|\mathbf{h}_2^{\natural}\|_2 = \|\mathbf{x}_2^{\natural}\|_2 = \kappa$ with $\kappa \in \{1, 2, 3\}$. The random initialization (8) is utilized. Fig. 1(b) shows the relative error versus the iteration count. As we can see, a larger κ yields slower convergence rate.

2 MAIN RESULTS

In this section, we shall verify the preceding numerical results with theoretical guarantees. Through theoretical analysis, we assume the design vectors \mathbf{a}_{ij} 's are $\mathbf{a}_{ij} \sim \mathcal{N}(\mathbf{0}, \frac{1}{2} \mathbf{I}_N) + i\mathcal{N}(\mathbf{0}, \frac{1}{2} \mathbf{I}_N)$. To present the main theorem, we first introduce several fundamental notions and definitions.

Throughout this paper, $f(n) = O(g(n))$ or $f(n) \lesssim g(n)$ denotes that there exists a constant $c > 0$ such that $|f(n)| \leq c|g(n)|$, whereas $f(n) \gtrsim g(n)$ means that there exists a constant $c > 0$ such that $|f(n)| \geq c|g(n)|$. $f(n) \gg g(n)$ denotes that there exists some sufficiently large constant $c > 0$ such that $|f(n)| \geq c|g(n)|$. In addition, the notation $f(n) \asymp g(n)$ means that there exists constants $c_1, c_2 > 0$ such that $c_1|g(n)| \leq |f(n)| \leq c_2|g(n)|$.

Furthermore, the incoherence parameter [9], which characterizes the incoherence between \mathbf{b}_j and \mathbf{h}_i for $1 \leq i \leq s, 1 \leq j \leq m$.

Definition 1 (Incoherence for blind demixing). *Let the incoherence parameter μ be the smallest number such that*

$$\max_{1 \leq i \leq s, 1 \leq j \leq m} \frac{|\mathbf{b}_j^* \mathbf{h}_i^{\natural}|}{\|\mathbf{h}_i^{\natural}\|_2} \leq \frac{\mu}{\sqrt{m}}. \quad (9)$$

Let $\tilde{\mathbf{h}}_i^t$ and $\tilde{\mathbf{x}}_i^t$ be

$$\tilde{\mathbf{h}}_i^t = \frac{1}{\omega_i^t} \mathbf{h}_i^t \quad \text{and} \quad \tilde{\mathbf{x}}_i^t = \omega_i^t \mathbf{x}_i^t, \quad (10)$$

for $i = 1, \dots, s$, respectively, where ω_i 's are alignment parameters. We further define the norm of signal component and the perpendicular component with respect to \mathbf{h}_i for $i = 1, \dots, s$, as

$$\alpha_{\mathbf{h}_i} := \langle \mathbf{h}_i^{\natural}, \tilde{\mathbf{h}}_i^t \rangle / \|\mathbf{h}_i^{\natural}\|_2, \quad (11)$$

$$\beta_{\mathbf{h}_i} := \left\| \tilde{\mathbf{h}}_i^t - \frac{\langle \mathbf{h}_i^{\natural}, \tilde{\mathbf{h}}_i^t \rangle}{\|\mathbf{h}_i^{\natural}\|_2^2} \mathbf{h}_i^{\natural} \right\|_2, \quad (12)$$

respectively. Here, ω_i 's are the alignment parameters. Similarly, the norms of the signal component and the perpendicular component with respect to \mathbf{x}_i for $i = 1, \dots, s$, can be represented as

$$\alpha_{\mathbf{x}_i} := \langle \mathbf{x}_i^{\natural}, \tilde{\mathbf{x}}_i^t \rangle / \|\mathbf{x}_i^{\natural}\|_2, \quad (13)$$

$$\beta_{\mathbf{x}_i} := \left\| \tilde{\mathbf{x}}_i^t - \frac{\langle \mathbf{x}_i^{\natural}, \tilde{\mathbf{x}}_i^t \rangle}{\|\mathbf{x}_i^{\natural}\|_2^2} \mathbf{x}_i^{\natural} \right\|_2, \quad (14)$$

respectively.

Theorem 1. *Assume that the initial points obey (8) for $i = 1, \dots, s$ and the stepsize $\eta > 0$ satisfies $\eta \asymp s^{-1}$. Suppose that the sample size satisfies $m \geq C\mu^2 s^2 \kappa^4 \max\{K, N\} \log^{12} m$ for some sufficiently large constant $C > 0$. Then with probability at least $1 - c_1 m^{-\nu} - c_2 m^{-c_2 N}$ for some constants $\nu, c_1, c_2 > 0$, there exists a sufficiently small constant $0 \leq \gamma \leq 1$ and $T_\gamma \lesssim s \log(\max\{K, N\})$ such that*

1. *The randomly initialized WF linearly converges to \mathbf{z}^{\natural} , i.e.,*

$$\text{dist}(\mathbf{z}^t, \mathbf{z}^{\natural}) \leq \gamma \left(1 - \frac{\eta}{16\kappa}\right)^{t-T_\gamma} \|\mathbf{z}^{\natural}\|_2, \quad t \geq T_\gamma,$$

2. *The magnitude ratios of the signal component to the perpendicular component with respect to \mathbf{h}_i^t and \mathbf{x}_i^t obey*

$$\max_{1 \leq i \leq s} \frac{\alpha_{\mathbf{h}_i^t}}{\beta_{\mathbf{h}_i^t}} \gtrsim \frac{1}{\sqrt{K \log K}} (1 + c_3 \eta)^t, \quad (15a)$$

$$\max_{1 \leq i \leq s} \frac{\alpha_{\mathbf{x}_i^t}}{\beta_{\mathbf{x}_i^t}} \gtrsim \frac{1}{\sqrt{N \log N}} (1 + c_4 \eta)^t, \quad (15b)$$

respectively, where $t = 0, 1, \dots$ for some constant $c_3, c_4 > 0$.

Theorem 1 provides precise statistical analysis on the computational efficiency of WF with random initialization. Specifically, in Stage I, it takes $T_\gamma = \mathcal{O}(s \log(\max\{K, N\}))$ iterations for randomly initialized WF to a local region near the ground truth that enjoys strong convexity and strong smoothness. The short duration of Stage I is own to the exponential growth of the magnitude ratio of the signal component to the perpendicular components (15). Moreover, in Stage II, it takes $\mathcal{O}(s \log(1/\varepsilon))$ iterations to reach ε -accurate solution at a linear convergence rate. Thus, the randomly initialized WF is guaranteed to converge to the ground truth with the iteration complexity $\mathcal{O}(s \log(\max\{K, N\}) + s \log(1/\varepsilon))$ given the sample size $m \gtrsim s^2 \max\{K, N\} \text{poly log}(m)$.

To further illustrate the relationship between the signal component $\alpha_{\mathbf{h}_i}$ (resp. $\alpha_{\mathbf{x}_i}$) and the perpendicular component $\beta_{\mathbf{h}_i}$ (resp. $\beta_{\mathbf{x}_i}$) for $i = 1, \dots, s$, we provide the simulation results under the setting of $K = N = 10$, $m = 50K$, $s = 4$ and $\eta = 0.1$ with $\|\mathbf{h}_i^{\natural}\|_2 = \|\mathbf{x}_i^{\natural}\|_2 = 1$ for $1 \leq i \leq s$. In particular, $\alpha_{\mathbf{h}_i}, \beta_{\mathbf{h}_i}$ versus iteration count (resp. $\alpha_{\mathbf{h}_i}, \beta_{\mathbf{h}_i}$ versus iteration count) for $i = 1, \dots, s$ is demonstrated in Fig. 1(c) (resp. Fig. 1(d)). Consider Fig. 1(a), Fig. 1(c) and Fig. 1(d) collectively, it shows that despite the rare decline of the relative error during Stage I, the sizes of the signal components, i.e., $\alpha_{\mathbf{h}_i}$ and $\alpha_{\mathbf{x}_i}$ for each $i = 1, \dots, s$, exponentially increase and the signal component becomes dominant component at the end of Stage I. Furthermore, the exponential growth of the ratio $\alpha_{\mathbf{h}_i}/\beta_{\mathbf{h}_i}$ (resp. $\alpha_{\mathbf{x}_i}/\beta_{\mathbf{x}_i}$) for each $i = 1, \dots, s$ is illustrated in Fig. 1(e) (resp. Fig. 1(f)).

3 DYNAMICS ANALYSIS

In this section, we shall briefly summarize the proof of the main theorem which is based on investigating the

dynamics of the iterates of WF with random initialization. The steps of proving Theorem 1 are summarized as follows.

1. Stage I:

- Dynamics of population-level state evolution.** Provide the population-level state evolution of $\alpha_{\mathbf{x}_i}$ (20a) and $\beta_{\mathbf{x}_i}$ (20b), $\alpha_{\mathbf{h}_i}$ (21a), $\beta_{\mathbf{h}_i}$ (21b) respectively, where the sample size approaches infinity. We then develop the approximate state evolution (23), which are remarkably close to the population-level state evolution, in the finite-sample regime. See details in Section 3.1.
- Dynamics of approximate state evolution.** Show that there exists some $T_\gamma = \mathcal{O}(s \log(\max\{K, N\}))$ such that $\text{dist}(\mathbf{z}^{T_\gamma}, \mathbf{z}^{\natural}) \leq \gamma$, if $\alpha_{\mathbf{h}_i}$ (11), $\beta_{\mathbf{h}_i}$ (12), $\alpha_{\mathbf{x}_i}$ (13) and $\beta_{\mathbf{x}_i}$ (14) satisfy the approximate state evolution (23). The exponential growth of the ratio $\alpha_{\mathbf{h}_i}/\beta_{\mathbf{h}_i}$ and $\alpha_{\mathbf{x}_i}/\beta_{\mathbf{x}_i}$ are further demonstrated under the same assumption. Please refer to Section 3.2.
- Leave-one-out arguments.** Prove that with high probability $\alpha_{\mathbf{h}_i}$, $\beta_{\mathbf{h}_i}$, $\alpha_{\mathbf{x}_i}$ and $\beta_{\mathbf{x}_i}$ satisfy the approximate state evolution (23) if the iterates $\{\mathbf{z}_i\}$ are independent with $\{\mathbf{a}_{ij}\}$. See details in Section A in the supplemental material. To achieve this, the ‘‘near-independence’’ between $\{\mathbf{z}_i\}$ and $\{\mathbf{a}_{ij}\}$ is established via exploiting leave-one-out arguments and some variants of the arguments. Specifically, the leave-one-out sequences and random-sign sequences are constructed in Section 3.3. The concentrations between the original and these auxiliary sequences are provided in the supplemental material.

2. Stage II: Local geometry in the region of incoherence and contraction. We invoke the prior theory provided in [2] to show local convergence of the random initialized WF in Stage II. Claim 15 in Stage II are further proven. Please refer to Section A.3 in the supplemental material.

3.1 Dynamics of Population-level State Evolution

In this subsection, we investigate the dynamics of population-level (where we have infinite samples) state evolution of $\alpha_{\mathbf{h}_i}$ (11), $\beta_{\mathbf{h}_i}$ (12), $\alpha_{\mathbf{x}_i}$ (13) and $\beta_{\mathbf{x}_i}$ (14).

Without loss the generality, we assume that $\mathbf{x}_i^{\natural} = q_i \mathbf{e}_1$ for $i = 1, \dots, s$, where $0 < q_i \leq 1, i = 1, \dots, s$ are some constants and $\kappa = \frac{\max_i q_i}{\min_i q_i}$, and \mathbf{e}_1 denotes the first standard basis vector. This assumption is

based on the rotational invariance of Gaussian distributions. Since the deterministic nature of $\{\mathbf{b}_j\}$, the ground truth signals $\{\mathbf{h}_i^{\natural}\}$ (channel vectors) cannot be transferred to a simple form, which yields more tedious analysis procedure. For simplification, for $i = 1, \dots, s$, we denote

$$\mathbf{x}_{i1}^t \quad \text{and} \quad \mathbf{x}_{i\perp}^t := [x_{ij}^t]_{2 \leq j \leq N}, \quad (16)$$

as the first entry and the second through the N^{th} entries of \mathbf{x}_i^t , respectively. Based on the assumption that $\mathbf{x}_i^{\natural} = q_i \mathbf{e}_1$ for $i = 1, \dots, s$, (13) and (14) can be reformulated as

$$\alpha_{\mathbf{x}_i} := \tilde{x}_{i1}^t \quad \text{and} \quad \beta_{\mathbf{x}_i} := \|\tilde{\mathbf{x}}_{i\perp}^t\|_2. \quad (17)$$

To study the population-level state evolution, we start by considering the case where the sequences $\{\mathbf{z}_i^t\}$ (refer to (3)) are established via the population gradient, i.e., for $i = 1, \dots, s$,

$$\begin{bmatrix} \mathbf{h}_i^{t+1} \\ \mathbf{x}_i^{t+1} \end{bmatrix} = \begin{bmatrix} \mathbf{h}_i^t \\ \mathbf{x}_i^t \end{bmatrix} - \eta \begin{bmatrix} \frac{1}{\|\mathbf{x}_i^t\|_2^2} \nabla_{\mathbf{h}_i} F(\mathbf{z}^t) \\ \frac{1}{\|\mathbf{h}_i^t\|_2^2} \nabla_{\mathbf{x}_i} F(\mathbf{z}^t) \end{bmatrix}, \quad (18)$$

where

$$\nabla_{\mathbf{h}_i} F(\mathbf{z}) := \mathbb{E}[\nabla_{\mathbf{h}_i} f(\mathbf{h}, \mathbf{x})] = \|\mathbf{x}_i\|_2^2 \mathbf{h}_i - (\mathbf{x}_i^{\natural*} \mathbf{x}_i) \mathbf{h}_i^{\natural},$$

$$\nabla_{\mathbf{x}_i} F(\mathbf{z}) := \mathbb{E}[\nabla_{\mathbf{x}_i} f(\mathbf{h}, \mathbf{x})] = \|\mathbf{h}_i\|_2^2 \mathbf{x}_i - (\mathbf{h}_i^{\natural*} \mathbf{h}_i) \mathbf{x}_i^{\natural}.$$

Here, the population gradients are computed based on the assumption that $\{\mathbf{x}_i\}$ (resp. $\{\mathbf{h}_i\}$) and $\{\mathbf{a}_{ij}\}$ (resp. $\{\mathbf{b}_j\}$) are independent with each other. With simple calculations, the dynamics for both the signal and the perpendicular components with respect to \mathbf{x}_i^t , $i = 1, \dots, s$ are given as:

$$\tilde{x}_{i1}^{t+1} = (1 - \eta) \tilde{x}_{i1}^t + \eta \frac{q_i^2}{\|\tilde{\mathbf{h}}_i^t\|_2^2} \mathbf{h}_i^{\natural*} \tilde{\mathbf{h}}_i^t, \quad (19a)$$

$$\tilde{\mathbf{x}}_{i\perp}^{t+1} = (1 - \eta) \tilde{\mathbf{x}}_{i\perp}^t. \quad (19b)$$

Assuming that $\eta > 0$ is sufficiently small and $\|\mathbf{x}_i^{\natural}\|_2 = \|\mathbf{x}_i^{\natural}\|_2 = q_i$ ($0 < q_i \leq 1$) for $i = 1, \dots, s$ and recognizing that $\|\tilde{\mathbf{h}}_i^t\|_2^2 = \alpha_{\mathbf{h}_i^t}^2 + \beta_{\mathbf{h}_i^t}^2$, we arrive at the following population-level state evolution for both $\alpha_{\mathbf{x}_i^t}$ and $\beta_{\mathbf{x}_i^t}$:

$$\alpha_{\mathbf{x}_i^{t+1}} = (1 - \eta) \alpha_{\mathbf{x}_i^t} + \eta \frac{q_i \alpha_{\mathbf{h}_i^t}}{\alpha_{\mathbf{h}_i^t}^2 + \beta_{\mathbf{h}_i^t}^2}, \quad (20a)$$

$$\beta_{\mathbf{x}_i^{t+1}} = (1 - \eta) \beta_{\mathbf{x}_i^t}. \quad (20b)$$

Likewise, the population-level state evolution for both $\alpha_{\mathbf{h}_i^t}$ and $\beta_{\mathbf{h}_i^t}$:

$$\alpha_{\mathbf{h}_i^{t+1}} = (1 - \eta) \alpha_{\mathbf{h}_i^t} + \eta \frac{q_i \alpha_{\mathbf{x}_i^t}}{\alpha_{\mathbf{x}_i^t}^2 + \beta_{\mathbf{x}_i^t}^2}, \quad (21a)$$

$$\beta_{\mathbf{h}_i^{t+1}} = (1 - \eta) \beta_{\mathbf{h}_i^t}. \quad (21b)$$

In finite-sample case, the dynamics of the randomly initialized WF iterates can be represented as

$$\mathbf{z}_i^{t+1} = \begin{bmatrix} \mathbf{h}_i^{t+1} \\ \mathbf{x}_i^{t+1} \end{bmatrix} = \begin{bmatrix} \mathbf{h}_i^t - \eta / \|\mathbf{x}_i^t\|_2^2 \cdot \nabla_{\mathbf{h}_i} F(\mathbf{z}) \\ \mathbf{x}_i^t - \eta / \|\mathbf{x}_i^t\|_2^2 \cdot \nabla_{\mathbf{x}_i} F(\mathbf{z}) \end{bmatrix} - \begin{bmatrix} \eta / \|\mathbf{x}_i^t\|_2^2 \cdot (\nabla_{\mathbf{h}_i} f(\mathbf{z}) - \nabla_{\mathbf{h}_i} F(\mathbf{z})) \\ \eta / \|\mathbf{h}_i^t\|_2^2 \cdot (\nabla_{\mathbf{x}_i} f(\mathbf{z}) - \nabla_{\mathbf{x}_i} F(\mathbf{z})) \end{bmatrix}. \quad (22)$$

Under the assumption that the last term in (22) is well-controlled, which will be justified in Section D in the supplemental material, we arrive at the approximate state evolution:

$$\alpha_{\mathbf{h}_i^{t+1}} = \left(1 - \eta + \frac{\eta q_i \psi_{\mathbf{h}_i^t}}{\alpha_{\mathbf{x}_i^t}^2 + \beta_{\mathbf{x}_i^t}^2}\right) \alpha_{\mathbf{h}_i^t} + \eta(1 - \rho_{\mathbf{h}_i^t}) \frac{q_i \alpha_{\mathbf{x}_i^t}}{\alpha_{\mathbf{x}_i^t}^2 + \beta_{\mathbf{x}_i^t}^2}, \quad (23a)$$

$$\beta_{\mathbf{h}_i^{t+1}} = \left(1 - \eta + \frac{\eta q_i \varphi_{\mathbf{h}_i^t}}{\alpha_{\mathbf{x}_i^t}^2 + \beta_{\mathbf{x}_i^t}^2}\right) \beta_{\mathbf{h}_i^t}, \quad (23b)$$

$$\alpha_{\mathbf{x}_i^{t+1}} = \left(1 - \eta + \frac{\eta q_i \psi_{\mathbf{x}_i^t}}{\alpha_{\mathbf{h}_i^t}^2 + \beta_{\mathbf{h}_i^t}^2}\right) \alpha_{\mathbf{x}_i^t} + \eta(1 - \rho_{\mathbf{x}_i^t}) \frac{q_i \alpha_{\mathbf{h}_i^t}}{\alpha_{\mathbf{h}_i^t}^2 + \beta_{\mathbf{h}_i^t}^2}, \quad (23c)$$

$$\beta_{\mathbf{x}_i^{t+1}} = \left(1 - \eta + \frac{\eta q_i \varphi_{\mathbf{x}_i^t}}{\alpha_{\mathbf{h}_i^t}^2 + \beta_{\mathbf{h}_i^t}^2}\right) \beta_{\mathbf{x}_i^t}, \quad (23d)$$

where $\{\psi_{\mathbf{h}_i^t}\}, \{\psi_{\mathbf{x}_i^t}\}, \{\varphi_{\mathbf{h}_i^t}\}, \{\varphi_{\mathbf{x}_i^t}\}, \{\rho_{\mathbf{h}_i^t}\}$ and $\{\rho_{\mathbf{x}_i^t}\}$ represent the perturbation terms.

3.2 Dynamics of Approximate State Evolution

It is easily seen that if $\alpha_{\mathbf{h}_i^t}$ (11), $\beta_{\mathbf{h}_i^t}$ (12), $\alpha_{\mathbf{x}_i^t}$ (13) and $\beta_{\mathbf{x}_i^t}$ (14) obey

$$\begin{aligned} |\alpha_{\mathbf{h}_i^t} - q_i| &\leq \frac{\sqrt{2}\gamma}{4\kappa\sqrt{s}} \quad \text{and} \quad \beta_{\mathbf{h}_i^t} \leq \frac{\sqrt{2}\gamma}{4\kappa\sqrt{s}} \quad \text{and} \\ |\alpha_{\mathbf{x}_i^t} - q_i| &\leq \frac{\sqrt{2}\gamma}{4\kappa\sqrt{s}} \quad \text{and} \quad \beta_{\mathbf{x}_i^t} \leq \frac{\sqrt{2}\gamma}{4\kappa\sqrt{s}}, \end{aligned} \quad (24)$$

for $i = 1, \dots, s$, then

$$\begin{aligned} \text{dist}(\mathbf{z}, \mathbf{z}^\dagger) &\leq \left[\frac{s\kappa^2}{2} \left(|\alpha_{\mathbf{h}_i^t} - q_i| + |\beta_{\mathbf{h}_i^t}| \right)^2 + \right. \\ &\quad \left. \frac{s\kappa^2}{2} \left(|\alpha_{\mathbf{x}_i^t} - q_i| + |\beta_{\mathbf{x}_i^t}| \right)^2 \right]^{1/2} \leq \gamma. \end{aligned} \quad (25)$$

In this subsection, we shall show that as long as the approximate state evolution (23) holds, there exists some constant $T_\gamma = \mathcal{O}(s \log \max\{K, N\})$ satisfying condition (24). This is demonstrated in the following Lemma. Prior to that, we first list several conditions and definitions that contribute to the lemma.

- The initial points obey

$$\alpha_{\mathbf{h}_i^0} \geq \frac{q_i}{K \log K} \quad \text{and} \quad \alpha_{\mathbf{x}_i^0} \geq \frac{q_i}{N \log N}, \quad (26a)$$

$$\sqrt{\alpha_{\mathbf{h}_i^0}^2 + \beta_{\mathbf{h}_i^0}^2} \in \left[1 - \frac{1}{\log K}, 1 + \frac{1}{\log K} \right] q_i, \quad (26b)$$

$$\sqrt{\alpha_{\mathbf{x}_i^0}^2 + \beta_{\mathbf{x}_i^0}^2} \in \left[1 - \frac{1}{\log N}, 1 + \frac{1}{\log N} \right] q_i, \quad (26c)$$

for $i = 1, \dots, s$.

- Define

$$T_\gamma := \min \{t : \text{satisfies (24)}\}, \quad (27)$$

where $\gamma > 0$ is some sufficiently small constant.

- Define

$$\begin{aligned} T_1 &:= \min \left\{ t : \min_i \frac{\alpha_{\mathbf{h}_i^t}}{q_i} \geq \frac{c_7}{\log^5 m}, \min_i \frac{\alpha_{\mathbf{x}_i^t}}{q_i} \right. \\ &\quad \left. \geq \frac{c'_7}{\log^5 m} \right\}, \quad (28) \\ T_2 &:= \min \left\{ t : \min_i \frac{\alpha_{\mathbf{h}_i^t}}{q_i} > c_8, \min_i \frac{\alpha_{\mathbf{x}_i^t}}{q_i} > c'_8 \right\}, \quad (29) \end{aligned}$$

for some small absolute positive constants $c_7, c'_7, c_8, c'_8 > 0$.

- For $0 \leq t \leq T_\gamma$, it has

$$\begin{aligned} \frac{1}{2\sqrt{K \log K}} &\leq \frac{\alpha_{\mathbf{h}_i^t}}{q_i} \leq 2, \quad c_5 \leq \frac{\beta_{\mathbf{h}_i^t}}{q_i} \leq 1.5 \quad \text{and} \\ \frac{\alpha_{\mathbf{h}_i^{t+1}} / \alpha_{\mathbf{h}_i^t}}{\beta_{\mathbf{h}_i^{t+1}} / \beta_{\mathbf{h}_i^t}} &\geq 1 + c_5 \eta, \quad i = 1, \dots, s, \end{aligned} \quad (30)$$

$$\begin{aligned} \frac{1}{2\sqrt{N \log N}} &\leq \frac{\alpha_{\mathbf{x}_i^t}}{q_i} \leq 2, \quad c_6 \leq \frac{\beta_{\mathbf{x}_i^t}}{q_i} \leq 1.5 \quad \text{and} \\ \frac{\alpha_{\mathbf{x}_i^{t+1}} / \alpha_{\mathbf{x}_i^t}}{\beta_{\mathbf{x}_i^{t+1}} / \beta_{\mathbf{x}_i^t}} &\geq 1 + c_6 \eta, \quad i = 1, \dots, s, \end{aligned} \quad (31)$$

for some constants $c_5, c_6 > 0$.

Lemma 1. *Assume that the initial points obey condition (26) and the perturbation terms in the approximate state evolution (23) obey $\max\{|\psi_{\mathbf{h}_i^t}|, |\psi_{\mathbf{x}_i^t}|, |\varphi_{\mathbf{h}_i^t}|, |\varphi_{\mathbf{x}_i^t}|, |\rho_{\mathbf{h}_i^t}|\} \leq \frac{c}{\log m}$, for $i = 1, \dots, s, t = 0, 1, \dots$ and some sufficiently small constant $c > 0$.*

1. Then for any sufficiently large K, N and the step-size $\eta > 0$ that obeys $\eta \asymp s^{-1}$, it follows

$$T_\gamma \lesssim s \log(\max\{K, N\}), \quad (32)$$

and (30), (31).

2. Then with the stepsize $\eta > 0$ following $\eta \asymp s^{-1}$, one has that $T_1 \leq T_2 \leq T_\gamma \lesssim s \log \max\{K, N\}$, $T_2 - T_1 \lesssim s \log \log m$, $T_\gamma - T_2 \lesssim s$.

Proof. Please refer to Appendix C in the supplemental material. \square

The random initialization (8) satisfies the condition (26) with probability at least $1 - \mathcal{O}(1/\sqrt{\log \min\{K, N\}})$ [16]. According to this fact, Lemma 1 ensures that under both random initialization (8) and approximate state evolution (23) with the stepsize $\eta \asymp s^{-1}$, Stage I only lasts a few iterations, i.e., $T_\gamma = \mathcal{O}(s \log \max\{K, N\})$. In addition, Lemma 1 demonstrates the exponential growth of the ratios, i.e., $\alpha_{\mathbf{h}_i^{t+1}}/\alpha_{\mathbf{h}_i^t}, \beta_{\mathbf{h}_i^{t+1}}/\beta_{\mathbf{h}_i^t}$, which contributes to the short duration of Stage I.

Moreover, Lemma 1 defines the midpoints T_1 when the sizes of the signal component, i.e., $\alpha_{\mathbf{h}_i^t}$ and $\alpha_{\mathbf{x}_i^t}$, $i = 1, \dots, s$, become sufficiently large, which is crucial to the following analysis. In particular, when establishing the approximate state evolution (23) in Stage I, we analyze two subphases of Stage I individually:

- Phase 1: consider the iterations in $0 \leq t \leq T_1$,
- Phase 2: consider the iterations in $T_1 < t \leq T_\gamma$,

where T_1 is defined in (28).

where T_1 is defined in (28).

3.3 Leave-One-Out Approach

According to Section 3.1 and Lemma 1, the unique challenge for establishing the approximate state evolution (23) is to bound the perturbation terms to certain order, i.e., $|\psi_{\mathbf{h}_i^t}|, |\psi_{\mathbf{x}_i^t}|, |\varphi_{\mathbf{h}_i^t}|, |\varphi_{\mathbf{x}_i^t}|, |\rho_{\mathbf{h}_i^t}|, |\rho_{\mathbf{x}_i^t}| \ll 1/\log m$ for $i = 1, \dots, s$. To achieve this goal, we exploit some variants of leave-one-out sequences [16] to establish the “near-independence” between $\{\mathbf{z}_i^t\}$ and $\{\mathbf{a}_i\}$. Hence, some terms can be approximated by a sum of independent variables with well-controlled weights, thereby be controlled via central limit theorem.

In the following, we define three sets of auxiliary sequences $\{\mathbf{z}^{t,(l)}\}$, $\{\mathbf{z}^{t,\text{sgn}}\}$ and $\{\mathbf{z}^{t,\text{sgn},(l)}\}$, respectively.

- *Leave-one-out sequences* $\{\mathbf{z}^{t,(l)}\}_{t \geq 0}$. For each $1 \leq l \leq m$, the auxiliary sequences $\{\mathbf{z}^{t,(l)}\}$ are established by dropping the l^{th} sample and runs randomly initialized WF with objective function

$$f^{(l)}(\mathbf{z}) = \sum_{j:j \neq l} \left| \sum_{i=1}^s \mathbf{b}_j^* \mathbf{h}_i \mathbf{x}_i^* \mathbf{a}_{ij} - y_j \right|^2. \quad (33)$$

Thus, the sequences $\{\mathbf{z}_i^{t,(l)}\}$ (recall the definition of \mathbf{z}_i (3)) are statistically independent of $\{\mathbf{a}_i\}$.

- *Random-sign sequences* $\{\mathbf{z}^{t,\text{sgn}}\}_{t \geq 0}$. Define the auxiliary design vectors $\{\mathbf{a}_{ij}^{\text{sgn}}\}$ as

$$\mathbf{a}_{ij}^{\text{sgn}} := \begin{bmatrix} \xi_{ij} a_{ij,1} \\ \mathbf{a}_{ij,\perp} \end{bmatrix}, \quad (34)$$

where $\{\xi_{ij}\}$ is a set of standard complex uniform random variables independent of $\{\mathbf{a}_{ij}\}$, i.e.,

$$\xi_{ij} \stackrel{\text{i.i.d.}}{=} u/|u|, \quad (35)$$

where $u \sim \mathcal{N}(0, \frac{1}{2}) + i\mathcal{N}(0, \frac{1}{2})$. Moreover, with the corresponding ξ_{ij} , the auxiliary design vector $\{\mathbf{b}_j^{\text{sgn}}\}$ is defined as $\mathbf{b}_j^{\text{sgn}} = \xi_{ij} \mathbf{b}_j$. With these auxiliary design vectors, the sequences $\{\mathbf{z}^{t,\text{sgn}}\}$ are generated by running randomly initialized WF with respect to the loss function $f^{\text{sgn}}(\mathbf{z})$

$$\sum_{j=1}^m \left| \sum_{i=1}^s \mathbf{b}_j^{\text{sgn}*} \mathbf{h}_i \mathbf{x}_i^* \mathbf{a}_{ij}^{\text{sgn}} - \mathbf{b}_j^{\text{sgn}*} \mathbf{h}_i^{\natural} \mathbf{x}_i^{\natural*} \mathbf{a}_{ij}^{\text{sgn}} \right|^2. \quad (36)$$

Note that these auxiliary design vectors, i.e., $\{\mathbf{a}_{ij}^{\text{sgn}}\}, \{\mathbf{b}_j^{\text{sgn}}\}$ produce the same measurements as $\{\mathbf{a}_{ij}\}, \{\mathbf{b}_j\}$, i.e., $\mathbf{b}_j^{\text{sgn}*} \mathbf{h}_i^{\natural} \mathbf{x}_i^{\natural*} \mathbf{a}_{ij}^{\text{sgn}} = \mathbf{b}_j^* \mathbf{h}_i^{\natural} \mathbf{x}_i^{\natural*} \mathbf{a}_{ij} = q_i a_{ij,1} \mathbf{b}_j^* \mathbf{h}_i^{\natural}$ for $1 \leq i \leq s, 1 \leq j \leq m$.

Note that all the auxiliary sequences are assumed to have the same initial point, namely, for $1 \leq l \leq m$, $\{\mathbf{z}^0\} = \{\mathbf{z}^{0,(l)}\} = \{\mathbf{z}^{0,\text{sgn}}\} = \{\mathbf{z}^{0,\text{sgn},(l)}\}$.

In view of the ambiguities, i.e., $\mathbf{h}_i^{\natural} \mathbf{x}_i^{\natural} = \frac{1}{\omega} \mathbf{h}_i^{\natural} (\omega \mathbf{x}_i^{\natural})^*$, several alignment parameters are further defined for the sequel analysis. Specifically, the alignment parameter between $\mathbf{z}_i^{t,(l)} = [\mathbf{h}_i^{t,(l)*} \mathbf{x}_i^{t,(l)*}]^*$ and $\tilde{\mathbf{z}}_i = [\tilde{\mathbf{h}}_i^{t*} \tilde{\mathbf{x}}_i^{t*}]^*$, where $\tilde{\mathbf{h}}_i^t = \frac{1}{\omega_i^t} \mathbf{h}_i^t$ and $\tilde{\mathbf{x}}_i^t = \omega_i^t \mathbf{x}_i^t$, is represented as

$$\omega_{i,\text{mutual}}^{t,(l)} := \arg \min_{\omega \in \mathbb{C}} \left\| \frac{1}{\omega} \mathbf{h}_i^{t,(l)} - \frac{1}{\omega_i^t} \mathbf{h}_i^t \right\|_2^2 + \left\| \omega \mathbf{x}_i^{t,(l)} - \omega_i^t \mathbf{x}_i^t \right\|_2^2, \quad (37)$$

for $i = 1, \dots, s$. In addition, we denote $\hat{\mathbf{z}}_i^{t,(l)} = [\hat{\mathbf{h}}_i^{t,(l)*} \hat{\mathbf{x}}_i^{t,(l)*}]^*$ where

$$\hat{\mathbf{h}}_i^{t,(l)} := \frac{1}{\omega_{i,\text{mutual}}^{t,(l)}} \mathbf{h}_i^{t,(l)} \quad \text{and} \quad \hat{\mathbf{x}}_i^{t,(l)} := \omega_{i,\text{mutual}}^{t,(l)} \mathbf{x}_i^{t,(l)}. \quad (38)$$

Define the alignment parameter between $\mathbf{z}_i^{t,\text{sgn}} =$

$[\mathbf{h}_i^{t,\text{sgn}*} \ \mathbf{x}_i^{t,\text{sgn}*}]^*$ and $\mathbf{z}_i^t = [\mathbf{h}_i^{t*} \ \mathbf{x}_i^{t*}]^*$ as

$$\omega_{i,\text{sgn}}^t := \arg \min_{\omega \in \mathbb{C}} \left\| \frac{1}{\omega} \mathbf{h}_i^{t,\text{sgn}} - \frac{1}{\omega^t} \mathbf{h}_i^t \right\|_2^2 + \|\omega \mathbf{x}_i^{t,\text{sgn}} - \omega_i^t \mathbf{x}_i^t\|_2^2, \quad (39)$$

for $i = 1, \dots, s$. In addition, we denote $\tilde{\mathbf{z}}_i^{t,\text{sgn}} = [\tilde{\mathbf{h}}_i^{t,\text{sgn}*} \ \tilde{\mathbf{x}}_i^{t,\text{sgn}*}]^*$ where

$$\tilde{\mathbf{h}}_i^{t,\text{sgn}} := \frac{1}{\omega_{i,\text{sgn}}^t} \mathbf{h}_i^{t,\text{sgn}} \quad \text{and} \quad \tilde{\mathbf{x}}_i^{t,\text{sgn}} := \omega_{i,\text{sgn}}^t \mathbf{x}_i^{t,\text{sgn}}. \quad (40)$$

3.4 Induction Hypotheses

In this subsection, we shall establish a collection of induction hypotheses which are crucial to the justification of approximate state evolution (23). We list all the induction hypotheses: for $1 \leq i \leq s$,

$$\begin{aligned} & \max_{1 \leq l \leq m} \text{dist} \left(\mathbf{z}_i^{t,(l)}, \tilde{\mathbf{z}}_i^t \right) \\ & \leq (\beta_{\mathbf{h}_i^t} + \beta_{\mathbf{x}_i^t}) \left(1 + \frac{1}{s \log m} \right)^t C_1 \frac{s \mu^2 \kappa \sqrt{\max\{K, N\} \log^8 m}}{m} \end{aligned} \quad (41a)$$

$$\begin{aligned} & \max_{1 \leq l \leq m} \text{dist} \left(\mathbf{h}_i^{l*} \mathbf{h}_i^{t,(l)}, \mathbf{h}_i^{l*} \tilde{\mathbf{h}}_i^t \right) \cdot \|\mathbf{h}_i^{l*}\|_2^{-1} \\ & \leq \alpha_{\mathbf{h}_i^t} \left(1 + \frac{1}{s \log m} \right)^t C_2 \frac{s \mu^2 \kappa \sqrt{K \log^{13} m}}{m} \end{aligned} \quad (41b)$$

$$\begin{aligned} & \max_{1 \leq l \leq m} \text{dist} \left(\mathbf{x}_i^{t,(l)}, \tilde{\mathbf{x}}_i^t \right) \\ & \leq \alpha_{\mathbf{x}_i^t} \left(1 + \frac{1}{s \log m} \right)^t C_2 \frac{s \mu^2 \kappa \sqrt{N \log^{13} m}}{m} \end{aligned} \quad (41c)$$

$$\begin{aligned} & \max_{1 \leq i \leq s} \text{dist} \left(\mathbf{h}_i^{t,\text{sgn}}, \tilde{\mathbf{h}}_i^t \right) \\ & \leq \alpha_{\mathbf{h}_i^t} \left(1 + \frac{1}{s \log m} \right)^t C_3 \sqrt{\frac{s \mu^2 \kappa^2 K \log^8 m}{m}} \end{aligned} \quad (41d)$$

$$\begin{aligned} & \max_{1 \leq i \leq s} \text{dist} \left(\mathbf{x}_i^{t,\text{sgn}}, \tilde{\mathbf{x}}_i^t \right) \\ & \leq \alpha_{\mathbf{x}_i^t} \left(1 + \frac{1}{s \log m} \right)^t C_3 \sqrt{\frac{s \mu^2 \kappa^2 N \log^8 m}{m}} \end{aligned} \quad (41e)$$

$$\begin{aligned} & \max_{1 \leq l \leq m} \left\| \tilde{\mathbf{h}}_i^t - \hat{\mathbf{h}}_i^{t,(l)} - \tilde{\mathbf{h}}_i^{t,\text{sgn}} + \hat{\mathbf{h}}_i^{t,\text{sgn},(l)} \right\|_2 \\ & \leq \alpha_{\mathbf{h}_i^t} \left(1 + \frac{1}{s \log m} \right)^t C_4 \frac{s \mu^2 \sqrt{K \log^{16} m}}{m}, \end{aligned} \quad (41f)$$

$$\begin{aligned} & \max_{1 \leq l \leq m} \left\| \tilde{\mathbf{x}}_i^t - \hat{\mathbf{x}}_i^{t,(l)} - \tilde{\mathbf{x}}_i^{t,\text{sgn}} + \hat{\mathbf{x}}_i^{t,\text{sgn},(l)} \right\|_2 \\ & \leq \alpha_{\mathbf{x}_i^t} \left(1 + \frac{1}{s \log m} \right)^t C_4 \frac{s \mu^2 \sqrt{N \log^{16} m}}{m}, \end{aligned} \quad (41g)$$

$$c_5 \leq \|\mathbf{h}_i^t\|_2, \|\mathbf{x}_i^t\|_2 \leq C_5, \quad (41h)$$

$$\|\mathbf{h}_i^t\|_2 \leq 5 \alpha_{\mathbf{h}_i^t} \sqrt{\log^5 m}, \quad (41i)$$

$$\|\mathbf{x}_i^t\|_2 \leq 5 \alpha_{\mathbf{x}_i^t} \sqrt{\log^5 m}, \quad (41j)$$

where C_1, \dots, C_5 and c_5 are some absolute positive constants and $\hat{\mathbf{x}}_i, \tilde{\mathbf{x}}_i, \hat{\mathbf{h}}_i, \tilde{\mathbf{h}}_i$ are defined in Section 3.3.

Specifically, (41a), (41c), (41d) and (41e) identify that the auxiliary sequences $\{\mathbf{z}^{t,(l)}\}$ and $\{\mathbf{z}^{t,\text{sgn}}\}$ are extremely close to the original sequences $\{\mathbf{z}^t\}$. In addition, as claimed in (41f) and (41g), $\tilde{\mathbf{h}}_i^t - \tilde{\mathbf{h}}_i^{t,\text{sgn}}$ (resp. $\tilde{\mathbf{x}}_i^t - \tilde{\mathbf{x}}_i^{t,\text{sgn}}$) and $\hat{\mathbf{h}}_i^{t,(l)} - \hat{\mathbf{h}}_i^{t,\text{sgn},(l)}$ (resp. $\hat{\mathbf{x}}_i^{t,(l)} - \hat{\mathbf{x}}_i^{t,\text{sgn},(l)}$) are also exceedingly close to each other. The hypotheses (41h) illustrates that the norm of the iterates $\{\mathbf{h}_i^t\}$ (resp. $\{\mathbf{x}_i^t\}$) is well-controlled in Phase 1. Moreover, (41i) (resp. (41j)) indicates that $\alpha_{\mathbf{h}_i^t}$ (resp. $\alpha_{\mathbf{x}_i^t}$) is comparable to $\|\mathbf{h}_i^t\|_2$ (resp. $\|\mathbf{x}_i^t\|_2$).

Lemma 4-Lemma 9 in the supplemental material demonstrate that if induction hypotheses (41) for $1 \leq i \leq s$ hold true up to the t^{th} iteration, then they are still satisfied in the $(t+1)^{\text{th}}$ iteration under the condition of sufficient samples. Thus auxiliary sequences are close sufficient to the original sequences, which contribute to bound the last term in (22) thereby validating the approximate evolution (23).

4 CONCLUSION

In this paper, we develop a least square estimation approach to solve blind demixing problem. To address the limitations of state-of-the-art algorithms, e.g., high computational complexity, low convergence rate and requirement of careful initialization, we developed the randomly initialized Wirtinger flow to solve the blind demixing problem. The global convergence guarantee of this algorithm is further provided in this paper. It shows that iterates of the randomly initialized Wirtinger flow can enter the local region that enjoys strong convexity and strong smoothness within a few iterations, followed by linearly converging to the ground truth.

Acknowledgements

This work was supported in part by the National Nature Science Foundation of China under Grant 61601290 and the Shanghai Sailing Program under Grant 16YF1407700.

References

- [1] S. Ling and T. Strohmer, "Blind deconvolution meets blind demixing: Algorithms and performance bounds," *IEEE Trans. Inf. Theory*, vol. 63, pp. 4497–4520, Jul. 2017.

-
- [2] J. Dong and Y. Shi, “Nonconvex demixing from bilinear measurements,” *IEEE Trans. Signal Process.*, vol. 66, pp. 5152–5166, Oct. 2018.
- [3] Y. Zhang, Y. Lau, H.-w. Kuo, S. Cheung, A. Papsupathy, and J. Wright, “On the global geometry of sphere-constrained sparse blind deconvolution,” in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, pp. 4894–4902, Jun. 2017.
- [4] E. A. Pnevmatikakis, D. Soudry, Y. Gao, T. A. Machado, J. Merel, D. Pfau, T. Reardon, Y. Mu, C. Lacefield, W. Yang, *et al.*, “Simultaneous denoising, deconvolution, and demixing of calcium imaging data,” *Neuron*, vol. 89, pp. 285–299, Jan. 2016.
- [5] H. Bristow, A. Eriksson, and S. Lucey, “Fast convolutional sparse coding,” in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, pp. 391–398, 2013.
- [6] M. D. Zeiler, D. Krishnan, G. W. Taylor, and R. Fergus, “Deconvolutional networks,” in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, pp. 2528–2535, Jun. 2010.
- [7] J. Dong and Y. Shi, “Blind demixing for low-latency communication,” *IEEE Trans. Wireless Commun.*, Dec. 2018.
- [8] X. Wang and H. V. Poor, “Blind equalization and multiuser detection in dispersive CDMA channels,” *IEEE Trans. Commun.*, vol. 46, pp. 91–103, Jan. 1998.
- [9] S. Ling and T. Strohmer, “Regularized gradient descent: A nonconvex recipe for fast joint blind deconvolution and demixing,” *Inf. Inference: J. IMA*, Mar. 2018.
- [10] R. Ge, F. Huang, C. Jin, and Y. Yuan, “Escaping from saddle points-online stochastic gradient for tensor decomposition,” in *28th Annu. Proc. Conf. Learn. Theory (COLT)*, pp. 797–842, Jul. 2015.
- [11] J. Sun, Q. Qu, and J. Wright, “A geometric analysis of phase retrieval,” *Found. Comput. Math.*, vol. 18, pp. 1131–1198, Oct. 2018.
- [12] C. Jin, R. Ge, P. Netrapalli, S. M. Kakade, and M. I. Jordan, “How to escape saddle points efficiently,” *Proc. Int. Conf. Mach. Learn. (ICML)*, 2017.
- [13] S. Bhojanapalli, B. Neyshabur, and N. Srebro, “Global optimality of local search for low rank matrix recovery,” in *Adv. Neural Inf. Process. Syst. (NIPS)*, pp. 3873–3881, 2016.
- [14] R. Ge, C. Jin, and Y. Zheng, “No spurious local minima in nonconvex low rank problems: A unified geometric analysis,” in *Proc. Int. Conf. Mach. Learn. (ICML)*, pp. 1233–1242, 2017.
- [15] M. Soltanolkotabi, A. Javanmard, and J. D. Lee, “Theoretical insights into the optimization landscape of over-parameterized shallow neural networks,” *IEEE Trans. Inf. Theory*, Jul. 2018.
- [16] Y. Chen, Y. Chi, J. Fan, and C. Ma, “Gradient descent with random initialization: Fast global convergence for nonconvex phase retrieval,” *Math. Program.*, to appear, 2018.
- [17] E. J. Candes, X. Li, and M. Soltanolkotabi, “Phase retrieval via Wirtinger flow: Theory and algorithms,” *IEEE Trans. Inf. Theory*, vol. 61, pp. 1985–2007, Apr. 2015.
- [18] C. Ma, K. Wang, Y. Chi, and Y. Chen, “Implicit regularization in nonconvex statistical estimation: Gradient descent converges linearly for phase retrieval, matrix completion and blind deconvolution,” in *Proc. Int. Conf. Mach. Learn. (ICML)*, pp. 3345–3354, Jul. 2018.
- [19] K. Pothoven, *Real and functional analysis*. Springer Science & Business Media, 2013.