

Yuanming Shi  
Jialin Dong  
Jun Zhang

# Low-overhead Communications in IoT Networks

Structured Signal Processing  
Approaches



Springer

# Low-overhead Communications in IoT Networks

Yuanming Shi • Jialin Dong • Jun Zhang

# Low-overhead Communications in IoT Networks

Structured Signal Processing Approaches

 Springer

Yuanming Shi  
School of Information Science  
and Technology  
Shanghai Tech University  
Shanghai, China

Jialin Dong  
School of Information Science  
and Technology  
ShanghaiTech University  
Shanghai, China

Jun Zhang  
Department of Electronic & Information  
Engineering  
Hong Kong Polytechnic University  
Kowloon, Hong Kong

ISBN 978-981-15-3869-8      ISBN 978-981-15-3870-4 (eBook)  
<https://doi.org/10.1007/978-981-15-3870-4>

© Springer Nature Singapore Pte Ltd. 2020

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors, and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, expressed or implied, with respect to the material contained herein or for any errors or omissions that may have been made. The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

This Springer imprint is published by the registered company Springer Nature Singapore Pte Ltd.  
The registered company address is: 152 Beach Road, #21-01/04 Gateway East, Singapore 189721, Singapore

# Preface

The past decades have witnessed a revolution in wireless communications and networking, which has profoundly changed our daily life. Particularly, it has enabled various innovative Internet of Things (IoT) applications, e.g., smart city, healthcare, and autonomous driving and drones. The IoT architecture is established by the proliferation of low-cost and small-size mobile devices. With the explosion of IoT devices, a heavy burden is placed on the wireless access. A key characteristic of IoT data traffic is the sporadic pattern, i.e., only a portion of all the devices are active at a given time instant. In particular, in many IoT applications, devices are designed to be inactive most of the time to save energy and only be activated by external events. Thus, with massive IoT devices, it is of vital importance to manage their random access procedures, detect the active ones, and decode their data at the access point. Massive IoT connectivity has been regarded as one of the key performance requirements of 5G and beyond networks.

The emerging IoT applications have stringent demands on low-latency communications and typically transmit short packets containing both the metadata and payload. The metadata may include packet initiation and termination information, logical addresses, security and synchronization information, etc. It also contains a channel estimation sequence that facilitates channel estimation at the access point. Additionally, the metadata includes various information about the packet structure, e.g., the pilot sequences used for random access and device identification information. Considering the typical small payload size of IoT applications, it is of critical importance to reduce the size of the overhead message, e.g., identification information, pilot symbols for channel estimation, control data, etc. Such low-overhead communications also help to improve the energy efficiency of IoT devices. Recently, structured signal processing approaches have been introduced and developed to reduce the overheads for key design problems in IoT networks, such as channel estimation, device identification, and message decoding. By exploiting underlying system and problem structures, including sparsity and low rank structures, these methods can achieve significant performance gains. Chapter 1 provides more background on low-overhead communications in IoT networks and introduces general structured signal processing techniques.

This monograph shall provide an overview of four structured signal processing models, i.e., a sparse linear model, a blind demixing model, a sparse blind demixing model, and a shuffled linear regression model. Chapter 2 introduces a sparse linear model for joint activity detection and channel estimation in IoT networks with grant-free random access. A convex relaxation approach based on  $\ell_p$ -norm minimization is firstly introduced, followed by a smoothed primal-dual first-order algorithm to solve it. For this convex relaxation approach, a trade-off between the computational cost and estimation accuracy is characterized by Proposition 2.1. The theoretical analysis of the convex relaxation approach is based on the conic integral geometry theory. This chapter only contains a brief introduction on the conic integral geometry theory. For more details, the interested reader can refer to Sect. 8.1 and other related mathematical literature enumerated in this monograph. Besides, an iterative threshold algorithm, namely approximate message passing (AMP), is introduced in Chap. 2, followed by the performance analysis based on the state evolution technique.

Blind demixing is introduced in Chap. 3, which facilitates joint data decoding and channel estimation without explicit pilot sequences. After presenting the basic convex relaxation approach for solving the blind demixing problem, we introduce three nonconvex approaches: the regularized Wirtinger flow, the regularization-free Wirtinger flow, and a Riemannian optimization algorithm. Theorems 3.1 and 3.2 provide the theoretical analysis of the convex relaxation approach and regularized Wirtinger flow, respectively. Furthermore, Theorem 3.3 presents the theoretical guarantees of the Wirtinger flow with the spectral initialization, which provides readers an easy access to well-round results. Readers who are interested in the intrinsic mechanism of the theoretical analysis can refer to Sect. 8.3 for more discussions. The theoretical analysis of the Wirtinger flow via random initialization is further provided in Sect. 8.4. Additionally, the basic concepts of Riemannian manifold optimization are presented in Sect. 8.5, which provide sufficient background for related algorithms in Chaps. 3 and 4. The extension of blind demixing, i.e., sparse blind demixing, is introduced in Chap. 4, which further takes device activity into consideration. The sparse blind demixing formulation is able to jointly consider device activity detection, data decoding, and channel estimation, for which three approaches are presented: a convex relaxation approach, a difference-of-convex-functions approach, and smoothed Riemannian optimization.

A further step to reduce the overhead is to remove the device identification information from the metadata. To support the joint data decoding and device identification, shuffled linear regression is introduced in Chap. 5. We first present maximum likelihood estimation (MLE) based approaches for solving the shuffled linear regression problem. Theorems 5.1 and 5.2 provide the statistical properties of the MLE, and both an upper bound and a lower bound on the probability of error of the permutation matrix estimator are introduced. To solve the MLE problem, two algorithms are presented: one is based on sorting, and the other algorithm returns an approximate solution to the MLE problem. Next, theoretical analysis of the shuffled linear regression problem based on the algebraic–geometric theory is presented. Based on the analysis, an algebraically initialized expectation–maximization algo-

rithm is introduced to solve the shuffled linear regression problem, which enjoys better algorithmic performance than previous works. To give a comprehensive introduction of the algebraic–geometric theory, besides the concepts mentioned in this chapter, we introduce several related definitions on the algebraic–geometric theory in Sect. 8.7, including the geometric characterization of dimension, algebraic characterization of dimension, homogenization, and regular sequences.

Furthermore, Chap. 6 provides some cutting-edge learning augmented techniques for structured signal processing on the aspects of structured signal model design (e.g., structured signal processing under a generative prior) and algorithm design (e.g., deep-learning-based algorithm). We begin with compressed sensing under a generative prior, and other structure signal processing techniques under a generative model are worth further investigating, e.g., blind deconvolution. We next consider the joint design of measurement matrix and sparse support recovery for the sparse linear model (e.g., compressed sensing). Some basic methods are firstly presented, i.e., sample scheduling and sensing matrix optimization, and then learning augmented techniques are introduced. Additionally, for estimating the sparse linear model, several deep-learning-based AMP methods are introduced in this chapter: learned AMP, learned Vector-AMP, and learned ISTA for group row sparsity. In Chap. 7, we summarize the book and discuss some potential extensions of the area of interest. Tables 7.1 and 7.2 list the main theorems, propositions, and algorithms presented in this monograph.

The monograph is not only suitable for beginners in structured signal processing for applications in IoT networks but also helpful to experienced researchers who intend to work in-depth on the theoretical analysis of structured signals. For beginners, the background of both low-overhead communications and structured signal processing in Chap. 1 is helpful, and the problem formulation section in each chapter may be referred for further details with respect to each model. Tables 1.1, 7.1, and 7.2 provide quick references for the main results. Readers who are more interested in the intrinsic mechanism of the theoretical analysis of the specific models can refer to Chap. 8.

Low-overhead communications supported by structured signal processing approaches have received significant attention in recent years. The main motivation of this monograph is to provide an overview of the major structured signal processing models, along with their applications in low-overhead communications in IoT networks. Practical algorithms, via both convex and nonconvex optimization approaches, and theoretical analysis, using various mathematical tools, will be introduced. While the structured signal models concerned in this monograph have certain limitations, we hope the presented results will galvanize researchers into investigating this influential and intriguing area.

## Acknowledgements

We express our gratitude to the support of National Nature Science Foundation of China under Grant 61601290. We also thank Xinyu Bian for proofreading an early version of the manuscript and Prof. Liang Liu for providing codes for part of the simulations in Sect. [2.4](#).

Shanghai, China  
Shanghai, China  
Kowloon, Hong Kong

Yuanming Shi  
Jialin Dong  
Jun Zhang

# Contents

<b>1</b>	<b>Introduction</b>	1
1.1	Low-Overhead Communications in IoT Networks	1
1.1.1	Grant-Free Random Access	2
1.1.2	Pilot-Free Communications	4
1.1.3	Identification-Free Communications	5
1.2	Structured Signal Processing	5
1.2.1	Example: Compressed Sensing	6
1.2.2	General Structured Signal Processing	7
1.3	Outline	8
	References	10
<b>2</b>	<b>Sparse Linear Model</b>	13
2.1	Joint Activity Detection and Channel Estimation	13
2.2	Problem Formulation	14
2.2.1	Single-Antenna Scenario	15
2.2.2	Multiple-Antenna Scenario	16
2.3	Convex Relaxation Approach	17
2.3.1	Method: $\ell_p$ -Norm Minimization	17
2.3.2	Algorithm: Smoothed Primal-Dual First-Order Methods	18
2.3.3	Analysis: Conic Integral Geometry	21
2.4	Iterative Thresholding Algorithm	27
2.4.1	Algorithm: Approximate Message Passing	28
2.4.2	Analysis: State Evolution	29
2.5	Summary	32
	References	32
<b>3</b>	<b>Blind Demixing</b>	35
3.1	Joint Data Decoding and Channel Estimation	35
3.2	Problem Formulation	37
3.2.1	Cyclic Convolution	37
3.2.2	System Model	38
3.2.3	Representation in the Fourier Domain	38

3.3	Convex Relaxation Approach .....	40
3.3.1	Method: Nuclear Norm Minimization .....	40
3.3.2	Theoretical Analysis .....	41
3.4	Nonconvex Approaches .....	42
3.4.1	Regularized Wirtinger Flow .....	43
3.4.2	Regularization-Free Wirtinger Flow .....	45
3.4.3	Riemannian Optimization Algorithm .....	49
3.4.4	Simulation Results .....	55
3.5	Summary .....	58
	References .....	58
<b>4</b>	<b>Sparse Blind Demixing .....</b>	<b>59</b>
4.1	Joint Device Activity Detection, Data Decoding, and Channel Estimation .....	59
4.2	Problem Formulation .....	60
4.2.1	Single-Antenna Scenario .....	60
4.2.2	Multiple-Antenna Scenario .....	61
4.3	Convex Relaxation Approach .....	61
4.4	Difference-of-Convex-Functions (DC) Programming Approach .....	63
4.4.1	Sparse and Low-Rank Optimization .....	64
4.4.2	A DC Formulation for Rank Constraint .....	65
4.4.3	DC Algorithm for Minimizing a DC Objective .....	66
4.4.4	Simulations .....	67
4.5	Smoothed Riemannian Optimization on Product Manifolds .....	69
4.5.1	Optimization on Product Manifolds .....	69
4.5.2	Smoothed Riemannian Optimization .....	70
4.5.3	Simulation Results .....	71
4.6	Summary .....	73
	References .....	73
<b>5</b>	<b>Shuffled Linear Regression .....</b>	<b>75</b>
5.1	Joint Data Decoding and Device Identification .....	75
5.2	Problem Formulation .....	77
5.3	Maximum Likelihood Estimation Based Approaches .....	78
5.3.1	Sorting Based Algorithms .....	78
5.3.2	Approximation Algorithm .....	80
5.4	Algebraic-Geometric Approach .....	82
5.4.1	Eliminating $\Pi$ via Symmetric Polynomials .....	83
5.4.2	Theoretical Analysis .....	85
5.4.3	Algebraically Initialized Expectation-Maximization .....	86
5.4.4	Simulation Results .....	88
5.5	Summary .....	89
	References .....	89

<b>6</b>	<b>Learning Augmented Methods</b> .....	91
6.1	Structured Signal Processing Under a Generative Prior .....	91
6.2	Joint Design of Measurement Matrix and Sparse Support Recovery .....	94
6.3	Deep-Learning-Based AMP .....	96
6.3.1	Learned AMP .....	97
6.3.2	Learned Vector-AMP .....	99
6.3.3	Learned ISTA for Group Row Sparsity .....	100
6.4	Summary .....	103
	References .....	105
<b>7</b>	<b>Conclusions and Discussions</b> .....	107
7.1	Summary .....	107
7.2	Discussions .....	109
	References .....	109
<b>8</b>	<b>Appendix</b> .....	111
8.1	Conic Integral Geometry .....	111
8.1.1	The Kinematic Formula for Convex Cones .....	111
8.1.2	Intrinsic Volumes and the Statistical Dimension .....	112
8.1.3	The Approximate Kinematic Formula .....	114
8.1.4	Computing the Statistical Dimension .....	114
8.2	Proof of Proposition 2.1 .....	115
8.3	Proof of Theorem 3.3 .....	116
8.3.1	Proof of Lemma 8.4 .....	129
8.4	Theoretical Analysis of Wirtinger Flow with Random Initialization for Blind Demixing .....	139
8.5	The Basic Concepts on Riemannian Optimization .....	142
8.6	Proof of Theorem 3.4 .....	147
8.7	Basic Concepts in Algebraic–Geometric Theory .....	149
8.7.1	Geometric Characterization of Dimension .....	149
	References .....	151

# Mathematical Notations

- The set of real numbers is denoted by  $\mathbb{R}$  and the set of positive real numbers is denoted by  $\mathbb{R}_+$ . The set of complex numbers is denoted by  $\mathbb{C}$ . Denote  $\mathbb{S}_+$  as the set of Hermitian positive semidefinite matrices. Moreover,  $\mathbb{N}$  represents the set of natural numbers.
- The boldface and lowercase alphabets, e.g.,  $\mathbf{x}$ ,  $\mathbf{y}$ , denote vectors. The zero vector is denoted by  $\mathbf{0}$ . A vector  $\mathbf{x} \in \mathbb{R}^d$  is in the column format. The transpose of a vector is denoted by  $\mathbf{x}^\top$ . The complex conjugate of  $\mathbf{x}$  is represented as  $\bar{\mathbf{x}}$ . The conjugate transpose of a vector is denoted by  $\mathbf{x}^H$  or  $\mathbf{x}^*$ .  $x_i$  denotes the  $i$ -th coordinate of a vector  $\mathbf{x}$ .
- For a complex vector  $\mathbf{x}$  or a complex matrix  $\mathbf{X}$ , the real parts of them are represented by  $\Re\{\mathbf{x}\}$  and  $\Re\{\mathbf{X}\}$ , respectively. Likewise, the imaginary parts are denoted as  $\Im\{\mathbf{x}\}$  and  $\Im\{\mathbf{X}\}$ .
- The boldface and uppercase alphabets, e.g.,  $\mathbf{A}$ ,  $\mathbf{B}$ , denote matrices.  $A_{ij}$  denotes the element at the  $i$ -th row and the  $j$ -th column.
- The support function of a vector  $\mathbf{x}$  is denoted as

$$\text{supp}(\mathbf{x}) := \{i : x_i \neq 0\}.$$

A vector  $\mathbf{x}$  such that  $|\text{supp}(\mathbf{x})| \leq s$  is defined as  $s$ -sparse.

- For a vector  $\mathbf{x} \in \mathbb{R}^d$  or  $\mathbf{x} \in \mathbb{C}^d$ , its  $\ell_p$ -norm is given by

$$\|\mathbf{x}\|_p = \sum_{i=1}^d |x_i|^p.$$

In certain cases, we define  $\ell_0$ -norm as  $\|\mathbf{x}\|_0 := |\text{supp}(\mathbf{x})|$ .

- For a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  or  $\mathbf{A} \in \mathbb{C}^{m \times n}$ , the Frobenius norm of  $\mathbf{A}$  is defined as

$$\|\mathbf{A}\|_F := \sqrt{\sum_{i,j} |A_{ij}|^2} = \sqrt{\sum_i \sigma_i(\mathbf{A})^2},$$

where  $\sigma_1(A) \geq \sigma_2(A) \geq \dots \geq \sigma_{\min\{m,n\}}(A)$  denote its singular values. The nuclear norm of  $\mathbf{A}$  is denoted as  $\|\mathbf{A}\|_* := \sum_i \sigma_i(\mathbf{A})$ . The spectral norm of a matrix  $\mathbf{A}$  is denoted as

$$\|\mathbf{A}\| := \max_i \sigma_i(\mathbf{A}).$$

- The cardinality of a set  $\mathcal{S}$  is denoted by  $|\mathcal{S}|$ .
- Random variables or events are denoted as uppercase letters, i.e.,  $X, Y, E$ .
- The indicator function of an event  $E$  is denoted by  $y = \mathbb{I}(E)$ , where  $y = 1$  if the event  $E$  is true, otherwise  $y = 0$ .
- Throughout this book,  $f(n) = \mathcal{O}(g(n))$  or  $f(n) \lesssim g(n)$  denotes that there exists a constant  $c > 0$  such that  $|f(n)| \leq c|g(n)|$ , whereas  $f(n) = \Omega(g(n))$  or  $f(n) \gtrsim g(n)$  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 exist constants  $c_1, c_2 > 0$  such that  $c_1|g(n)| \leq |f(n)| \leq c_2|g(n)|$ .
- For a general cone  $C \subset \mathbb{R}^d$ , the *polar cone*  $C^\circ$  is the set of outward normals of  $C$ :

$$C^\circ := \{\mathbf{u} \in \mathbb{R}^d : \langle \mathbf{u}, \mathbf{x} \rangle \leq 0 \text{ for all } \mathbf{x} \in C\}.$$

The polar cone  $C^\circ$  is always closed and convex.

# Chapter 1

## Introduction

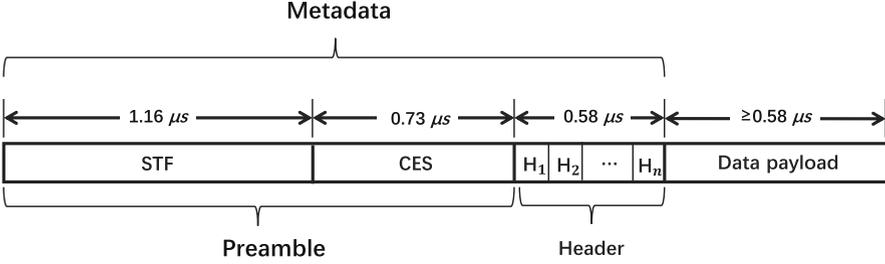


**Abstract** This chapter presents a background on low-overhead communications in IoT networks and structured signal processing. It starts with introducing three key techniques for low-overhead communications: grant-free random access, pilot-free communications, and identification-free communications. Then different models for structured signal processing to support low-overhead communications are presented, which form the main theme of this monograph. A classical exemplary of structure signal processing, i.e., compressed sensing, is provided to illustrate the main principles of algorithm design and theoretical analysis. Finally, the outline of the monograph is presented.

### 1.1 Low-Overhead Communications in IoT Networks

The proliferation of low-cost and small-size computing devices endowed with communication and sensing capabilities is paving the way for the era of IoT. These devices can support various innovative applications, including smart city, health-care [1], and autonomous driving [22] and drones [27]. The explosion of IoT devices places a heavy burden on the wireless network, as they demand scalable wireless access, which has been put forward as a key challenge of 5G and beyond networks [23]. A key characteristic of IoT data traffic is the sporadic pattern where only a small part of all devices are active at any time instant. In particular, in many IoT applications, devices are designed to be inactive most of the time to save energy and only be activated by external events [26]. Thus, with massive IoT devices, it is of vital importance to manage their random access procedure, detect the active ones, and decode their data at the access point.

Moreover, the emerging IoT applications have stringent demands on low-latency communications, and typically transmit short packets containing both the metadata and payload [16]. An exemplary packet structure is illustrated in Fig. 1.1 (please refer to [29] for more details). The metadata may include packet initiation and termination information, logical addresses, security and synchronization information, etc. [16]. In the example showed in Fig. 1.1, we simply illustrate the metadata that contains a preamble and a header coming from the media-access-control



- STF: short training field
- CES: channel estimation sequence
- H<sub>1</sub>: pilot sequence
- H<sub>2</sub>: device identification information

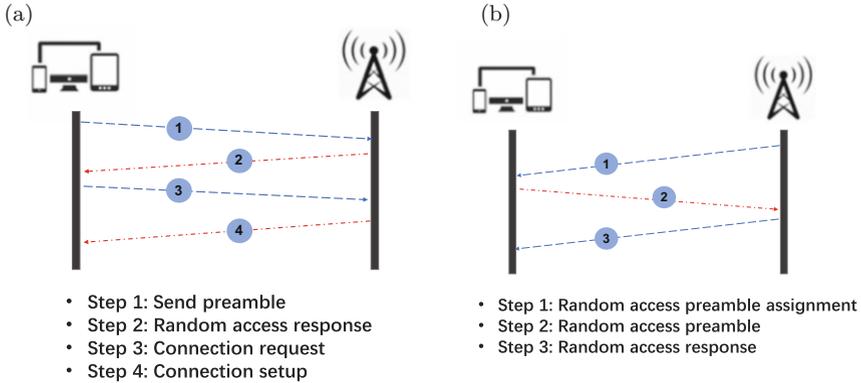
**Fig. 1.1** An exemplary packet structure

(MAC) layer and physical (PHY) layer. Specifically, the preamble contains a short training field (STF), which will be used for packet detection, indication of the modulation type, frequency offset estimation, synchronization, etc. It also contains a channel estimation sequence (CES) that facilitates channel estimation at the access point. Additionally, the header includes various information about the packet structure, e.g., the pilot sequences used for random access and device identification information. The header also includes the modulation and coding scheme adopted for transmitting the data payload. Furthermore, it may include the length of the payload and a header checksum field [16].

From the packet structure in Fig. 1.1, we see that the efficiency of short-packet transmissions, in terms of energy, latency, and bandwidth cost, critically depends on the size of the metadata, which is comparable to the payload size in many cases. To improve the communication efficiency, plenty of efforts have been made to reduce the size of the metadata, which result in *low-overhead communications*. Reducing overheads will not only improve spectral efficiency, reduce latency, but also achieve significant energy saving, which is especially important for resource-constrained IoT devices. In the sequel, we shall introduce three representative methods for reducing overheads.

### 1.1.1 Grant-Free Random Access

Conventionally, the grant-based random access scheme (illustrated in Fig. 1.2a) is applied to allow multiple users to access the network, e.g., in 4G LTE networks [3, 19, 26]. Under this scheme, each active device randomly chooses a pilot sequence from a predefined set of orthogonal preamble sequences to inform the base station (BS) of the device's active state. A connection between the BS and the active device



**Fig. 1.2** Random access schemes. Note that for the grant-based scheme, steps 1–3 may need to repeat multiple times to establish a connection due to contention. (a) Grant-based. (b) Grant-free

will be established if the pilot sequence of this active device is not occupied by others. In this case, the BS will send a contention-resolution message to inform the device of the radio resources reserved for its data transmission. If two or more devices have selected the same pilot sequence, their connection requests collide. Once the BS detects this collision, it will not reply with a contention-resolution message. Instead, the affected devices have to restart the random access procedure again, which leads to high latency. Note that the messages sent by the active devices in the first and third phases correspond to metadata, as they are control information for establishing the connection without carrying any payload. Besides the overhead, a major drawback of the grant-based random access scheme is the limited number of active devices that can receive the grant to access the network. For example, as shown in [26], for a network with one BS and 2000 devices, a minimum length of the pilot sequence of 470, out of the total 1000 symbols, is needed to guarantee a 90% success rate. Even equipped with advanced contention-resolution strategies [6], 930 out of 1000 symbols are still required for transmitting the pilot sequence.

To address the collision issue of the random access scheme caused by a massive number of devices, the grant-free random access scheme illustrated in Fig. 1.2b has been proposed. With this new scheme, the devices do not need to wait for any grant to access the network and can directly transmit the coupled metadata and data to the BS. In this way, the BS can perform user activity detection, channel estimation, and/or data decoding simultaneously [33, 36–39]. The essential idea underlying this line of studies is to connect with sparse signal processing and leverage the compressed sensing technique. In particular, a compressed sensing problem is established by exploiting the sparsity in the user activity pattern. The received signal at the BS equipped with a single antenna is given by a *sparse linear model*:

$$\mathbf{y} = \mathbf{A}\mathbf{x}, \quad (1.1)$$

where  $\mathbf{x}$ , denoting the activity of devices, is a sparse vector due to the sporadic traffic pattern. Then compressed sensing techniques can be applied to recover the sparse vector. Such grant-free random access has received lots of attention recently. To be specific, when the channel state sequences (recall CES in Fig. 1.1) are contained in the metadata, a joint device activity and data detection problem was studied in [39]. Regarding the sparse linear model proposed in [39], the matrix  $\mathbf{A}$  in (1.1) captures the channel estimation sequences and pilot sequences, and the vector  $\mathbf{x}$  in (1.1) represents the information symbols of all devices, where the value is 0 for each inactive device. To improve the efficiency, the overhead caused by metadata has been further reduced. When the CES in Fig. 1.1 is eliminated from the metadata during the packet transmission [33, 36, 37], performed joint channel and data estimation based on various compressed sensing techniques with  $\mathbf{A}$  in (1.1) capturing the data for all (active and inactive) devices and pilot sequences. Moreover, device activity detection and channel estimation were jointly achieved in the work [38]. In this scenario, the matrix  $\mathbf{A}$  in (1.1) characterizes the pilot sequences, and the vector  $\mathbf{x}$  in (1.1) contains the device activity and channel information.

### 1.1.2 Pilot-Free Communications

In the grant-free random access scheme, the pilot sequence is needed for activity detection, which requires extra bandwidth and induces additional overhead. A more aggressive approach is the pilot-free communication scheme that removes both the fields  $H_1$  and CES in Fig. 1.1 from the metadata. To elude the pilot sequences, more powerful signal processing techniques are needed for data detection. Specifically, a blind demixing based approach has been developed in [11, 14, 24, 25]. Consider an IoT network with one BS and  $s$  devices. Each device transmits an encoded signal  $\mathbf{f} = \mathbf{A}\mathbf{x}$  to the BS through the channel  $\mathbf{g}$ , where  $\mathbf{x}$  is the message and  $\mathbf{A}$  is the encoding matrix, and the received signal at the BS is represented by the cyclic convolution operator  $\circledast$ ,

$$\mathbf{y} = \sum_{i=1}^s \mathbf{f}_i \circledast \mathbf{g}_i, \quad (1.2)$$

which is a *blind demixing model* that facilitates to demix the original signals  $\{\mathbf{f}_i\}$  from the observation  $\mathbf{y}$  without the knowledge of the channel states  $\{\mathbf{g}_i\}$ . The blind demixing based approach can achieve joint data decoding (i.e., recover data  $\mathbf{x}$  for each device) and channel estimation. With pilot-free communication supported by blind demixing, the overhead during the transmission is effectively reduced via waiving both the pilot sequences and channel estimation sequences from the metadata.

Considering the sporadic traffic pattern in the IoT network where only part (denoted as  $\mathcal{S}$ ) of the devices are active, a sparse blind demixing model is further developed in [12, 18], given by a *sparse blind demixing model*:

$$\mathbf{y} = \sum_{i \in \mathcal{S}} \mathbf{f}_i \otimes \mathbf{g}_i. \quad (1.3)$$

The estimation for the sparse blind demixing model aims to achieve joint device activity detection and data decoding without the channel state information. Similar to the blind demixing model, it can facilitate to reduce the overhead caused by the channel state information and pilot sequence in the metadata, using more sophisticated detection algorithms.

### 1.1.3 Identification-Free Communications

Besides the above methods of reducing overhead, excluding the identification information is an important consideration in some IoT applications. Specifically, the identification-free communication scheme eliminates the field  $H_2$  in Fig. 1.1 from the metadata. As an example, suppose that multiple sensors are deployed to take measurements of an unknown parameter vector  $\mathbf{x}$ . In this case, the overhead is mainly dominated by the identity information contained in the metadata [21]. To reduce the overhead, a shuffled linear regression model has been developed. It is established by introducing an unknown permutation matrix  $\mathbf{\Pi}$  of which the  $i$ -th row is the canonical vector  $\mathbf{e}_{\pi(i)}^\top$  of all zeros except a 1 at position  $\pi(i)$ :

$$\mathbf{y} = \mathbf{\Pi} \mathbf{A} \mathbf{x}. \quad (1.4)$$

The goal of the data fusion is to recover  $\mathbf{x}$  from the permuted data  $\mathbf{y}$  based on the known sensing matrix  $\mathbf{A}$ . That is, the identities of the signals sent by the sensors are not accessible to the fusion. To address this challenging problem, a line of literatures have developed advanced algorithms from theoretical and practical points of view [30, 31, 34, 35].

## 1.2 Structured Signal Processing

The techniques mentioned above to achieve low-overhead communications rely on structured signal processing, which exploits underlying structures of the signals or systems, e.g., sparsity, low-rankness, group sparsity or permutation, for effective signal estimation and detection. In this section, we first take a basic structured signal processing problem, i.e., compressed sensing, as an example, to illustrate the main design principles. Then general structured signal processing techniques are introduced.

### 1.2.1 Example: Compressed Sensing

The key point of compressed sensing is to recover a sparse signal from very few linear measurements. Mathematically, given a *sensing matrix*, i.e.,  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , the *compressed sensing problem* can be formulated as recovering  $\mathbf{x} \in \mathbb{R}^n$  from the observation of

$$\mathbf{y} = \mathbf{A}\mathbf{x} \in \mathbb{R}^m, \quad (1.5)$$

based on the assumption that  $\mathbf{x}$  has very few nonzero elements, i.e., the  $\ell_0$ -norm  $\|\mathbf{x}\|_0$  is small. In the sequel, we introduce three key ingredients of a compressed sensing problem: recovery algorithms, measurement mechanisms, and theoretical guarantees.

It is intuitive to recover  $\mathbf{x}$  from the observation  $\mathbf{y}$  via solving

$$\underset{\mathbf{x} \in \mathbb{C}^n}{\text{minimize}} \quad \|\mathbf{x}\|_0 \quad \text{subject to} \quad \mathbf{y} = \mathbf{A}\mathbf{x}. \quad (1.6)$$

The paper [5] showed that problem (1.6) enables to recover a  $k$ -sparse signal exactly with a high probability with only  $m = k + 1$  random measurements from a Gaussian distributed sensing matrix. Unfortunately, problem (1.6) is a combinatorial optimization problem with an excessive complexity if solved by enumeration [28]. Thus, the tightest convex norm of  $\ell_0$ -norm, i.e., the  $\ell_1$ -norm, is proposed to relax  $\ell_0$ -norm [8], which leads to

$$\underset{\mathbf{x} \in \mathbb{C}^n}{\text{minimize}} \quad \|\mathbf{x}\|_1 \quad \text{subject to} \quad \mathbf{y} = \mathbf{A}\mathbf{x}. \quad (1.7)$$

Intuitively, this  $\ell_1$  minimization formulation facilitates to induce sparsity due to the shape of the  $\ell_1$  ball.

There have been various types of algorithms developed for different formulations of sparse recovery. The most commonly used formulation is the convex relaxation based on (1.7). Given a certain parameter  $\lambda > 0$ , problem (1.7) can also be represented as an unconstrained optimization problem,

$$\underset{\mathbf{x} \in \mathbb{C}^n}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{x}\|_1. \quad (1.8)$$

Various algorithms have been developed to solve problem (1.8), including interior-point methods [7], projected gradient methods [17], iterative thresholding [10], and the approximate message passing algorithm [15].

Besides effective recovery algorithms, there exist rigorous theoretical guarantees on the recovery of sparse signals, based on specific conditions of the measurements matrix. In particular, the restricted isometry property (RIP) of a sensing matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  was introduced in [7] that measures the degree to which each subset of  $k$

column vectors of  $\mathbf{A} \in \mathbb{R}^{m \times n}$  is close to being an isometry. A typical theoretical result based on RIP analysis is stated as follows:

*Example 1.1* If there exists a  $\delta_k \in (0, 1)$  such that the sensing matrix  $\mathbf{A}$  satisfies

$$(1 - \delta_k) \|\mathbf{x}\|_2^2 \leq \|\mathbf{Ax}\|_2^2 \leq (1 + \delta_k) \|\mathbf{x}\|_2^2 \quad (1.9)$$

for any  $\mathbf{x}$  that belongs to the set of  $k$ -sparse vectors, then problem (1.7) can facilitate to exactly recover the sparse vector  $\mathbf{x}$  with high probability, provided the number of measurements  $m \gtrsim \delta_k^{-2} k \log(n/k)$ .

In addition, the exact location of the phase transition for problem (1.7) can be obtained based on conic geometry theory, where a parameter, called the statistical dimension, is introduced to capture the dimension of a linear subspace to the set of convex cones [2]. It demonstrates that under the assumption of i.i.d. standard normal measurements, the transition occurs where the number of measurements, i.e.,  $m$ , equals the statistical dimension of the descent cone. The shift from failure to success occurs over a range of about  $\mathcal{O}(\sqrt{n})$  measurements.

## 1.2.2 General Structured Signal Processing

Compressive sensing techniques have been successfully applied in many application domains, which have inspired lots of interest in exploiting structures other than sparsity [4, 9, 13, 14, 20, 26, 32, 36]. In the following, we give a brief introduction of the structured signal processing approaches that will be applied for low-overhead communications in this monograph. A general structured signal processing problem with a vector variable is given by

$$\underset{\mathbf{x} \in \mathcal{D}_v}{\text{minimize}} f(\mathcal{A}\mathbf{x}), \quad (1.10)$$

where  $\mathcal{A}$  is a linear operator representing the measurement mechanism,  $f$  is a loss function, and  $\mathcal{D}_v$  is a space of structured vectors (e.g., sparse vectors). For example, in the sparse linear model in (1.1), the operator  $\mathcal{A}$  captures the set of pilot matrices and  $\mathbf{x}$  is a sparse vector. In the shuffled linear regression problem (1.4), the operator  $\mathcal{A}$  indicates the permuted sensing matrix. Likewise, a general structured signal processing problem with a matrix variable is given by

$$\underset{\mathbf{X} \in \mathcal{D}_m}{\text{minimize}} f(\mathcal{A}\mathbf{X}), \quad (1.11)$$

where  $\mathcal{D}_m$  is a space of structured matrices (e.g., low-rank matrices, sparse matrices, low-rank and sparse matrices, etc.). Specifically, in the blind demixing model (1.2),  $\mathbf{X}$  is a collection of rank-1 matrices, while  $\mathbf{X}$  in the sparse blind demixing (1.3) is a collection of low-rank matrices endowed with group sparsity.

The convex program based on the norm operator is typically an effective way to solve problems (1.10) and (1.11) with  $\mathbf{x}$  and  $\mathbf{X}$  enjoying certain structures such as sparsity, low-rankness, and group sparsity. Specifically, the convex program with a vector variable can be represented as

$$\underset{\mathbf{x}}{\text{minimize}} f(\mathcal{A}\mathbf{x}), \quad \text{subject to} \quad \|\mathbf{x}\|_1 \leq \alpha. \quad (1.12)$$

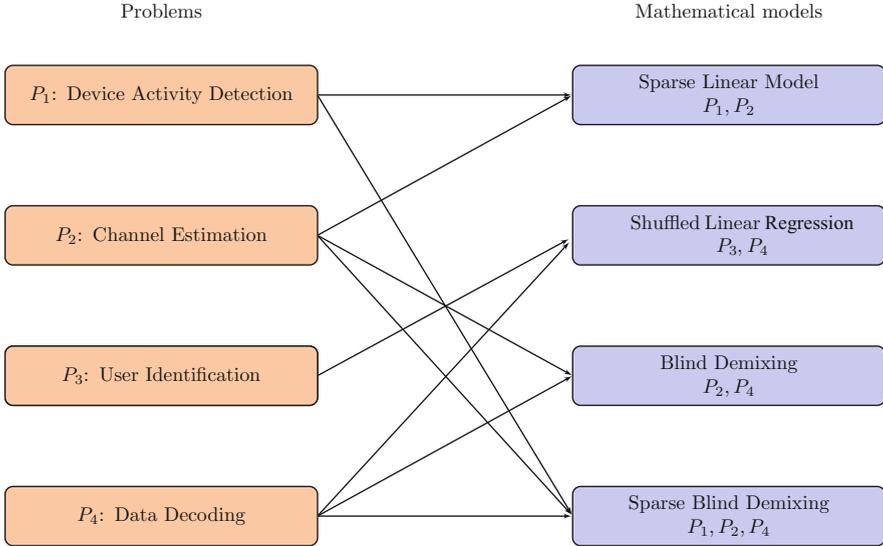
Here, the  $\ell_1$ -norm can be used for inducing sparsity of a vector. Moreover, the convex program with a matrix variable is given by

$$\underset{\mathbf{X}}{\text{minimize}} f(\mathcal{A}\mathbf{X}), \quad \text{subject to} \quad \mathcal{M}(\mathbf{X}) \leq \alpha, \quad (1.13)$$

where the operator  $\mathcal{M}(\cdot)$  indicates a norm operator to induce low-rankness, group sparsity, or simultaneous low-rankness and group sparsity. The convex programs (1.12) and (1.13) can be solved via semidefinite programs. Considering the computational complexity and the scalability of the convex program, it motivates to develop nonconvex algorithms that enjoy lower computational complexity. This monograph provides a comprehensive discussion on various algorithms for solving structured signal estimation problems for low-overhead communications from both computational and theoretical points of view.

### 1.3 Outline

This monograph aims at providing an introduction to key models, algorithms, and theoretical results of structured signal processing in achieving low-overhead communications in IoT networks. Specifically, the content is organized according to four clearly defined categories, i.e., the sparse linear model, blind demixing, sparse blind demixing, and shuffled linear regression, which are summarized in Fig. 1.3. Key problems of low-overhead communications to be considered are also shown in Fig. 1.3. Detailed discussions are provided on methods for solving the above mentioned structured signal processing problems, including convex relaxation approaches, nonconvex approaches, and other optimization algorithms. Moreover, a significant part in each chapter is devoted to statistical theory, demonstrating how to set the corresponding algorithms on solid theoretical foundations, which includes conic integral geometry, algebraic geometric, and Riemannian optimization theory. The proofs of some key results are also included in order to illustrate the theoretical building blocks. For the ease of reference, a brief summary of different models introduced in this monograph and corresponding theory and algorithms is provided in Table 1.1. Moreover, Chap. 6 will introduce the latest developments in learning augmented methods for structured signal processing.



**Fig. 1.3** A schematic plot showing the mathematical models and corresponding problems

**Table 1.1** Summary of different models, applications, and corresponding theory and algorithms

Model	Application	Formulation	Method (M), theory (T), and algorithm (A)
Sparse linear model	Device activity detection and channel estimation	Model: (2.3) Problem: (2.9)	M: convex relaxation (2.10) T: conic integral geometry
			M: iterative thresholding A: approximate message passing
Blind demixing	Data decoding and channel estimation	Model: (3.13) Problem: (3.20)	M: convex relaxation: (3.20) T: restricted isometry property
			M: nonconvex A: Riemannian trust-region (3.41) Wirtinger flow (regularized (3.22), regularization-free (3.28))
Sparse blind demixing	Device activity detection and data decoding and channel estimation	Model: (4.2) Problem: (4.6)	M: convex relaxation (4.8)
			M: difference-of-convex-functions approach (4.18) A: majorization minimization
			M: smoothed Riemannian optimization (4.30)
Shuffled linear regression	Device activity detection and data decoding	Model: (5.8) Problem: (5.9)	M: maximum likelihood estimation (5.9) A: algorithm based on sorting, approximation algorithm
			M: algebraic–geometric approach (5.31) T: algebraic geometry

## References

1. Al-Fuqaha, A., Guizani, M., Mohammadi, M., Aledhari, M., Ayyash, M.: Internet of things: a survey on enabling technologies, protocols, and applications. *IEEE Commun. Surv. Tutorials* **17**(4), 2347–2376 (2015)
2. Amelunxen, D., Lotz, M., McCoy, M.B., Tropp, J.A.: Living on the edge: phase transitions in convex programs with random data. *Inf. Inference* **3**(3), 224–294 (2014)
3. Arunabha, G., Zhang, J., Andrews, J.G., Muhamed, R.: *Fundamentals of LTE*. Prentice-Hall, Englewood Cliffs (2010)
4. Bajwa, W.U., Haupt, J., Sayeed, A.M., Nowak, R.: Compressed channel sensing: a new approach to estimating sparse multipath channels. *Proc. IEEE* **98**(6), 1058–1076 (2010)
5. Baraniuk, R.G.: Compressive sensing. *IEEE Signal Process. Mag.* **24**(4), 118–121 (2007)
6. Björnson, E., De Carvalho, E., Sørensen, J.H., Larsson, E.G., Popovski, P.: A random access protocol for pilot allocation in crowded massive MIMO systems. *IEEE Trans. Wirel. Commun.* **16**(4), 2220–2234 (2017)
7. Candès, E.J., Romberg, J., Tao, T.: Robust uncertainty principles: exact signal reconstruction from highly incomplete frequency information. *IEEE Trans. Inf. Theory* **52**(2), 489–509 (2006)
8. Chen, S.S., Donoho, D.L., Saunders, M.A.: Atomic decomposition by basis pursuit. *SIAM Rev.* **43**(1), 129–159 (2001)
9. Choi, J.W., Shim, B., Ding, Y., Rao, B., Kim, D.I.: Compressed sensing for wireless communications: useful tips and tricks. *IEEE Commun. Surv. Tutorials* **19**(3), 1527–1550 (2017)
10. Daubechies, I., Defrise, M., De Mol, C.: An iterative thresholding algorithm for linear inverse problems with a sparsity constraint. *Commun. Pure Appl. Math.* **57**(11), 1413–1457 (2004)
11. Dong, J., Shi, Y.: Nonconvex demixing from bilinear measurements. *IEEE Trans. Signal Process.* **66**(19), 5152–5166 (2018)
12. Dong, J., Shi, Y., Ding, Z.: Sparse blind demixing for low-latency signal recovery in massive IoT connectivity. In: *Proceedings of the IEEE International Conference on Acoustics Speech Signal Process (ICASSP)*, pp. 4764–4768. IEEE, Piscataway (2019)
13. Dong, J., Shi, Y., Ding, Z.: Sparse blind demixing for low-latency signal recovery in massive IoT connectivity. In: *Proceedings of the IEEE International Conference on Acoustics Speech Signal Process (ICASSP)*, pp. 4764–4768 (2019)
14. Dong, J., Yang, K., Shi, Y.: Blind demixing for low-latency communication. *IEEE Trans. Wirel. Commun.* **18**(2), 897–911 (2019)
15. Donoho, D.L., Maleki, A., Montanari, A.: Message-passing algorithms for compressed sensing. *Proc. Natl. Acad. Sci.* **106**(45), 18914–18919 (2009)
16. Durisi, G., Koch, T., Popovski, P.: Toward massive, ultrareliable, and low-latency wireless communication with short packets. *Proc. IEEE* **104**(9), 1711–1726 (2016)
17. Figueiredo, M.A., Nowak, R.D., Wright, S.J.: Gradient projection for sparse reconstruction: application to compressed sensing and other inverse problems. *IEEE J. Sel. Top. Sign. Proces.* **1**(4), 586–597 (2007)
18. Fu, M., Dong, J., Shi, Y.: Sparse blind demixing for low-latency wireless random access with massive connectivity. In: *Proceedings of the IEEE Vehicular Technology Conference (VTC)*, pp. 4764–4768. IEEE, Piscataway (2019)
19. Hasan, M., Hossain, E., Niyato, D.: Random access for machine-to-machine communication in LTE-advanced networks: issues and approaches. *IEEE Commun. Mag.* **51**(6), 86–93 (2013)
20. Jiang, T., Shi, Y., Zhang, J., Letaief, K.B.: Joint activity detection and channel estimation for IoT networks: phase transition and computation-estimation tradeoff. *IEEE Internet Things J.* **6**(4), 6212–6225 (2018)
21. Keller, L., Siavoshani, M.J., Fragouli, C., Argyraki, K., Diggavi, S.: Identity aware sensor networks. In: *IEEE INFOCOM*, pp. 2177–2185. IEEE, Piscataway (2009)

22. Kong, L., Khan, M.K., Wu, F., Chen, G., Zeng, P.: Millimeter-wave wireless communications for IoT-cloud supported autonomous vehicles: overview, design, and challenges. *IEEE Commun. Mag.* **55**(1), 62–68 (2017)
23. Letaief, K.B., Chen, W., Shi, Y., Zhang, J., Zhang, Y.A.: The roadmap to 6G: AI empowered wireless networks. *IEEE Commun. Mag.* **57**(8), 84–90 (2019)
24. Ling, S., Strohmer, T.: Blind deconvolution meets blind demixing: algorithms and performance bounds. *IEEE Trans. Inf. Theory* **63**(7), 4497–4520 (2017)
25. Ling, S., Strohmer, T.: Regularized gradient descent: a nonconvex recipe for fast joint blind deconvolution and demixing. *Inf. Inference J. IMA* **8**(1), 1–49 (2019)
26. Liu, L., Larsson, E.G., Yu, W., Popovski, P., Stefanovic, C., De Carvalho, E.: Sparse signal processing for grant-free massive connectivity: a future paradigm for random access protocols in the Internet of Things. *IEEE Signal Process. Mag.* **35**(5), 88–99 (2018)
27. Motlagh, N.H., Bagaa, M., Taleb, T.: UAV-based IoT platform: a crowd surveillance use case. *IEEE Commun. Mag.* **55**(2), 128–134 (2017)
28. Muthukrishnan, S., et al.: Data streams: algorithms and applications. *Found. Trends Theor. Comput. Sci.* **1**(2), 117–236 (2005)
29. Nitsche, T., Cordeiro, C., Flores, A.B., Knightly, E.W., Perahia, E., Widmer, J.: IEEE 802.11 ad: directional 60 GHz communication for multi-Gigabit-per-second Wi-Fi. *IEEE Commun. Mag.* **52**(12), 132–141 (2014)
30. Pananjady, A., Wainwright, M.J., Courtade, T.A.: Linear regression with shuffled data: statistical and computational limits of permutation recovery. *IEEE Trans. Inf. Theory* **64**(5), 3286–3300 (2018)
31. Peng, L., Song, X., Tsakiris, M.C., Choi, H., Kneip, L., Shi, Y.: Algebraically-initialized expectation maximization for header-free communication. In: *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 5182–5186. IEEE, Piscataway (2019)
32. Qin, Z., Fan, J., Liu, Y., Gao, Y., Li, G.Y.: Sparse representation for wireless communications: a compressive sensing approach. *IEEE Signal Process. Mag.* **35**(3), 40–58 (2018)
33. Schepker, H.F., Bockelmann, C., Dekorsy, A.: Exploiting sparsity in channel and data estimation for sporadic multi-user communication. In: *Proceedings of the International Symposium on Wireless Communication Systems*, pp. 1–5. VDE, Frankfurt (2013)
34. Tsakiris, M.C., Peng, L.: Homomorphic sensing. In: *Proceedings of the International Conference on Machine Learning (ICML)*, pp. 6335–6344 (2019)
35. Tsakiris, M.C., Peng, L., Conca, A., Kneip, L., Shi, Y., Choi, H.: An algebraic-geometric approach to shuffled linear regression (2018). arXiv:1810.05440
36. Wunder, G., Boche, H., Strohmer, T., Jung, P.: Sparse signal processing concepts for efficient 5G system design. *IEEE Access* **3**, 195–208 (2015)
37. Wunder, G., Jung, P., Wang, C.: Compressive random access for post-LTE systems. In: *Proceedings of the IEEE International Conference on Communications Workshops (ICC)*, pp. 539–544. IEEE, Piscataway (2014)
38. Xu, X., Rao, X., Lau, V.K.: Active user detection and channel estimation in uplink CRAN systems. In: *Proceedings of the IEEE International Conference on Communications (ICC)*, pp. 2727–2732. IEEE, Piscataway (2015)
39. Zhu, H., Giannakis, G.B.: Exploiting sparse user activity in multiuser detection. *IEEE Trans. Commun.* **59**(2), 454–465 (2010)

# Chapter 2

## Sparse Linear Model



**Abstract** In this chapter, a sparse linear model for joint activity detection and channel estimation in IoT networks is introduced. We present the problem formulation for both the cases of single-antenna and multiple-antenna BSs. A convex relaxation approach based on  $\ell_p$ -norm minimization is firstly introduced, followed by a smoothed primal-dual first-order algorithm to solve it. The theoretical analysis of the convex relaxation approach based on the conic integral geometry theory is further presented. Furthermore, an iterative threshold algorithm, namely approximate message passing (AMP), is introduced, followed by the performance analysis based on the state evolution technique. Simulation results are also presented to demonstrate the performance of different algorithms.

### 2.1 Joint Activity Detection and Channel Estimation

Under the grant-free random access scheme, the metadata contains control information, e.g., the user identifier and pilot for channel estimation, and payload data that are transmitted together to the BS [24, 25]. Due to the finite channel coherence time and the massive number of devices in the IoT network, it is impossible to assign orthogonal pilot sequences to different devices [26]. Moreover, incorporating a separate pilot sequence for channel estimation in the metadata would bring redundant overheads. Considering the typical small payload size of IoT applications, it is of vital importance to reduce the overheads.

One unique characteristic of massive IoT connectivity is the sporadic data traffic, i.e., only a part of devices in the network are active at each time slot. This is because IoT devices are often designed to be in the sleep mode most of the time to conserve energy and are only triggered by external events to transmit data [26]. Exploiting this fact, a *sparse linear model* can well capture the problem of massive connectivity, which enables joint device activity detection and channel estimation [26]. This structured model describes an underdetermined linear system with more

unknown variables than equations. Considering an IoT network consisting of one BS and  $N$  devices, a sparse linear model can be established as

$$\mathbf{y} = \mathbf{A}\mathbf{x}, \quad (2.1)$$

where  $\mathbf{y} \in \mathbb{R}^L$  is the received signal at the BS,  $\mathbf{A} \in \mathbb{R}^{L \times N}$  is the set of pilot sequences, and  $\mathbf{x} \in \mathbb{R}^N$  is a sparse vector containing the information of the activity states of devices and the channel states. Particularly,  $\mathbf{A}$  is chosen from a set of non-orthogonal preamble sequences. The corresponding element of  $\mathbf{x}$  is 0 for an inactive device; otherwise, it denotes the channel coefficient for an active device. Therefore, by recovering the sparse vector  $\mathbf{x}$  from the observation  $\mathbf{y}$ , device activity detection and channel estimation can be simultaneously achieved.

To solve the estimation problem (2.1), the work [36] proposed a modified Bayesian compressed sensing algorithm. To further improve the performance of the algorithm, the works [31, 34, 35] developed the AMP algorithm with the performance analysis based on the fading coefficients and statistical channel information. The rigorous analysis has been recently investigated in a line of literatures. It shows that a state evolution analysis [3, 24, 25] of the AMP algorithm enables to characterize the false alarm and miss detection probabilities for activity detection. Recently, the paper [21] has developed a structured group sparsity estimation approach to achieve joint device activity detection and channel estimation. To increase the convergence rate and guarantee the accuracy, a smoothing method has been proposed in [21] to solve the group sparsity estimation problem, and sharp computation and estimation trade-offs of this method were further provided [21].

In the following, we first illustrate how the sparse linear model helps to formulate the joint activity detection and channel estimation problem. Then, effective algorithms and rigorous analysis are provided, and both convex and nonconvex approaches are considered.

## 2.2 Problem Formulation

This section presents the problem formulation for joint activity detection and channel estimation, for both single-antenna and multi-antenna BSs. Assume there is one BS along with  $N$  devices in an IoT network. Due to sporadic traffic, only a part of the devices are active in each time slot. For each coherent block in a synchronized wireless system with block fading, the indicator function that implies the device activity is defined as:

$$\alpha_i = \begin{cases} 1, & \text{if device } i \text{ is active,} \\ 0, & \text{otherwise,} \end{cases} \quad \forall i \in \{1, \dots, N\}. \quad (2.2)$$

Hence,  $\mathcal{S} = \{i \mid \alpha_i = 1, i = 1, \dots, N\}$  denotes the set of active devices within a coherence block, with the number of active devices being  $|\mathcal{S}|$ .

### 2.2.1 Single-Antenna Scenario

Assume that the BS is equipped with a single antenna, and denote the channel coefficient from device  $i$  to the BS as  $h_i$  for  $i = 1, \dots, N$ . Define  $\mathbf{q}_i \in \mathbb{C}^L$  as the pilot sequence transmitted from device  $i$ , where  $L$  is the length of the pilot sequence which is much smaller than the number of devices, i.e.,  $L \ll N$ , due to the finite coherence time. The received signal over  $L$  symbols at the BS is given by

$$\mathbf{y} = \sum_{i=1}^N \alpha_i h_i \mathbf{q}_i + \mathbf{n} = \sum_{i \in \mathcal{S}} h_i \mathbf{q}_i + \mathbf{n} = \mathbf{A} \mathbf{x} + \mathbf{n}, \quad (2.3)$$

where  $\mathbf{y} = [y_1, \dots, y_L]^\top \in \mathbb{C}^L$  is the received signal,  $q_{i,\ell} \sim \mathcal{CN}(0, 1) \in \mathbb{C}$  for  $i = 1, \dots, N, \ell = 1, \dots, L$  are pilot symbols, and  $\mathbf{n} \in \mathbb{C}^L \sim \mathcal{CN}(\mathbf{0}, \sigma^2 \mathbf{I})$  is the additive white Gaussian noise. Moreover,

$$\mathbf{A} = [\mathbf{q}_1, \dots, \mathbf{q}_N] \in \mathbb{C}^{L \times N}$$

is the collection of pilot sequences of all the devices, and

$$\mathbf{x} = [x_1, \dots, x_N]^\top \in \mathbb{C}^N$$

with  $x_i = \alpha_i h_i$  for  $i = 1, \dots, N$  contain device activity indicators and channel states. Here, Eq. (2.3) gives a sparse linear model. The task for the BS is to jointly detect the active devices and estimate the channel coefficients by recovering  $\mathbf{x}$  from the observation  $\mathbf{y}$ , which can then be used for data detection. An example of the sparse linear model is illustrated in Example 2.1.

*Example 2.1* Consider a network with two devices and one BS equipped with a single antenna. Assume that the pilot sequences  $\mathbf{A}$  are predefined as:

$$\mathbf{A} = \begin{bmatrix} 1 & -\sqrt{3} \\ 1 & \sqrt{3} \\ 2 & 0 \end{bmatrix}. \quad (2.4)$$

Assuming that the second device is inactive and the channel state of the active device (i.e., device 1) is  $h_1 = 1$ , we have

$$\mathbf{x} = \begin{bmatrix} 1 \cdot h_1 \\ 0 \cdot h_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}. \quad (2.5)$$

It yields a sparse linear model:

$$\mathbf{y} = \mathbf{A} \mathbf{x} = \begin{bmatrix} 1 & -\sqrt{3} \\ 1 & \sqrt{3} \\ 2 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 2 \end{bmatrix}. \quad (2.6)$$

### 2.2.2 Multiple-Antenna Scenario

Inspired by the successful application of the sparse linear model for device activity detection, multi-antenna technologies have been applied to enhance the detection performance. It generalizes the sparse signal-recovery problem to the case with a group of measurement vectors. These signal vectors are assumed to be sparse and share a common support, corresponding to the active devices. This induces a group sparsity structure, which helps to improve the performance of device activity detection and channel estimation.

Assume the BS is equipped with  $M$  antennas. The  $\ell$ -th received signal at the BS is denoted as  $\mathbf{y}(\ell) \in \mathbb{C}^M$  for all  $\ell = 1, \dots, L$ , which is given by

$$\mathbf{y}(\ell) = \sum_{i=1}^N \mathbf{h}_i \alpha_i q_i(\ell) + \mathbf{n}(\ell) = \sum_{i \in \mathcal{S}} \mathbf{h}_i q_i(\ell) + \mathbf{n}(\ell), \quad (2.7)$$

for all  $\ell = 1, \dots, L$ . Here,  $q_i(\ell) \sim \mathcal{CN}(0, 1) \in \mathbb{C}$  is the pilot symbol transmitted from device  $i$  at time slot  $\ell$ ,  $\mathbf{h}_i \in \mathbb{C}^M$  denotes the channel vector from device  $i$  to the BS antennas, and  $\mathbf{n}(\ell) \in \mathbb{C}^M \sim \mathcal{CN}(\mathbf{0}, \sigma^2 \mathbf{I})$  is the independent additive white Gaussian noise.

By accumulating the signal vectors over  $L$  time slots, we get the aggregated received signal matrix

$$\mathbf{Y} = [\mathbf{y}(1), \dots, \mathbf{y}(L)]^\top \in \mathbb{C}^{L \times M},$$

the channel matrix

$$\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_N]^\top \in \mathbb{C}^{N \times M},$$

the additive noise matrix

$$\mathbf{N} = [\mathbf{n}(1), \dots, \mathbf{n}(L)] \in \mathbb{C}^{L \times M},$$

and pilot matrix

$$\mathbf{Q} = [\mathbf{q}(1), \dots, \mathbf{q}(L)]^\top \in \mathbb{C}^{L \times N},$$

where  $\mathbf{q}(\ell) = [q_1(\ell), \dots, q_N(\ell)]^\top \in \mathbb{C}^N$ . Thus, (2.7) can be rewritten as

$$\mathbf{Y} = \mathbf{Q}\mathbf{\Theta} + \mathbf{N}, \quad (2.8)$$

where the matrix  $\mathbf{\Theta}$  is given by  $\mathbf{\Theta} = \mathbf{D}\mathbf{H} \in \mathbb{C}^{N \times M}$  with  $\mathbf{D} = \text{diag}(\alpha_1, \dots, \alpha_n) \in \mathbb{R}^{N \times N}$  being the diagonal activity matrix. Hence, the matrix  $\mathbf{\Theta}$  endows with a group sparse structure. The task for the multi-antenna BS is to detect the active devices and estimate the channel matrix by recovering  $\mathbf{\Theta}$  from the observation  $\mathbf{Y}$ .

## 2.3 Convex Relaxation Approach

### 2.3.1 Method: $\ell_p$ -Norm Minimization

In this section, we present a convex relaxation approach for joint device activity detection and channel estimation. For the noiseless case, the straightforward idea of recovering a sparse signal  $\mathbf{x}$  of which most elements are zeros is to find the sparsest signal among all those that generate the observation  $\mathbf{y} = \mathbf{A}\mathbf{x}$ . It results in the following problem:

$$\begin{aligned} & \underset{\mathbf{x} \in \mathbb{C}^N}{\text{minimize}} && \|\mathbf{x}\|_0 \\ & \text{subject to} && \mathbf{y} = \mathbf{A}\mathbf{x}, \end{aligned} \quad (2.9)$$

where the  $\ell_0$ -norm describes the number of nonzeros in  $\mathbf{x}$ . However, the problem is NP-hard due to the inevitable combinatorial search [27]. A convex relaxation approach can be applied by replacing the  $\ell_0$ -norm by the  $\ell_1$ -norm. The method of  $\ell_1$ -norm minimization [8, 10, 14], which exploits  $\ell_1$ -norm to induce the sparsity of the signal  $\mathbf{x}$ , is a well-established approach to solve compressed sensing problems.

The optimization problem that recovers  $\mathbf{x}$  from the noisy observation  $\mathbf{y}$  in (2.3) is formulated as

$$\begin{aligned} & \underset{\mathbf{x} \in \mathbb{C}^N}{\text{minimize}} && \|\mathbf{x}\|_1 \\ & \text{subject to} && \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2 \leq \epsilon, \end{aligned} \quad (2.10)$$

where the parameter  $\epsilon > 0$  is a prior threshold such that  $\mathbf{n}$  in (2.3) obeys  $\|\mathbf{n}\|_2 \leq \epsilon$ . Given the estimate vector  $\hat{\mathbf{x}}$ , the activity matrix can be recovered as

$$\hat{\mathbf{C}} = \text{diag}(\hat{a}_1, \dots, \hat{a}_n),$$

where  $\hat{a}_i = 1$  if  $|\hat{x}_i| \geq \gamma_0$  for a small enough threshold  $\gamma_0$  ( $\gamma_0 \geq 0$ ); otherwise,  $\hat{a}_i = 0$ . The estimated channel vector for the active devices is thus given by  $\hat{\mathbf{h}}$  with its  $i$ -th element as  $\hat{h}_i = \hat{x}_i$ , where  $i \in \{j | \hat{a}_j = 1\}$ .

Likewise, the optimization problem in the multiple-antenna scenario can be presented as

$$\begin{aligned} & \underset{\Theta \in \mathbb{C}^{N \times M}}{\text{minimize}} && \mathcal{R}(\Theta) := \sum_{i=1}^N \|\theta^i\|_2 \\ & \text{subject to} && \|\mathbf{Q}\Theta - \mathbf{Y}\|_F \leq \epsilon, \end{aligned} \quad (2.11)$$

where  $\epsilon > 0$  is a priori such that  $\mathbf{N}$  in (2.8) obeys  $\|\mathbf{N}\|_F \leq \epsilon$ , and  $\theta^i$  is the  $i$ -th row of matrix  $\Theta$ . Here the function  $\mathcal{R}(\Theta)$  induces the group sparsity via mixed  $\ell_1/\ell_2$ -

norm, where  $\ell_2$ -norm  $\|\boldsymbol{\theta}^i\|_2$  bounds the magnitude of the elements of  $\boldsymbol{\theta}^i$ , while  $\ell_1$ -norm induces the sparsity of  $[\|\boldsymbol{\theta}^1\|_2, \dots, \|\boldsymbol{\theta}^N\|_2]$ . Given the estimated matrix  $\hat{\boldsymbol{\Theta}}$ , the activity matrix can be recovered as  $\hat{\mathbf{C}} = \text{diag}(\hat{a}_1, \dots, \hat{a}_n)$ , where  $\hat{a}_i = 1$  if  $\|\hat{\boldsymbol{\theta}}^i\|_2 \geq \gamma_0$  for a small enough threshold  $\gamma_0$  ( $\gamma_0 \geq 0$ ); otherwise,  $\hat{a}_i = 0$ . The estimated channel matrix for the active devices is thus given by  $\hat{\mathbf{H}}$  with its  $i$ -th row as  $\hat{\mathbf{h}}^i = \hat{\boldsymbol{\theta}}^i$ , where  $i \in \{j | \hat{a}_j = 1\}$ .

The convex relaxation approaches can be applied to solve problems (2.10) and (2.11) in polynomial time. However, the general interior point solvers that are typically used to deal with SDP are impractical to be applied in large-scale problems, due to the high computational complexity. It motivates to develop fast, first-order algorithms with reduced computational complexity.

### 2.3.2 Algorithm: Smoothed Primal-Dual First-Order Methods

The first-order methods, e.g., gradient methods, proximal methods [30], alternating direction method of multipliers (ADMM) algorithm [6, 33], fast ADMM algorithm [19], and Nesterov-type algorithms [4], can efficiently solve large-scale problems. Furthermore, one way to lower the computational complexity is to accelerate the convergence rate without increasing the computational cost of each iteration. It was shown in [29] that with a large data size it is possible to increase the step size in the projected gradient method, thereby achieving a faster convergence rate. The paper [18] showed that via adjusting the original iterations, it is possible to achieve faster convergence rates and maintain the estimation accuracy without greatly increasing the computational cost of each iteration. Furthermore, the acceleration of convergence rates can be achieved via smoothing techniques such as convex relaxation [9], or simply adding a smooth function to smooth the non-differentiable objective function [4, 7, 22]. However, the quantity of smoothing should be chosen thoughtfully to guarantee the performance of sporadic device activity detection in IoT networks. To address the limitations above, the paper [21] proposed a smoothed primal-dual first-order method to solve the high-dimensional group sparsity estimation problem. The sharp trade-offs between the computational cost and estimation accuracy are rigorously characterized in [21], which is further discussed in Sect. 2.3.3.2. The smoothing algorithm is first presented in the following.

By adding a smoothing function  $\frac{\mu}{2} \|\boldsymbol{\Theta}\|_F^2$ , where  $\mu$  is a positive scalar and called as the smoothing parameter, problem (2.11) is reformulated as

$$\begin{aligned} & \underset{\boldsymbol{\Theta} \in \mathbb{C}^{N \times M}}{\text{minimize}} && \tilde{\mathcal{R}}(\boldsymbol{\Theta}) := \mathcal{R}(\boldsymbol{\Theta}) + \frac{\mu}{2} \|\boldsymbol{\Theta}\|_F^2 \\ & \text{subject to} && \|\mathbf{Q}\boldsymbol{\Theta} - \mathbf{Y}\|_F \leq \epsilon. \end{aligned} \quad (2.12)$$

To facilitate algorithm design, the sparse linear observation is represented in the real domain as follows:

$$\begin{aligned}\tilde{Y} &= \tilde{Q}\tilde{\Theta}_0 + \tilde{N} \\ &= \begin{bmatrix} \Re\{\tilde{Q}\} & -\Im\{\tilde{Q}\} \\ \Im\{\tilde{Q}\} & \Re\{\tilde{Q}\} \end{bmatrix} \begin{bmatrix} \Re\{\Theta_0\} \\ \Im\{\Theta_0\} \end{bmatrix} + \begin{bmatrix} \Re\{N\} \\ \Im\{N\} \end{bmatrix}.\end{aligned}\quad (2.13)$$

The function  $\tilde{\mathcal{H}}(\Theta)$  with respect to the complex matrix  $\Theta \in \mathbb{C}^{N \times M}$  can be further converted to the function  $\tilde{\mathcal{H}}_G(\tilde{\Theta})$  with respect to the real matrix  $\tilde{\Theta} \in \mathbb{R}^{2N \times M}$  as

$$\tilde{\mathcal{H}}_G(\tilde{\Theta}) = \sum_{i=1}^N \|\tilde{\Theta}_{\mathcal{Y}_i}\|_F + \frac{\mu}{2} \|\tilde{\Theta}_{\mathcal{Y}_i}\|_F^2. \quad (2.14)$$

Here

$$\tilde{\Theta}_{\mathcal{Y}_i} = [(\tilde{\theta}^i)^\top, (\tilde{\theta}^{i+N})^\top]^\top$$

is the row submatrix of  $\tilde{\Theta}$  consisting of the rows indexed by  $\mathcal{Y}_i = \{i, i + N\}$ . Hence, problem (2.12) can be approximated as the following structured group sparse estimation problem

$$\begin{aligned}\text{minimize } & \tilde{\mathcal{H}}_G(\tilde{\Theta}) \\ & \tilde{\Theta} \in \mathbb{R}^{2N \times M} \\ \text{subject to } & \|\tilde{Q}\tilde{\Theta} - \tilde{Y}\|_F \leq \epsilon,\end{aligned}\quad (2.15)$$

where  $\tilde{Q} \in \mathbb{R}^{2L \times 2N} \sim \mathcal{N}(\mathbf{0}, 0.5I)$  is designed as a Gaussian random matrix. Due to the indifferentiability of problem (2.15), it yields a slow coverage rate when solved by the subgradient method. Fortunately, the dual formulation of problem (2.15) leverages the benefits from smoothing techniques. In particular, the smoothed dual problem can be transferred to an unconstrained problem with the composite objective function consisting of a convex, nonsmooth function and a convex, smooth function. The dual problem of (2.15) is represented as

$$\begin{aligned}\text{maximize } & \mathcal{D}(\mathbf{Z}, t) := \inf_{\tilde{\Theta}} \left\{ \tilde{\mathcal{H}}(\tilde{\Theta}) - \langle \mathbf{Z}, \tilde{Q}\tilde{\Theta} - \tilde{Y} \rangle - t\epsilon \right\} \\ \text{subject to } & \|\mathbf{Z}\|_F \leq t,\end{aligned}$$

where  $\mathbf{Z} \in \mathbb{R}^{2N \times M}$  and  $t > 0$ . Since the parameter  $\epsilon \geq 0$  eludes the dual variable  $t$ , it yields the unconstrained problem:

$$\text{minimize } \mathcal{D}(\mathbf{Z}) := -\inf_{\tilde{\Theta}} \left\{ \tilde{\mathcal{H}}(\tilde{\Theta}) - \langle \mathbf{Z}, \tilde{Q}\tilde{\Theta} - \tilde{Y} \rangle - \epsilon \|\mathbf{Z}\|_F \right\}. \quad (2.16)$$

The dual objective function  $\tilde{\mathcal{D}}(\mathbf{Z})$  can be further represented as a composite function

$$\mathcal{D}(\mathbf{Z}) = \tilde{\mathcal{D}}(\mathbf{Z}) + \mathcal{H}(\mathbf{Z}), \quad (2.17)$$

where

$$\tilde{\mathcal{D}}(\mathbf{Z}) = -\inf_{\tilde{\Theta}} \left\{ \tilde{\mathcal{R}}(\tilde{\Theta}) - \langle \mathbf{Z}, \tilde{\mathbf{Q}}\tilde{\Theta} \rangle \right\} - \langle \mathbf{Z}, \tilde{\mathbf{Y}} \rangle$$

and  $\mathcal{H}(\mathbf{Z}) = \epsilon \|\mathbf{Z}\|_F$ . The gradient of the function  $\tilde{\mathcal{D}}(\mathbf{Z})$  is

$$\nabla \tilde{\mathcal{D}}(\mathbf{Z}) = -\tilde{\mathbf{Y}} + \tilde{\mathbf{Q}}\tilde{\Theta}_{\mathbf{Z}},$$

where

$$\tilde{\Theta}_{\mathbf{Z}} := \arg \min_{\tilde{\Theta}} \left\{ \tilde{\mathcal{R}}(\tilde{\Theta}) - \langle \mathbf{Z}, \tilde{\mathbf{Q}}\tilde{\Theta} \rangle \right\}. \quad (2.18)$$

In addition,  $\nabla \tilde{\mathcal{D}}(\mathbf{Z})$  is Lipschitz continuous with the Lipschitz constant being bounded by  $L_s := \mu^{-1} \|\tilde{\mathbf{Q}}\|_2^2$ . The composite form in (2.17) can be solved by a set of first-order approaches [4]. These methods are exceptionally sensitive to the smoothing parameter  $\mu$ , which means that a larger value of the smoothing parameter  $\mu$  induces a faster convergence rate. For instance, the Lan, Lu, and Monteiro's algorithm [23] is illustrated in Algorithm 2.1 as a typical example to show the benefits of smoothing.

---

#### Algorithm 2.1: Lan, Lu, and Monteiro's algorithm

---

Input : Pilot matrix  $\tilde{\mathbf{Q}} \in \mathbb{R}^{2L \times 2N}$ , Lipschitz constant  $L_s := \mu^{-1} \|\tilde{\mathbf{Q}}\|_2^2$ , observation matrix  $\tilde{\mathbf{Y}} \in \mathbb{R}^{2L \times M}$ , and parameter  $\epsilon$ .

- 1  $\mathbf{Z}_0 \leftarrow \mathbf{0}, \bar{\mathbf{Z}}_0 \leftarrow \mathbf{Z}_0, t_0 \leftarrow 1$
- 2 for  $k = 0, 1, 2, \dots$  do
- 3      $\mathbf{B}_k \leftarrow (1 - t_k)\mathbf{Z}_k + t_k \bar{\mathbf{Z}}_k$
- 4      $\tilde{\Theta}_k \leftarrow \mu^{-1} \text{SoftThreshold}(\tilde{\mathbf{Q}}^T \mathbf{B}_k, 1)$
- 5      $\bar{\mathbf{Z}}_{k+1} \leftarrow \text{Shrink}(\bar{\mathbf{Z}}_k - (\tilde{\mathbf{Q}}\tilde{\Theta}_k - \tilde{\mathbf{Y}})/L_s/t_k, \epsilon/L_s/t_k)$
- 6      $\mathbf{Z}_{k+1} \leftarrow \text{Shrink}(\mathbf{B}_k - (\tilde{\mathbf{Q}}\tilde{\Theta}_k - \tilde{\mathbf{Y}})/L_s, \epsilon/t_k)$
- 7      $t_{k+1} \leftarrow 2/(1 + (1 + 4/t_k^2)^{1/2})$
- 8 end

---

In Algorithm 2.1, Line 4 is the solution to (2.18), Lines 5 and 6 are the solutions to the following gradient mapping, respectively,

$$\begin{aligned}\bar{\mathbf{Z}}_{k+1} &\leftarrow \arg \min_{\mathbf{Z} \in \mathbb{R}^{2N \times M}} \left\{ \langle \nabla \tilde{\mathcal{D}}(\mathbf{Z}), \mathbf{Z} \rangle + \frac{1}{2} t_k L_s \|\mathbf{Z} - \bar{\mathbf{Z}}_k\|_F + \mathcal{H}(\mathbf{Z}) \right\}, \\ \mathbf{Z}_{k+1} &\leftarrow \arg \min_{\mathbf{Z} \in \mathbb{R}^{2N \times M}} \left\{ \langle \nabla \tilde{\mathcal{D}}(\mathbf{Z}), \mathbf{Z} \rangle + \frac{1}{2} L_s \|\mathbf{Z} - \mathbf{B}_k\|_F + \mathcal{H}(\mathbf{Z}) \right\}.\end{aligned}$$

Denote  $\mathbf{Z}^*$  as an optimal solution for (2.16), then the convergence behavior of Algorithm 2.1 is demonstrated as [4]

$$\mathcal{D}(\mathbf{Z}_{k+1}) - \mathcal{D}(\mathbf{Z}^*) \leq \frac{2\|\tilde{\mathbf{Q}}\|_2^2 \|\mathbf{Z}_0 - \mathbf{Z}^*\|_F^2}{\mu k^2}. \quad (2.19)$$

Based on (2.19), the number of iterations

$$\left\lceil \sqrt{2\|\tilde{\mathbf{Q}}\|_2^2 / (\mu \epsilon_0) \|\mathbf{Z}_0 - \mathbf{Z}^*\|_F} \right\rceil$$

is required to reach the accuracy of  $\epsilon_0$ . That is, a larger  $\mu$  would lead to a faster convergence rate.

### 2.3.3 Analysis: Conic Integral Geometry

The paper [21] discussed the trade-off between the estimation accuracy and computational cost in terms of the smoothing method described in Sect. 2.3.2, which is achieved by characterizing the convergence rate in terms of the smoothing parameter, pilot sequence length, and estimation accuracy. The analysis is based on the theory of conic integral geometry [1, 28, 32]. Prior to focusing on conic integral geometry for the sparse linear model, you may refer to Sect. 8.1 to have a basic overview of conic integral geometry.

#### 2.3.3.1 Conic Integral Geometry for the Sparse Linear Model

Considering the smoothing method illustrated in Sect. 2.3.2, it is critical to find a proper smoothing parameter  $\mu$ , which can be achieved by analyzing the trade-off between the estimation accuracy and computational cost of the convex optimization problem (2.11). *Conic integral geometry* theory turns out to be a promising and powerful tool to predict phase transitions (including the location and width of the transition region) for random cone programs in the real field case [1, 28, 32]. Based on the conic integral geometry, the paper [21] proposed to approximate the

original complex estimation problem (2.11) by a real estimation problem, followed by analyzing on the performance of the proposed smoothing method concerning the smoothing parameter  $\mu$ .

In the noiseless scenario, we consider the following approximated problem:

$$\begin{aligned} & \underset{\tilde{\Theta} \in \mathbb{R}^{2N \times M}}{\text{minimize}} \mathcal{R}_G(\tilde{\Theta}) \\ & \text{subject to } \tilde{Y} = \bar{Q}\tilde{\Theta}, \end{aligned} \quad (2.20)$$

where

$$\mathcal{R}_G(\tilde{\Theta}) = \sum_{i=1}^N \|\tilde{\Theta} \gamma_i\|_F$$

and  $\tilde{\Theta}$ ,  $\bar{Q}$ , and  $\tilde{Y}$  are defined in (2.15). To deal with problem (2.20), several definitions and facts in convex analysis [1] are introduced first.

**Definition 2.1 (Descent Cone)** The descent cone  $\mathcal{D}(\mathcal{R}, \mathbf{x})$  of a proper convex function  $\mathcal{R} : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\pm\infty\}$  at point  $\mathbf{x} \in \mathbb{R}^d$  is the conic hull of the perturbations that do not increase  $\mathcal{R}$  near  $\mathbf{x}$ , i.e.,

$$\mathcal{D}(\mathcal{R}, \mathbf{x}) = \bigcup_{\tau > 0} \left\{ \mathbf{y} \in \mathbb{R}^d : \mathcal{R}(\mathbf{x} + \tau \mathbf{y}) \leq \mathcal{R}(\mathbf{x}) \right\}.$$

**Fact 2.1 (Optimality Condition)** Let  $\mathcal{R}$  be a proper convex function. Matrix  $\tilde{\Theta}_0$  is the unique optimal solution to problem (2.20) if and only if

$$\mathcal{D}(\mathcal{R}_G, \tilde{\Theta}_0) \cap \text{null}(\bar{Q}, M) = \{\mathbf{0}\},$$

where

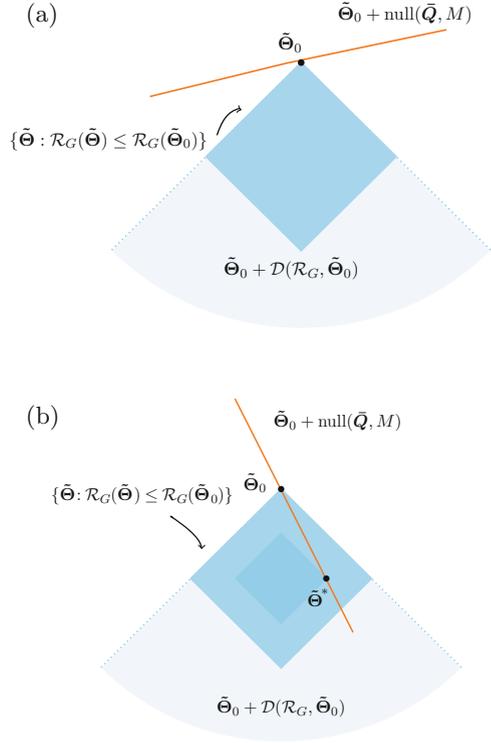
$$\text{null}(\bar{Q}, M) = \{\mathbf{Z} \in \mathbb{R}^{2N \times M} : \bar{Q}\mathbf{Z} = \mathbf{0}_{2L \times M}\}$$

denotes the null space of the operator  $\bar{Q} \in \mathbb{R}^{2L \times 2N}$ .

Figure 2.1 illustrates the geometry of the optimality condition described in Fact 2.1. Specifically, problem (2.20) succeeds to yield optimal solution if and only if the null space of  $\bar{Q}$  misses the cone of descent directions of  $\mathcal{R}_G$  at the ground truth  $\tilde{\Theta}_0$ , which is illustrated in Fig. 2.1a; otherwise, it fails to obtain the optimal solution follows  $\tilde{\Theta}^* \neq \tilde{\Theta}_0$ , which is illustrated in Fig. 2.1b.

To characterize the phase transition in two intersection cones, the concept of statistical dimension is proposed in [1] that is the generalization of the dimension of linear subspaces.

**Fig. 2.1** Optimality condition for problem (2.20).  
**(a)** Problem succeeds. **(b)** Problem fails



**Definition 2.2 (Statistical Dimension)** The statistical dimension  $\delta(C)$  of a closed convex cone  $C$  in  $\mathbb{R}^d$  is defined as:

$$\delta(C) = \mathbb{E}[\|\Pi_C(\mathbf{g})\|_2^2], \quad (2.21)$$

where  $\mathbf{g} \in \mathbb{R}^d$  is a standard normal vector, and

$$\Pi_C(\mathbf{x}) = \arg \min\{\|\mathbf{x} - \mathbf{y}\|_2 : \mathbf{y} \in C\}$$

denotes the Euclidean projection onto  $C$ .

The statistical dimension enables to measure the size of convex cones. Based on the statistical dimensions of general convex cones, the approximated conic kinematic formula can be presented as follows [2].

**Theorem 2.1 (Approximate Kinematic Formula)** Fix a tolerance  $\eta \in (0, 1)$ . Let  $C$  and  $K$  be convex cones in  $\mathbb{R}^d$ , but one of them is not a subspace. Draw a random orthogonal basis  $U$ . Then

$$\begin{aligned} \delta(C) + \delta(K) \leq d - a_\eta \sqrt{d} &\implies \mathbb{P}\{C \cap UK \neq \{\mathbf{0}\}\} \leq \eta \\ \delta(C) + \delta(K) \geq d + a_\eta \sqrt{d} &\implies \mathbb{P}\{C \cap UK \neq \{\mathbf{0}\}\} \geq 1 - \eta, \end{aligned}$$

where  $a_\eta := \sqrt{8 \log(4/\eta)}$ .

Theorem 2.1 captures a phase transition on whether the two randomly rotated cones share a ray. In particular, the two randomly rotated cones share a ray with high probability, if the total statistical dimension of the two cones exceeds the ambient dimension  $d$ ; otherwise, they fail to share a ray.

### 2.3.3.2 Computation and Estimation Trade-Offs

For the smoothing method introduced in Sect. 2.3.2, a trade-off between the computational cost and estimation accuracy is characterized based on the general results in Theorem 2.1. This trade-off plays a vital role in massive connectivity with a finite time budget and a modest requirement on estimation accuracy.

The basis of the trade-off is introduced in the sequel. From the geometric point of view, the smoothing term in  $\tilde{\mathcal{R}}(\tilde{\Theta})$  (with  $\mu > 0$ ) increases the sublevel set of  $\mathcal{R}(\Theta)$ , which derives a problem that can be solved via computationally efficient algorithms with an accelerated convergence rate. However, this geometric modification leads to a loss of the estimation accuracy. Thus, it leads to a trade-off between the computational time and estimation accuracy. The trade-off can be identified by Theorem 2.1 based on the statistical dimension of the decent cone of the smoothed regularizer in (2.20), i.e.,

$$\tilde{\mathcal{R}}_G(\tilde{\Theta}) = \mathcal{R}_G(\tilde{\Theta}) + \frac{\mu}{2} \|\tilde{\Theta}\|_F^2. \quad (2.22)$$

We begin with the basic notation used in Proposition 2.1, for some  $\tilde{\Theta} \in \mathbb{R}^{2N \times M}$  satisfying  $\tilde{\Theta} \gamma_j = \mathbf{0}$  for  $j \neq i$ , we have

$$\forall \tilde{\Theta} \gamma_i \in \mathbb{R}^{2 \times M} : \quad \|\tilde{\Theta} \gamma_i\|_F \geq \|(\tilde{\Theta}_0) \gamma_i\|_F + \langle \mathbf{Z} \gamma_i, \tilde{\Theta} \gamma_i - (\tilde{\Theta}_0) \gamma_i \rangle, \quad (2.23)$$

which implies  $\mathbf{Z} \gamma_i \in \partial \|(\tilde{\Theta}_0) \gamma_j\|_F$ . In particular, the statistical dimension  $\delta(\mathcal{D}(\tilde{\mathcal{R}}_G, \tilde{\Theta}_0))$  can be exactly computed by the following result.

**Proposition 2.1 (Statistical Dimension Bound for  $\tilde{\mathcal{R}}_G$ )** *Let  $\Theta_0 \in \mathbb{C}^{N \times M}$  be with  $K$  nonzero rows, and define the normalized sparsity as  $\rho := K/N$ . An upper bound of the statistical dimension of the descent cone of  $\tilde{\mathcal{R}}_G$  at*

$$\tilde{\Theta}_0 = [(\Re\{\Theta_0\})^T, (\Im\{\Theta_0\})^T]^T \in \mathbb{R}^{2N \times M}$$

is given by

$$\frac{\delta(\mathcal{D}(\tilde{\mathcal{R}}_G; \tilde{\Theta}_0))}{N} \leq \inf_{\tau \geq 0} \left\{ \rho(2M + \tau^2(1 + 2\mu\bar{a} + \mu^2\bar{b})) + (1 - \rho) \frac{2^{1-M}}{\Gamma(M)} \int_{\tau}^{\infty} (u - \tau)^2 u^{2M-1} e^{-\frac{u^2}{2}} du \right\}, \quad (2.24)$$

where  $\Gamma(\cdot)$  denotes the Gamma function. The unique optimum  $\tau^*$  which minimizes the right-hand side of (2.24) is the solution of

$$\frac{2^{1-M}}{\Gamma(M)} \int_{\tau}^{\infty} \left( \frac{u}{\tau} - 1 \right) u^{2M-1} e^{-\frac{u^2}{2}} du = \frac{\rho(1 + 2\mu\bar{a} + \mu^2\bar{b})}{1 - \rho}, \quad (2.25)$$

where  $\bar{a} = \frac{1}{S} \sum_{i=1}^S \|(\tilde{\Theta}_0)_{\mathcal{Y}_i}\|_F$ ,  $\bar{b} = \frac{1}{S} \sum_{i=1}^S \|(\tilde{\Theta}_0)_{\mathcal{Y}_i}\|_F^2$ .

**Proof** Please refer to Sect. 8.2 for details.

Although the convergence rate of proposed smoothing algorithm, i.e., Algorithm 2.1, can be accelerated by increasing the smoothing parameter, Proposition 2.1 shows that a larger smoothing parameter leads to a larger statistical dimension  $\delta(\mathcal{D}(\tilde{\mathcal{R}}_G, \tilde{\Theta}_0))$  since the bound in (2.24) increases with  $\mu$ .

### 2.3.3.3 Simulation Results

Proposition 2.1 is verified in Fig. 2.2 with the BS equipped with 2 antennas, the total number of devices being 100, and the channel matrix and pilot matrix generated as

$$\mathbf{H} \sim \mathcal{CN}(\mathbf{0}, \mathbf{I}) \quad \text{and} \quad \mathbf{Q} \sim \mathcal{CN}(\mathbf{0}, \mathbf{I}),$$

respectively. The recovery is considered to be successful if  $\|\hat{\Theta} - \Theta_0\|_F \leq 10^{-5}$ . The number of active devices is fixed as  $|\mathcal{S}| = 10$ . Figure 2.2 shows the impact on the exact recovery when changing the smoothing parameter  $\mu$ . It shows that a larger smoothing parameter will induce a larger statistical dimension of the descent cone of  $\tilde{\mathcal{R}}(\Theta)$ . In other words, longer pilot sequences are required for exact signal recovery.

The effectiveness of the smoothing method illustrated in Algorithm 2.1 is evaluated under the scenario where the base station is equipped with 10 antennas, and the total number of devices is set to be 2000. The number of active devices is fixed as  $|\mathcal{S}| = 100$ . Considering problem (2.15), the channel matrix follows  $\mathbf{H} \sim \mathcal{CN}(\mathbf{0}, \mathbf{I})$ , the pilot matrix follows  $\mathbf{Q} \sim \mathcal{CN}(\mathbf{0}, \mathbf{I})$  and the additive noise matrix follows  $\mathbf{N} \sim \mathcal{CN}(\mathbf{0}, 0.01\mathbf{I})$ . Figure 2.3 demonstrates the convergence rate of Algorithm 2.1 under different smoothing parameters with a fixed pilot sequence length  $L = 500$ . It shows that increasing the smooth parameter enables to accelerate the convergence rate significantly.

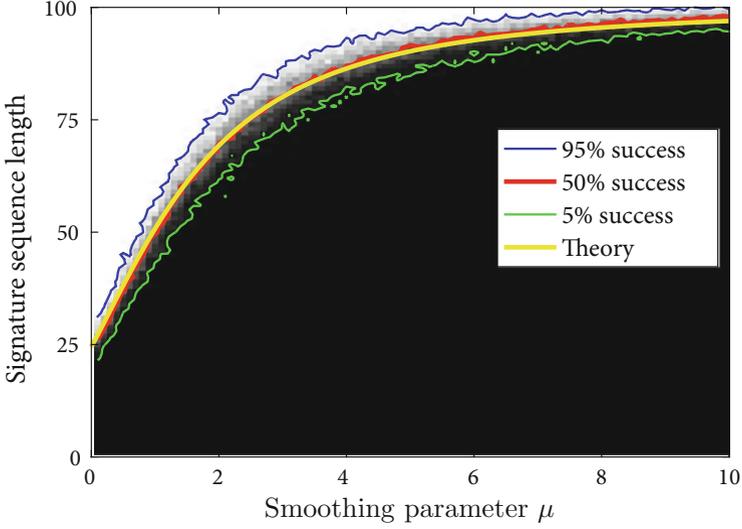


Fig. 2.2 Phase transitions in massive device connectivity via smoothing

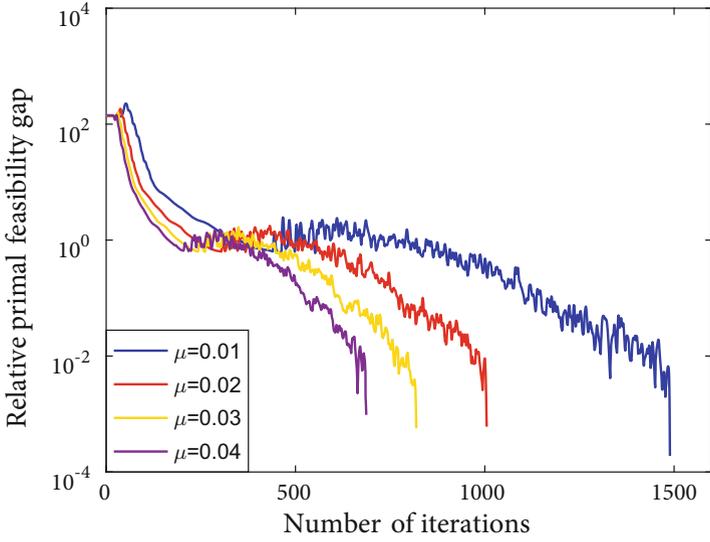
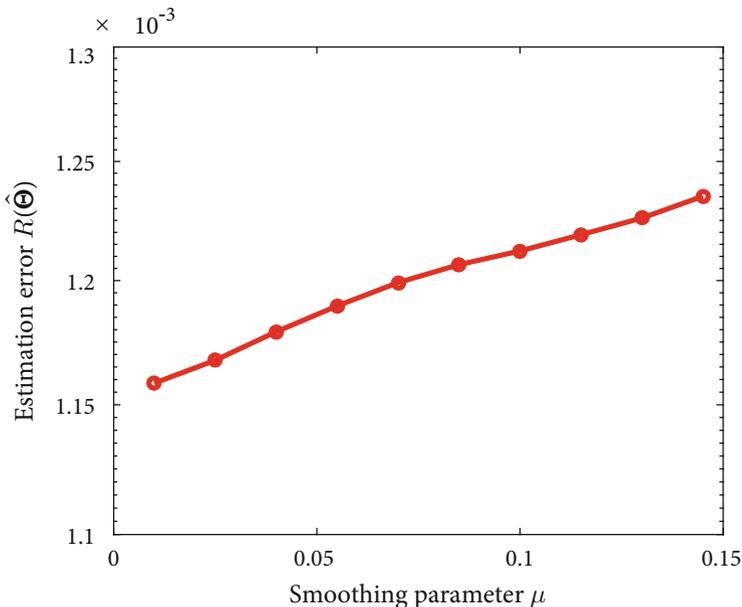


Fig. 2.3 Convergence rate of Algorithm 2.1

Furthermore, with a fixed pilot sequence length  $L = 500$ , problem (2.15) is solved by Algorithm 2.1 under different smoothing parameters  $\mu$ . Algorithm 2.1 stops when

$$\left| \|\tilde{\mathbf{Q}}\tilde{\boldsymbol{\Theta}} - \tilde{\mathbf{Y}}\|_F - \epsilon \right| / \epsilon \leq 10^{-3},$$



**Fig. 2.4** Estimation error versus smoothing parameter  $\mu$

where the parameter  $\epsilon$  is given by

$$\epsilon = \sigma \sqrt{2LM - \delta(\mathcal{D}(\tilde{\mathcal{R}}_G, \tilde{\Theta}_0))}.$$

The simulation result illustrated in Fig. 2.4 is obtained by averaging over 300 channel realizations. It shows that the average squared estimation error becomes large as the smoothing parameter  $\mu$  increases. This can be justified by Proposition 2.1 that the increase of smoothing parameter results in the increase of statistical dimension  $\delta(\mathcal{D}(\tilde{\mathcal{R}}_G, \tilde{\Theta}_0))$ .

## 2.4 Iterative Thresholding Algorithm

Despite attractive theoretical guarantees for the sparse linear model, convex relaxation methods that are solved via a second-order cone program (SOCP) fail in the high-dimensional data setting due to the high computational cost. One way to improve the computational efficiency is the smoothed primal-dual first-order method introduced in the previous section. Another line of literatures that aim to reduce the computational complexity for solving the sparse linear model estimation problem focus on iterative thresholding algorithms [13]. Unfortunately, such fast

iterative thresholding algorithms suffer from worse sparsity-undersampling trade-offs than convex optimization [15], and the sparsity-undersampling trade-off is precisely controlled by the sampling ratio  $\delta = L/N$  and sparsity ratio  $\rho = |\mathcal{S}|/N$  with  $L$ ,  $N$ ,  $\mathcal{S}$  defined in the model (2.3). To resolve this issue, approximate message passing [12, 15, 24] has been proposed for sparse recovery.

### 2.4.1 Algorithm: Approximate Message Passing

The approximate message passing (AMP) algorithm was proposed and developed in a line of literatures [12, 15, 24], which is an efficient iterative thresholding method for solving the linear model estimation problem (2.3). For simplicity, we take the single-antenna scenario for example.

The goal of the AMP algorithm is to evaluate an estimator  $\hat{\mathbf{x}}(\mathbf{y})$  from the observation  $\mathbf{y}$  (2.3) that minimizes the mean-squared error (MSE)

$$\text{MSE} = \mathbb{E}_{\mathbf{x}, \mathbf{y}} \|\hat{\mathbf{x}}(\mathbf{y}) - \mathbf{x}\|_2^2, \quad (2.26)$$

where the signals  $x_i = \alpha_i h_i$  for  $i = 1, \dots, n$  are assumed to follow a Bernoulli-Gaussian distribution. Starting from  $\mathbf{x}^0 = \mathbf{0}$  and  $\mathbf{r}^0 = \mathbf{y}$ , the iterative update of the AMP algorithm at the  $t$ -th iteration is given by Donoho et al. [15]

$$x_i^{t+1} = \eta_{t,i}((\mathbf{r}^t)^H \mathbf{a}_i + x_i^t), \quad (2.27a)$$

$$\mathbf{r}^{t+1} = \mathbf{y} - \mathbf{A}\mathbf{x}^{t+1} + \frac{N}{L} \mathbf{r}^t \sum_{n=1}^N \frac{\eta'_{t,i}((\mathbf{r}^t)^H \mathbf{a}_i + x_i^t)}{N}, \quad (2.27b)$$

where  $\mathbf{x}^t = [x_1^t, \dots, x_N^t]^\top \in \mathbb{C}^N$  is the estimate of  $\mathbf{x}$  at the  $t$ -th iteration,  $\mathbf{r}^t = [r_1^t, \dots, r_L^t]^\top \in \mathbb{C}^L$  denotes the residual,

$$\eta_{t,i}(\cdot) : \mathbb{C} \rightarrow \mathbb{C}$$

is the denoiser which facilitates to induce the sparsity, and  $\eta'_{t,i}(\cdot)$  is the first-order derivative of  $\eta_{t,i}(\cdot)$ . The performance of the AMP algorithm highly depends on the design of the denoiser  $\eta_{t,i}(\cdot)$ , which will be discussed in the sequel.

## 2.4.2 Analysis: State Evolution

### 2.4.2.1 State Evolution

In order to precisely capture the dynamic property of the AMP algorithm, thereby facilitating the design of the denoiser  $\eta_{t,i}(\cdot)$ , a state evolution formalism was first proposed in the paper [15]. In this formalism, the MSE (2.26) is a state variable and its variation from iteration to iteration can be represented by a plain iterative function, i.e.,  $\tau_t$ .

Define a set of random variables  $\hat{X}_i^t$  at the  $t$ -th iteration of the AMP algorithm as

$$\hat{X}_i^t = X_i + \tau_t V_i, \quad i = 1, \dots, n, \quad (2.28)$$

where the distributions of  $X_i$ 's are characterized by the random variables  $X_i$ 's, and  $V_i$  obeys the normal distribution, i.e.,  $V_i \in \mathcal{C}\mathcal{N}(0, 1)$ . In addition,  $V_i$  is independent of  $X_i$  and  $V_j$  for  $\forall j \neq i$ , and  $\tau_t$  is the state variable represented as

$$\tau_{t+1}^2 = \frac{\sigma^2}{\xi} + \frac{N}{L} \mathbb{E} \left[ |\eta_{t,i}(X_i + \tau_t V_i) - X_i|^2 \right], \quad (2.29)$$

where the expectation is over the random variables  $V_i$ 's and  $X_i$ 's for  $i = 1, \dots, N$ .

The theoretical analysis of AMP is based on the state evolution in the asymptotic regime, when  $L$  (i.e., the length of pilot sequences),  $K$  (i.e., the average number of active devices in each time slot),  $N$  (i.e., the total number of devices)  $\rightarrow \infty$ , while their ratios converge to some positive values  $N/L \rightarrow \omega$  and

$$K/N \rightarrow \epsilon = \lim_{N \rightarrow \infty} \sum_i \epsilon_i / N$$

with  $\omega, \epsilon \in (0, \infty)$ . In massive IoT connectivity, these assumptions imply that the length of the pilot sequence, i.e.,  $L$ , is in the same order of the number of active users, i.e.,  $K$ , or total users, i.e.,  $N$ .

### 2.4.2.2 Denoiser Designs

In general, the prior distribution of  $\mathbf{x}$  is assumed to be unknown. In this case, a soft-thresholding denoiser is designed to induce sparsity for  $\mathbf{x}$ , which is given by Donoho et al. [16]:

$$\eta_{t,i}(\hat{x}_i^t) = \left( \hat{x}_i^t - \frac{\theta_i^t \hat{x}_i^t}{|\hat{x}_i^t|} \right) \mathbb{I}(|\hat{x}_i^t| > \theta_i^t), \quad (2.30)$$

where the parameter  $\theta_i^t$  is the threshold for the  $i$ -th device activity detection at the  $t$ -th iteration of the AMP algorithm. Based on the state evolution (2.29), the parameter  $\theta_i^t$  can be optimized to minimize the MSE (2.26). After the  $t$ -th iteration proceeded by the AMP algorithm with the denoiser (2.30), device  $i$  is evaluated to be active if

$$|(\mathbf{r}^t)^H \mathbf{a}_i + x_i^t| > \theta_i^t,$$

otherwise it is evaluated to be inactive.

If the prior distribution of  $\mathbf{x}$  in (2.3) is known, the minimum mean-squared error (MMSE) denoiser via the Bayesian approach can be developed for the AMP algorithm [16]. Based on the random variables defined in (2.28) and assuming the channel signal  $h_i \sim \mathcal{CN}(0, 1)$  for  $i = 1, \dots, N$ , the MMSE denoiser is given in the form of a conditional expectation [16],

$$\begin{aligned} \eta_{t,i}(\hat{x}_i^t) &= \mathbb{E}[X_i | \hat{X}_i^t = \hat{x}_i^t] \\ &= \phi_{t,i}(1 + \tau_t)^{-1} \hat{x}_i^t, \quad \forall t, i, \end{aligned} \quad (2.31)$$

where

$$\phi_{t,i} = \frac{1}{1 + \frac{1-\epsilon}{\epsilon} \exp(-(\pi_{t,i} - \psi_{t,i}))}, \quad (2.32)$$

$$\pi_{t,i} = (\tau_t^{-2} - (\tau_t^2 + 1)^{-1}) |\hat{x}_i^t|^2, \quad (2.33)$$

$$\psi_{t,i} = \log \det(1 + \tau_t^{-2}). \quad (2.34)$$

Note that the above MMSE denoiser is a nonlinear function of  $\hat{x}_i^t$  due to the functional form of  $\phi_{t,i}$ .

### 2.4.2.3 Asymptotic Performance of Device Activity Detection

Based on the soft thresholding (2.30) and MMSE denoisers (2.31), a miss detection occurs when

$$|(\mathbf{r}^t)^H \mathbf{a}_i + x_i^t| < \theta_i^t$$

with device  $i$  actually being active, while a false alarm occurs when

$$|(\mathbf{r}^t)^H \mathbf{a}_i + x_i^t| > \theta_i^t$$

with device  $i$  actually being inactive. Since the statistical distribution of the thresholding term, i.e.,  $(\mathbf{r}^t)^H \mathbf{a}_i + x_i^t$ , can be identified by  $\hat{x}_i^t$  defined in (2.28), the

probabilities of miss detection and false alarm for device  $i$  at the  $t$ -th iteration of the AMP algorithm can be given by Donoho et al. [15]

$$P_{t,i}^{\text{MD}} = \Pr(\hat{x}_i^t < \theta_i^t | \alpha_i = 1), \quad (2.35)$$

$$P_{t,i}^{\text{FA}} = \Pr(\hat{x}_i^t > \theta_i^t | \alpha_i = 0), \quad (2.36)$$

respectively. The probabilities of missed detection (2.35) and false alarm (2.36) depend on the values of  $\tau_t$ 's (2.29) which can be tracked over iterations based on the state evolution (2.29).

Considering a general multiple-antenna scenario, the theorem in [24] characterizes  $P_{t,i}^{\text{MD}}(M)$  and  $P_{t,i}^{\text{FA}}$  analytically in terms of  $\tau_t^2$  and the number of antennas  $M$ . In particular, the miss detection and false alarm probabilities of AMP algorithm with  $M$  antennas are denoted by  $P_{t,i}^{\text{MD}}(M)$  and  $P_{t,i}^{\text{FA}}(M)$ . The paper [24] demonstrates that with proper thresholds for device detection, i.e.,  $\theta_i^t$ 's, highly accurate device activity detection can be achieved in the asymptotic regime of  $M \rightarrow \infty$ :

$$\lim_{M \rightarrow \infty} P_{t,i}^{\text{MD}}(M) = \lim_{M \rightarrow \infty} P_{t,i}^{\text{FA}}(M) = 0, \quad \forall t, i. \quad (2.37)$$

It thus indicates that the AMP-based algorithm can accomplish perfect device activity detection in the massive MIMO connectivity systems.

#### 2.4.2.4 Simulation Results

To further illustrate the performance of AMP algorithm for solving the sparse linear model estimation problem, the probabilities of missed detection and false alarm versus the length of the pilot sequences,  $L$ , with different numbers of antennas at the BS, i.e.,  $M = 4, 8, \text{ or } 16$ , are illustrated in Fig. 2.5. In particular, with a given value of  $M$ , the average probabilities of missed detection and false alarm over all devices are denoted as

$$P^{\text{MD}}(M) = \sum_{n=1}^N P_{\infty,n}^{\text{MD}}(M)/N$$

and

$$P^{\text{FA}}(M) = \sum_{n=1}^N P_{\infty,n}^{\text{FA}}(M)/N,$$

respectively, where  $P_{\infty,n}^{\text{MD}}(M)$  and  $P_{\infty,n}^{\text{FA}}(M)$  are defined in (2.35) and (2.36). Figure 2.5 demonstrates that both  $P^{\text{MD}}$  and  $P^{\text{FA}}$  decrease as the pilot sequence length  $L$  increases and when  $M$  increases.

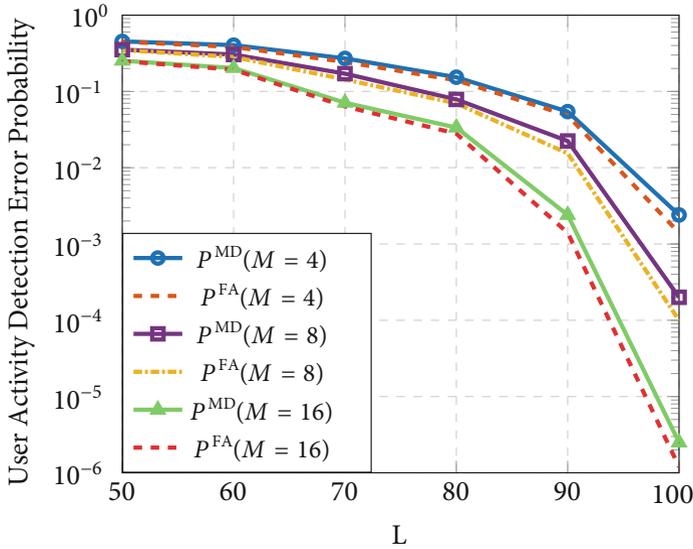


Fig. 2.5 Probabilities of missed detection and false alarm versus pilot sequence lengths

## 2.5 Summary

This chapter introduced a sparse linear model for joint device activity detection and channel estimation in grant-free random access. Such an access scheme reduces the overhead by removing dedicated channel estimation sequences in IoT networks. To solve the estimation problem, both convex relaxation approach and nonconvex approach have been investigated from the practical and theoretical points of view. Recently, there is a line of studies focusing on solving the sparse linear model via deep-learning-based methods from both empirical and theoretical points of view [5, 11, 17, 20, 37], which provide an also interesting direction for future study.

## References

1. Amelunxen, D., Lotz, M., McCoy, M.B., Tropp, J.A.: Living on the edge: phase transitions in convex programs with random data. *Inf. Inference* **3**(3), 224–294 (2014)
2. Bastug, E., Bennis, M., Debbah, M.: Living on the edge: the role of proactive caching in 5G wireless networks. *IEEE Commun. Mag.* **52**(8), 82–89 (2014)
3. Bayati, M., Montanari, A.: The dynamics of message passing on dense graphs, with applications to compressed sensing. *IEEE Trans. Inf. Theory* **57**(2), 764–785 (2011)
4. Becker, S.R., Candès, E.J., Grant, M.C.: Templates for convex cone problems with applications to sparse signal recovery. *Math. Program. Comput.* **3**(3), 165–218 (2011)
5. Borgerding, M., Schniter, P., Rangan, S.: AMP-inspired deep networks for sparse linear inverse problems. *IEEE Trans. Signal Process.* **65**(16), 4293–4308 (2017)

6. Boyd, S., Parikh, N., Chu, E., Peleato, B., Eckstein, J.: Distributed optimization and statistical learning via the alternating direction method of multipliers. *Found. Trends Mach. Learn.* **3**(1), 1–122 (2011)
7. Bruer, J.J., Tropp, J.A., Cevher, V., Becker, S.R.: Designing statistical estimators that balance sample size, risk, and computational cost. *IEEE J. Sel. Top. Sign. Process.* **9**(4), 612–624 (2015)
8. Candes, E., Tao, T.: Near optimal signal recovery from random projections: universal encoding strategies? *IEEE Trans. Inf. Theory* **52**(12), 5406–5425 (2006)
9. Chandrasekaran, V., Jordan, M.I.: Computational and statistical tradeoffs via convex relaxation. *Proc. Natl. Acad. Sci.* **110**(13), E1181–E1190 (2013)
10. Chen, S.S., Donoho, D.L., Saunders, M.A.: Atomic decomposition by basis pursuit. *SIAM Rev.* **43**(1), 129–159 (2001)
11. Chen, X., Liu, J., Wang, Z., Yin, W.: Theoretical linear convergence of unfolded ISTA and its practical weights and thresholds. In: *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 9061–9071 (2018)
12. Chen, Z., Sohrabi, F., Yu, W.: Sparse activity detection for massive connectivity. *IEEE Trans. Signal Process.* **66**(7), 1890–1904 (2018)
13. Daubechies, I., Defrise, M., De Mol, C.: An iterative thresholding algorithm for linear inverse problems with a sparsity constraint. *Commun. Pure Appl. Math.* **57**(11), 1413–1457 (2004)
14. Donoho, D.L., et al.: Compressed sensing. *IEEE Trans. Inf. Theory* **52**(4), 1289–1306 (2006)
15. Donoho, D.L., Maleki, A., Montanari, A.: Message-passing algorithms for compressed sensing. *Proc. Nat. Acad. Sci.* **106**(45), 18914–18919 (2009)
16. Donoho, D.L., Maleki, A., Montanari, A.: Message passing algorithms for compressed sensing: I. motivation and construction. In: *Proceedings of the IEEE Information Theory Workshop on Information Theory*, pp. 1–5. IEEE, Piscataway (2010)
17. Fletcher, A.K., Pandit, P., Rangan, S., Sarkar, S., Schniter, P.: Plug-in estimation in high-dimensional linear inverse problems: a rigorous analysis. In: *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 7440–7449 (2018)
18. Giryes, R., Eldar, Y.C., Bronstein, A.M., Sapiro, G.: Tradeoffs between convergence speed and reconstruction accuracy in inverse problems. *IEEE Trans. Signal Process.* **66**(7), 1676–1690 (2018)
19. Goldstein, T., O’Donoghue, B., Setzer, S., Baraniuk, R.: Fast alternating direction optimization methods. *SIAM J. Imaging Sci.* **7**(3), 1588–1623 (2014)
20. Ito, D., Takabe, S., Wadayama, T.: Trainable ISTA for sparse signal recovery. *IEEE Trans. Signal Process.* **67**(12), 3113–3125 (2019)
21. Jiang, T., Shi, Y., Zhang, J., Letaief, K.B.: Joint activity detection and channel estimation for IoT networks: phase transition and computation-estimation tradeoff. *IEEE Internet Things J.* **6**(4), 6212–6225 (2018)
22. Lai, M.J., Yin, W.: Augmented  $\ell_1$  and nuclear-norm models with a globally linearly convergent algorithm. *SIAM J. Imaging Sci.* **6**(2), 1059–1091 (2013)
23. Lan, G., Lu, Z., Monteiro, R.D.: Primal-dual first-order methods with  $\mathcal{O}(1/\epsilon)$  iteration-complexity for cone programming. *Math. Program.* **126**(1), 1–29 (2011)
24. Liu, L., Yu, W.: Massive connectivity with massive MIMO—part I: device activity detection and channel estimation. *IEEE Trans. Signal Process.* **66**(11), 2933–2946 (2018)
25. Liu, L., Yu, W.: Massive connectivity with massive MIMO—part II: achievable rate characterization. *IEEE Trans. Signal Process.* **66**(11), 2947–2959 (2018)
26. Liu, L., Larsson, E.G., Yu, W., Popovski, P., Stefanovic, C., De Carvalho, E.: Sparse signal processing for grant-free massive connectivity: a future paradigm for random access protocols in the Internet of Things. *IEEE Signal Process. Mag.* **35**(5), 88–99 (2018)
27. Muthukrishnan, S., et al.: Data streams: algorithms and applications. *Found. Trends Theor. Comput. Sci.* **1**(2), 117–236 (2005)
28. Oymak, S., Hassibi, B.: Sharp MSE bounds for proximal denoising. *Found. Comput. Math.* **16**(4), 965–1029 (2016)
29. Oymak, S., Recht, B., Soltanolkotabi, M.: Sharp time–data tradeoffs for linear inverse problems. *IEEE Trans. Inf. Theory* **64**(6), 4129–4158 (2018)

30. Parikh, N., Boyd, S.: Proximal algorithms. *Found. Trends Optim.* **1**(3), 127–239 (2014)
31. Schepker, H.F., Bockelmann, C., Dekorsy, A.: Exploiting sparsity in channel and data estimation for sporadic multi-user communication. In: *Proceedings of the International Symposium on Wireless Communication Systems*, pp. 1–5. VDE, Frankfurt (2013)
32. Schneider, R., Weil, W.: *Stochastic and Integral Geometry*. Springer, Berlin (2008)
33. Shi, Y., Zhang, J., Letaief, K.B., Bai, B., Chen, W.: Large-scale convex optimization for ultradense Cloud-RAN. *IEEE Trans. Wirel. Commun.* **22**(3), 84–91 (2015)
34. Wunder, G., Boche, H., Strohmer, T., Jung, P.: Sparse signal processing concepts for efficient 5G system design. *IEEE Access* **3**, 195–208 (2015)
35. Wunder, G., Jung, P., Wang, C.: Compressive random access for post-LTE systems. In: *Proceedings of the IEEE International Conference on Communications Workshops (ICC) Workshops*, pp. 539–544. IEEE, Piscataway (2014)
36. Xu, X., Rao, X., Lau, V.K.: Active user detection and channel estimation in uplink CRAN systems. In: *Proceedings of the International Conference on Communications (ICC)*, pp. 2727–2732. IEEE, Piscataway (2015)
37. Zhang, J., Ghanem, B.: ISTA-net: Interpretable optimization-inspired deep network for image compressive sensing. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1828–1837 (2018)

# Chapter 3

## Blind Demixing



**Abstract** This chapter presents a blind demixing model for joint data decoding and channel estimation in IoT networks, without transmitting pilot sequences. The problem formulation based on the cyclic convolution in the time domain is first introduced, which is then reformulated in the Fourier domain for the ease of algorithm design. A convex relaxation approach based on nuclear norm minimization is first presented as a basic solution. Next, several nonconvex approaches are introduced, including both regularized and regularization-free Wirtinger flow and the Riemannian optimization algorithm. The mathematical tools for analyzing nonconvex approaches are also provided.

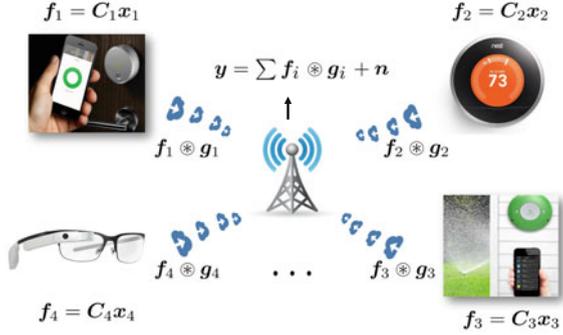
### 3.1 Joint Data Decoding and Channel Estimation

For data transmission in IoT networks, as the blocklength of packets is typically very short, the channel estimation sequences (CES) (illustrated in Fig. 1.1) occupy the primary part of the packet [11]. Thus, CES overhead reduction becomes critical to achieve low-overhead communications. To exclude the CES overhead, the BS may jointly decode data and estimate channel states, which can be established as a blind demixing model (3.1) [8].

For an IoT network containing one BS and  $s$  devices, as shown in Fig. 3.1, the observation signal vector is the mixture of the encoding signals generated from  $s$  devices and passed through the corresponding channels. The goal of the BS is to jointly decode data and estimate the channel states, which can be captured via a *blind demixing model* consisting of both summation operation and convolution operation. For ease of algorithm design, the measurements in the blind demixing are represented in the Fourier domain, which are given by

$$y_j = \sum_{i=1}^s \mathbf{b}_j^H \mathbf{h}_i \mathbf{x}_i^H \mathbf{a}_{ij}, \quad 1 \leq j \leq L. \quad (3.1)$$

**Fig. 3.1** The blind demixing model in an IoT network



Denote

$$\mathbf{y} = [y_1, \dots, y_L]^T \in \mathbb{C}^L$$

as the received signal at the BS represented in the Fourier domain,  $\{\mathbf{b}_j\}, \{\mathbf{a}_{ij}\}$  are design vectors, and  $\{\mathbf{h}_i\}, \{\mathbf{x}_i\}$  are channel states and data signals, respectively. Particularly, the design vectors  $\{\mathbf{b}_j\}$  indicate the Fourier transform operation and the design vectors  $\{\mathbf{a}_{ij}\}$  indicate the encoding procedure. By evaluating the vectors  $\{\mathbf{h}_i\}, \{\mathbf{x}_i\}$  from the observation  $\mathbf{y}$ , data decoding and channel estimation can be simultaneously accomplished.

There is a growing body of recent works paying attention on the blind demixing model (3.1). In particular, semidefinite programming has been developed in [9] to solve the blind demixing problems by lifting the bilinear model into the rank-one matrix model. However, it is computationally expensive to deal with large-scale problems. To address this issue, nonconvex algorithms, e.g., regularized Wirtinger flow with spectral initialization [10], have been developed to optimize the variables in the vector space. The Riemannian trust-region optimization algorithm without regularization was further developed in [8] to improve the convergence rate compared to the regularized Wirtinger flow algorithm [10]. Recently, concerning the blind demixing problem, theoretical guarantees for regularization-free Wirtinger flow with spectral initialization were established in [6]. To further find a natural initialization for the practitioners that works equally well as spectral initialization, the paper [7] established the global convergence guarantee of Wirtinger flow with random initialization for blind demixing.

In the sequel, the procedure of establishing the blind demixing model based on the convolution operations in IoT networks will be first illustrated. We further clarify the vital role that the blind demixing model plays in joint data decoding and channel estimation. Then, effective algorithms and rigorous analysis are provided for both convex and nonconvex approaches.

## 3.2 Problem Formulation

In this section, the basic concept of the cyclic convolution is first introduced, followed by a detailed description of the blind demixing model based on the cyclic convolution.

### 3.2.1 Cyclic Convolution

The elementary concept of the cyclic convolution is first introduced to characterize the connection among the channel state, received signal, and transmitted signal, thereby assisting the presentation of the blind demixing model.

Denote  $p[n]$  and  $\theta[n]$  as the transmitted signal and received signal in the  $n$ -th time slot, respectively. Define  $q_\ell$  as the  $\ell$ -th tap channel impulse response which is constant with  $n$ . Hence, the channel is assumed to be linear time-invariant. Thus, the discrete-time model is represented as

$$\theta[n] = \sum_{\ell=0}^{L_t-1} q_\ell p[n - \ell], \quad (3.2)$$

where  $L_t$  is the number of nonzero taps. A *cyclic prefix* is added to  $\mathbf{p}$ , which yields the symbol vector, i.e.,  $\mathbf{d} \in \mathbb{C}^{N_p+L_t-1}$ :

$$\mathbf{d} = [p[N_p - L_t + 1], \dots, p[N_p - 1], p[0], p[1], \dots, p[N_p - 1]]^\top. \quad (3.3)$$

The output over the time interval  $n \in [L_t, N_p + L_t - 1]$  is represented as

$$\theta[n] = \sum_{\ell=0}^{L_t-1} q_\ell d[(n - L_t - \ell) \text{ modulo } N_p]. \quad (3.4)$$

Denote the output of length  $N_p$  as

$$\boldsymbol{\theta} = [\theta[L_t], \dots, \theta[N_p + L_t - 1]]^\top, \quad (3.5)$$

and the channel impulse as

$$\mathbf{q} = [q_0, q_1, \dots, q_{L_t-1}, 0, \dots, 0]^\top \in \mathbb{C}^{N_p},$$

and then (3.4) can be reformulated as

$$\boldsymbol{\theta} = \mathbf{q} \circledast \mathbf{p},$$

where the notion  $\circledast$  denotes the *cyclic convolution*.

### 3.2.2 System Model

Consider a network with one BS and  $s$  devices. Denote the original signals of length  $N$  from the  $i$ -th user as  $\mathbf{x}_i \in \mathbb{C}^N$ . The transmitted signals over  $L$  time slots from the  $i$ -th user are represented as

$$\mathbf{f}_i = \mathbf{C}_i \mathbf{x}_i, \quad (3.6)$$

where  $\mathbf{C}_i \in \mathbb{C}^{L \times N}$  with  $L > N$  as the encoding matrix and known to the BS. The signals  $\mathbf{f}_i$ 's are transmitted through individual time-invariant channels endowed with impulse responses  $\mathbf{h}_i$ 's where a maximum delay of at most  $K$  samples is contained in  $\mathbf{h}_i \in \mathbb{C}^K$ . The zero-padded channel vector  $\mathbf{g}_i \in \mathbb{C}^L$  is given as

$$\mathbf{g}_i = [\mathbf{h}_i^\top, 0, \dots, 0]^\top. \quad (3.7)$$

Hence, based on the cyclic convolution operation, the received signal is given as

$$\mathbf{z} = \sum_{i=1}^s \mathbf{f}_i \circledast \mathbf{g}_i + \mathbf{n}, \quad (3.8)$$

where  $\mathbf{n}$  is the additive white complex Gaussian noise. The BS needs to recover the data signals  $\{\mathbf{x}_i\}_{i=1}^s$  from the observation  $\mathbf{z}$  without knowing channel states  $\{\mathbf{g}_i\}_{i=1}^s$ . This model is called a *blind demixing* model.

### 3.2.3 Representation in the Fourier Domain

For the ease of algorithm design and theoretical analysis, the blind demixing model based on cyclic convolution is represented in the Fourier domain. This is achieved by left multiplying the signals in the time domain with the unitary discrete Fourier transform (DFT) matrix and converting the convolution operation in the time domain to the componentwise production operation in the Fourier domain [8, 9]:

$$\mathbf{y} = \mathbf{F} \mathbf{z} = \sum_i (\mathbf{F} \mathbf{C}_i \mathbf{x}_i) \odot \mathbf{B} \mathbf{h}_i + \mathbf{F} \mathbf{n}, \quad (3.9)$$

where the operation  $\odot$  is the componentwise product. Here, the first  $K$  columns of the unitary discrete Fourier transform (DFT) matrix  $\mathbf{F} \in \mathbb{C}^{L \times L}$  with  $\mathbf{F} \mathbf{F}^\mathbf{H} = \mathbf{I}_L$  form the known matrix

$$\mathbf{B} := [\mathbf{b}_1, \dots, \mathbf{b}_L]^\mathbf{H} \in \mathbb{C}^{L \times K}$$

with  $\mathbf{b}_j \in \mathbb{C}^K$  for  $1 \leq j \leq L$ . An example of the sparse linear model is illustrated in Example 3.1.

*Example 3.1* Consider a network with two devices and one single-antenna BS. We assume that  $K = N = 1$  and the data signals  $\{\mathbf{x}_i\}_{i=1}^2$  as  $\mathbf{x}_1 = 1, \mathbf{x}_2 = 2$ , and channel signals  $\{\mathbf{h}_i\}_{i=1}^2$  as  $\mathbf{h}_1 = 3, \mathbf{h}_2 = 4$ . In addition, three time slots are considered in this example such that the encoding matrices are given by

$$\mathbf{C}_1 = \begin{bmatrix} 1 \\ 1 \\ 2 \end{bmatrix} \quad \text{and} \quad \mathbf{C}_2 = \begin{bmatrix} 4 \\ 2 \\ 1 \end{bmatrix}. \quad (3.10)$$

Based on the unitary discrete Fourier transform (DFT) matrix  $\mathbf{F} \in \mathbb{C}^{3 \times 3}$ :

$$\mathbf{F} = \begin{bmatrix} 0.5774 & 0.5774 & 0.5774 \\ 0.5774 & -0.2887 - 0.5i & -0.2887 + 0.5i \\ 0.5774 & -0.2887 + 0.5i & -0.2887 - 0.5i \end{bmatrix}, \quad (3.11)$$

it yields the blind demixing model:

$$\begin{aligned} \mathbf{y} &= \sum_i (\mathbf{F} \mathbf{C}_i \mathbf{x}_i) \odot \mathbf{B} \mathbf{h}_i \\ &= \begin{bmatrix} 2.3096 \\ -0.2887 + 0.5i \\ -0.2887 - 0.5i \end{bmatrix} \odot \begin{bmatrix} 1.7322 \\ 1.7322 \\ 1.7322 \end{bmatrix} + \begin{bmatrix} 8.8036 \\ 2.8870 - 0.5i \\ 2.8870 + 0.5i \end{bmatrix} \odot \begin{bmatrix} 2.3096 \\ 2.3096 \\ 2.3096 \end{bmatrix} \\ &= \begin{bmatrix} 22.6706 \\ 6.1677 - 1.4435i \\ 6.1677 + 1.4435i \end{bmatrix}. \end{aligned} \quad (3.12)$$

Generally, the blind demixing model can be formulated as the sum of bilinear measurements of vectors  $\mathbf{x}_i^\natural \in \mathbb{C}^N$ ,  $\mathbf{h}_i^\natural \in \mathbb{C}^K$ ,  $i = 1, \dots, s$ , i.e.,

$$y_j = \sum_{i=1}^s \mathbf{b}_j^H \mathbf{h}_i^\natural \mathbf{x}_i^{\natural H} \mathbf{a}_{ij} + e_j, \quad 1 \leq j \leq L, \quad (3.13)$$

where  $y_j$  is the  $j$ -th entry of  $\mathbf{y}$  in (3.9),  $\mathbf{b}_j \in \mathbb{C}^K$  denotes the  $j$ -th column of  $\mathbf{B}^H$ , and  $\mathbf{a}_{ij} \in \mathbb{C}^N$  denotes the  $j$ -th column of  $(\mathbf{F} \mathbf{C}_i)^H$ . Furthermore, the noise  $e_j$  obeys

$$e_j \sim \mathcal{CN} \left( 0, \frac{\sigma^2 d_0^2}{2L} \right) \quad (3.14)$$

with

$$d_0 = \sqrt{\sum_{i=1}^s \|\mathbf{h}_i^{\flat}\|_2^2 \|\mathbf{x}_i^{\flat}\|_2^2} \quad (3.15)$$

and  $\sigma^2$  as the measurement of noise variance.

### 3.3 Convex Relaxation Approach

In this section, a convex relaxation approach for estimating the blind demixing model is introduced, followed by theoretical analysis.

#### 3.3.1 Method: Nuclear Norm Minimization

To begin with, a low-rank matrix optimization problem is established via lifting the bilinear model in (3.13) to the rank-one matrices space. Based on (3.9), the  $j$ -th entry of the first term in (3.9) can be formulated as

$$[(\mathbf{F}\mathbf{C}_i\mathbf{x}_i) \odot \mathbf{B}\mathbf{h}_i]_j = (\mathbf{c}_{ij}^{\text{H}}\mathbf{x}_i)(\mathbf{b}_j^{\text{H}}\mathbf{h}_i) = \langle \mathbf{c}_{ij}\bar{\mathbf{b}}_j^{\text{H}}, \mathbf{X}_i \rangle,$$

where  $\mathbf{c}_{ij}^{\text{H}}$  is the  $j$ -th row of  $\mathbf{F}\mathbf{C}_i$ ,  $\mathbf{b}_j^{\text{H}}$  is the  $j$ -th row of  $\mathbf{B}$ , and  $\mathbf{X}_i = \mathbf{x}_i\bar{\mathbf{h}}_i^{\text{H}} \in \mathbb{C}^{N \times K}$  is a rank-one matrix. We have

$$y_j = \left\langle \begin{bmatrix} \mathbf{x}_1\bar{\mathbf{h}}_1^{\text{H}} & 0 & \cdots & 0 \\ 0 & \mathbf{x}_2\bar{\mathbf{h}}_2^{\text{H}} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \mathbf{x}_s\bar{\mathbf{h}}_s^{\text{H}} \end{bmatrix}, \begin{bmatrix} \mathbf{c}_{1j}\bar{\mathbf{b}}_j^{\text{H}} & 0 & \cdots & 0 \\ 0 & \mathbf{c}_{2j}\bar{\mathbf{b}}_j^{\text{H}} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \mathbf{c}_{sj}\bar{\mathbf{b}}_j^{\text{H}} \end{bmatrix} \right\rangle + e_l. \quad (3.16)$$

Thus, the received signal at the BS in the Fourier domain is given by

$$\mathbf{y} = \sum_{k=1}^s \mathcal{A}_i(\mathbf{X}_i) + \mathbf{e}, \quad (3.17)$$

where the vector  $\mathbf{e}$  denotes the additive Gaussian noise and the linear operator  $\mathcal{A}_i : \mathbb{C}^{N \times K} \rightarrow \mathbb{C}^L$  is represented as

$$\mathcal{A}_i(\mathbf{X}_i) := \left\{ \langle \mathbf{c}_{ij}\bar{\mathbf{b}}_j^{\text{H}}, \mathbf{X}_i \rangle \right\}_{i=1}^L = \{ \langle \mathbf{A}_{ij}, \mathbf{X}_i \rangle \}_{i=1}^L, \quad (3.18)$$

with

$$\mathbf{A}_{ij} = \mathbf{c}_{ij} \bar{\mathbf{b}}_j^H.$$

In addition, the operation  $\mathcal{A}_i^* : \mathbb{C}^L \rightarrow \mathbb{C}^{N \times K}$  can be represented as

$$\mathcal{A}_i^*(\mathbf{y}) = \sum_{j=1}^L y_j \mathbf{b}_j \mathbf{c}_{ij}^H. \quad (3.19)$$

The goal of blind demixing problem is to find the rank-one matrices that match the observation, formulated as

$$\begin{aligned} & \text{find } \text{rank}(\mathbf{W}_i) = 1, \quad i = 1, \dots, s \\ & \text{subject to } \left\| \sum_{i=1}^s \mathcal{A}_i(\mathbf{W}_i) - \mathbf{y} \right\|_2 \leq \varepsilon, \end{aligned} \quad (3.20)$$

where the parameter  $\varepsilon$  is a bound for  $\|\mathbf{e}\|_2$  (recall that  $\mathbf{e}$  appeared in (3.17)). Nevertheless, due to the nonconvexity of the rank function, problem (3.20) is NP-hard and thus intractable. The nuclear norm minimization approach has been proposed to relax the rank function [9], which gives the following formulation:

$$\begin{aligned} & \underset{\mathbf{W}_i, i=1, \dots, s}{\text{minimize}} \quad \sum_{i=1}^s \|\mathbf{W}_i\|_* \\ & \text{subject to } \left\| \sum_{i=1}^s \mathcal{A}_i(\mathbf{W}_i) - \mathbf{y} \right\|_2 \leq \varepsilon. \end{aligned} \quad (3.21)$$

The problem (3.21) can be solved via semidefinite program. Based on the estimated  $\hat{\mathbf{X}}_i$ , the corresponding  $\hat{\mathbf{h}}_i$  and  $\hat{\mathbf{x}}_i$  can be set as the right and left singular values of  $\hat{\mathbf{X}}_i$ , respectively.

### 3.3.2 Theoretical Analysis

To present theoretical analysis for methods that solve the blind demixing problem, several notions are first introduced. For simplicity, we summarize the parameters involved in the analysis of solving the blind demixing problem (3.21) via semidefinite program in Table 3.1, and the detailed formulations of these parameters can be found in [9].

The paper [9] demonstrates that the method (3.21) provides an effective way to solve the blind demixing problem and is also robust to noise, as illustrated in Theorem 3.1.

**Table 3.1** Conditions involved in Theorem 3.1 and corresponding section mentioned in [9]

Condition	Parameter	Reference
Joint incoherent pattern on the matrices $\mathbf{B}$	$\mu_{\max}, \mu_{\min}$	Sect. II-C
Incoherence between $\mathbf{b}_j$ and $\mathbf{h}_i$	$\mu$	Sect. II-D
Upper bound on $\ \mathcal{A}_i\  := \sup_{\mathbf{X} \neq \mathbf{0}} \ \mathcal{A}_i(\mathbf{X})\ _F / \ \mathbf{X}\ _F$	$\gamma$	Sect. II-E

**Theorem 3.1** *Considering the blind demixing model (3.17) in the noiseless scenario, if*

$$L \geq Cs^2 \max\{\mu_{\max}^2 K, \mu_h^2 N\} \log^2 L \log \gamma,$$

where  $C > 0$  is sufficiently large, and  $\mu_{\max}, \mu, \gamma$  are summarized in Table 3.1, then the convex relaxation approach (3.21) recovers the ground truth rank-one matrices exactly with high probability.

**Proof** The proof details of Theorem 3.1 can be referred to the paper [9] and the proof architecture of Theorem 3.1 is briefly summarized. We further present a sufficient condition and an approximate dual certificate condition for the minimizer of (3.21) to be the unique solution to (3.20). These conditions stipulate that matrices  $\mathcal{A}_i$  need to satisfy two key properties. The first property can be regarded as a modification of the celebrated *Restricted Isometry Property (RIP)* [3], as it requires  $\mathcal{A}_i$  to act in a certain sense as “local” approximate isometries [4]. The second property requires the two operators  $\mathcal{A}_i$  and  $\mathcal{A}_j$  to satisfy a “local” *mutual incoherence property*. With these two key properties in place, an approximate dual certificate can be established that fulfills the sufficient condition. With all these tools in place, the proof of Theorem 3.1 can be completed.

*Remark 3.1* Theorem 3.1 demonstrates that the successful recovery of  $\{\mathbf{W}_i\}_{i=1}^s$  in problem (3.21) in the noiseless scenario via semidefinite programming can be achieved with high probability as long as the number of measurements satisfies

$$L \gtrsim s^2 \max\{\mu_{\max}^2 K, \mu_h^2 N\} \log^2 L.$$

The paper [9] further considers problem (3.21) in the noisy scenario and provides the performance guarantee of recovering  $\{\mathbf{W}_i\}_{i=1}^s$  in problem (3.21) under the same conditions as in Theorem 3.1.

### 3.4 Nonconvex Approaches

While convex techniques can be exploited to solve the blind demixing problem provably and robustly under certain assumptions, the resulting algorithms are computationally expensive for large-scale problems. This motivates the develop-

ment of efficient nonconvex approaches, which are introduced in this section. The nonconvex approaches introduced in this section can be separated into two types: the Wirtinger flow based approach, which is an iterative algorithm based on the gradients derived in the complex space, and the Riemannian optimization based approach, which is developed on the Riemannian manifold search space.

### 3.4.1 Regularized Wirtinger Flow

Considering the rank-one structure of the blind demixing model, matrix factorization provides an efficient method to address the low-rank optimization problem. Specifically, Ling and Strohmer [10] developed an algorithm to solve the blind demixing problem based on matrix factorization and the regularized Wirtinger flow. The regularized optimization problem is established as

$$\underset{\mathbf{u}_k, \mathbf{v}_k, k=1, \dots, s}{\text{minimize}} \quad F(\mathbf{u}, \mathbf{v}) := g(\mathbf{u}, \mathbf{v}) + \lambda R(\mathbf{u}, \mathbf{v}), \quad (3.22)$$

where

$$g(\mathbf{u}, \mathbf{v}) := \left\| \sum_{k=1}^s \mathcal{A}_k(\mathbf{u}_k \mathbf{v}_k^H) - \mathbf{y} \right\|^2$$

with  $\mathbf{u}_k \in \mathbb{C}^N$ ,  $\mathbf{v}_k \in \mathbb{C}^K$  and the aim of the regularizer  $R(\mathbf{u}, \mathbf{v})$  is to enable the iterates to lie in the *basin of attraction* [10]. The algorithm begins with a spectral initialization point and updates the iterates as:

$$\mathbf{u}_k^{[t+1]} = \mathbf{u}_k^{[t]} - \eta \nabla F_{\mathbf{u}_k}(\mathbf{u}_k^{[t]}, \mathbf{v}_k^{[t]}), \quad (3.23)$$

$$\mathbf{v}_k^{[t+1]} = \mathbf{v}_k^{[t]} - \eta \nabla F_{\mathbf{v}_k}(\mathbf{u}_k^{[t]}, \mathbf{v}_k^{[t]}), \quad (3.24)$$

where  $\nabla F_{\mathbf{u}_k}$  is the derivative of the objective function (3.22) with respect to  $\mathbf{u}_k$ . The following theorem provided in [10] demonstrates that the regularized Wirtinger gradient descent will guarantee the linear convergence of the iterates, and the recovery is exact in the noiseless scenario and stable in the presence of noise.

Denote the condition number as

$$\kappa := \frac{\max_i \|\mathbf{x}_i^{\natural}\|_2}{\min_i \|\mathbf{x}_i^{\natural}\|_2}, \quad (3.25)$$

and recall that  $d_0$  in (3.15). Furthermore, for simplicity, we summarize the parameters involved in the analysis of solving the blind demixing problem (3.22) via regularized Wirtinger flow in Table 3.2, and the detailed formulations of these parameters can be referred to the references presented in the table.

**Table 3.2** Conditions involved in Theorem 3.2 and corresponding section mentioned in [10]

Condition	Parameter	Reference
Local regularity condition	$\omega$	Sect. 5.1
Robustness condition on $\ \mathcal{A}^*(\mathbf{e})\ $	$\gamma_e$	Sect. 6.5

---

**Algorithm 3.1: 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  $1 \leq i \leq s$ :

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

- 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  $(\mathbf{u}^{(0)}, \mathbf{v}^{(0)}, \{d_i\}_{i=1}^s)$ .
- 

**Theorem 3.2** Starting from the initial point generated via Algorithm 3.1, the regularized Wirtinger flow algorithm derives a sequence of iterates  $(\mathbf{u}^{[t]}, \mathbf{v}^{[t]})$  which converges to the global minimum linearly,

$$\sum_{k=1}^s \left\| \mathbf{u}_k^{[t]} (\mathbf{v}_k^{[t]})^H - \mathbf{h}_k^\natural \mathbf{x}_k^{\natural H} \right\|_F \leq \frac{d_0}{\sqrt{2s\kappa^2}} (1 - \eta\omega)^{t/2} + 60\sqrt{s}\gamma_e \quad (3.26)$$

with high probability if the number of measurements  $L$  satisfies

$$L \geq C(\mu^2 + \sigma^2)s^2\kappa^4 \max\{K, N\} \log^2 L, \quad (3.27)$$

where  $C > 0$  is sufficiently large. Here, the parameter  $\sigma$  and  $d_0$  are defined in (3.14), and  $\omega, \gamma_e$  are summarized in Table 3.2.

**Proof** The convergence analysis provided in Theorem 3.2 relies on four conditions: local regularity condition of the objective function  $F(\mathbf{u}, \mathbf{v})$  (3.22), local smoothness condition of the objective function  $F(\mathbf{u}, \mathbf{v})$  (3.22), local restricted isometry property, and robustness condition. Under the assumptions mentioned in Theorem 3.2, with the spectral initialization being in the basin of attraction, the four conditions can be guaranteed, which yield the results of Theorem 3.2.

The performance of regularized Wirtinger flow is further illustrated in Sect. 3.4.4.

*Remark 3.2* Even though Theorem 3.2 demonstrates that the regularized Wirtinger flow endows with a linear convergence rate, it requires extra regularization added on the objective function, and the step size, i.e.,  $\eta \lesssim \frac{1}{s\kappa m}$  [10], lacks of aggressiveness. To exclude the regularization and achieve more aggressive step size, the papers [6, 7]

have recently investigated regularization-free Wirtinger flow which yields a more aggressive step size, i.e.,  $\eta \lesssim s^{-1}$ .

### 3.4.2 Regularization-Free Wirtinger Flow

Another line of studies has focused on the blind demixing model that is in the form of the bilinear model (3.13). In this section, we would formulate an optimization problem concerning the bilinear formulation of blind demixing and introduce efficient regularizer-free algorithms. The theoretical analysis on these algorithms will also be discussed.

A least-squares optimization problem under the scheme of the bilinear formulation of blind demixing is given by

$$\underset{\{\mathbf{h}_i\}, \{\mathbf{x}_i\}}{\text{minimize}} f(\mathbf{h}, \mathbf{x}) := \sum_{j=1}^m \left| \sum_{i=1}^s \mathbf{b}_j^H \mathbf{h}_i \mathbf{x}_i^H \mathbf{a}_{ij} - y_j \right|^2. \quad (3.28)$$

For simplification, the objective function in (3.28) is denoted as

$$f(\mathbf{z}) := f(\mathbf{h}, \mathbf{x}),$$

where

$$\mathbf{z} = [\mathbf{z}_1^H \cdots \mathbf{z}_s^H]^H \in \mathbb{C}^{2sK} \text{ with } \mathbf{z}_i = [\mathbf{h}_i^H \ \mathbf{x}_i^H]^H \in \mathbb{C}^{2K}.$$

A line of literatures, e.g., [5–7, 10], have developed effective algorithms to solve problem (3.28). The blind demixing problem can be generally solved via two procedures [6, 10], i.e., Stage I: find an initial point that is in the neighborhood of the ground truth, which can be accomplished via spectral initialization; Stage II: optimize the initial estimate via an iterative algorithm, e.g., Wirtinger flow:

$$\begin{bmatrix} \mathbf{h}_i^{t+1} \\ \mathbf{x}_i^{t+1} \end{bmatrix} = \begin{bmatrix} \mathbf{h}_i^t \\ \mathbf{x}_i^t \end{bmatrix} - \eta \begin{bmatrix} \frac{1}{\|\mathbf{x}_i^t\|_2^2} \nabla_{\mathbf{h}_i} f(\mathbf{z}^t) \\ \frac{1}{\|\mathbf{h}_i^t\|_2^2} \nabla_{\mathbf{x}_i} f(\mathbf{z}^t) \end{bmatrix}, \quad i = 1, \dots, s, \quad (3.29)$$

where  $\eta > 0$  is the step size,  $\nabla_{\mathbf{h}_i} f(\mathbf{z})$  and  $\nabla_{\mathbf{x}_i} f(\mathbf{z})$  represent the Wirtinger gradient of  $f(\mathbf{z})$  with respect to  $\mathbf{h}_i$  and  $\mathbf{x}_i$ , respectively.

The discrepancy between the estimate  $\mathbf{z}$  and the ground truth  $\mathbf{z}^\natural$  is defined as the distance function:

$$\text{dist}(\mathbf{z}, \mathbf{z}^\natural) = \left( \sum_{i=1}^s \text{dist}^2(\mathbf{z}_i, \mathbf{z}_i^\natural) \right)^{1/2}, \quad (3.30)$$

**Table 3.3** Conditions involved in Theorem 3.3 and corresponding references

Conditions	Reference
Incoherence between $\mathbf{a}_j$ and $\mathbf{x}_i$	(6b) in [6]
Incoherence between $\mathbf{b}_j$ and $\mathbf{h}_i$	(6c) in [6]
Robustness condition: $\ \mathcal{A}_i^*(\mathbf{e})\  \leq \gamma_e$	Sect. 6.5 in [10]

where

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

for  $i = 1, \dots, s$ . Here,  $d_i = \|\mathbf{h}_i^{\natural}\|_2^2 + \|\mathbf{x}_i^{\natural}\|_2^2$  and each  $\alpha_i$  is an alignment parameter. Without loss of generality, the ground-truth vectors are assumed to obey  $\|\mathbf{h}_i^{\natural}\|_2 = \|\mathbf{x}_i^{\natural}\|_2$  for  $i = 1, \dots, s$ . Recall the operator in (3.19) such that

$$\mathcal{A}_i^*(\mathbf{e}) = \sum_{j=1}^m e_j \mathbf{b}_j \mathbf{a}_{ij}^{\text{H}}, \quad i = 1, \dots, s,$$

and the condition number  $\kappa$  in (3.25), then the theorem of Wirtinger flow with the spectral initialization for solving the blind demixing problem is presented in Theorem 3.3.

For simplicity, we summarize the conditions involved in the analysis of solving the blind demixing problem (3.28) via Wirtinger flow with spectral initialization in Table 3.3, and the detailed formulations of these parameters can be referred to the references presented in the table.

**Theorem 3.3** *Assuming that the step size obeys  $\eta > 0$ ,  $\eta \asymp s^{-1}$ , and the conditions in Table 3.3 are satisfied, the iterates (including the spectral initialization point) in Wirtinger flow satisfy*

$$\text{dist}(\mathbf{z}^t, \mathbf{z}^{\natural}) \leq C_1 \left(1 - \frac{\eta}{16\kappa}\right)^t \left( \frac{1}{\log^2 m} - \frac{48\sqrt{s}\kappa^2}{\eta} \cdot \gamma_e \right) + \frac{48C_1\sqrt{s}\kappa^2}{\eta} \gamma_e,$$

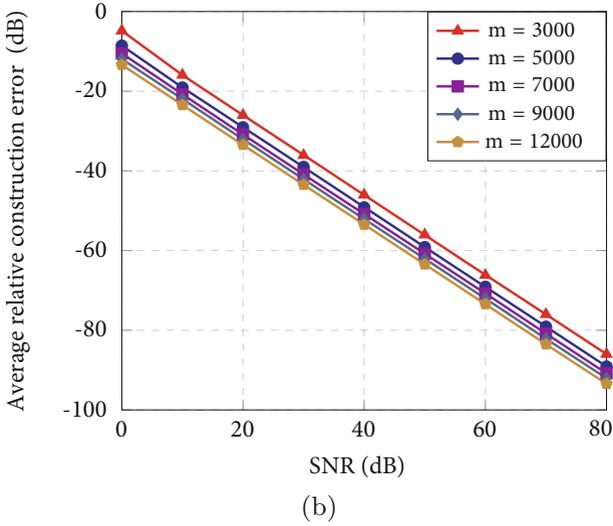
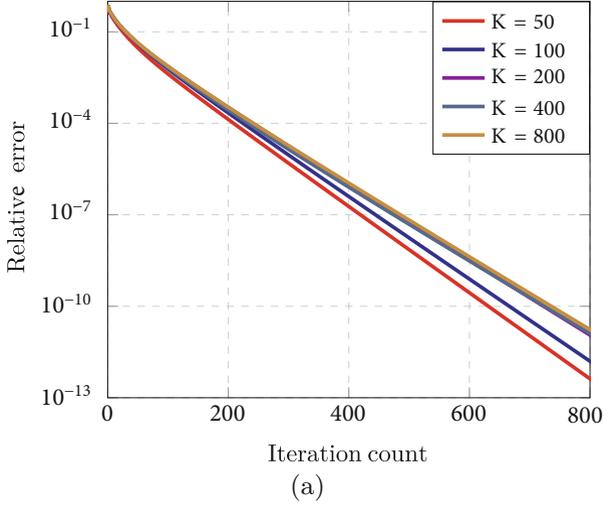
for all  $t \geq 0$ , with high probability if the number of measurements satisfies

$$m \geq C(\mu^2 + \sigma^2)s^2\kappa^4 K \log^8 m$$

for some constants  $C_1 > 0$  and adequately large constant  $C > 0$ .

**Proof** Please refer to Sect. 8.3 for details.

Theorem 3.3 provides the justification for a more aggressive step size (i.e.,  $\eta \asymp s^{-1}$ ) even without regularization, compared to the step size (i.e.,  $\eta \lesssim \frac{1}{s\kappa m}$ ) given in [10] for regularized Wirtinger flow. In addition, the performance of the Wirtinger flow algorithm with spectral initialization is illustrated in Fig. 3.2a, b, which is endorsed



**Fig. 3.2** Numerical results of Wirtinger flow with spectral initialization

by Theorem 3.3. To be specific, for each  $K \in \{50, 100, 200, 400, 800\}$ ,  $s = 10$ , and  $m = 50K$ , the design vectors  $\mathbf{a}_{ij}$ 's and  $\mathbf{b}_j$ 's for each  $1 \leq i \leq s, 1 \leq j \leq m$  are generated based on the instructions in Sect. 3.2.3. The underlying signals  $\mathbf{h}_i, \mathbf{x}_i \in \mathbb{C}^K, 1 \leq i \leq s$ , are generated as random vectors with unit norm. With the chosen step size  $\eta = 0.1$  in all settings, Fig. 3.2a shows the relative error, i.e.,

$$\frac{\sum_{i=1}^s \|\mathbf{h}_i^t \mathbf{x}_i^{tH} - \mathbf{h}_i^{\natural} \mathbf{x}_i^{\natural H}\|_F}{\sum_{i=1}^s \|\mathbf{h}_i^{\natural} \mathbf{x}_i^{\natural H}\|_F}, \quad (3.31)$$

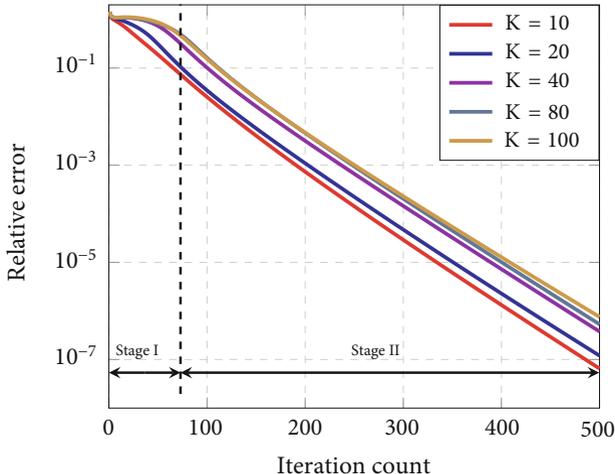


Fig. 3.3 Numerical result of Wirtinger flow with random initialization

versus the iteration count. Figure 3.2a shows that in the noiseless case, Wirtinger flow with a constant step size enjoys a linear convergence rate, which barely changes as the problem scale changes. Additionally, Fig. 3.2b shows the relative error (3.31) versus the signal-to-noise ratio (SNR), where the SNR is defined as  $\text{SNR} := \|\mathbf{y}\|_2 / \|\mathbf{e}\|_2$ . Both the relative error and the SNR are represented in the dB scale.

The random initialization strategy has recently been proven in [7] to be good enough for Wirtinger flow to guarantee linear converge rate when solving blind demixing problems. Specifically, in Stage I, it takes  $\mathcal{O}(s \log(\max\{K, N\}))$  iterations for randomly initialized Wirtinger flow to reach a local region near the ground truth. Furthermore, in Stage II, it takes  $\mathcal{O}(s \log(1/\varepsilon))$  iterations to attain an  $\varepsilon$ -accurate estimator, i.e.,  $\text{dist}(\mathbf{z}, \mathbf{z}^\natural) \leq \varepsilon$ , at a linear convergence rate. Please refer to Sect. 8.4 for details on theoretical guarantees for this case. Figure 3.3 shows the performance of the Wirtinger flow algorithm with random initialization, showing the relative error (3.31) versus the iteration count. In the simulation, the ground truth signals and initial points are randomly generated as

$$\mathbf{h}_i^\natural \sim \mathcal{CN}(\mathbf{0}, K^{-1} \mathbf{I}_K), \quad \mathbf{x}_i^\natural \sim \mathcal{CN}(\mathbf{0}, N^{-1} \mathbf{I}_N), \quad (3.32)$$

$$\mathbf{h}_i^0 \sim \mathcal{CN}(\mathbf{0}, K^{-1} \mathbf{I}_K), \quad \mathbf{x}_i^0 \sim \mathcal{CN}(\mathbf{0}, N^{-1} \mathbf{I}_N), \quad (3.33)$$

for  $i = 1, \dots, s$ . In all simulations, we set  $K = N$  for each  $K \in \{10, 20, 40, 80, 100\}$ ,  $s = 10$ , and  $m = 50K$ , and with the chosen step size  $\eta = 0.1$ .

The above nonconvex algorithm has a low iteration cost, and the overall computational complexity can be further decreased via reducing the iteration complexity,

i.e., accelerating the convergence rate. This motivates to develop the Riemannian optimization algorithm which will be introduced in the next section.

### 3.4.3 Riemannian Optimization Algorithm

The paper [8] developed a Riemannian trust-region algorithm on the complex product manifolds to solve the blind demixing problem, which enjoys a fast convergence rate. Prior to introducing this algorithm for solving the blind demixing problem, we start with some basic concepts of Riemannian manifold optimization, and readers can refer to Sect. 8.5 in the book [1] for more details.

#### 3.4.3.1 An Example on Riemannian Optimization

In order to optimize a smooth function on a manifold, several geometric concepts in terms of manifolds are required. To be specific, tangent vectors on manifolds generalize the notion of a direction, and an inner product of tangent vectors generalizes a notion of length that applies to these tangent vectors. A *Riemannian manifold*, generally denoted as  $\mathcal{M}$ , is the manifold of which tangent spaces  $T_x\mathcal{M}$  are endowed with a smoothly varying inner product. The smoothly varying inner product is called the *Riemannian metric*, generally denoted as

$$g_x(\eta_x, \zeta_x),$$

where  $\mathbf{x} \in \mathcal{M}$  and  $\eta_x, \zeta_x \in T_x\mathcal{M}$ . Some examples of Riemannian manifold can be enumerated as: sphere, orthogonal Stiefel manifold, Grassmann manifold, rotation group, positive definite matrices, fixed-rank matrices, etc.

Consider minimizing a smooth function on the sphere  $\mathbb{S}^{n-1} = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{x}^\top \mathbf{x} = 1\}$ :

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} f(\mathbf{x}) = -\mathbf{x}^\top \mathbf{A} \mathbf{x} \quad \text{subject to} \quad \mathbf{x}^\top \mathbf{x} = 1, \quad (3.34)$$

where  $\mathbf{A}$  is a symmetric matrix. As illustrated in Fig. 3.4, the Riemannian optimization procedure on the sphere can be separated into three steps:

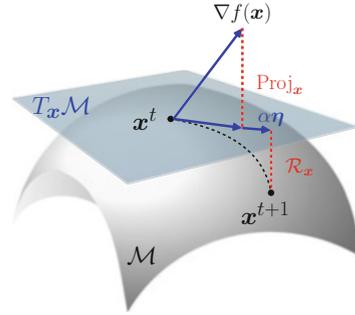
1. Compute the Euclidean gradient in  $\mathbb{R}^n$ :

$$\nabla f(\mathbf{x}) = -2\mathbf{A}\mathbf{x}. \quad (3.35)$$

2. Compute the Riemannian gradient on the sphere  $\mathbb{S}^{n-1}$  via projecting  $\nabla f(\mathbf{x})$  to the tangent space  $T_x\mathcal{M}$ :

$$\text{grad} f(\mathbf{x}) = \text{Proj}_x \nabla f(\mathbf{x}) = (\mathbf{I} - \mathbf{x}\mathbf{x}^\top) \nabla f(\mathbf{x}). \quad (3.36)$$

**Fig. 3.4** Schematic viewpoint of Riemannian optimization on the Riemannian manifold



3. Move along the descent direction  $\eta = \text{grad}f(x)$  and retract the directional vector  $\alpha\eta$  to the sphere, where  $\alpha > 0$  is the step size. The retraction operator on  $\mathbb{S}^{n-1}$  is given by

$$\mathcal{R}_x(\alpha\eta) = \text{qf}(x + \alpha\eta), \quad (3.37)$$

where  $\text{qf}(\cdot)$  denotes the mapping that maps a matrix to the Q factor of its QR decomposition.

Furthermore, for the ease of implementing the optimization scheme on manifolds, a powerful Matlab toolbox, namely Manopt [2], has been developed, which contains a larger library of manifolds (e.g., sphere, orthogonal Stiefel manifold, Grassmann manifold, rotation group, positive definite matrices, fixed-rank matrices, etc.) and various Riemannian optimization algorithms (e.g., steepest descent, conjugate gradient, stochastic gradient descent, trust-regions algorithm, etc.).

### 3.4.3.2 Riemannian Optimization on Product Manifolds for Blind Demixing

Due to the multiple rank-one matrices in the blind demixing problem (3.20), problem (3.20) can be reformulated as minimizing a smooth function on the product of multiple fixed-rank matrices. The product of multiple fixed-rank matrices is a product manifold and is also a Riemannian manifold [1]. The example mentioned above paves the way for dealing with more complicated Riemannian optimization algorithms on product manifolds for solving the blind demixing problem.

Firstly, a linear map is developed to handle complex asymmetric matrices. The linear map facilitates to convert the optimization variables to a Hermitian positive semidefinite matrix. Define a linear map

$$\mathcal{J}_i : \mathbb{S}_+^{(N+K)} \rightarrow \mathbb{C}^m$$

with respect to a Hermitian positive semidefinite (PSD) matrix  $\mathbf{Y}_i$  that obeys

$$[\mathcal{J}_i(\mathbf{Y}_i)]_i = \langle \mathbf{J}_{ij}, \mathbf{Y}_i \rangle \quad (3.38)$$

with  $\mathbf{Y}_i \in \mathbb{S}_+^{(N+K)}$  and  $\mathbf{J}_{ij}$  as

$$\mathbf{J}_{ij} = \begin{bmatrix} \mathbf{0}_{N \times N} & \mathbf{A}_{ij} \\ \mathbf{0}_{K \times N} & \mathbf{0}_{K \times K} \end{bmatrix} \in \mathbb{C}^{(N+K) \times (N+K)}, \quad (3.39)$$

where  $\mathbf{A}_{ij} = \mathbf{a}_{ij} \bar{\mathbf{b}}_j^H$ . Note that based on (3.38), we have

$$[\mathcal{J}_i(\mathbf{M}_i)]_i = \langle \mathbf{J}_{ij}, \mathbf{M}_i \rangle = \langle \mathbf{A}_{ij}, \mathbf{x}_i \bar{\mathbf{h}}_i^H \rangle, \quad (3.40)$$

where  $\mathbf{M}_i = \mathbf{w}_i \mathbf{w}_i^H$  with  $\mathbf{w}_i = [\mathbf{x}_i^H \bar{\mathbf{h}}_i^H]^H \in \mathbb{C}^{N+K}$ . Based on the matrix factorization, a manifold optimization problem with respect to Hermitian positive semidefinite (PSD) matrices can be established as:

$$\underset{\mathbf{v}=\{\mathbf{w}_k\}_{k=1}^s}{\text{minimize}} f(\mathbf{v}) := \left\| \sum_{k=1}^s \mathcal{J}_k(\mathbf{w}_k \mathbf{w}_k^H) - \mathbf{y} \right\|^2, \quad (3.41)$$

where  $\mathbf{v} \in \mathcal{M}^s$  with  $\mathbf{w}_k \in \mathcal{M} := \mathbb{C}_*^{N+K}$  for  $k = 1, \dots, s$ , where the space  $\mathbb{C}_*^n$  is the complex Euclidean space  $\mathbb{C}^n$  without the origin. According to (3.40), the data signal estimation  $\hat{\mathbf{x}}_k$  can be represented by the first  $N$  rows of the estimation  $\hat{\mathbf{w}}_k$ .

Since the quotient manifold is abstract, the matrix representations in the computational space  $\mathcal{M}^s$  of the geometric concepts in the quotient space are required. In particular, to develop the Riemannian optimization algorithm over the product manifolds, various geometric concepts need to be derived, such as the notion of length (i.e., Riemannian metric  $g_{\mathbf{w}_k}$ ), set of directional derivatives (i.e., horizontal space  $\mathcal{H}_{\mathbf{w}_k}$ ), and motion along geodesics (i.e., retraction  $\mathcal{R}_{\mathbf{w}_k}$ ) [1]. The concrete optimization-related ingredients are shown in Table 3.4. Based on these ingredients, we develop a Riemannian algorithm to solve the blind demixing problem (3.41).

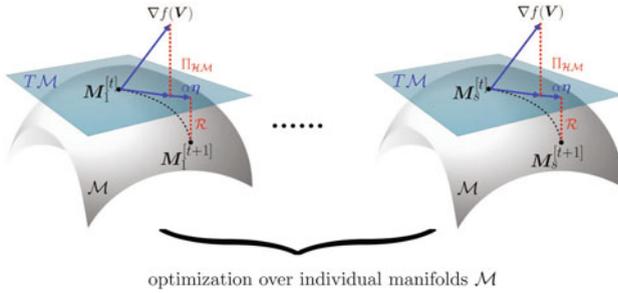
Based on the geometry of the product manifolds, the Riemannian optimization algorithm operated on the product manifolds  $\mathcal{M}^s$  can be elementwisely developed on the individual manifold  $\mathcal{M}$ . To be specific, for each  $k = 1, 2, \dots, s$ , the descent direction  $\boldsymbol{\eta}$  is detected on the horizontal space  $\mathcal{H}_{\mathbf{w}_k} \mathcal{M}$  parallelly, and  $\boldsymbol{\eta}$  is parallelly retracted on the individual manifold  $\mathcal{M}$  via the retraction mapping  $\mathcal{R}_{\mathbf{w}_k}$ . In addition, Fig. 3.5 shows the schematic viewpoint of Algorithm 3.2.

**Riemannian Gradient Descent with Spectral Initialization** In the Riemannian gradient descent algorithm, i.e., Algorithm 3.3, the search direction is given by

$$\boldsymbol{\eta} = -\text{grad}_{\mathbf{w}_k^{[r]}} f/g_{\mathbf{w}_k^{[r]}}(\mathbf{w}_k^{[r]}, \mathbf{w}_k^{[r]}),$$

**Table 3.4** Elementwise optimization-related ingredients for Problem (3.41)

	Minimize $w_k \in \mathcal{M} \quad \left\  \sum_{k=1}^S \mathcal{F}_k(w_k w_k^H) - \mathbf{y} \right\ ^2$
Computational space: $\mathcal{M}$	$\mathbb{C}_*^{N+K}$
Quotient space: $\mathcal{M} / \sim$	$\mathbb{C}_*^{N+K} / \text{SU}(1)$
Riemannian metric: $g_{w_k}$	$g_{w_k}(\zeta_{w_k}, \eta_{w_k}) = \text{Tr}(\zeta_{w_k}^H \eta_{w_k} + \eta_{w_k}^H \zeta_{w_k})$
Horizontal space: $\mathcal{H}_{w_k} \mathcal{M}$	$\eta_{w_k} \in \mathbb{C}^{N+K} : \eta_{w_k}^H w_k = w_k^H \eta_{w_k}$
Horizontal space projection	$\Pi_{\mathcal{H}_{w_k} \mathcal{M}}(\eta_{w_k}) = \eta_{w_k} - a w_k,$ $a = (w_k^H \eta_{w_k} - \eta_{w_k}^H w_k) / 2 w_k^H w_k$
Riemannian gradient: $\text{grad}_{w_k} f$	$\text{grad}_{w_k} f = \Pi_{\mathcal{H}_{w_k} \mathcal{M}}(\frac{1}{2} \nabla_{w_k} f(v))$
Riemannian Hessian: $\text{Hess}_{w_k} f[\eta_{w_k}]$	$\text{Hess}_{w_k} f[\eta_{w_k}] = \Pi_{\mathcal{H}_{w_k} \mathcal{M}}(\frac{1}{2} \nabla_{w_k}^2 f(v)[\eta_{w_k}])$
Retraction: $\mathcal{R}_{w_k} : T_{w_k} \mathcal{M} \rightarrow \mathcal{M}$	$\mathcal{R}_{w_k}(\eta_{w_k}) = w_k + \eta_{w_k}$

**Fig. 3.5** Schematic viewpoint of Riemannian optimization on the product manifolds

where  $g_{w_k^{[t]}}$  is the Riemannian metric and

$$\text{grad}_{w_k^{[t]}} f \in \mathcal{H}_{w_k} \mathcal{M}$$

is the Riemannian gradient. Therefore, the sequence of the iterates is given by

$$w_k^{[t+1]} = \mathcal{R}_{w_k^{[t]}}(\alpha_t \eta),$$

---

**Algorithm 3.2: Riemannian optimization on product manifolds**


---

Given: Riemannian manifold  $\mathcal{M}^s$  with Riemannian metric  $g_v$ , retraction mapping

$\mathcal{R}_v = \{\mathcal{R}_{\mathbf{w}_k}\}_{k=1}^s$ , objective function  $f$  and the step size  $\alpha$ .

Output:  $\mathbf{v} = \{\mathbf{w}_k\}_{k=1}^s$

- 1: Initialize: initial point  $\mathbf{v}^{[0]} = \{\mathbf{w}_k^{[0]}\}_{k=1}^s, t = 0$
  - 2: while not converged do
  - 3:   for all  $k = 1, \dots, s$  do in parallel
  - 4:     Compute a descent direction  $\eta$ . (e.g., via implementing trust-region method)
  - 5:     Update  $\mathbf{w}_k^{[t+1]} = \mathcal{R}_{\mathbf{w}_k^{[t]}}(\alpha\eta)$
  - 6:      $t = t + 1$ .
  - 7:   end for
  - 8: end while
- 

where the step size  $\alpha_t > 0$  and

$$\mathcal{R}_{\mathbf{w}_k}(\xi) = \mathbf{w}_k + \xi, \quad (3.42)$$

with  $\xi \in \mathcal{H}_{\mathbf{w}_k} \mathcal{M}$ . Here, the retraction map

$$\mathcal{R}_{\mathbf{w}_k} : \mathcal{H}_{\mathbf{w}_k} \mathcal{M} \rightarrow \mathcal{M}$$

is an approximation of the exponential map that characterizes the motion of “moving along geodesics on the Riemannian manifold.” More details on computing the retraction are available in [1, Section 4.1.2]. The statistical analysis of the Riemannian gradient descent algorithm will be provided in the sequel, which demonstrates the linear rate of the proposed algorithm for converging to the ground truth signals.

**Theorem 3.4** *Suppose the rows of the encoding matrices, i.e.,  $\mathbf{c}_{ij}$ 's, follow the i.i.d. complex Gaussian distribution, i.e.,*

$$\mathbf{c}_{ij} \sim \mathcal{N}\left(0, \frac{1}{2}\mathbf{I}_N\right) + i\mathcal{N}\left(0, \frac{1}{2}\mathbf{I}_N\right)$$

*and the step size obeys  $\alpha_t > 0$  and  $\alpha_t \equiv \alpha \asymp s^{-1}$ , then the iterates (including the spectral initialization point) in Algorithm 3.3 satisfy*

$$\text{dist}(\mathbf{v}^t, \mathbf{v}^\natural) \leq C_1 \left(1 - \frac{\alpha}{16\kappa}\right)^t \frac{1}{\log^2 L} \quad (3.43)$$

*for all  $t \geq 0$  and some constant  $C_1 > 0$ , with probability at least  $1 - c_1 L^{-\gamma} - c_1 L e^{-c_2 K}$  if the number of measurements*

$$L \geq C\mu^2 s^2 \kappa^4 \max\{K, N\} \log^8 L$$

*for some constants  $\gamma, c_1, c_2 > 0$  and sufficiently large constant  $C > 0$ .*

---

**Algorithm 3.3: Riemannian gradient descent with spectral initialization**


---

Given: Riemannian manifold  $\mathcal{M}^s$  with optimization-related ingredients, objective function  $f$ ,  $\{\mathbf{c}_{ij}\}$ ,  $\{\mathbf{b}_j\}$ ,  $\{y_j\}$  and the stepsize  $\alpha$ .

Output:  $\mathbf{v} = \{\mathbf{w}_k\}_{k=1}^s$

1: Spectral Initialization:

2: for all  $i = 1, \dots, s$  do in parallel

3: Let  $\sigma_1(N_i)$ ,  $\check{\mathbf{h}}_i^0$  and  $\check{\mathbf{x}}_i^0$  be the leading singular value, left singular vector and right singular vector of matrix  $N_i := \sum_{j=1}^m y_j \mathbf{b}_j \mathbf{c}_{ij}^H$ , respectively.

4: Set  $\mathbf{w}_i^{[0]} = \begin{bmatrix} \mathbf{x}_i^0 \\ \mathbf{h}_i^0 \end{bmatrix}$  where  $\mathbf{x}_i^0 = \sqrt{\sigma_1(N_i)} \check{\mathbf{x}}_i^0$  and  $\mathbf{h}_i^0 = \sqrt{\sigma_1(N_i)} \check{\mathbf{h}}_i^0$ .

5: end for

6: for all  $t = 1, \dots, T$

7: for all  $i = 1, \dots, s$  do in parallel

8:  $\boldsymbol{\eta} = -\frac{1}{g_{\mathbf{w}_k^{[t]}(\mathbf{w}_k^{[t]}, \mathbf{w}_k^{[t]})}} \text{grad}_{\mathbf{w}_k^{[t]}} f$

9: Update  $\mathbf{w}_k^{[t+1]} = \mathcal{R}_{\mathbf{w}_k^{[t]}}(\alpha_t \boldsymbol{\eta})$

10: end for

11: end for

---

**Proof** Please refer to Sect. 8.6 for details.

Theorem 3.4 demonstrates that the number of measurements  $\mathcal{O}(s^2 \kappa^4 \max\{K, N\} \log^8 L)$  are sufficient for the Riemannian gradient descent algorithm (with spectral initialization), i.e., Algorithm 3.3, to linearly converge to the ground truth signals.

**Riemannian Trust-Region Algorithm** A scalable algorithm that enjoys superlinear convergence rate, i.e., the Riemannian trust-region algorithm, can be developed on the product manifolds to detect the descent direction  $\boldsymbol{\eta}$  [1, Section 7]. In order to parallelly search the descent direction on the horizontal space  $\mathcal{H}_{\mathbf{v}} \mathcal{M}^s$ , the method of searching the direction  $\boldsymbol{\eta}_{\mathbf{w}_k}$  on the horizontal space  $\mathcal{H}_{\mathbf{w}_k} \mathcal{M}$  is developed. At each iteration, define the point on the manifold as  $\mathbf{w}_k \in \mathcal{M}$ , a trust-region subproblem is described as follows [1]:

$$\begin{aligned} & \underset{\boldsymbol{\eta}_{\mathbf{w}_k}}{\text{minimize}} \quad m(\boldsymbol{\eta}_{\mathbf{w}_k}) \\ & \text{subject to} \quad g_{\mathbf{w}_k}(\boldsymbol{\eta}_{\mathbf{w}_k}, \boldsymbol{\eta}_{\mathbf{w}_k}) \leq \delta^2, \end{aligned} \quad (3.44)$$

where  $\boldsymbol{\eta}_{\mathbf{w}_k} \in \mathcal{H}_{\mathbf{w}_k} \mathcal{M}$ ,  $\delta$  is the trust-region radius, and the objective function is represented as

$$m(\boldsymbol{\eta}_{\mathbf{w}_k}) = g_{\mathbf{w}_k}(\boldsymbol{\eta}_{\mathbf{w}_k}, \text{grad}_{\mathbf{w}_k} f) + \frac{1}{2} g_{\mathbf{w}_k}(\boldsymbol{\eta}_{\mathbf{w}_k}, \text{Hess}_{\mathbf{w}_k} f[\boldsymbol{\eta}_{\mathbf{w}_k}]), \quad (3.45)$$

and  $\text{Hess}_{\mathbf{w}_k} f[\boldsymbol{\eta}_{\mathbf{w}_k}]$  and  $\text{grad}_{\mathbf{w}_k} f$  are the matrix representations of Riemannian Hessian and Riemannian gradient in the quotient space, respectively. In addition, the iterate being updated or maintained depends on whether the decrease of the

function  $m(\boldsymbol{\eta}_{\mathbf{w}_k})$  is satisfied or not [1, Section 7]. If the decrease is sufficient, the iterate is updated as

$$\mathcal{R}_{\mathbf{w}_k}(\boldsymbol{\eta}_{\mathbf{w}_k}) = \mathbf{w}_k + \boldsymbol{\eta}_{\mathbf{w}_k}. \quad (3.46)$$

Under the above framework, the Riemannian trust-region algorithm is parallelly developed on individual manifolds to solve problem (3.41).

### 3.4.4 Simulation Results

In this section, we compare three algorithms for estimating the blind demixing model: nuclear norm minimization (NNM) in Sect. 3.3.1, regularized Wirtinger flow (RGD) in Sect. 3.4.1, and Riemannian trust-region algorithm (RTR) in Sect. 3.4.3.

The ground-truth vectors, i.e.,  $\mathbf{x}_k \in \mathbb{C}^N$  and  $\mathbf{h}_k \in \mathbb{C}^K$  for  $k = 1, \dots, s$ , are generated as standard complex Gaussian vectors whose entries are drawn i.i.d. from the standard normal distribution. In addition, the relative construction error with respect to the rank-one matrices, i.e.,  $\mathbf{X}_i = \mathbf{h}_i \mathbf{x}_i^H$ , is adopted to evaluate the performance of the algorithms, given as

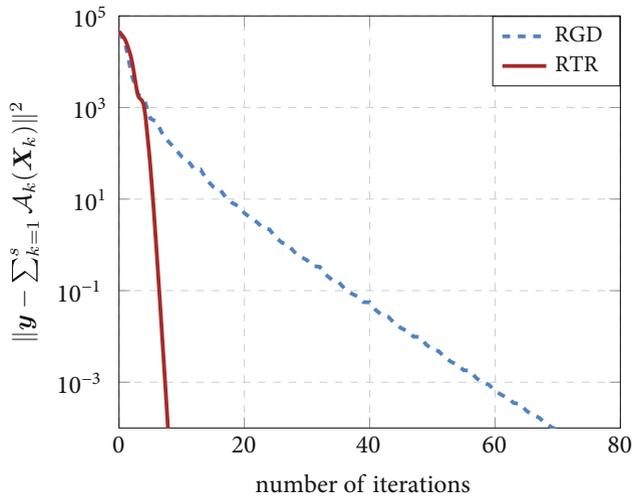
$$\text{err}(\mathbf{X}) = \frac{\sqrt{\sum_{k=1}^s \|\mathbf{X}_k - \hat{\mathbf{X}}_k\|_F^2}}{\sqrt{\sum_{k=1}^s \|\hat{\mathbf{X}}_k\|_F^2}}, \quad (3.47)$$

where  $\{\mathbf{X}_k\}$  are estimated matrices and  $\{\hat{\mathbf{X}}_k\}$  are ground truth matrices. The initialization strategy, i.e., Algorithm 3.1, is adopted for all the nonconvex optimization algorithms, i.e., RGD and RTR. The RTR algorithm stops when the norm of Riemannian gradient is less than  $10^{-8}$  or the number of iterations exceeds 500. The stopping criteria of RGD is adopted from the paper [10].

In the noiseless scenario, two nonconvex algorithms are compared under the setting of  $N = K = 50$ ,  $L = 1250$ , and  $s = 5$ . The convergence rates of nonconvex algorithms are illustrated in Fig. 3.6. In the noisy scenario, assume the additive noise term in (3.17) obeys

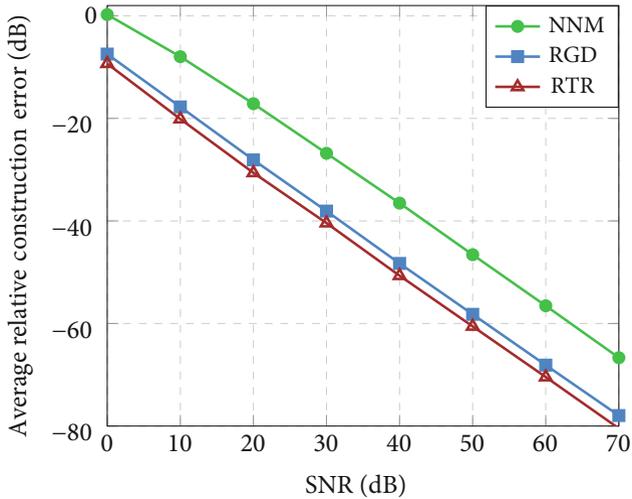
$$\mathbf{e} = \sigma \cdot \|\mathbf{y}\| \cdot \frac{\boldsymbol{\omega}}{\|\boldsymbol{\omega}\|}, \quad (3.48)$$

where  $\boldsymbol{\omega} \in \mathbb{C}^L$  denotes a standard complex Gaussian vector. Three algorithms with respect to different signal-to-noise ratios (SNR)  $\sigma$  are compared under the setting of  $L = 1500$ ,  $N = K = 50$ , and  $s = 2$ . In each circumstance, ten independent trails are simulated. Figure 3.7 shows the average relative construction error in dB against the signal-to-noise ratio (SNR). It concludes that the average relative construction



**Fig. 3.6** Convergence rate of nonconvex algorithms

error decreases as SNR increases, which demonstrates that RTR is robust to the noise.



**Fig. 3.7** Average relative construction error versus SNR (dB)

**Table 3.5** Summary of approaches for solving blind demixing problem

	Convex relaxation	Nonconvex approach			
Algorithm	NNM Sect. 3.3.1	RGD Sect. 3.4.1	RTR Sect. 3.4.3	WF Sect. 3.4.2	WF Sect. 3.4.2
Regularizer	×	√	×	×	×
Initialization	×	Spectral	Spectral	Spectral	Random
Condition number	1	$\kappa$	$\kappa$	$\kappa$	$\kappa$
Sample complexity	$m \geq Cs^2\mu^2$ $K \log^2 m$	$m \geq Cs^2\mu^2$ $\kappa^4 K \log^2 m$	$m \geq Cs^2\mu^2$ $\kappa^4 K \log^8 m$	$m \geq Cs^2\mu^2$ $\kappa^4 K \log^8 m$	$m \geq Cs^2\mu^2$ $\kappa^4 K \log^8 m$
Computational complexity	—	$\mathcal{O}\left(sm \log \frac{1}{\epsilon}\right)$	Not provided	$\mathcal{O}\left(s \log \frac{1}{\epsilon}\right)$	$\mathcal{O}\left(s \log K + s \log \frac{1}{\epsilon}\right)$

### 3.5 Summary

This chapter introduced a blind demixing model that facilitates to jointly decode data and estimate the channel state in IoT networks. The low-overhead communications can be achieved via the blind demixing model since it excludes the channel estimation sequence in the metadata. The convex relaxation method is introduced to solve the blind demixing problem based on its low-rank property. To further reduce the computational complexity, first-order algorithms, e.g., Wirtinger flow and regularized Wirtinger flow, have been developed. In addition, a Riemannian trust-region algorithm that enjoys faster convergence than the first-order algorithm has also been presented. The summary of both convex and nonconvex approaches for solving the blind demixing problem in the noiseless scenario is provided in Table 3.5.<sup>1</sup> State-of-the-art theoretical analysis is developed under the assumption of the Gaussian encoding matrices. It is intriguing to explore more general types of encoding matrices, e.g., sub-Gaussian matrices, in future works.

### References

1. Absil, P.A., Mahony, R., Sepulchre, R.: Optimization Algorithms on Matrix Manifolds. Princeton University Press, Princeton (2009)
2. Boumal, N., Mishra, B., Absil, P.A., Sepulchre, R.: Manopt, a Matlab toolbox for optimization on manifolds. *J. Mach. Learn. Res.* **15**, 1455–1459 (2014). <http://www.manopt.org>
3. Candes, E.J., Romberg, J.K., Tao, T.: Stable signal recovery from incomplete and inaccurate measurements. *Commun. Pure Appl. Math.* **59**(8), 1207–1223 (2006)
4. Candes, E.J., Strohmer, T., Vershynina, V.: PhaseLift: exact and stable signal recovery from magnitude measurements via convex programming. *Commun. Pure Appl. Math.* **66**(8), 1241–1274 (2013)
5. Chi, Y., Lu, Y.M., Chen, Y.: Nonconvex optimization meets low-rank matrix factorization: an overview. arXiv preprint. arXiv:1809.09573 (2018)
6. Dong, J., Shi, Y.: Nonconvex demixing from bilinear measurements. *IEEE Trans. Signal Process.* **66**(19), 5152–5166 (2018)
7. Dong, J., Shi, Y.: Blind demixing via Wirtinger flow with random initialization. In: Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS), vol. 89, pp. 362–370 (2019)
8. Dong, J., Yang, K., Shi, Y.: Blind demixing for low-latency communication. *IEEE Trans. Wireless Commun.* **18**(2), 897–911 (2019)
9. Ling, S., Strohmer, T.: Blind deconvolution meets blind demixing: algorithms and performance bounds. *IEEE Trans. Inf. Theory* **63**(7), 4497–4520 (2017)
10. Ling, S., Strohmer, T.: Regularized gradient descent: a nonconvex recipe for fast joint blind deconvolution and demixing. *Inf. Inference: J. IMA* **8**(1), 1–49 (2018)
11. Parvez, I., Rahmati, A., Guvenc, I., Sarwat, A.I., Dai, H.: A survey on low latency towards 5G: RAN, core network and caching solutions. arXiv preprint. arXiv:1708.02562 (2017)

---

<sup>1</sup>In Table 3.5, the parameters are consistent with the blind demixing model (3.13). The incoherence parameter  $\mu$  is mentioned in Table 3.1. We assume that  $K = N$  and there is some sufficient large constant  $C > 0$ .

# Chapter 4

## Sparse Blind Demixing



**Abstract** This chapter extends the models presented in Chaps. 2 and 3 to the scenario involving device activity detection. The new setting induces a sparse blind demixing model for developing methods for joint device activity detection, data decoding, and channel estimation in IoT networks. The signal model is first presented, in the scenario with either a single-antenna or multi-antenna BS. A convex relaxation approach is first introduced as a basic method to solve the nonconvex estimation problem. We further present a difference-of-convex-functions (DC) approach which turns out to be a powerful tool to solve the resulting sparse and low-rank optimization problem with matrix lifting. Furthermore, a smooth Riemannian optimization algorithm operating on the product manifold is introduced for solving the sparse blind demixing problem directly.

### 4.1 Joint Device Activity Detection, Data Decoding, and Channel Estimation

In Chap. 2, a sparse linear model has been developed for grant-free random access to jointly detect device activity and estimate channel state information. Under this scheme, pilot sequences are needed for activity detection, which lead to excess overhead for short packet communications. To avoid the transmission of pilot sequences, more powerful signal processing techniques are needed for data detection. Assuming the active device set is known, the blind demixing model has been introduced in Chap. 3 to achieve pilot-free communications in the massive IoT network via joint data decoding and channel estimation. To further account for the sporadic activity pattern in massive IoT networks, a *sparse blind demixing model* was proposed in [4, 6] to reduce the overhead during the transmission via joint device activity detection, data decoding, and channel estimation in a unified way.

Considering an IoT network containing one BS equipped with a single antenna, where only part (denoted as the set  $\mathcal{S}$ ) of the devices are active, the sparse blind demixing model represented in the Fourier domain is given as

$$y_j = \sum_{k \in \mathcal{S}} \mathbf{b}_j^H \mathbf{h}_k \mathbf{x}_k^H \mathbf{a}_{kj}, \quad 1 \leq j \leq m, \quad (4.1)$$

where

$$\mathbf{y} = [y_1, \dots, y_m]^T \in \mathbb{C}^m$$

is the received signal at the BS represented in the Fourier domain,  $\{\mathbf{b}_j\}$  are design vectors that indicate the Fourier transform operation,  $\{\mathbf{a}_{kj}\}$  are design vectors that indicate the encoding procedure, and  $\{\mathbf{h}_k\}$ ,  $\{\mathbf{x}_k\}$  are channel signals and data signals, respectively. By detecting the active set  $\mathcal{S}$  and the vectors  $\{\mathbf{h}_k\}$ ,  $\{\mathbf{x}_k\}$  for  $k \in \mathcal{S}$  from the observation  $\mathbf{y}$ , device activity detection, data decoding, and channel estimation can be simultaneously achieved. This is a highly challenging problem.

In the sequel, we first introduce the problem formulation of the sparse blind demixing model. Various approaches for solving the corresponding nonconvex estimation problem are then introduced: (1) a convex relaxation approach based on the minimization of nuclear norms and  $\ell_1/\ell_2$ -norms, (2) a difference-of-convex (DC) function approach based on the minimization of DC objective functions. Along the discussion, we also identify theoretical analysis for the sparse blind demixing model as future research directions.

## 4.2 Problem Formulation

In this section, we present problem formulation for joint activity detection, data decoding, and channel estimation for both scenarios of single-antenna and multi-antenna BSs. Considering an IoT network consisting of one BS and  $s$  single-antenna devices with sporadic traffic, in each coherence block, only an unknown subset of devices are active, defined as  $\mathcal{S} \subseteq \{1, 2, \dots, s\}$ .

### 4.2.1 Single-Antenna Scenario

In the single-antenna BS scenario, the problem formulation of the sparse blind demixing model can be derived from the blind demixing model mentioned in Sect. 3.2 with an additional consideration of the sparse activity pattern. The data signal transmitted by the  $k$ -th user is denoted as  $\mathbf{x}_k^{\dagger} \in \mathbb{C}^N$ . Assume that an encoding matrix over the  $m$  time slots is assigned to each device  $k$ . Over  $m$  time slots, the

received signals at the BS in the frequency domain are presented as [3, 9]

$$y_j = \sum_{k \in \mathcal{S}} \mathbf{b}_j^H \mathbf{h}_k^{\natural} \mathbf{x}_k^{\natural*} \mathbf{a}_{kj} + e_j, \quad 1 \leq j \leq m, \quad (4.2)$$

which resembles the blind demixing model defined in (3.13) as presented in Sect. 3.2.3. From the observation  $y_j$  for  $1 \leq j \leq m$ , the active set  $\mathcal{S}$ , data information  $\{\mathbf{x}_k\}$ , and the channel state information  $\{\mathbf{h}\}$  can be recovered. Hence, joint device activity detection, data decoding, and channel estimation can be achieved.

## 4.2.2 Multiple-Antenna Scenario

Considering an IoT network consisting of a BS equipped with  $r$  antennas and  $s$  single-antenna devices with sporadic traffic. Denote  $\mathbf{g}_{ij}^{\natural} \in \mathbb{C}^m$  as the channel impulse response from the  $j$ -th device to the  $i$ -th antenna of the BS and recall the transmitted signal at the  $j$ -th device defined in (3.6). Thus, the observations  $\mathbf{z}_i \in \mathbb{C}^m$  at the  $i$ -th antenna of the BS are represented as

$$\mathbf{z}_i = \sum_{j \in \mathcal{S}} \mathbf{f}_j^{\natural} * \mathbf{g}_{ij}^{\natural} + \mathbf{n}_i, \quad \forall i = 1 \dots r, \quad (4.3)$$

where  $\mathbf{n}_i \in \mathbb{C}^m$  is additive white complex Gaussian noise. The sparse blind demixing model in the single-antenna scenario is the specific case of (4.3) when  $r = 1$ . Given the observations  $\{\mathbf{z}_i\}$ , our goal is to detect the active device set  $\mathcal{S}$  and recover the associated  $\{\mathbf{f}_j^{\natural}\}$  and  $\{\mathbf{g}_{ij}^{\natural}\}$  simultaneously.

Similar to the model in the single-antenna scenario, i.e., (4.2), the  $l$ -th entry of  $y_i$  is given by

$$y_i[l] = \sum_{j \in \mathcal{S}} \mathbf{b}_l^H \mathbf{h}_{ij}^{\natural} \mathbf{x}_j^{\natural*} \mathbf{c}_{jl} + \xi_i[l], \quad l = 1, \dots, m, \quad i = 1, \dots, r. \quad (4.4)$$

The goal is to simultaneously detect active device set  $\mathcal{S}$  and recover both  $\{\mathbf{x}_j^{\natural}\}$  and  $\{\mathbf{h}_{ij}^{\natural}\}$  from the observations  $\{\mathbf{z}_i\}$ .

## 4.3 Convex Relaxation Approach

In this section, we present a convex relaxation approach to solve the sparse blind demixing problem. Taking the single-antenna scenario as an example, the optimization problem is firstly established for the sparse blind demixing model (4.2). Then a convex relaxation approach is further presented to solve this resulting nonconvex optimization problem.

Define a collection of groups as

$$\mathcal{G} = \{\mathcal{G}_1, \mathcal{G}_2, \dots, \mathcal{G}_s\} \quad (4.5)$$

with

$$\mathcal{G}_k = \{N(k-1) + 1, \dots, Nk\}$$

and

$$\mathcal{G}_i \cap \mathcal{G}_j = \emptyset$$

for  $i \neq j$ , and denote an aggregative vector as

$$\mathbf{x} = [\mathbf{x}_1^\top, \dots, \mathbf{x}_s^\top]^\top \in \mathbb{C}^{Ns},$$

where the index set is

$$\mathcal{V} = \{1, 2, \dots, Ns\}.$$

With the support of the data vector defined as

$$\text{Supp}(\mathbf{x}) = \{i | x_i \neq 0, \forall i \in \mathcal{V}\},$$

the sparse blind demixing problem can be formulated as

$$\begin{aligned} & \underset{\{\mathbf{x}_k\}, \{\mathbf{h}_k\}}{\text{minimize}} && \sum_{k=1}^s \mathbb{I}(\text{Supp}(\mathbf{x}) \cap \mathcal{G}_k \neq \emptyset) \\ & \text{subject to} && \sum_{j=1}^m \left| \sum_{k=1}^s \mathbf{b}_j^H \mathbf{h}_k \mathbf{x}_k^H \mathbf{a}_{kj} - y_j \right|^2 \leq \epsilon, \end{aligned} \quad (4.6)$$

where parameter  $\epsilon > 0$  is known a priori. Denoting  $\mathbf{x}^*$  as a solution of problem (4.6), the set of active devices is given as

$$\mathcal{S}^* = \{k : \text{Supp}(\mathbf{x}) \cap \mathcal{G}_k \neq \emptyset\}.$$

Due to the nonconvex bilinear constraint and the combinatorial objective function, problem (4.6) is highly intractable, which motivates to develop efficient algorithms with good performance.

A natural way is to lift the bilinear model into the linear model with a low-rank matrix [9], i.e.,

$$\mathbf{b}_j^H \mathbf{h}_k \mathbf{x}_k^H \mathbf{a}_{kj} = \mathbf{b}_j^H \mathbf{W}_k \mathbf{a}_{kj} \quad (4.7)$$

with  $\mathbf{W}_k \in \mathbb{C}^{K \times N}$  and

$$\text{rank}(\mathbf{W}_k) = 1, \forall k = 1, \dots, s.$$

The natural idea is to exploit a convex relaxation method to deal with the sparsity and low-rankness in matrices  $\mathbf{W}_k$ 's of problem (4.6):

$$\begin{aligned} & \underset{\{\mathbf{W}_k\}}{\text{minimize}} && \lambda_1 \sum_{k=1}^s \|\mathbf{W}_k\|_* + \lambda_2 \sum_{k=1}^s \|\mathbf{W}_k\|_F \\ & \text{subject to} && \sum_{j=1}^m \left| \sum_{k=1}^s \mathbf{b}_j^H \mathbf{W}_k \mathbf{a}_{kj} - y_j \right|^2 \leq \epsilon, \end{aligned} \quad (4.8)$$

where  $\lambda_1 \geq 0$  and  $\lambda_2 \geq 0$  are the regularization parameters. The group sparsity structure in the aggregated data signals  $\mathbf{x}$  induces a group sparsity structure in the lifting vector

$$\text{vec}(\mathbf{W}) = [\text{vec}(\mathbf{W}_1)^H, \dots, \text{vec}(\mathbf{W}_s)^H]^H \in \mathbb{C}^{KNs},$$

where  $\text{vec}(\mathbf{M})$  is the vectorization of matrix  $\mathbf{M}$ . Furthermore, the  $\ell_1/\ell_2$ -norm is adopted to induce the group sparsity in the vector  $\text{vec}(\mathbf{W})$ , i.e.,

$$\|\text{vec}(\mathbf{W})\|_{1,2} = \sum_{k=1}^s \|\text{vec}(\mathbf{W}_k)\|_2 = \sum_{k=1}^s \|\mathbf{W}_k\|_F.$$

#### 4.4 Difference-of-Convex-Functions (DC) Programming Approach

Although the convex relaxation approach (4.8) provides a natural way to solve problem (4.6), the results obtained from norm relaxation are usually suboptimal to the original nonconvex optimization problem [10]. Moreover, two regularization parameters are introduced by the combination of norms, which are difficult to tune. Additionally, there is no efficient convex relaxation approach to simultaneously induce low-rankness and sparsity [2]. To address these issues, the paper [6] developed a difference-of-convex-functions (DC) representation for the rank function in order to satisfy the fixed-rank constraint.

In the sequel, we consider the sparse blind demixing model under the multiple-antenna BS scenario. Specifically, the sparse blind demixing problem is reformulated as a sparse and low-rank matrix recovery problem via lifting the bilinear model into the linear model. Based on the linear model, an exact DC formulation for the rank constraint is further established, followed by developing an efficient DC algorithm (DCA) for minimizing the DC objective.

#### 4.4.1 Sparse and Low-Rank Optimization

The estimation problem for sparse blind demixing with a multiple-antenna BS can be established in the similar form of the optimization problem (4.6). To facilitate the design of the DC algorithm, a sparse and low-rank optimization problem is first established. Denote

$$\mathbf{h}_j^{\natural} = [\mathbf{h}_{1j}^{\natural\text{H}}, \dots, \mathbf{h}_{rj}^{\natural\text{H}}]^{\text{H}}, \forall j = 1, \dots, s, \quad (4.9)$$

where  $\mathbf{h}_j^{\natural} \in \mathbb{C}^{rk}$ . Define a set of matrices

$$\mathbf{X}_{ij}^{\natural} = \mathbf{h}_{ij}^{\natural} \mathbf{x}_j^{\natural\text{H}}, \quad (4.10)$$

where  $\mathbf{X}_{ij}^{\natural} \in \mathbb{C}^{k \times d}$ . Here,  $\mathbf{X}_{ij}^{\natural}$  is a rank-one matrix when  $j \in \mathcal{S}$ , otherwise a zero matrix. Define

$$\mathbf{X}_j^{\natural} = [\mathbf{X}_{1j}^{\natural\text{H}}, \mathbf{X}_{2j}^{\natural\text{H}}, \dots, \mathbf{X}_{rj}^{\natural\text{H}}]^{\text{H}} = \mathbf{h}_j^{\natural} \mathbf{x}_j^{\natural\text{H}}, \quad (4.11)$$

where  $\mathbf{X}_j^{\natural} \in \mathbb{C}^{rk \times d}$ .  $\mathbf{X}_j^{\natural}$  is a rank-one matrix when  $j \in \mathcal{S}$ ; otherwise it is a zero matrix. With a matrix defined as  $\mathbf{E}_i \in \mathbb{R}^{k \times rk}$

$$\mathbf{E}_i = [\mathbf{e}_{k(i-1)+1}, \mathbf{e}_{k(i-1)+2}, \dots, \mathbf{e}_{ki}]^{\text{H}}, \forall i = 1, \dots, r, \quad (4.12)$$

where  $\mathbf{e}_l$  denotes the  $rk$ -dimensional standard basis vector, a linear map  $\mathcal{A}_{ij} : \mathbb{C}^{rk \times d} \rightarrow \mathbb{C}^m$  for  $1 \leq i \leq r, 1 \leq j \leq s$  is given by

$$\mathcal{A}_{ij}(\mathbf{Z}) := \{\langle \mathbf{b}_l \mathbf{c}_{jl}^{\text{H}}, \mathbf{E}_i \mathbf{Z} \rangle\}_{l=1}^m, \quad (4.13)$$

where  $\mathbf{Z} \in \mathbb{C}^{rk \times d}$  and  $\mathbf{E}_i \mathbf{X}_j^{\natural} = \mathbf{X}_{ij}^{\natural}$ . Thus, the model (4.4) can be transformed into

$$\mathbf{y}_i = \sum_{j=1}^s \mathcal{A}_{ij}(\mathbf{X}_j^{\natural}) + \boldsymbol{\xi}_i, \forall i = 1, \dots, r. \quad (4.14)$$

The measurements  $\{y_i\}$  are the linear combinations of the corresponding entries of every column block in the lifted matrix  $\mathbf{X}^\natural \in \mathbb{C}^{rk \times ds}$ , which is given by

$$\mathbf{X}^\natural = \begin{bmatrix} \mathbf{h}_{11}^\natural \mathbf{x}_1^{\natural H} & \mathbf{h}_{12}^\natural \mathbf{x}_2^{\natural H} & \cdots & \mathbf{h}_{1s}^\natural \mathbf{x}_s^{\natural H} \\ \mathbf{h}_{21}^\natural \mathbf{x}_1^{\natural H} & \mathbf{h}_{22}^\natural \mathbf{x}_2^{\natural H} & \cdots & \mathbf{h}_{2s}^\natural \mathbf{x}_s^{\natural H} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{h}_{r1}^\natural \mathbf{x}_1^{\natural H} & \mathbf{h}_{r2}^\natural \mathbf{x}_2^{\natural H} & \cdots & \mathbf{h}_{rs}^\natural \mathbf{x}_s^{\natural H} \end{bmatrix} = [\mathbf{X}_1^\natural, \dots, \mathbf{X}_s^\natural].$$

Instead of recovering both  $\{\mathbf{h}_{ij}^\natural\}$  and  $\{\mathbf{x}_j^\natural\}$ , problem  $\mathcal{P}$  is solved with the recovery of the matrix  $\mathbf{X}^\natural$ . Notice that  $\mathbf{X}^\natural$  has block-low-rank and column sparse structures. The goal is to recover  $\mathbf{X}^\natural$  from the observation  $y_i$  for  $i = 1, \dots, r$ . Since  $\mathbf{X}^\natural$  has block-low-rank and column sparse structures, we can establish a sparse and low-rank optimization problem as follows:

$$\begin{aligned} & \underset{\{\mathbf{X}_j\}}{\text{minimize}} \quad \left\| [\|\text{vec}(\mathbf{X}_1)\|_2, \dots, \|\text{vec}(\mathbf{X}_j)\|_2] \right\|_0 \\ & \text{subject to} \quad \sum_{i=1}^r \left\| \mathbf{y}_i - \sum_{j=1}^s \mathcal{A}_{ij}(\mathbf{X}_j) \right\|_2^2 \leq \epsilon \\ & \quad \text{rank}(\mathbf{X}_j) \leq 1, \forall j = 1, \dots, s, \end{aligned} \quad (4.15)$$

where  $\{\mathbf{X}_j\} \in \mathbb{C}^{rk \times d}$ .

#### 4.4.2 A DC Formulation for Rank Constraint

Before giving an exact DC formulation for the rank constraint, we introduce the definition of Ky Fan  $k$ -norm.

**Definition 4.1** Ky Fan  $k$ -norm [7]: the Ky Fan  $k$ -norm of a matrix  $\mathbf{X} \in \mathbb{C}^{m \times n}$  is defined as the sum of its largest- $k$  singular values, i.e.,

$$\|\|\mathbf{X}\|_k = \sum_{i=1}^k \sigma_i(\mathbf{X}), \quad (4.16)$$

where  $k \leq \min\{m, n\}$ .

Since the rank of a matrix is equal to the number of its nonzero singular values, for any matrix  $\mathbf{X} \in \mathbb{C}^{m \times n}$  whose rank is less than  $k$ , it can yield from Definition 4.1 that [7]:

$$\text{rank}(\mathbf{X}) \leq k \Leftrightarrow \|\mathbf{X}\|_* - \|\mathbf{X}\|_k = 0. \quad (4.17)$$

Instead of using the discontinuous rank function, a continuous DC function  $\|\mathbf{X}\|_* - \|\mathbf{X}\|_k$  can be adopted for inducing low-rankness property of a matrix.

### 4.4.3 DC Algorithm for Minimizing a DC Objective

Based on (4.17), problem (4.15) can be further formulated as the minimization problem with a DC objective function:

$$\begin{aligned} & \underset{\{\mathbf{X}_j\}}{\text{minimize}} \quad \sum_{j=1}^s \left( \|\mathbf{X}_j\|_* - \|\mathbf{X}_j\|_1 \right) \\ & \text{subject to} \quad \sum_{i=1}^r \left\| \mathbf{y}_i - \sum_{j=1}^s \mathcal{A}_{ij}(\mathbf{X}_j) \right\|_2^2 \leq \epsilon. \end{aligned} \quad (4.18)$$

To address the nonconvexity of the DC objective function, a DC algorithm based on majorization-minimization (MM) has been proposed in [12]. At each iteration, the DC algorithm solves a convex subproblem, given by

$$\begin{aligned} & \underset{\{\mathbf{X}_j\}}{\text{minimize}} \quad \sum_{j=1}^s \left( \|\mathbf{X}_j\|_* - \langle \mathbf{X}_j, \mathbf{Y}_j^{t-1} \rangle \right) \\ & \text{subject to} \quad \sum_{i=1}^r \left\| \mathbf{y}_i - \sum_{j=1}^s \mathcal{A}_{ij}(\mathbf{X}_j) \right\|_2^2 \leq \epsilon, \end{aligned} \quad (4.19)$$

where  $\mathbf{Y}_j^{t-1} \in \mathbb{C}^{r \times d}$  is a subgradient of  $\|\mathbf{X}_j\|_1$  at  $\mathbf{X}_j^{t-1}$  and can be efficiently derived from the singular value decomposition, given by

$$\partial \|\mathbf{X}_j^t\|_1 = \{ \mathbf{U} \text{diag}(\mathbf{q}) \mathbf{V}^H : \mathbf{q} = [1, 0, \dots, 0] \}. \quad (4.20)$$

The DC algorithm is illustrated in Algorithm 4.1.

---

**Algorithm 4.1: DC algorithm for problem (4.18)**


---

Input:  $\{\mathcal{A}_{ij}\}, \{\mathbf{y}_i\}$ , upper bound  $\epsilon$ , a small value  $\eta$ Output:  $\{\mathbf{X}_j^t\}$ 

```

Initialisation :  $\{\mathbf{X}_j^0\}$ 
1:  $k = 1$ 
   LOOP Process
2: for  $t = 1, 2, \dots$  do
3:   Select  $\{\mathbf{Y}_j^{t-1} \in \partial \|\mathbf{X}_j^{t-1}\|_k\}$ 
4:   Solve the convex problem (4.19), and obtain the
      optimal solution  $\{\mathbf{X}_j^t\}$ 
5:   if  $\sum_{j=1}^S (\|\mathbf{X}_j^t\|_* - \|\|\mathbf{X}_j^t\|\|_1) < \eta$  then
6:     break
7:   end if
8: end for
9: return  $\{\mathbf{X}_j^t\}$ 

```

---

#### 4.4.4 Simulations

In this section, we conduct numerical experiments to compare the proposed DC approach with the convex relaxation methods for empirical recovery performance and test the robustness against noise.

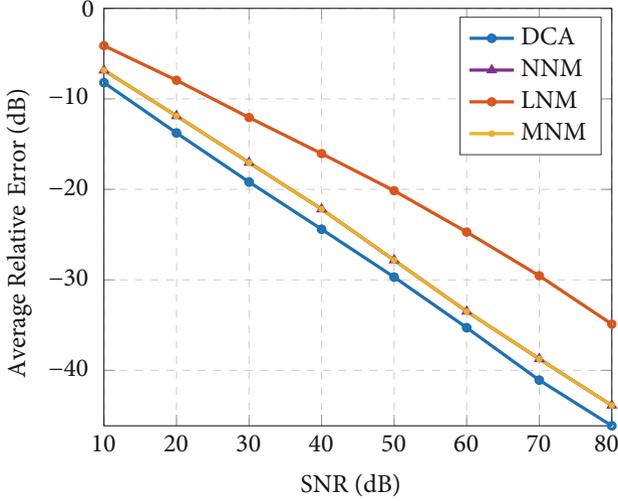
For  $j \notin \mathcal{S}$ , set the ground truth data signal as  $\mathbf{x}_j^{\natural} = \mathbf{0}_d$ , and for  $j \in \mathcal{S}$ ,  $\mathbf{x}_j^{\natural}$  is drawn i.i.d. from the standard complex Gaussian distribution. Both the channel states  $\{\mathbf{h}_{ij}^{\natural}\}$  and matrices  $\{\mathbf{C}_j\}$  are drawn i.i.d. from the standard complex Gaussian distribution. To measure the accuracy of estimation, the relative construction error is defined as

$$\text{error}(\mathbf{X}) = \frac{\sqrt{\sum_{j=1}^S \|\mathbf{X}_j^{\natural} - \mathbf{X}_j\|_F^2}}{\sqrt{\sum_{j=1}^S \|\mathbf{X}_j^{\natural}\|_F^2}}, \quad (4.21)$$

where the ground truth matrices are denoted as  $\{\mathbf{X}_j^{\natural}\}$ , and  $\{\mathbf{X}_j\}$  are the estimated matrices.

We compare the empirical recovery and robustness performance of the following four algorithms:

- **DC algorithm (DCA)**: The termination criterion is either the iteration number exceeding 200 or  $\sum_{j=1}^S (\|\mathbf{X}_j^t\|_* - \|\|\mathbf{X}_j^t\|\|_1) < 10^{-6}$ .
- **mixed norm minimization(MNM)**: The regularization terms  $\lambda_1$  and  $\lambda_1$  are chosen via cross validation.



**Fig. 4.1** Probability of successful recovery with different sample sizes  $m$

- **nuclear norm minimization (NNM)**: The algorithm is similar to **MNM** except for  $\lambda_1 = 1, \lambda_2 = 0$ .
- **$\ell_1/\ell_2$ -norm minimization(LNM)**: The algorithm is similar to **MNM** except for  $\lambda_1 = 0, \lambda_2 = 1$ .

The empirical recovery performance of the above four algorithms is investigated in the noiseless scenario under the setting of  $k = 5, d = 20, s = 10, |\mathcal{S}| = 4$ , and  $r = 3$ . For each setting, 20 independent trails are performed and the recovery is regarded as a success if the  $\text{error}(\mathbf{X}) < 10^{-2}$ . Figure 4.1 shows the performance of recovery with varying the number of measurements.

The robustness of the four algorithms with respect to noise is further investigated. The noise  $\xi_i$  is generated as

$$\xi_i = \sigma \cdot \|y_i\| \cdot \frac{z_i}{\|z_i\|}, \quad \forall i = 1 \dots r, \quad (4.22)$$

where  $z_i \in \mathbb{C}^m$  is the normalized standard Gaussian vector. Under the setting of  $k = 5, d = 20, s = 10, |\mathcal{S}| = 4, r = 3$ , and  $m = 670$ , 20 independent trails are performed with respect to different  $\sigma$ . Figure 4.2 illustrates the average relative error in dB against the signal-to-noise-ratio (SNR). It shows that DCA enjoys a higher accuracy of reconstruction than other algorithms.

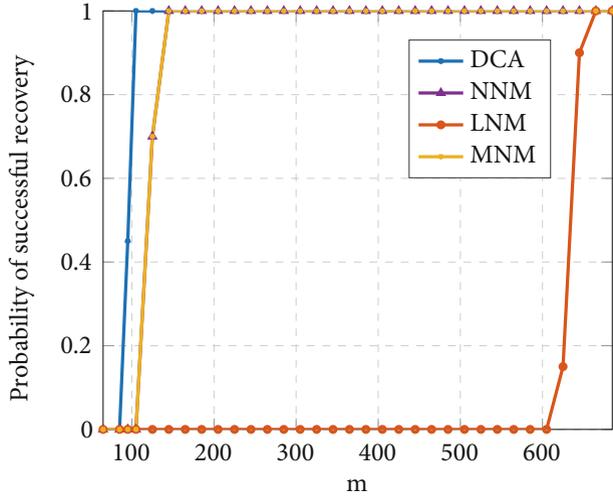


Fig. 4.2 Robustness under different SNR(dB)

## 4.5 Smoothed Riemannian Optimization on Product Manifolds

Another line of literatures have developed efficient nonconvex algorithms to solve the sparse and low-rank optimization problem [8, 13] in the natural vector space via matrix factorization. For instance, an alternating minimization approach was developed in [8] for solving the sparse blind deconvolution problem. However, the additional group sparsity structure of the sparse blind demixing problem (4.6) brings unique challenges to develop the nonconvex optimization paradigm. To address this challenge, a smoothed Riemannian optimization approach is introduced to solve sparse blind demixing problem [4], thereby achieving better performance with low computational complexity. More details on the manifold optimization can be referred to Sect. 3.4.3.

### 4.5.1 Optimization on Product Manifolds

To begin with, problem (4.6) is formulated as a regularized optimization problem under fixed-rank constraints. For  $k = 1, \dots, s$ ,  $j = 1, \dots, m$ , define

$$\mathbf{c}_j = [\mathbf{b}_j^H, \mathbf{0}_N^H]^H \in \mathbb{C}^{N+K}, \quad \mathbf{d}_{kj} = [\mathbf{0}_K^H, \mathbf{a}_{kj}^H]^H \in \mathbb{C}^{N+K},$$

it yields

$$\mathbf{c}_j^H \mathbf{M}_k \mathbf{d}_{kj} = \mathbf{b}_j^H \mathbf{h}_k \mathbf{x}_k^H \mathbf{a}_{kj}, \quad (4.23)$$

where

$$\mathbf{M}_k = \mathbf{w}_k \mathbf{w}_k^H \in \mathbb{S}_+^{N+K} \quad (4.24)$$

is a Hermitian positive semidefinite matrix with

$$\mathbf{w}_k = [\mathbf{h}_k^H, \mathbf{x}_k^H]^H \in \mathbb{C}^{N+K}. \quad (4.25)$$

Hence, problem (4.6) can be represented as the optimization problem on the product of Hermitian positive semidefinite matrices:

$$\begin{aligned} & \underset{\mathbf{M}}{\text{minimize}} \quad \sum_{j=1}^m \left| \sum_{k=1}^s \mathbf{c}_j^H \mathbf{M}_k \mathbf{d}_{kj} - y_j \right|^2 + \lambda f(\mathbf{M}) \\ & \text{subject to} \quad \text{rank}(\mathbf{M}_k) = 1, \quad k = 1, \dots, s, \end{aligned} \quad (4.26)$$

where  $\mathbf{M} = \{\mathbf{M}_k\}_{k=1}^s$  with  $\mathbf{M}_k \in \mathbb{S}_+^{N+K}$ ,  $\lambda > 0$  is the regularization parameter, and  $f(\mathbf{M})$  is the function to induce the sparsity structure. Here,  $\mathbf{M}_k$  is in the space of the manifold encoded by complex symmetric rank-one matrices, i.e.,  $\mathbf{M}_k \in \mathcal{M}_k$  [5]. It yields that  $\mathbf{M} \in \mathcal{M}^s$ , where

$$\mathcal{M}^s := \mathcal{M}_1 \times \mathcal{M}_2 \times \dots \times \mathcal{M}_s \quad (4.27)$$

represents the product of manifolds  $\mathcal{M}_k$ . By exploiting the quotient manifold geometry of the product of complex symmetric rank-one matrices, computationally efficient Riemannian optimization algorithms can be developed on product manifolds.

## 4.5.2 Smoothed Riemannian Optimization

The smooth objective function is normally required [1, 11] in order to develop Riemannian optimization algorithm for solving problem (4.26). To achieve this goal, the smoothed  $\ell_1/\ell_2$ -norm is introduced, represented as

$$f_\epsilon(\mathbf{M}) = \sum_{k=1}^s \left( \|\mathbf{M}_k\|_F^2 + \epsilon^2 \right)^{1/2} \quad (4.28)$$

with  $\epsilon > 0$  as the smoothing parameter with a small value. This can be used for inducing the group sparsity structure in vector

$$\text{vec}(\mathbf{M}) = [\text{vec}(\mathbf{M}_1)^H, \dots, \text{vec}(\mathbf{M}_s)^H]^H. \quad (4.29)$$

Therefore, the proposed smoothed Riemannian optimization approach over the product manifold  $\mathcal{M}^s$  for sparse blind demixing problem (4.6) is given by

$$\underset{\mathbf{M} \in \mathcal{M}^s}{\text{minimize}} \sum_{j=1}^m \left| \sum_{k=1}^s \mathbf{c}_j^H \mathbf{M}_k \mathbf{d}_{kj} - y_j \right|^2 + \lambda f_\epsilon(\mathbf{M}), \quad (4.30)$$

where the objective function is smooth and the constraint is a manifold.

Due to the geometry of the product manifolds, the Riemannian optimization algorithms developed on the product manifold  $\mathcal{M}^s$  can be elementwisely operated over the individual manifolds  $\mathcal{M}_k$  [5]. For individual manifold  $\mathcal{M}_k$ , the descent direction is detected on the horizontal space of the manifold and then retract it on the manifold via retraction operation. Therein, the detection of the descent direction can be achieved by the Riemannian optimization algorithms, e.g., conjugate gradient descent algorithm [1].

### 4.5.3 Simulation Results

In this section, to illustrate the advantages of the smoothed Riemannian optimization for solving the sparse blind demixing problem (4.30), the Riemannian conjugate-gradient descent algorithm (RCGD) is compared with the other three algorithms mentioned in Sect. 4.4.4, i.e., nuclear norm minimization (NNM),  $\ell_1/\ell_2$ -norm minimization (LNM), and mixed norm minimization (MNM). Here, the RCGD algorithm adopts the initialization strategy in [5] and stops when the norm of Riemannian gradient falls below  $10^{-8}$  or the number of iterations exceeds 500.

The empirical recovery performance of the above four algorithms, i.e., RCGD, NNM, LMN, and MNM, are investigated under the setting of  $N = K = 10$ ,  $s = 10$ ,  $|\mathcal{A}| = 3$ . For each setting, 30 independent trails are performed and the recovery is considered as a success if  $\text{err}(\mathbf{x}) \leq 10^{-2}$ . Figure 4.3 illustrates the probability of successful recovery with respect to different sample sizes  $m$ . It shows that the smoothed Riemannian optimization algorithm achieves much better performance than other algorithms. That is, it exactly recovers the ground truth signals with less samples.

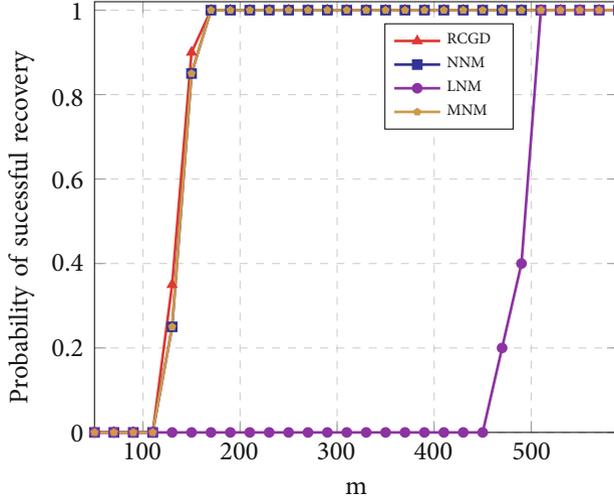


Fig. 4.3 Probability of successful recovery with different sample sizes  $m$

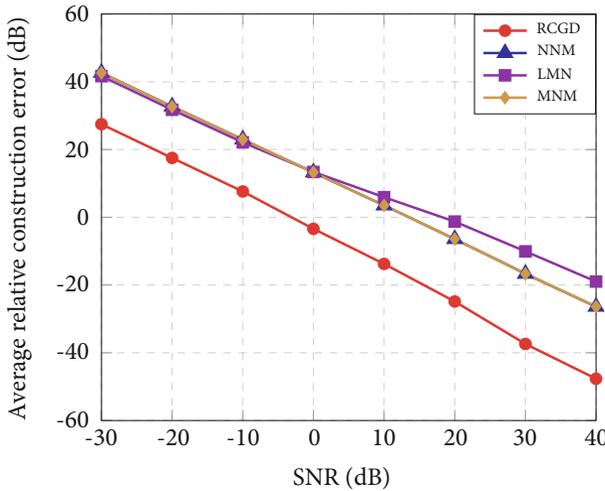


Fig. 4.4 Average relative construction error vs. SNR (dB)

The average relative construction error of the four algorithms is further investigated to explore the robustness of the proposed smoothed Riemannian optimization algorithm against additive noise. The four algorithms for each level of signal-to-noise ratio (SNR)  $1/\sigma$  are compared in the setting of  $m = 550$ ,  $N = K = 10$ ,  $s = 10$ ,  $|\mathcal{A}| = 3$ . For each setting, 20 independent trials are performed. The average relative construction error in dB against the SNR is showed in Fig. 4.4, which demonstrates that RCGD is robust to the noise and can achieve better performance than other algorithms.

## 4.6 Summary

This chapter introduced a sparse blind demixing model with both single-antenna and multi-antenna BSs for joint device activity detection, data decoding, and channel estimation in IoT networks with the grant-free random access scheme. It enjoys attractive advantages by removing the overhead caused by channel estimation sequence and device activity information. According to the simultaneous group sparse and low-rank variables in the sparse blind demixing model, the convex relaxation approach based on the norm minimization was first introduced. To further pursue higher accuracy of signal reconstruction compared to the convex relaxation approach, the approach that minimizes the difference-of-convex (DC) objective functions was developed. Another line of works has been focused on establishing Riemannian manifold to characterize the structured variables in the sparse blind demixing model. It is also interesting to further investigate the geometry property, i.e., group sparsity and low-rankness, of the sparse blind demixing model, thereby facilitating to design efficient algorithms with satisfactory performance, i.e., low sample complexity or high accuracy of estimation. A rigorous theoretical analysis on the sparse blind demixing problem is also of interest for future study, to characterize the number of measurements required for exact recovery.

## References

1. Absil, P.A., Mahony, R., Sepulchre, R.: Optimization Algorithms on Matrix Manifolds. Princeton University Press, Princeton (2009)
2. Aghasi, A., Bahmani, S., Romberg, J.: A tightest convex envelope heuristic to row sparse and rank one matrices. In: Proceedings of the IEEE Global Conference on Signal and Information Processing (GlobalSIP), p. 627. IEEE, Piscataway (2013)
3. Dong, J., Shi, Y.: Nonconvex demixing from bilinear measurements. *IEEE Trans. Signal Process.* **66**(19), 5152–5166 (2018)
4. Dong, J., Shi, Y., Ding, Z.: Sparse blind demixing for low-latency signal recovery in massive IoT connectivity. In: Proceedings of the IEEE International Conference on Acoustics Speech Signal Processing (ICASSP), pp. 4764–4768. IEEE, Piscataway (2019)
5. Dong, J., Yang, K., Shi, Y.: Blind demixing for low-latency communication. *IEEE Trans. Wireless Commun.* **18**(2), 897–911 (2019)
6. Fu, M., Dong, J., Shi, Y.: Sparse blind demixing for low-latency wireless random access with massive connectivity. In: Proceedings of the IEEE Vehicular Technology Conference (VTC), pp. 4764–4768. IEEE, Piscataway (2019)
7. Gotoh, J.Y., Takeda, A., Tono, K.: DC formulations and algorithms for sparse optimization problems. *Math. Program.* **169**(1), 141–176 (2018)
8. Lee, K., Wu, Y., Bresler, Y.: Near-optimal compressed sensing of a class of sparse low-rank matrices via sparse power factorization. *IEEE Trans. Inf. Theory* **64**(3), 1666–1698 (2018)
9. Ling, S., Strohmer, T.: Blind deconvolution meets blind demixing: algorithms and performance bounds. *IEEE Trans. Inf. Theory* **63**(7), 4497–4520 (2017)
10. Lu, C., Tang, J., Yan, S., Lin, Z.: Nonconvex nonsmooth low rank minimization via iteratively reweighted nuclear norm. *IEEE Trans. Image Process.* **25**(2), 829–839 (2016)

11. Shi, Y., Mishra, B., Chen, W.: Topological interference management with user admission control via Riemannian optimization. *IEEE Trans. Wireless Commun.* **16**(11), 7362–7375 (2017)
12. Tao, P.D., An, L.T.H.: Convex analysis approach to DC programming: theory, algorithms and applications. *Acta Math. Vietnam.* **22**(1), 289–355 (1997)
13. Zhang, Y., Kuo, H.W., Wright, J.: Structured local optima in sparse blind deconvolution (2018). Preprint. arXiv: 1806.00338

# Chapter 5

## Shuffled Linear Regression



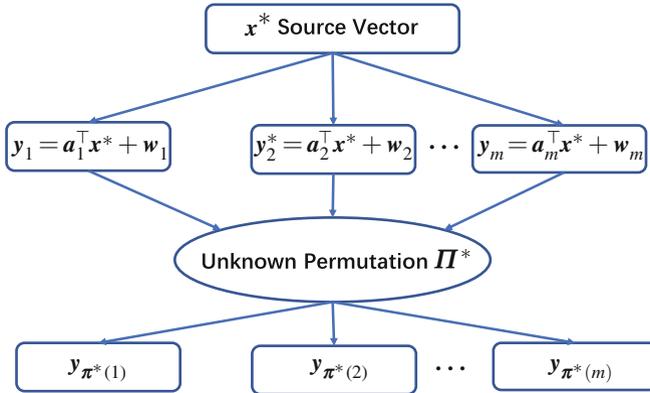
**Abstract** In this chapter, we shall introduce a shuffled linear regression model for joint data decoding and device identification in IoT networks. It is first formulated as a maximum likelihood estimation (MLE) problem. To solve this MLE problem, two algorithms are presented: one is based on sorting, and the other algorithm returns an approximate solution to the MLE problem. Next, theoretical analysis on the shuffled linear regression based on the algebraic-geometric theory is presented. Based on the analysis, an algebraically initialized expectation-maximization algorithm is introduced to solve the problem.

### 5.1 Joint Data Decoding and Device Identification

In the massive IoT scenario, the device identity information plays a vital role in differential updates, spatial correlation [12], and multi-stage collection [11], for which sensors are used to reconstruct the spatial field. It would take excess time if the identity information has to be collected regularly. Hence, a significant gain in the efficiency of communication procedure can be obtained by excluding the identification information in the header of the packet structure. This yields a joint data decoding and device identification problem at the BS, which may also act as a data fusion center. To achieve this goal, a *shuffled linear regression* has been recently investigated in a line of literature [9, 10, 14, 15] that can be exploited to remove the metadata used for device identification. The shuffled linear regression for identification-free communication is illustrated in Fig. 5.1. Considering a massive sensor network that contains  $m$  sensor nodes to capture the parameter data  $\mathbf{x} \in \mathbb{R}^n$  generated from  $n$  devices, a shuffled linear regression can be represented as

$$\mathbf{y} = \mathbf{\Pi} \mathbf{A} \mathbf{x}, \tag{5.1}$$

where  $\mathbf{y} \in \mathbb{R}^m$  is the permuted signal received at the BS,  $\mathbf{A} \in \mathbb{R}^{m \times n}$  is an encoding matrix, and  $\mathbf{\Pi}$  is an unknown  $m \times m$  permutation matrix whose  $i$ -th row is the canonical vector  $\mathbf{e}_{\pi(i)}^\top$  of all zeros except a 1 at position  $\pi(i)$ . The recovery of the shuffled linear regression (5.1) enables the BS to decode the signal



**Fig. 5.1** An example to illustrate the shuffled linear regression for identification-free communication

$\mathbf{x} = [x_1, \dots, x_n]^\top$  corresponding to each device from subsampled and permuted measurements  $\mathbf{y}$ . A simple linear shuffled model is illustrated in Example 5.1.

*Example 5.1* Consider a sensor network with three sensor nodes and two devices. Assume that the parameter data  $\mathbf{x}$  and the encoding matrix are given by

$$\mathbf{A} = \begin{bmatrix} 1 & 2 \\ -2 & 4 \\ 0 & -5 \end{bmatrix} \quad \text{and} \quad \mathbf{x} = \begin{bmatrix} 3 \\ 4 \end{bmatrix}. \quad (5.2)$$

Based on a permutation matrix  $\mathbf{\Pi}$ , i.e.,

$$\mathbf{\Pi} = \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \quad (5.3)$$

it yields a shuffled linear model:

$$\mathbf{y} = \mathbf{\Pi} \cdot \mathbf{A} \mathbf{x} = \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 11 \\ 10 \\ -20 \end{bmatrix} = \begin{bmatrix} -20 \\ 11 \\ 10 \end{bmatrix}. \quad (5.4)$$

Recently, important theoretical advances have been made to understand this problem, which can be mainly separated into three types: statistical approaches, algebraic geometry approaches, and alternating minimization approaches. For statistical approaches, the works [2, 4, 13] have developed algorithms based on

the maximum likelihood estimator  $\mathbf{x}_{\text{ML}}$  given by

$$(\mathbf{\Pi}_{\text{ML}}, \mathbf{x}_{\text{ML}}) = \underset{\mathbf{\Pi}^*, \mathbf{x}^*}{\text{argmin}} \|\mathbf{y} - \mathbf{\Pi}^* \mathbf{A} \mathbf{x}^*\|_2.$$

In addition, the conditions when the estimator recovers the ground truth in (5.1), i.e.,  $\mathbf{\Pi}_{\text{ML}} = \mathbf{\Pi}$ , have been established in the works [9, 13].

When the ratio of shuffled entries to all of the data entries is small, one may apply alternating minimization or multi-start gradient descent to solve (5.9) [2], which is an NP-hard problem for  $n > 1$  [15]. Due to high nonconvexity, such methods are very sensitive to initialization. This issue is addressed by the algebraically initialized expectation-maximization method proposed in [15], which uses the solution to the polynomial system of equations mentioned above to obtain a high-quality initialization. This approach is robust to small levels of noise, efficient for  $n \leq 5$ , and is able to handle fully shuffled data.

In the following, we will first demonstrate that the shuffled linear regression provides a way to achieve joint data decoding and device identification. Furthermore, two types of methods for solving the estimation problem in the shuffled linear regression will be introduced along with theoretical analysis, which include a maximum likelihood estimation based approach and an algebraic-geometric approach.

## 5.2 Problem Formulation

Consider a massive sensor network that contains  $m$  sensor nodes. Based on the correspondence pairs  $\{\mathbf{u}_j, y_j\}_{j=1}^m$ , the aim is to find the parameter vector

$$\mathbf{x} = [x_1, \dots, x_n]^\top \in \mathbb{R}^n$$

that characterizes the environment information, e.g., temperature, humidity, and pressure. The measurements are given by

$$y_j = \mathbf{a}_j^\top \mathbf{x}, \quad \forall j = 1, \dots, m, \quad (5.5)$$

$$\mathbf{a}_j := [a_1(\mathbf{u}_j), \dots, a_n(\mathbf{u}_j)]^\top, \quad (5.6)$$

where  $a_i : \mathbb{R}^s \rightarrow \mathbb{R}$  are known functions. The shuffled linear regression (5.5) can be considered as a special type of data corruption where the correspondences are missing. It can support identification-free communications, where the identification information, i.e., correspondences in the model (5.5), is excluded from the packet structure.

Given functions of the input samples represented in (5.6), i.e.,

$$\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_m]^\top \in \mathbb{R}^{m \times n}, \quad (5.7)$$

and a shuffled signal  $\mathbf{y} = [y_{j_1}, \dots, y_{j_m}]^\top \in \mathbb{R}^m$  with the unknown shuffling indices  $j_1, \dots, j_m$ , we have

$$\mathbf{y} = (\mathbf{\Pi})^\top \mathbf{A} \mathbf{x} + \mathbf{w} \in \mathbb{R}^m, \quad (5.8)$$

where  $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{\Pi}$  is an  $m \times m$  permutation matrix, and the vector

$$\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_m)$$

represents the additive Gaussian noise. The goal of the shuffled linear regression is to efficiently estimate both the signal  $\mathbf{x}$  and the permutation matrix  $\mathbf{\Pi}$  from  $\mathbf{y}$ . Thus, joint data decoding and device identification in the IoT network is achieved.

### 5.3 Maximum Likelihood Estimation Based Approaches

Several methods have recently been developed to solve the shuffled linear regression problem. Specifically, the estimation of shuffled linear regression can be achieved via the *maximum likelihood estimator* (MLE) [2, 9, 13]:

$$(\hat{\mathbf{\Pi}}_{\text{ML}}, \hat{\mathbf{x}}_{\text{ML}}) = \underset{\mathbf{\Pi}^*, \mathbf{x}^*}{\operatorname{argmin}} \|\mathbf{\Pi}^* \mathbf{y} - \mathbf{A} \mathbf{x}^*\|_2. \quad (5.9)$$

Based on this estimator, several algorithms have been developed and the theoretical analyses have been established, e.g., [2, 4, 16]. It shows that  $\hat{\mathbf{\Pi}}_{\text{ML}}$  is divergent from  $\mathbf{\Pi}$  in (5.8) with high probability if the SNR is not large enough [4]. The detailed guarantees will be discussed in Theorems 5.1 and 5.2. Furthermore, [16] shows that if the SNR approaches infinity,  $\hat{\mathbf{x}}_{\text{ML}}$  approaches  $\mathbf{x}$  in (5.8). Moreover, advanced algorithms based on algebraic-geometric approaches proposed recently to address the computational issue will also be introduced.

#### 5.3.1 Sorting Based Algorithms

The paper [9] analyzed the shuffled linear regression problem under the assumption that the entries of the matrix  $\mathbf{A}$  are drawn i.i.d. from a standard Gaussian distribution. The paper [9] established sharp conditions on the sample size  $m$ , dimension  $n$ , and SNR, under which  $\mathbf{\Pi}$  is exactly recoverable. From the computational point of view, the paper [9] demonstrated that the maximum likelihood estimate of  $\mathbf{\Pi}$  is NP-hard

to compute, and it proposed a polynomial-time algorithm based on sorting, which is introduced in the sequel.

Theorems 5.1 and 5.2 in the following provide the statistical properties of the MLE (5.9). Based on the maximum likelihood estimator (5.9), an upper bound on the probability of error of  $\hat{\boldsymbol{\Pi}}_{\text{ML}}$  is provided in the following theorem given by [9] with  $c_1, c_2 > 0$ .

**Theorem 5.1** *For any  $n < m$  and  $\epsilon < \sqrt{m}$ , if*

$$\log \left( \frac{\|\mathbf{x}\|_2^2}{\sigma^2} \right) \geq \left( c_1 \frac{m}{m-n} + \epsilon \right) \log m, \quad (5.10)$$

then  $\mathbb{P}\{\hat{\boldsymbol{\Pi}}_{\text{ML}} \neq \boldsymbol{\Pi}\} \leq c_2 m^{-2\epsilon}$ .

Furthermore, the lower bound on the probability of error of  $\hat{\boldsymbol{\Pi}}_{\text{ML}}$  is provided as follows.

**Theorem 5.2** *For any  $\delta \in (0, 2)$ , if*

$$2 + \log \left( 1 + \frac{\|\mathbf{x}\|_2^2}{\sigma^2} \right) \leq (2 - \delta) \log m, \quad (5.11)$$

then  $\mathbb{P}\{\hat{\boldsymbol{\Pi}} \neq \boldsymbol{\Pi}\} \geq 1 - c_3 e^{-c_4 m^\delta}$  for any estimator  $\hat{\boldsymbol{\Pi}}$ .

We can conclude from Theorem 5.2 that if condition (5.11) is satisfied, the recovery probability approaches 1 when  $m$  tends to infinity.

Since Eq. (5.9) requires a combinatorial minimization over  $n!$  permutations, advanced algorithms are needed to compute  $\hat{\boldsymbol{\Pi}}_{\text{ML}}$  efficiently. To begin with, the maximum likelihood estimate of the permutation is represented as [9]

$$\hat{\boldsymbol{\Pi}}_{\text{ML}} = \arg \min_{\boldsymbol{\Pi}} \|P_{\boldsymbol{\Pi}}^\perp \mathbf{y}\|_2^2, \quad (5.12)$$

where

$$P_{\boldsymbol{\Pi}}^\perp = \mathbf{I} - \boldsymbol{\Pi} \mathbf{A} (\mathbf{A}^\top \mathbf{A})^{-1} (\boldsymbol{\Pi} \mathbf{A})^\top.$$

When  $n = 1$  and representing the design vector as  $\mathbf{a}$ , Eq. (5.12) can be represented as

$$\begin{aligned} \hat{\boldsymbol{\Pi}}_{\text{ML}} &= \arg \max_{\boldsymbol{\Pi}} \|\mathbf{a}_{\boldsymbol{\Pi}}^\top \mathbf{y}\|^2 \\ &= \arg \max_{\boldsymbol{\Pi}} \max \left\{ \mathbf{a}_{\boldsymbol{\Pi}}^\top \mathbf{y}, -\mathbf{a}_{\boldsymbol{\Pi}}^\top \mathbf{y} \right\} \\ &= \arg \min_{\boldsymbol{\Pi}} \max \left\{ \|\mathbf{a}_{\boldsymbol{\Pi}} - \mathbf{y}\|_2^2, \|\mathbf{a}_{\boldsymbol{\Pi}} + \mathbf{y}\|_2^2 \right\}. \end{aligned} \quad (5.13)$$

The polynomial-time algorithm illustrated in Algorithm 5.1 is developed based on (5.13). This is achieved based on the fact that for fixed vectors  $\mathbf{x}$  and  $\mathbf{y}$ ,

$$\|\mathbf{x}_{\Pi} - \mathbf{y}\|$$

can be minimized for  $\Pi$  by sorting  $\mathbf{x}$  according to the order of  $\mathbf{y}$ . The theoretical analysis of Algorithm 5.1 from the computational points of view is illustrated in Theorem 5.3.

---

Algorithm 5.1: Exact algorithm for implementing Eq. (5.12)

---

Input: design vector  $\mathbf{a}$ , observation vector  $\mathbf{y}$   
 1  $\Pi_1 \leftarrow$  permutation that sorts  $\mathbf{a}$  according to  $\mathbf{y}$   
 2  $\Pi_2 \leftarrow$  permutation that sorts  $-\mathbf{a}$  according to  $\mathbf{y}$   
 3  $\hat{\Pi}_{\text{ML}} \leftarrow \arg \max\{|\mathbf{a}_{\Pi_1}^\top \mathbf{y}|, |\mathbf{a}_{\Pi_2}^\top \mathbf{y}|\}$   
 Output:  $\hat{\Pi}_{\text{ML}}$

---

**Theorem 5.3** *For  $n = 1$ , the MLE estimator  $\hat{\Pi}_{\text{ML}}$  can be computed via Algorithm 5.1 in time  $\mathcal{O}(m \log m)$  for any choice of the measurement matrix  $\mathbf{A}$ . In contrast, if  $n > 1$ , then  $\hat{\Pi}_{\text{ML}}$  is NP-hard to compute.*

Theorem 5.3 shows that the algorithmic advantages enjoyed in the case of  $n = 1$  cannot extend to general cases of  $n > 1$ . For  $n > 1$  a natural method is brute force search: for each permutation  $\Pi$  of the  $m!$  permutations, check whether the linear system

$$\Pi \mathbf{y} = \mathbf{A} \mathbf{x}$$

is consistent, followed by solving it if it is consistent. This algorithm ends with the complexity of  $\mathcal{O}(n^2(m+1)!)$ . An approximate algorithm that is more efficient than the brute force has been proposed in [4], which makes progress on both computational and statistical aspects. It is introduced in the next section.

### 5.3.2 Approximation Algorithm

Considering the least squares problem (5.9), an approximation approach [4] is proposed that for any  $\epsilon \in (0, 1)$ , it returns an  $(1 + \epsilon)$ -approximation in time  $\mathcal{O}((m/\epsilon)^n)$ .

The approximation algorithm, shown as Algorithm 5.3, uses a careful enumeration to beat the naive brute force running time of  $\Omega(n!)$ . The ‘‘Row Sampling’’

algorithm [3] is exploited in the beginning of Algorithm 5.3 in order to narrow the search space. The details of the ‘‘Row Sampling’’ algorithm [3] is presented in the following. The ‘‘Row Sampling’’ is illustrated in Algorithm 5.2 with the following notations:

- For each  $i \in [n]$ ,  $e_i$  is the  $i$ -th coordinate basis vector in  $\mathbb{R}^n$ .
- $L(\mathbf{x}, \delta_L, \mathbf{A}, \ell) := \frac{\mathbf{x}^\top (\mathbf{A} - (\ell + \delta_L) \mathbf{I}_k)^{-2} \mathbf{x}}{\phi(\ell + \delta_L, \mathbf{A}) - \Phi(\ell, \mathbf{A})} - (\ell + \delta_L) \mathbf{I}_k)^{-1} \mathbf{x}$ , where  $\phi(\ell, \mathbf{A}) := \sum_{i=1}^k \frac{1}{\lambda_i(\mathbf{A}) - \ell}$  and  $(\lambda_i(\mathbf{A}))_{i=1}^k$  are the eigenvalues of  $\mathbf{A}$ .
- $\hat{U}(\mathbf{x}, \delta, \mathbf{B}, u) := \frac{\mathbf{x}^\top (\mathbf{B} - u' \mathbf{I}_r)^{-2} \mathbf{x}}{\phi'(u, \mathbf{B}) - \phi'(u', \mathbf{B})} - \mathbf{x}^\top (\mathbf{B} - u' \mathbf{I}_r)^{-1} \mathbf{x}$ , where  $u' = u + \delta$ ,  $\phi'(u, \mathbf{B}) := \sum_{i=1}^r \frac{1}{u - \lambda_i(\mathbf{B})}$ , and  $(\lambda_i(\mathbf{B}))_{i=1}^k$  are the eigenvalues of  $\mathbf{B}$ .

---

**Algorithm 5.2: ‘‘Row Sampling’’ algorithm [3]**

---

input Matrix  $\mathbf{A} = [\mathbf{A}_1 \cdots \mathbf{A}_n]^\top \in \mathbb{R}^{m \times n}$  such that  $\mathbf{A}^\top \mathbf{A} = \mathbf{I}_n$ ; integer  $r \geq n$ .

output Matrix  $\mathbf{S} = (S_{i,j})_{(i,j) \in \times [m]} \in \mathbb{R}^{r \times m}$ .

- 1: Set  $\mathbf{Q}_0 = \mathbf{0}_{n \times n}$ ,  $\mathbf{B}_0 = \mathbf{0}_{m \times m}$ ,  $\mathbf{S} = \mathbf{0}_{r \times m}$ ,  $\delta = (1 + m/r)(1 - \sqrt{n/r})^{-1}$  and  $\delta_L = 1$ .
  - 2: for  $\tau = 0$  to  $r - 1$  do
  - 3:   Let  $\ell_\tau = \tau - \sqrt{rk}$  and  $u_\tau = \delta(\tau + \sqrt{mr})$ .
  - 4:   Select  $i_\tau \in [m]$  and number  $t_\tau > 0$  such that  $\hat{U}(e_{i_\tau}, \delta, \mathbf{B}_\tau, u_\tau) \leq \frac{1}{t_\tau} \leq L(\mathbf{A}_{i_\tau}, \delta_L, \mathbf{Q}_\tau, \ell_\tau)$ .
  - 5:   Set  $\mathbf{Q}_{\tau+1} = \mathbf{Q}_\tau + t_\tau \mathbf{A}_{i_\tau} \mathbf{A}_{i_\tau}^\top$ ,  $\mathbf{B}_{\tau+1} = \mathbf{B}_\tau + t_\tau e_{i_\tau} e_{i_\tau}^\top$  and  $S_{\tau+1, i_\tau} = \sqrt{r^{-1}(1 - \sqrt{n/r})} / \sqrt{t_\tau}$ .
  - 6: end for
  - 7: return  $\mathbf{S}$ .
- 

The theoretical guarantee of Algorithm 5.3 is given in the following Theorem 5.4 [4].

**Theorem 5.4** *Algorithm 5.3 returns  $\hat{\mathbf{x}} \in \mathbb{R}^n$  and  $\hat{\mathbf{\Pi}}$  satisfying*

$$\|\hat{\mathbf{\Pi}} \mathbf{y} - \mathbf{A} \hat{\mathbf{x}}\|_2^2 \leq (1 + \epsilon) \min_{\mathbf{x}, \mathbf{\Pi}} \|\mathbf{\Pi} \mathbf{y} - \mathbf{A} \mathbf{x}\|_2^2.$$

It shows that Algorithm 5.3 enjoys recovery guarantees for  $\mathbf{x}$  and  $\mathbf{\Pi}$  when the data come from the Gaussian measurement model (5.8). Moreover, the overall running time is  $\mathcal{O}((m/\epsilon)^{(k)})$  which is remarkably lower than that of naive brute force search, i.e.,  $\Omega(n!)$ .

---

**Algorithm 5.3: Approximation algorithm for least squares problem (5.9)**


---

Input Sample matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ; observation  $\mathbf{y} \in \mathbb{R}^m$ ; approximation parameter  $\epsilon \in (0, 1)$ .

Assume  $\mathbf{A}^\top \mathbf{A} = \mathbf{I}_n$ .

Output Parameter vector  $\hat{\mathbf{x}} \in \mathbb{R}^n$  and permutation matrix  $\hat{\Pi}$ .

- 1: Run Algorithm 5.2 with input matrix  $\mathbf{A}$  to obtain a matrix  $\mathbf{S} \in \mathbb{R}^{r \times m}$  with  $r = 4n$ .
  - 2: Let  $\mathcal{B}$  be the set of vectors  $\mathbf{b} = (b_1, b_2, \dots, b_n)^\top \in \mathbb{R}^n$  satisfying the following: for each  $i \in [n]$ ,
    - if the  $i$ -th column of  $\mathbf{S}$  is all zeros, then  $b_i = 0$ ;
    - otherwise,  $b_i \in \{y_1, y_2, \dots, y_n\}$ .
  - 3: Let  $c := 1 + 4(1 + \sqrt{m/(4n)})^2$ .
  - 4: for each  $\mathbf{b} \in \mathcal{B}$  do
    - 5: Compute  $\tilde{\mathbf{x}}_{\mathbf{b}} \in \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{S}(\mathbf{b} - \mathbf{A}\mathbf{x})\|_2^2$ , and let  $r_{\mathbf{b}} := \min_{\Pi} \|\Pi \mathbf{y} - \mathbf{A}\tilde{\mathbf{x}}_{\mathbf{b}}\|_2^2$ .
    - 6: Construct a  $\sqrt{\epsilon r_{\mathbf{b}}/c}$ -net  $\mathcal{N}_{\mathbf{b}}$  for the Euclidean ball of radius  $\sqrt{c r_{\mathbf{b}}}$  around  $\tilde{\mathbf{x}}_{\mathbf{b}}$ , so that for each  $\mathbf{v} \in \mathbb{R}^k$  with  $\|\mathbf{v} - \tilde{\mathbf{x}}_{\mathbf{b}}\|_2 \leq \sqrt{c r_{\mathbf{b}}}$ , there exists  $\mathbf{v}' \in \mathcal{N}_{\mathbf{b}}$  such that  $\|\mathbf{v} - \mathbf{v}'\|_2 \leq \sqrt{\epsilon r_{\mathbf{b}}/c}$ .
  - 7: end for
  - 8: return  $\hat{\mathbf{x}} \in \operatorname{argmin}_{\mathbf{x} \in \bigcup_{\mathbf{b} \in \mathcal{B}} \mathcal{N}_{\mathbf{b}}} \min_{\Pi} \|\Pi \mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$  and  $\hat{\Pi} \in \operatorname{argmin}_{\Pi} \|\Pi \mathbf{y} - \mathbf{A}\hat{\mathbf{x}}\|_2^2$ .
- 

However, the approximation guarantee is not robust to even mild levels of noise. Thus, it motivates other advanced algorithms, e.g., alternative minimization approaches [1] and algebraic geometric approaches [15].

## 5.4 Algebraic-Geometric Approach

Recently, the paper [1] proposed a practical algorithm for solving shuffled linear regression (5.9) via alternating minimization: estimating  $\Pi^*$  via sorting an estimate  $\xi^*$  and estimating  $\xi^*$  via least-squares given an estimate  $\Pi^*$ . Nevertheless, this approach is very sensitive to initialization and generally works only when the observation data is partially shuffled. To address the limitations of alternating minimization approach, the paper [15] proposed an algebraic geometric approach, which uses symmetric polynomials and leads to a polynomial system of  $n$  equations in  $n$  unknowns, containing  $\mathbf{x}$  in its root locus.

Based on algebraic geometry theory, the paper [15] proved that this polynomial system is consistent with at most  $n!$  complex roots as long as the independent samples are generic. This fact implies that this polynomial system can always be solved, and its most suitable root can be used as initialization to the expectation maximization (EM) algorithm.

### 5.4.1 Eliminating $\Pi$ via Symmetric Polynomials

Prior to introducing the algebraically initialized expectation-maximization, we first describe the main idea of the algebraic-geometric approach to solve the shuffled linear regression estimation problem (5.8). Denote the ring of polynomials with real coefficients over variables

$$\mathbf{z} := [z_1, \dots, z_m]^\top$$

as

$$\mathbb{R}[\mathbf{z}] := \mathbb{R}[z_1, \dots, z_m].$$

A symmetric polynomial<sup>1</sup>  $p \in \mathbb{R}[\mathbf{z}]$  means that it is invariant to any permutation of the variables  $\mathbf{z}$ , given by

$$p(\mathbf{z}) := p(z_1, \dots, z_m) = p(z_{\pi(1)}, \dots, z_{\pi(m)}) =: p(\mathbf{\Pi z}), \quad (5.14)$$

where  $\pi$  is a permutation on  $\{1, \dots, m\}$  and  $\mathbf{\Pi}$  is an  $m \times m$  permutation matrix.

Recall the shuffled linear regression problem (5.8) in the noiseless scenario and let  $(\mathbf{\Pi}^*, \mathbf{x}^*)$  with

$$\mathbf{x}^* = [x_1^*, \dots, x_n^*]^\top$$

being a solution. Based on a symmetric polynomial  $p \in \mathbb{R}[\mathbf{z}]$ , we get

$$\mathbf{\Pi}^* \mathbf{y} = \mathbf{A} \mathbf{x}^* \xrightarrow{p: \text{symmetric}} p(\mathbf{y}) = p(\mathbf{\Pi}^* \mathbf{y}) = p(\mathbf{A} \mathbf{x}^*). \quad (5.15)$$

In (5.15), the symmetric polynomial  $p$  plays a vital role in eliminating the unknown permutation  $\mathbf{\Pi}^*$  and providing a constraint that only be relative with the known  $\mathbf{A}$ ,  $\mathbf{y}$ ,

$$\hat{p}(\mathbf{x}) := p(\mathbf{A} \mathbf{x}) - p(\mathbf{y}) = 0. \quad (5.16)$$

We aim to find the solution  $\mathbf{x}^*$  that satisfies (5.16), thereby finding all solutions to the estimation problem (5.8).

To achieve this goal, we first introduce the concept of *algebraic variety* which is used to characterize the solutions of (5.16). Recall that the polynomial  $\hat{p}$  in (5.16) is an element of the polynomial ring  $\mathbb{R}[\mathbf{x}]$  in  $n$  variables  $\mathbf{x} := [x_1, \dots, x_n]^\top$ , and the set of its roots, denoted as

$$\mathcal{V}(\hat{p}) := \{\mathbf{x} \in \mathbb{R}^n : \hat{p}(\mathbf{x}) = 0\},$$

---

<sup>1</sup>We do not distinguish between  $p$  and  $p(\mathbf{z})$ .

is called an algebraic variety. In particular,  $\mathcal{V}(\hat{p})$  defines a hypersurface of  $\mathbb{R}^n$ . Geometrically, the solutions to (5.16) are the intersection points of the corresponding  $n$  hypersurfaces

$$\mathcal{V}(\hat{p}_1), \dots, \mathcal{V}(\hat{p}_n),$$

which include all solutions to problem (5.8), as well as potentially irrelevant points. Theorems provided in Sect. 5.4.2 investigate a system of  $n$  equations in  $n$  unknowns and the method of filtering its roots of interest is introduced in Sect. 5.4.3.

In addition, an example is provided in the following to illustrate the symmetric polynomial.

*Example 5.2* Consider the data

$$\mathbf{A} = \begin{bmatrix} 1 & 2 \\ -2 & 4 \\ 0 & -5 \end{bmatrix}, \mathbf{y} = \begin{bmatrix} -20 \\ 11 \\ 10 \end{bmatrix}. \quad (5.17)$$

It is simple to find that there is a unique permutation

$$\mathbf{\Pi}^* = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix} \quad (5.18)$$

that results in a consistent linear system of equations

$$\begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} -20 \\ 11 \\ 10 \end{bmatrix} = \begin{bmatrix} 1 & 2 \\ -2 & 4 \\ 0 & -5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad (5.19)$$

with solution  $\xi_1^* = 3$ ,  $\xi_2^* = 4$ . Now consider the symmetric polynomial

$$p_1(z_1, z_2, z_3) = z_1 + z_2 + z_3, \quad (5.20)$$

and based on (5.16) it yields the constraint

$$(x_1 + 2x_2) + (4x_2 - 2x_1) - 5x_2 = -20 + 11 + 10, \quad (5.21)$$

$$\Leftrightarrow x_2 - x_1 = 1, \quad (5.22)$$

$$\Leftrightarrow \hat{p}_1(\mathbf{x}) := p_1(\mathbf{A}\mathbf{x}) - p_1(\mathbf{y}) = 0. \quad (5.23)$$

Indeed, we see that the solution  $\xi^* = [3, 4]^\top$  satisfies (5.23).

### 5.4.2 Theoretical Analysis

The theoretical analysis on symmetric polynomials for both exact and corrupted data are provided in the following.

#### 5.4.2.1 Exact Data

As Example 5.2 suggests, a natural choice for  $n$  symmetric polynomials are the first  $n$  power sums

$$p_k(\mathbf{z}) \in \mathbb{R}[\mathbf{z}] := \mathbb{R}[z_1, \dots, z_m], \quad k \in [n] := \{1, \dots, n\},$$

denoted as

$$p_k(\mathbf{z}) := z_1^k + \dots + z_m^k. \quad (5.24)$$

Based on (5.16), we conclude that any solution  $\xi^*$  of (5.8) in the noiseless scenario must obey the polynomial constraints

$$\hat{p}_k(\mathbf{x}) = 0, \quad k \in [n], \quad \text{where} \quad (5.25)$$

$$\hat{p}_k(\mathbf{x}) := p_k(\mathbf{Ax}) - p_k(\mathbf{y}) = \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{x})^k - \sum_{j=1}^m y_j^k, \quad (5.26)$$

and  $\mathbf{a}_i^\top$  presents the  $i$ th row of  $\mathbf{A}$ . The following theorem provides a theoretical guarantee that the number of other irrelevant solutions must be finite.

**Theorem 5.5 ([15])** *Assuming  $\mathbf{A}$  is generic and  $\mathbf{y}$  is some permutation of a vector. The algebraic variety*

$$\mathcal{V}(\hat{p}_1, \dots, \hat{p}_n)$$

*consists of all*

$$\xi_1^*, \dots, \xi_\ell^* \in \mathbb{R}^n$$

*such that there exists permutations*

$$\Pi_1^*, \dots, \Pi_\ell^*$$

*with*

$$\Pi_i^* \mathbf{y} = \mathbf{A} \xi_i^*, \quad \forall i \in [\ell],$$

*while it may include at most  $n! - \ell$  other points of  $\mathbb{C}^n$ .*

Theorem 5.5 demonstrates that the system of polynomial equations

$$\hat{p}_1(\mathbf{x}) = \cdots = \hat{p}_n(\mathbf{x}) = 0, \quad (5.27)$$

always has a finite number of solutions in  $\mathbb{C}^n$  (at most  $n!$ ), which contain all possible solutions  $\xi_1^*, \dots, \xi_\ell^* \in \mathbb{R}^n$  of problem (5.8).

### 5.4.2.2 Corrupted Data

The following theorem addresses the issue of corrupted data which is common in practical applications. Considering the corrupted data which is denoted as  $\tilde{\mathbf{A}}, \tilde{\mathbf{y}}$ , the linear system can be represented as

$$\mathbf{\Pi} \tilde{\mathbf{y}} = \tilde{\mathbf{A}} \mathbf{x}. \quad (5.28)$$

There exists a permutation  $\mathbf{\Pi} = \tilde{\mathbf{\Pi}}^*$  such that (5.28) is approximately consistent, if the degree of corruption is sufficiently small. In order to get an approximate solution of (5.28), a *corrupted* power-sum polynomial is defined as

$$\tilde{p}_k(\mathbf{x}) := p_k(\tilde{\mathbf{A}} \mathbf{x}) - p_k(\tilde{\mathbf{y}}), \quad k \in [n], \quad (5.29)$$

and the polynomial system  $\tilde{\mathcal{P}}$  is considered, given by

$$\tilde{p}_1 = \cdots = \tilde{p}_n = 0. \quad (5.30)$$

These are  $n$  equations of degrees  $1, 2, \dots, n$  with  $n$  unknowns. The system of polynomial equations (5.30) with respect to corrupted data is investigated in the following theorem.

**Theorem 5.6 ([15])** *If  $\tilde{\mathbf{A}}$  is generic and  $\tilde{\mathbf{y}} \in \mathbb{R}^m$  is any vector, then  $\mathcal{V}(\tilde{p}_1, \dots, \tilde{p}_n)$  is non-empty containing at most  $n!$  points of  $\mathbb{C}^n$ .*

Theorem 5.6 demonstrates that the system of polynomial equations (5.30) always has at least one solution. We can conclude that an approximate solution to the shuffled linear system (5.28) lies in a finite number of solutions of the system (5.30). Theorems 5.5 and 5.6 provide theoretical guarantees for developing algebraical method to solve the shuffled linear regression problem. The algorithm, called algebraically initialized expectation-maximization, is introduced in the next section.

### 5.4.3 Algebraically Initialized Expectation-Maximization

If there is a unique solution  $\xi^*$  to the shuffled linear regression problem (5.8), Theorem 5.5 ensures that  $\xi^*$  is one of the finitely many complex roots of the polynomial system (5.25) of  $n$  equations in  $n$  unknowns. Moreover, in the case of

corrupted data, Theorem 5.6 ensures that the system is consistent with  $L \leq n!$  complex roots, and if the corruption degree is modest, one of the roots can be a good approximation to the maximum likelihood estimator (MLE)  $\hat{\xi}_{\text{ML}}$  (5.9). Thus, the goal is to filter that root and refine it.

Particularly, several state-of-the-art polynomial system solvers [5] can be exploited to solve the polynomial system of equations (5.25). With the computed roots

$$\hat{\xi}_1, \dots, \hat{\xi}_L \in \mathbb{C}^n, L \leq n!,$$

of the polynomial system, only their real parts

$$(\hat{\xi}_1)_{\mathbb{R}}, \dots, (\hat{\xi}_L)_{\mathbb{R}}$$

are retained, which can be used for obtaining an approximation to the ML estimator  $\hat{\xi}_{\text{ML}}$ . This is achieved by selecting the root that yields the smallest  $\ell_2$  error among all possible permutations  $\Pi$ :

$$\hat{\xi}_{\text{AI}} := \operatorname{argmin}_{i \in [L]} \left\{ \min_{\Pi} \left\| \Pi \mathbf{y} - \mathbf{A}(\hat{\xi}_i)_{\mathbb{R}} \right\|_2 \right\}. \quad (5.31)$$

Furthermore, the *algebraic initialization*  $\hat{\xi}_{\text{AI}}$  is utilized as an initialization to the expectation maximization algorithm [1] which implements alternating minimization to solve (5.9). This method is called *Algebraically Initialized Expectation-Maximization (AI-EM)* [15] illustrated in Algorithm 5.4.

---

#### Algorithm 5.4: Algebraically initialized expectation-maximization

---

procedure AI-EM( $\mathbf{y} \in \mathbb{R}^m$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $T \in \mathbb{N}$ ,  $\epsilon \in \mathbb{R}_+$ )

- 1:  $p_k(\mathbf{z}) := \sum_{j=1}^m z_j^k$ ,  $\hat{p}_k := p_k(\mathbf{A}\mathbf{x}) - p_k(\mathbf{y})$ ,  $k \in [n]$ ;
  - 2: Compute roots  $\{\hat{\xi}_i\}_{i=1}^L \subset \mathbb{C}^n$  of  $\{\hat{p}_k = 0, k \in [n]\}$ ;
  - 3: Extract the real parts  $\{(\hat{\xi}_i)_{\mathbb{R}}\}_{i=1}^L$  of  $\{\hat{\xi}_i\}_{i=1}^L \subset \mathbb{C}^n$ ;
  - 4:  $\{\xi_0, \Pi_0\} \leftarrow \operatorname{argmin}_{\xi \in \{(\hat{\xi}_i)_{\mathbb{R}}\}_{i=1}^L, \Pi} \|\Pi \mathbf{y} - \mathbf{A}\xi\|_2$ ;
  - 5:  $t \leftarrow 0$ ,  $\Delta \mathcal{J} \leftarrow \infty$ ,  $\mathcal{J} \leftarrow \|\Pi_0 \mathbf{y} - \mathbf{A}\xi_0\|_2$ ;
  - 6: while  $t < T$  and  $\Delta \mathcal{J} > \epsilon \mathcal{J}$  do
  - 7:    $t \leftarrow t + 1$ ;
  - 8:    $\xi_t \leftarrow \operatorname{argmin}_{\xi \in \mathbb{R}^n} \|\Pi_{t-1} \mathbf{y} - \mathbf{A}\xi\|_2$ ;
  - 9:    $\Pi_t \leftarrow \operatorname{argmin}_{\Pi} \|\Pi \mathbf{y} - \mathbf{A}\xi_t\|_2$ ;
  - 10:    $\Delta \mathcal{J} \leftarrow \mathcal{J} - \|\Pi_t \mathbf{y} - \mathbf{A}\xi_t\|_2$ ;
  - 11:    $\mathcal{J} \leftarrow \|\Pi_t \mathbf{y} - \mathbf{A}\xi_t\|_2$ ;
  - 12: end while
  - 13: Return  $\xi_t, \Pi_t$ .
- end procedure
-

### 5.4.4 Simulation Results

To further illustrate the advantage of Algorithm 5.4, we compare it with two variations of EM algorithms that were proposed in [1]: (1) *LS-EM* that computes the MLE via alternating minimization with the initialization satisfying

$$\xi_{0,LS} := \operatorname{argmin}_{\xi \in \mathfrak{R}^n} \|\mathbf{y} - \mathbf{A}\xi\|_2; \quad (5.32)$$

(2) *Soft-EM* that uses the same initialization as LS-EM, but exploits a dynamic empirical average of permutation matrices drawn from a suitable Markov chain to optimize the permutation operation.

All methods are evaluated by measuring the relative estimation error between the estimator  $\hat{\xi}$  and the ground truth  $\xi^*$ , given by

$$100 \frac{\|\xi^* - \hat{\xi}\|_2}{\|\xi^*\|_2} \%. \quad (5.33)$$

For AI-EM, the estimation error between the best root  $\xi_{AI}^*$  of the polynomial system is defined as

$$\xi_{AI}^* := \operatorname{argmin}_{\hat{\xi}_i, i \in [L]} \|\xi^* - (\hat{\xi}_i)_{\mathfrak{R}}\|_2, \quad (5.34)$$

and the estimator  $\hat{\xi}_{AI}$  is computed as in (5.31).

Figure 5.2 illustrates the estimation error of the three methods with fully shuffled data under the setting of  $n = 3$ ,  $\sigma = 0:0.01:0.1$  and  $m = 500$ . The simulation

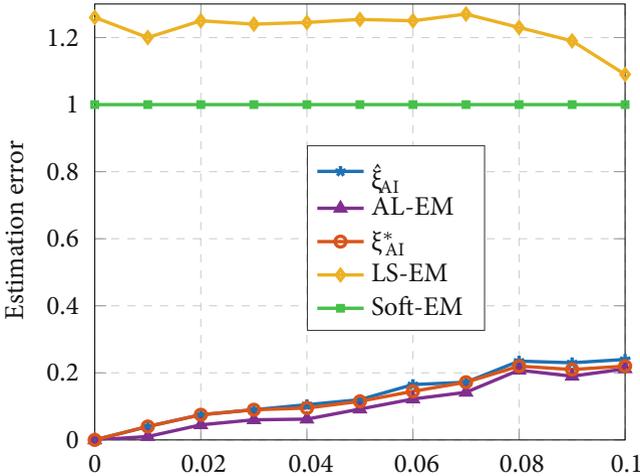


Fig. 5.2 Estimation error for fully shuffled data

results are averaged over 100 independent trials. It shows LS-EM and Soft-EM fail. It can be explained that when the data are fully shuffled, the least-squares initialization (5.32) exploited by both LS-EM and Soft-EM rather deviates from the ground truth  $\xi^*$ .

## 5.5 Summary

This chapter summarized a shuffled linear regression model to support joint data decoding and device identification, thereby reducing the overhead in massive connectivity systems. The methods developed for solving the shuffled linear regression estimation problem are presented in this chapter from the numerical and theoretical points of view. The methods can be mainly categorized into two types: maximum likelihood estimation based approach and algebraic geometric approach. Besides the application introduced in this chapter, the shuffled linear regression method and its variations arise in many applications, e.g., image processing [6], user de-anonymization [8], and correspondence estimation [7]. Recently, an abstraction of shuffled linear problems which is called homomorphic sensing has been studied in [14], and an algebraic theory for homomorphic sensing has been developed. The paper [14] provides the first working solutions for the unlabeled sensing problem for small dimensions. It is still a principle direction of study to develop more efficient algorithms and corresponding theoretical guarantees for homomorphic sensing.

## References

1. Abid, A., Zou, J.: A stochastic expectation-maximization approach to shuffled linear regression. In: Proceedings of the 56th Annual Allerton Conference on Communication, Control, and Computing, pp. 470–477. IEEE, Piscataway (2018)
2. Abid, A., Poon, A., Zou, J.: Linear regression with shuffled labels (2017). Preprint. arXiv:1705.01342
3. Boutsidis, C., Drineas, P., Magdon-Ismail, M.: Near-optimal coresets for least-squares regression. *IEEE Trans. Inf. Theory*. **59**(10), 6880–6892 (2013)
4. Hsu, D.J., Shi, K., Sun, X.: Linear regression without correspondence. In: Advances in Neural Information Processing Systems (NeurIPS), pp. 1531–1540 (2017)
5. Lazard, D.: Thirty years of polynomial system solving, and now? *J. Symb. Comput.* **44**(3), 222–231 (2009)
6. Lian, W., Zhang, L., Yang, M.H.: An efficient globally optimal algorithm for asymmetric point matching. *IEEE Trans. Pattern Anal. Mach. Intell.* **39**(7), 1281–1293 (2016)
7. Marques, M., Stosić, M., Costeira, J.: Subspace matching: unique solution to point matching with geometric constraints. In: Proceedings of the International Conference on Computer Vision, pp. 1288–1294. IEEE, Piscataway (2009)
8. Narayanan, A., Shmatikov, V.: Robust de-anonymization of large sparse datasets. In: Proceedings of the IEEE Symposium on Security and Privacy, pp. 111–125 (2008)

9. Pananjady, A., Wainwright, M.J., Courtade, T.A.: Linear regression with shuffled data: statistical and computational limits of permutation recovery. *IEEE Trans. Inf. Theory* **64**(5), 3286–3300 (2018)
10. Peng, L., Song, X., Tsakiris, M.C., Choi, H., Kneip, L., Shi, Y.: Algebraically-initialized expectation maximization for header-free communication. In: *Proceedings of the IEEE International Conference on Acoustics Speech Signal Processing (ICASSP)*, pp. 5182–5186 (2019)
11. Pradhan, S.S., Kusuma, J., Ramchandran, K.: Distributed compression in a dense microsensor network. *IEEE Signal Process. Manag.* **19**(2), 51–60 (2002)
12. Scaglione, A., Servetto, S.: On the interdependence of routing and data compression in multi-hop sensor networks. *Wirel. Netw.* **11**(1–2), 149–160 (2005)
13. Slawski, M., Ben-David, E., et al.: Linear regression with sparsely permuted data. *Electron. J. Stat.* **13**(1), 1–36 (2019)
14. Tsakiris, M.C., Peng, L.: Homomorphic sensing. In: *Proceedings of the International Conference on Machine Learning (ICML)*, pp. 6335–6344 (2019)
15. Tsakiris, M.C., Peng, L., Conca, A., Kneip, L., Shi, Y., Choi, H.: An algebraic-geometric approach to shuffled linear regression (2018). Preprint. arXiv: 1810.05440
16. Unnikrishnan, J., Haghghatshoar, S., Vetterli, M.: Unlabeled sensing with random linear measurements. *IEEE Trans. Inf. Theory.* **64**(5), 3237–3253 (2018)

# Chapter 6

## Learning Augmented Methods



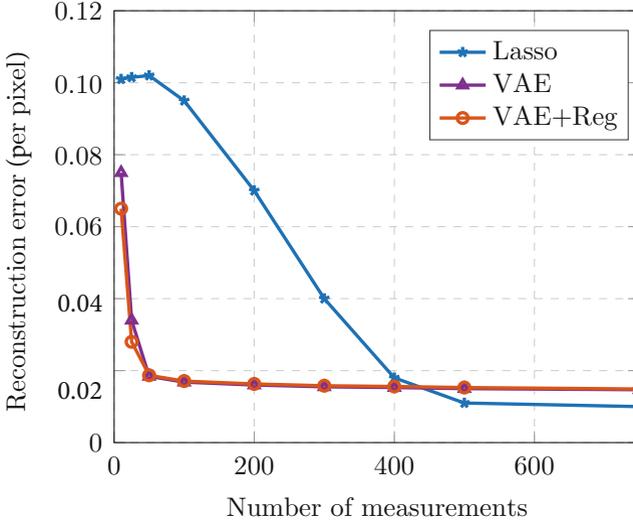
**Abstract** In this chapter, we introduce some cutting-edge learning augmented techniques to enhance the performance of structured signal processing. We start with compressed sensing under a generative prior, which can better capture the underlying signal structure than the traditional sparse prior. We then present learning augmented techniques for the joint design of measurement matrix and sparse support recovery for the sparse linear model (e.g., compressed sensing). Furthermore, several deep-learning-based AMP methods for the sparse linear model are introduced, including learned AMP, learned Vector-AMP, and learned ISTA for group row sparsity.

### 6.1 Structured Signal Processing Under a Generative Prior

Recall the sparse linear model defined in (2.3). The sparse signal  $\mathbf{x}$  can be recovered via solving a convex optimization problem known as Lasso [21]:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1, \quad (6.1)$$

where the parameter  $\lambda > 0$  controls the sparsity level. Instead of focusing on the sparsity of  $\mathbf{x}$  in (2.3), the paper [3] has recently estimated  $\hat{\mathbf{x}}$  based on the structure derived from a *generative model*. It demonstrates that the data distribution can be identified by neural network based generative models, e.g., variational auto-encoders (VAEs) [12] and generative adversarial networks (GANs) [8]. The neural network based generative model learns a generator  $G(\mathbf{z}) : \mathbf{z} \in \mathbb{R}^k \rightarrow G(\mathbf{z}) \in \mathbb{R}^n$  that maps a low dimensional space  $\mathbf{z}$  to the high-dimensional sample space. This generator is trained to generate vectors that approximate the vectors in the training dataset. Here, the generator characterizes a probability distribution over vectors in the sample space, and based on the training dataset, the generator is trained to allocate higher probabilities to more likely vectors. Thus, the notion of a vector in certain space can be generally captured by a pre-trained generator. It was shown in



**Fig. 6.1** The performance of compressed sensing under generative models

[3] that a vector in the space of a given system setting is close to some point in the range of  $G$ .

The paper [3] also proposed an algorithm that exploits generative models to solve the compressed sensing problem. This algorithm optimizes the variable  $\mathbf{z} \in \mathbb{R}^k$  via gradient descent such that the corresponding generator  $G(\mathbf{z})$  yields a small measurement error, i.e.,

$$\|AG(\mathbf{z}) - \mathbf{y}\|_2^2. \quad (6.2)$$

Even though the objective function (6.2) is nonconvex, it was empirically demonstrated in [3] that gradient descent works well, and can yield significantly better performance than Lasso with relatively few measurements. Figure 6.1 illustrates the comparison of signal reconstruction from compressed linear measurements with the sparse linear model and a generative model, i.e., VAE. The experiment is run on the MNIST dataset, a classic 10-class hand-written digit classification dataset. Each pixel value of an image is either 0 (background) or 1 (foreground), so the digit images are reasonably sparse in the pixel space. Three different algorithms are considered: the classic Lasso algorithm for the sparse linear model in the pixel space, and two algorithms, “VAE” and “VAE+Reg,” for the generative model with loss functions

$$\|AG(\mathbf{z}) - \mathbf{y}\|_2^2,$$

and

$$\|AG(\mathbf{z}) - \mathbf{y}\|_2^2 + \lambda \|\mathbf{z}\|_2,$$

where  $G(\mathbf{z})$  is the generative model trained on the MNIST dataset. This result shows that signal reconstruction with generative models requires  $10\times$  fewer measurements than a conventional sparse linear model to achieve 10% error.

The paper [3] provided the theoretical performance guarantee. Specifically, it was demonstrated that, as long as a good approximate solution to the objective (6.2) is found by gradient descent, the generator  $G(\mathbf{z})$ , which yields the closest possible point in the range of  $G$ , will be sufficiently close to the ground truth  $\mathbf{x}^*$ . The proof provided in [3] relies on the set-restricted eigenvalue condition which is a generalization of the restricted eigenvalue condition (*REC*). Moreover, [3] shows that for some generators, e.g., VAEs and GANs, random Gaussian measurement matrices can satisfy the set-restricted eigenvalue condition with high probability. That is, for  $d$ -layer neural networks,  $O(kd \log n)$  Gaussian measurements sufficiently guarantee high-accuracy reconstruction with high probability. Specifically, the result states as below.

**Theorem 6.1 ([3])** *Assuming that  $G : \mathbb{R}^k \rightarrow \mathbb{R}^n$  is an  $L$ -Lipschitz function, and  $A \in \mathbb{R}^{m \times n}$  is a random Gaussian matrix for  $m = O(k \log \frac{Lr}{\delta})$ , obeying  $A_{ij} \sim N(0, 1/m)$ . For any  $\mathbf{x}^* \in \mathbb{R}^n$  and observation  $\mathbf{y} = A\mathbf{x}^* + \boldsymbol{\eta}$ , let the estimator  $\hat{\mathbf{z}}$  minimize (6.2) to within additive  $\varepsilon$  of the optimum over vectors with  $\|\hat{\mathbf{z}}\|_2 \leq r$ . Then with  $1 - e^{-\Omega(m)}$  probability, there is*

$$\|G(\hat{\mathbf{z}}) - \mathbf{x}^*\|_2 \leq 6 \min_{\substack{\mathbf{z}^* \in \mathbb{R}^k \\ \|\mathbf{z}^*\|_2 \leq r}} \|G(\mathbf{z}^*) - \mathbf{x}^*\|_2 + 3\|\boldsymbol{\eta}\|_2 + 2\varepsilon + 2\delta. \quad (6.3)$$

The first two terms on the right-hand side of (6.3) identify the minimum error of any vector in the range of the generator and the norm of the noise, respectively. The third term  $\varepsilon$  comes from the distance between the global optimum and the convergence result generated by gradient descent.

The results of [3] have inspired lots of follow-up studies. Recently, the paper [23] proposed a novel framework that significantly improves both the performance and speed of signal recovery by jointly training a generator and the optimization process for reconstruction via meta-learning. The paper explored training the measurements with different objectives, and derived a family of models based on minimizing measurement errors. We will provide an overview of this work in Sect. 6.2.

Besides the compressed sensing problem, the generative prior has also been applied to the blind image deconvolution problem [1]. It can also provide some insights on applying the generative prior in the blind demixing problem introduced in Chap. 3 and the sparse blind demixing problem introduced in Chap. 4.

## 6.2 Joint Design of Measurement Matrix and Sparse Support Recovery

The design of the measurement matrix in compressed sensing is of critical importance for both practical implementation and performance enhancement (e.g., achieving a better compression or allowing a higher signal reconstruction quality). Thus, it has received intensive attentions, and good progresses have been made. Learning augmented techniques have recently exploited in joint design of measurement matrix and sparse support recovery. This section first introduces basic methods for solving this problem, followed by the learning augmented methods.

**Sample Scheduling** The paper [10] proposed an adaptive CS based sample scheduling mechanism (ACS) with respect to different per-sampling-window bases for wireless sensor networks. For each basis, given a sensing quality, ACS estimates the minimum required sample rate, thereby correspondingly adjusting sensors' sample rates.

**Sensing Matrix Optimization** To optimize sensing matrices, some techniques known as mutual coherence minimization [6, 7, 14, 24] are developed, without additional assumption on the class of acquired signals. Another line of research shows that better results can be achieved with some priors on the input signal. For instance, when the energy of the signals to be acquired is not evenly distributed, i.e., when they are both sparse and localized. Mathematically, for the sparse signal  $\mathbf{x}$  in compressed sensing (1.5), it holds that  $\mathbb{E}(\mathbf{x}\mathbf{x}^\top)$  is not a multiplier of the  $n \times n$  identity matrix  $\mathbf{I}_n$ . To characterize this property, the paper [14] introduced a design criterion, which is called *rakeness*, to identify the amount of energy that the measurements seize from the acquired signal. The proposed *rakeness* approach [14] aligns statistical properties of the compressed sensing stage with that of the input signal  $\mathbf{x}$ , while simultaneously preserving conditions for a correct signal reconstruction required by the standard compressed sensing theory.

Following the idea proposed in [14], the paper [15] proposed sensing matrix optimization techniques that exploit statistical properties of the process generating  $\mathbf{x}$  in compressed sensing (1.5). One method is nearly orthogonal CS, which is based on a geometric constraint enforcing diversity between different compressed measurements. Another method, named, maximum-energy CS, is a heuristic screening of candidate measurements that relies on a self-adapted optimization procedure.

**Learning Augmented Methods** Recent works [13, 17, 18, 22, 23] consider joint design of signal compression and recovery methods using auto-encoder [13, 17, 18, 22] and generative adversarial networks (GAN) [23] in deep learning. In particular, linear compression for real signals was considered in [18, 22]; nonlinear compression for real signals was considered in [17, 23]. The paper [13] studied linear compression for complex signals.

The fundamental idea of joint design of signal compression and recovery methods using auto-encoder [13, 17, 18] can be summarized as follows.

- Collect an input signal  $\mathbf{x}^{(i)}$  from a training set  $\mathcal{D}_{\text{train}} = \{\mathbf{x}^{(i)}\}_{i=1}^s$ .
- Reconstruct input's components and reduce its dimensionality via  $d$  layer operations such as convolutional layer [17], linear reduction mapping [18], etc. Denote the set of parameters at  $d$  layers as  $\Omega = \{\mathbf{W}_j, \mathbf{b}_j\}_{j=1}^d$ .
- Take undersampled measurements.
- Increase measurements dimensionality via operations such as convolutional layers [17], nonlinear inverse mapping [18], etc.
- Convert the output to a reconstructed signal. Denote this mapping from original signals to reconstructed signals as  $\hat{\mathbf{x}} = \mathcal{F}(\mathbf{x}, \Omega)$ .

The mean-squared error (MSE) can be adopted as a loss function over the training set  $\mathcal{D}_{\text{train}}$

$$\mathcal{L}(\Omega) = \frac{1}{s} \sum_{i=1}^s \left\| \mathcal{F}(\mathbf{x}, \Omega) - \mathbf{x}^{(i)} \right\|_2^2. \quad (6.4)$$

The stochastic gradient descent (SGD) or ADMM optimizer [11] can be applied for minimizing  $\mathcal{L}(\Omega)$  (6.4) and learning parameters. Recently, the work [13] extended the methods in [17, 18] to joint linear compression and recovery methods for complex signal estimation, which is more challenging. The proposed architecture includes two components, an auto-encoder and a hard thresholding module. The proposed auto-encoder successfully deals with complex signals via exploiting standard auto-encoder for real numbers. The key technique is to establish the encoder which mimics the noisy linear measurement process. The model for complex numbers in compressed sensing, i.e.,

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{z},$$

can be equivalently expressed via the following two expressions of real numbers:

$$\Re(\mathbf{y}) = \Re(\mathbf{A})\Re(\mathbf{x}) - \Im(\mathbf{A})\Im(\mathbf{x}) + \Re(\mathbf{z}), \quad (6.5)$$

$$\Im(\mathbf{y}) = \Im(\mathbf{A})\Re(\mathbf{x}) + \Re(\mathbf{A})\Im(\mathbf{x}) + \Im(\mathbf{z}). \quad (6.6)$$

Besides the decoders mentioned above, a recent paper [22] presented a  $\ell_1$  decoder to learn linear encoders that adjust to data. The convex and non-smooth  $\ell_1$  decoder cannot be trained via standard gradient-based, i.e., gradient propagation. To address this issue, the paper [22] relies on the idea of unrolling the convex decoder into  $T$  projected subgradient steps. Denote the  $\ell_1$ -minimization as:

$$L(\mathbf{A}, \mathbf{y}) := \arg \min_{\mathbf{x} \in \mathbb{R}^d} \|\mathbf{x}\|_1 \quad \text{s.t. } \mathbf{A}\mathbf{x} = \mathbf{y}. \quad (6.7)$$

Mathematically, given a training set  $\mathcal{D}_{\text{train}}$ , the problem of finding the best  $\mathbf{A}$  can be formulated as

$$\min_{\mathbf{A} \in \mathbb{R}^{m \times d}} f(\mathbf{A}) := \sum_{i=1}^s \|\mathbf{x}^{(i)} - L(\mathbf{A}, \mathbf{A}\mathbf{x}^{(i)})\|_2^2.$$

Here  $L(\cdot, \cdot)$  is the  $\ell_1$  decoder defined in (6.7). Unfortunately, it is difficult to compute the gradient of  $f(\mathbf{A})$  with respect to  $\mathbf{A}$ , due to the optimization problem defined in (6.7). The paper [22] addressed this issue by replacing the  $\ell_1$ -minimization with the iterations of  $T$ -step projected subgradient, which approximately computes the gradients. Define an approximate function  $\tilde{f}(\mathbf{A}) : \mathbb{R}^{m \times d} \mapsto \mathbb{R}$ , and this procedure can be represented as

$$\begin{aligned} \tilde{f}(\mathbf{A}) &:= \sum_{i=1}^s \|\mathbf{x}^{(i)} - \hat{\mathbf{x}}^{(i)}\|_2^2, \quad \text{where} \\ \hat{\mathbf{x}}^{(i)} &= T\text{-step projected subgradient of} \\ &L(\mathbf{A}, \mathbf{A}\mathbf{x}^{(i)}), \text{ for } i = 1, \dots, s, \end{aligned} \tag{6.8}$$

which is called *unrolling*.

Another type of data-driven approaches is using learning augmented generative adversarial networks (GAN) [23] for joint design of signal compression and recovery methods. The paper [23] generalized the measurement matrix  $\mathbf{A}$  in compressed sensing under a generative prior  $G_\phi(\cdot)$ , i.e., (6.2). To achieve this, the paper [23] defines a measurement function  $\mathbf{y} \leftarrow F_\phi(\mathbf{x})$  where  $\mathbf{x} = G_\phi(\mathbf{z})$ , thus both  $F_\phi$  and  $G_\phi$  can be trained via deep neural networks. The key point of this generalized setting is recovering the signal  $\mathbf{x}$  from inverting the measurement function  $\mathbf{x} \leftarrow F_\phi^{-1}(\mathbf{y})$  via minimizing the measurement error:

$$E_\phi(\mathbf{y}, \mathbf{z}) = \|\mathbf{y} - F_\phi(G_\phi(\mathbf{z}))\|_2^2. \tag{6.9}$$

### 6.3 Deep-Learning-Based AMP

Deep learning recently has achieved great successes in many applications, which has inspired recent developments of deep-learning-based methods for structured signal processing. In this section, we introduce two neural network architectures proposed in [4]. Similar to the approximate message passing (AMP) algorithms that decouple prediction errors across iterations, the deep-learning-based AMP in [4] decouples prediction errors across layers. In [4], the proposed methods were applied to solve the compressive random access and massive-MIMO channel estimation in 5G networks. These methods also bring some insights on developing deep-learning-

based AMP to solve the joint device activity detection and channel estimation problem in massive IoT networks.

We first introduce the “learned AMP” network proposed in [4], followed by the “learned VAMP” network [4] that offers increased robustness to deviations in the i.i.d. Gaussian measurement matrix. In both cases, the linear transforms and scalar nonlinearities of the network are simultaneously learned. An straightforward interpretation of learned VAMP is demonstrated in [4] that with i.i.d. measurements, the linear transforms and scalar nonlinearities established by the VAMP algorithm follow the values learned through back-propagation. Furthermore, we introduce a learned-based algorithm for group row sparsity (LISTA-GS) to estimate the sparse linear model (2.8) in the multiple-antenna scenario in Sect. 6.3.3.

$$\mathbf{Y} = \mathbf{Q}\Theta + \mathbf{N}. \quad (6.10)$$

Besides the approaches introduced in this section, the authors in [9, 16, 20, 25] exploited properties of sparsity patterns of real signals [9, 16, 25] and complex signals [20] from training samples using data-driven approaches based on deep learning, which also brings some insights for future study.

### 6.3.1 Learned AMP

For compressed sensing (2.3), i.e.,

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{n},$$

where  $\mathbf{y} \in \mathbb{C}^L$ ,  $\mathbf{x} \in \mathbb{C}^N$ ,

$$\mathbf{n} \in \mathbb{C}^L \sim \mathcal{CN}(\mathbf{0}, \sigma^2 \mathbf{I}) \quad (6.11)$$

is the additive white Gaussian noise, the sparse signal  $\mathbf{x}$  can be estimated by Lasso (6.1) given by

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1.$$

It can be solved by the approximate message passing algorithm (2.27) introduced in Sect. 2.4.1:

$$\mathbf{r}_t = \mathbf{y} - \mathbf{A}\hat{\mathbf{x}}_t + \frac{1}{M} \|\hat{\mathbf{x}}_t\|_0 \mathbf{r}_{t-1} \quad (6.12a)$$

$$\hat{\mathbf{x}}_{t+1} = \eta_{st} \left( \hat{\mathbf{x}}_t + \mathbf{A}^\top \mathbf{r}_t; \frac{\alpha}{\sqrt{M}} \|\mathbf{r}_t\|_2 \right), \quad (6.12b)$$

where the initial points are set as  $\hat{\mathbf{x}}_0 = \mathbf{0}$ ,  $\mathbf{r}_{-1} = \mathbf{0}$ ,  $t \in \{0, 1, 2, \dots\}$ , and  $\eta_{st}(\cdot; \lambda) : \mathbb{R}^N \rightarrow \mathbb{R}^N$  is the ‘‘soft thresholding’’ shrinkage function, component wisely defined as

$$[\eta_{st}(\mathbf{r}; \lambda)]_i \triangleq \text{sgn}(r_i) \max\{|r_i| - \lambda, 0\}. \quad (6.13)$$

In (6.12b),  $\alpha$  is a tuning parameter that correlates with  $\lambda$  in (6.1). The AMP-inspired deep networks for solving sparse linear problems have been proposed in [4], which are introduced in the sequel.

The paper [4] established a neural network via unfolding the iterations of AMP- $\ell_1$  from (6.12), followed by learning the MSE-minimal values of the network parameters, which is called ‘‘LAMP- $\ell_1$ .’’ The  $t$ -th layer of the LAMP- $\ell_1$  network is represented by

$$\hat{\mathbf{x}}_{t+1} = \beta_t \eta_{st} \left( \hat{\mathbf{x}}_t + \mathbf{B}_t \mathbf{r}_t; \frac{\alpha_t}{\sqrt{M}} \|\mathbf{r}_t\|_2 \right) \quad (6.14a)$$

$$\mathbf{r}_{t+1} = \mathbf{y} - \mathbf{A} \hat{\mathbf{x}}_{t+1} + \frac{\beta_t}{M} \|\hat{\mathbf{x}}_{t+1}\|_0 \mathbf{r}_t, \quad (6.14b)$$

where the first-layer inputs are set as  $\hat{\mathbf{x}}_0 = \mathbf{0}$  and  $\mathbf{r}_0 = \mathbf{y}$ . The paper [4] refers to networks that use fixed  $\mathbf{B}$  over all layers  $t$  as ‘‘tied,’’ where the LAMP- $\ell_1$  parameters are  $\boldsymbol{\Omega} = \{\mathbf{B}, \{\alpha_t, \beta_t\}_{t=0}^{T-1}\}$ . Those that depend on  $t$ , i.e.,  $\mathbf{B}_t$ , are named as ‘‘untied,’’ where the LAMP- $\ell_1$  parameters are  $\boldsymbol{\Omega} = \{\mathbf{B}_t, \alpha_t, \beta_t\}_{t=0}^{T-1}$ . The parameters  $\boldsymbol{\Omega}$  of ‘‘tied’’ case and ‘‘untied’’ case can be further learned by minimizing the MSE on the training data, which are illustrated in Algorithms 6.1 and 6.2, respectively. In [4], it was demonstrated that LAMP- $\ell_1$  yields a faster convergence rate than AMP from both empirical and theoretical points of view.

---

Algorithm 6.1: Tied LAMP- $\ell_1$  parameter learning [4]

---

- 1: Input  $\mathbf{B} = \mathbf{I}$ ,  $\alpha_0 = 1$ ,  $\beta_0 = 1$
  - 2: Learn  $\boldsymbol{\Omega}_0^{\text{tied}} = \{\mathbf{B}, \alpha_0\}$
  - 3: for  $t = 1$  to  $T - 1$  do
  - 4:   Initialize  $\alpha_t = \alpha_{t-1}$ ,  $\beta_t = \beta_{t-1}$
  - 5:   Learn  $\{\alpha_t, \beta_t\}$  with fixed  $\boldsymbol{\Omega}_{t-1}^{\text{tied}}$
  - 6:   Re-learn  $\boldsymbol{\Omega}_t^{\text{tied}} = \{\mathbf{B}, \{\alpha_i, \beta_i\}_{i=1}^t, \alpha_0\}$
  - 7: end for
  - 8: Return  $\boldsymbol{\Omega}_{T-1}^{\text{tied}}$
-

---

**Algorithm 6.2: Untied LAMP- $\ell_1$  parameter learning [4]**


---

- 1: Learn  $\{\mathbf{\Omega}_t^{\text{tied}}\}_{t=1}^{T-1}$  via Algorithm 6.1
  - 2: Initialize  $\mathbf{B}_0 = \mathbf{I}$ ,  $\alpha_0 = 1$ ,  $\beta_0 = 1$
  - 3: Learn  $\mathbf{\Omega}_0^{\text{untied}} = \{\mathbf{B}_0, \alpha_0\}$
  - 4: for  $t = 1$  to  $T - 1$  do
  - 5:   Initialize  $\mathbf{B}_t = \mathbf{B}_{t-1}$ ,  $\alpha_t = \alpha_{t-1}$ ,  $\beta_t = \beta_{t-1}$
  - 6:   Learn  $\{\mathbf{B}_t, \alpha_t, \beta_t\}$  with fixed  $\mathbf{\Omega}_{t-1}^{\text{untied}}$
  - 7:   Set  $\mathbf{\Omega}_t^{\text{untied}} = \{\mathbf{B}_t, \alpha_t, \beta_t\}_{t=0}^t \setminus \beta_0$  (“ $\setminus$ ” denotes the set difference operation)
  - 8:   if  $\mathbf{\Omega}_t^{\text{tied}}$  enjoys better performance than  $\mathbf{\Omega}_t^{\text{untied}}$  then
  - 9:     Replace  $\mathbf{\Omega}_t^{\text{untied}}$  with  $\mathbf{\Omega}_t^{\text{tied}}$
  - 10:   end if
  - 11:   Re-learn  $\mathbf{\Omega}_t^{\text{untied}}$
  - 12: end for
  - 13: Return  $\mathbf{\Omega}_{T-1}^{\text{untied}}$
- 

### 6.3.2 Learned Vector-AMP

The VAMP algorithm illustrated in Algorithm 6.3 has been recently proposed in [19] to address AMP’s fragility concerning the matrix  $\mathbf{A}$ . Compared to the original AMP, the VAMP algorithm enjoys lower per-iteration complexity and fewer iterations required to convergence. The procedure of the VAMP algorithm is elaborated in the following.

We begin with the definition of the right-rotationally invariant matrices. For the matrix  $\mathbf{A} \in \mathbb{R}^{L \times N}$  in compressed sensing (2.3), suppose that

$$\mathbf{A} = \mathbf{U} \mathbf{\Lambda} \mathbf{V}^T \quad (6.15)$$

satisfies that  $\mathbf{s} \in \mathbb{R}_+^r$  where  $r \triangleq \text{rank}(\mathbf{A})$  contains the positive singular values of  $\mathbf{A}$ , then  $\mathbf{\Lambda} = \text{diag}(\mathbf{s}) \in \mathbb{R}^{r \times r}$ ,  $\mathbf{U}^T \mathbf{U} = \mathbf{I}_r$ , and  $\mathbf{V}^T \mathbf{V} = \mathbf{I}_r$ . The matrix  $\mathbf{A}$  is *right-rotationally invariant* if  $\mathbf{V}$  consists of the first  $r$  columns of a random matrix uniformly distributed on the group of  $n \times n$  orthogonal matrices. With any random orthogonal  $\mathbf{U}$  and a particular distribution on  $\mathbf{s}$ , i.i.d. Gaussian matrices are right-rotationally invariant. The paper [19] demonstrates that with large enough dimensions  $m, n$ , VAMP behaves well when the sensing matrix  $\mathbf{A}$  in compressed sensing is an i.i.d. Gaussian matrix.

The VAMP algorithm consists of two stages which endow with different estimators: LMMSE stage with the estimator

$$\tilde{\boldsymbol{\eta}}(\tilde{\mathbf{r}}_t; \tilde{\sigma}_t, \hat{\theta}) := \mathbf{V} \left( \text{diag}(\mathbf{s})^2 + \frac{\sigma^2}{\tilde{\sigma}_t^2} \mathbf{I}_R \right)^{-1} \left( \text{diag}(\mathbf{s}) \mathbf{U}^T \mathbf{y} + \frac{\sigma^2}{\tilde{\sigma}_t^2} \mathbf{V}^T \tilde{\mathbf{r}}_t \right), \quad (6.16)$$

where  $\sigma$  is the standard deviation, and the parameter  $\hat{\theta}$  is given by

$$\hat{\theta} := \{\mathbf{U}, \mathbf{s}, \mathbf{V}, \sigma\}, \quad (6.17)$$

---

**Algorithm 6.3: Vector-AMP [19]**


---

Require: LMMSE estimator  $\tilde{\eta}(\cdot; \tilde{\sigma}, \hat{\theta})$  (6.16), shrinkage  $\eta(\cdot; \sigma, \omega)$  (6.18), max iteration  $T$ , parameters  $\{\omega_t\}_{t=1}^T$  and  $\hat{\theta}$ .

- 1: Set initial  $\tilde{\mathbf{r}}_1$  and  $\tilde{\sigma}_1 > 0$ .
- 2: for  $t = 1, 2, \dots, T$  do
- 3:   // LMMSE stage:
- 4:    $\tilde{\mathbf{x}}_t = \tilde{\eta}(\tilde{\mathbf{r}}_t; \tilde{\sigma}_t, \hat{\theta})$    // estimation
- 5:    $\tilde{\mathbf{u}}_t = \langle \tilde{\eta}'(\tilde{\mathbf{r}}_t; \tilde{\sigma}_t, \hat{\theta}) \rangle$    // divergence computation
- 6:    $\mathbf{r}_t = (\tilde{\mathbf{x}}_t - \tilde{\mathbf{u}}_t \tilde{\mathbf{r}}_t) / (1 - \tilde{\mathbf{u}}_t)$    // Onsager correction
- 7:    $\tilde{\sigma}_t^2 = \tilde{\sigma}_t^2 \tilde{\mathbf{u}}_t / (1 - \tilde{\mathbf{u}}_t)$    // variance computation
- 8:   // Shrinkage stage:
- 9:    $\hat{\mathbf{x}}_t = \eta(\mathbf{r}_t; \sigma_t, \omega_t)$    // estimation
- 10:    $\mathbf{u}_t = \langle \eta'(\mathbf{r}_t, \sigma_t, \omega_t) \rangle$    // divergence computation
- 11:    $\tilde{\mathbf{r}}_{t+1} = (\hat{\mathbf{x}}_t - \mathbf{u}_t \mathbf{r}_t) / (1 - \mathbf{u}_t)$    // Onsager correction
- 12:    $\tilde{\sigma}_{t+1}^2 = \sigma_t^2 \mathbf{u}_t / (1 - \mathbf{u}_t)$    // variance computation
- 13: end for
- 14: Return  $\hat{\mathbf{x}}_T$ .

---

and a shrinkage stage with the estimator

$$\eta(\mathbf{r}_t; \sigma_t, \alpha) = \eta_{st}(\mathbf{r}_t; \alpha \sigma_t), \quad (6.18)$$

where  $\eta_{st}(\cdot, \cdot)$  is given by (6.12b). Lines 5 and 10 in Algorithm 6.3 compute the average of the diagonal entries of the Jacobian of  $\tilde{\eta}(\cdot; \tilde{\sigma}_t, \hat{\theta})$  and  $\eta(\cdot; \sigma_t, \omega_t)$ , respectively, which can be referred to [4] for detailed representation.

Based on VAMP illustrated in Algorithm 6.3, we move to the learned VAMP (LVAMP) algorithm proposed in [4]. The  $t$ -th layer of the learned VAMP (LVAMP) network consists of four stages: (1) LMMSE optimization, (2) decoupling, (3) shrinkage, and (4) decoupling. The learnable parameters in the  $t$ -th layer are the LMMSE stage parameters  $\hat{\omega}_t = \{\mathbf{U}_t, \mathbf{s}_t, \mathbf{V}_t, \sigma_t^2\}$  (6.17) and the shrinkage parameters  $\omega_t$ . Similar to VAMP, the network parameters of LVAMP are concerned in two cases: “tied” and “untied.” In the tied case, the network parameters are  $\{\hat{\theta}, \{\omega_t\}_{t=1}^T\}$ , while in the untied case, it is  $\{\hat{\theta}_t, \omega_t\}_{t=1}^T$ . Algorithms 6.1 and 6.2 can be exploited to learn the LVAMP parameters for the tied case and the untied case, respectively (with  $\hat{\theta}_t$  replacing  $\mathbf{B}_t$  and with  $\theta_t$  replacing  $\{\alpha_t, \beta_t\}$ ).

### 6.3.3 Learned ISTA for Group Row Sparsity

Recall the sparse linear model (2.8) in the multiple-antenna scenario discussed in Sect. 2.2.2, represented by

$$\mathbf{Y} = \mathbf{S}\mathbf{X} + \mathbf{Z}, \quad (6.19)$$

where  $\mathbf{Y} \in \mathbb{C}^{L \times M}$ ,  $\mathbf{X} \in \mathbb{C}^{N \times M}$  endowed with group row sparsity.

---

**Algorithm 6.4: ISTA-GS**


---

 Require:  $\tilde{\mathbf{S}}, \tilde{\mathbf{Y}}, \lambda$ , Iterations.

 ISTA-GS( $\tilde{\mathbf{S}}, \tilde{\mathbf{Y}}, \tilde{\mathbf{X}}, \lambda$ , Iterations):

- 1: Initialize:  $\tilde{\mathbf{X}} \leftarrow \mathbf{0}$ ,  $C \leftarrow$  largest eigenvalue of  $\tilde{\mathbf{S}}^T \tilde{\mathbf{S}}$
  - 2: for  $i = 1$  to Iterations do
  - 3:    $\tilde{\mathbf{X}} = \eta_{\lambda/C}(\tilde{\mathbf{X}} + \frac{1}{C} \tilde{\mathbf{S}}^T (\tilde{\mathbf{Y}} - \tilde{\mathbf{S}} \tilde{\mathbf{X}}))$
  - 4: end for
  - 5: return  $\tilde{\mathbf{X}}$ .
- 

Recall the real-valued counterpart (2.13), and the real-valued counterpart of (6.19) can be represented as

$$\begin{aligned} \tilde{\mathbf{Y}} &= \tilde{\mathbf{S}} \tilde{\mathbf{X}} + \tilde{\mathbf{Z}} \\ &= \begin{bmatrix} \Re\{\tilde{\mathbf{S}}\} - \Im\{\tilde{\mathbf{S}}\} \\ \Im\{\tilde{\mathbf{S}}\} \quad \Re\{\tilde{\mathbf{S}}\} \end{bmatrix} \begin{bmatrix} \Re\{\tilde{\mathbf{X}}\} \\ \Im\{\tilde{\mathbf{X}}\} \end{bmatrix} + \begin{bmatrix} \Re\{\tilde{\mathbf{Z}}\} \\ \Im\{\tilde{\mathbf{Z}}\} \end{bmatrix}. \end{aligned} \quad (6.20)$$

The following problem can be established to estimate the group sparse  $\tilde{\mathbf{X}}$ :

$$\underset{\tilde{\mathbf{X}} \in \mathbb{R}^{2N \times M}}{\text{minimize}} \|\tilde{\mathbf{Y}} - \tilde{\mathbf{S}} \tilde{\mathbf{X}}\|_F^2 + \lambda \mathcal{R}(\tilde{\mathbf{X}}). \quad (6.21)$$

To solve problem (6.21), we start with the ISTA for group row sparse (*ISTA-GS*) illustrated in Algorithm 6.4. Specifically, in the  $k$ -th iteration, the update rule is represented as

$$\tilde{\mathbf{X}}^{k+1} = \eta_{\lambda/C} \left( \tilde{\mathbf{X}}^k + \frac{1}{C} \tilde{\mathbf{S}}^T (\tilde{\mathbf{Y}} - \tilde{\mathbf{S}} \tilde{\mathbf{X}}^k) \right), \quad (6.22)$$

where  $C$  is the largest eigenvalue of matrix  $\tilde{\mathbf{S}}^T \tilde{\mathbf{S}}$ , and  $\eta_\theta(\mathbf{x}^n)$  denotes the group soft-thresholding function for the  $n$ -th row in matrix  $\mathbf{X}$  (i.e.,  $\mathbf{x}^n$ ) [2]. Specifically,  $\eta_\theta(\mathbf{x}^n)$  is defined as

$$\eta_\theta(\mathbf{x}^n) = \max \left\{ 0, \frac{\|\mathbf{x}\|_2 - \theta}{\|\mathbf{x}\|_2} \right\} \mathbf{x}^n. \quad (6.23)$$

Such an iterative algorithm takes a large number of iterations to converge. To address this issue, we present the LISTA-GS method, which parameterizes the iterative method.

Inspired by [5, 9] and by denoting  $\mathbf{W}_1 = \frac{1}{C} \tilde{\mathbf{S}}$ ,  $\mathbf{W}_2 = \mathbf{I} - \frac{1}{C} \tilde{\mathbf{S}}^T \tilde{\mathbf{S}}$ , and  $\theta = \frac{\lambda}{C}$ , we rewrite (6.22) as

$$\tilde{\mathbf{X}}^{k+1} = \eta_{\theta^k}(\mathbf{W}_1 \tilde{\mathbf{Y}} + \mathbf{W}_2 \tilde{\mathbf{X}}^k). \quad (6.24)$$

The key idea of the proposed LISTA-GS method is to view matrix  $\mathbf{W}_1$ , matrix  $\mathbf{W}_2$ , and scalar  $\theta^k$  in (6.24) as trainable parameters. As a result, (6.24) can be modeled as a one-layer RNN. Moreover, the unrolled RNN with  $K$  iterations for group row sparse can be expressed as

$$\tilde{\mathbf{X}}^{k+1} = \eta_{\theta^k} (\mathbf{W}_1^k \tilde{\mathbf{Y}} + \mathbf{W}_2^k \tilde{\mathbf{X}}^k), k = 0, 1, \dots, K - 1, \quad (6.25)$$

where all parameters  $\Theta = \{\mathbf{W}_1^k, \mathbf{W}_2^k, \theta^k\}_{k=0}^{K-1}$  are trainable. This is a main difference from the problem formulation in (6.22).

### 6.3.3.1 Simulations Results

In this section, we present the simulation results of the LISTA-GS method for the joint device activity detection and channel estimation, and compare the results with that of the ISTA-GS method.

In simulations, we generate the signature matrix according to the complex Gaussian distribution, i.e.,  $\mathbf{S} \sim \mathcal{CN}(\mathbf{0}, \mathbf{I})$ . The channels are assumed to suffer from independent Rayleigh fading, i.e.,  $\mathbf{H} \sim \mathcal{CN}(\mathbf{0}, \mathbf{I})$ . In addition, we set the length of the signature sequence ( $L$ ), the total number of devices ( $N$ ), and the number of antennas at the BS ( $M$ ) as 100, 200, and 2, respectively. Each entry of the activity sequence  $\{a_1, \dots, a_N\}$  follows the Bernoulli distribution with mean  $p = 0.1$ , i.e.,  $\mathbb{P}(a_n = 1) = 0.1$  and  $\mathbb{P}(a_n = 0) = 0.9, \forall n \in \mathcal{N}$ . After normalizing  $\mathbf{S}$ , we transform all these complex-valued matrices into real-value matrices according to (6.20). Hence, we obtain the training data set  $\{\tilde{\mathbf{X}}_i^*, \tilde{\mathbf{Y}}_i\}_{i=1}^N$ . In the training stage, the batch size is set to be 64 and the validation set contains 1000 samples. The learning rate is set to be  $10^{-3}$ . In the testing stage, 1000 samples are generated to test the trained LISTA-GS model. As for ISTA-GS, we set  $\lambda = 0.2, 0.1$ , and  $0.05$ . We choose  $K = 16$  layers for the LISTA-GS method in all the simulations. Furthermore, we initialize the parameters as  $W_1^0 = \frac{1}{C} \tilde{\mathbf{S}}, W_2^0 = \mathbf{I} - \frac{1}{C} \tilde{\mathbf{S}}^T \tilde{\mathbf{S}}$ , and  $\theta = \frac{0.1}{C}$ .

We adopt the normalized mean square error (NMSE) to evaluate the performance of LISTA-GS and ISTA-GS in recovering the real-valued  $\tilde{\mathbf{X}}$ , defined as

$$\text{NMSE}(\tilde{\mathbf{X}}, \tilde{\mathbf{X}}^*) = 10 \log_{10} \left( \frac{\mathbb{E} \|\tilde{\mathbf{X}} - \tilde{\mathbf{X}}^*\|_F^2}{\mathbb{E} \|\tilde{\mathbf{X}}^*\|_F^2} \right), \quad (6.26)$$

where  $\tilde{\mathbf{X}}^*$  represents the ground truth and  $\tilde{\mathbf{X}}$  is the estimate obtained by the ISTA-GS and LISTA-GS methods.

As suggested in [5], we train the LISTA-GS model by adopting the layer-wise training strategy, which has been widely used in the previous ISTA models. To stabilize the training process, we add two decayed learning rates, i.e.,  $\beta_1 = 0.2\beta_0$  and  $\beta_2 = 0.02\beta_0$ , where  $\beta_0$  is the initial learning rate. Note that  $\Theta^i =$

$\{\mathbf{W}_1^k, \mathbf{W}_2^k, \theta^k\}_{k=0}^i$  are all the weights from layer 0 to layer  $i$  and  $m(\cdot)$  is the learning multiplier. We train the RNN layer by layer and the training process of each layer is described as follows:

- Suppose that  $\Theta^{i-1}$  is pre-trained for layer  $i$ . Initialize the learning multipliers  $m(W_1^i), m(W_2^i), m(\theta^i) = 1$ .
- Train  $\{\mathbf{W}_1^i, \mathbf{W}_2^i, \theta^i\}$  with  $\beta_0$ .
- Multiply the learning multiplier to their weights and train  $\Theta^i = \Theta^{i-1} \cup \{\mathbf{W}_1^i, \mathbf{W}_2^i, \theta^i\}$  with  $\beta_1$  and  $\beta_2$ .
- Multiply a decay rate to each learning multiplier.
- Move to train the next layer.

Figure 6.2a shows the NMSE of the proposed LISTA-GS and the baseline ISTA-GS methods over iterations in a noiseless scenario. For the baseline ISTA-GS method, there exists an inherent tradeoff between the convergence rate and the NMSE. In particular, a smaller value of  $\lambda$  results in a more accurate solution but leads to a lower convergence rate, and vice versa. Besides, we observe that LISTA-GS method achieves a much faster convergence rate as well as a much lower NMSE than ISTA-GS for different values of  $\lambda$ . This is because LISTA-GS treats  $\lambda$  as a weight in the training process, yielding a good solution that balances the tradeoff.

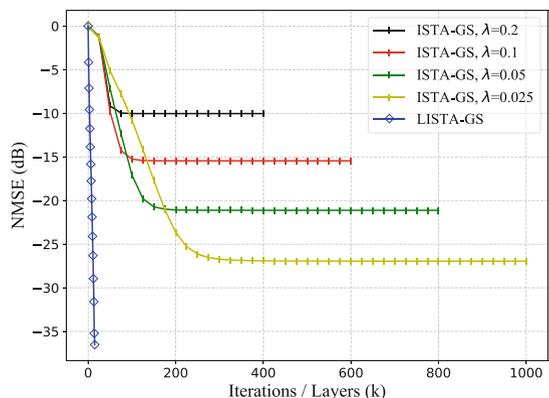
Figure 6.2b illustrates the NMSE of the proposed LISTA-GS method over iterations in a noisy scenario with different values of signal-to-noise-ratio (SNR). It can be observed that LISTA-GS can also reach convergence in a few iterations (e.g., less than 16) in the noisy case. As the value of SNR increases, the received power of the pilot sequence increases, which in turn decreases the achievable NMSE.

Figure 6.2c plots the impact of SNR on the NMSE of LISTA-GS and ISTA-GS. The proposed LISTA-GS method achieves a much lower NMSE than ISTA-GS for different values of SNR. In addition, the NMSEs of both methods decrease as the value of SNR increases.

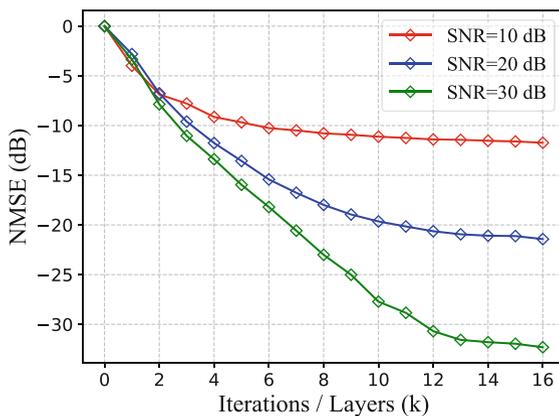
## 6.4 Summary

In this chapter, we introduced some cutting-edge learning augmented techniques for both structured signal modeling (e.g., structured signal processing under a generative prior [1, 3]) and algorithm design (e.g., learning augmented algorithms [4]). We also introduced Learned ISTA for group row sparsity to solve the sparse linear model in the multiple-antenna scenario. These techniques are summarized in Table 6.1. We hope that this basic idea on learning augmented techniques will provide an intriguing direction for future investigations.

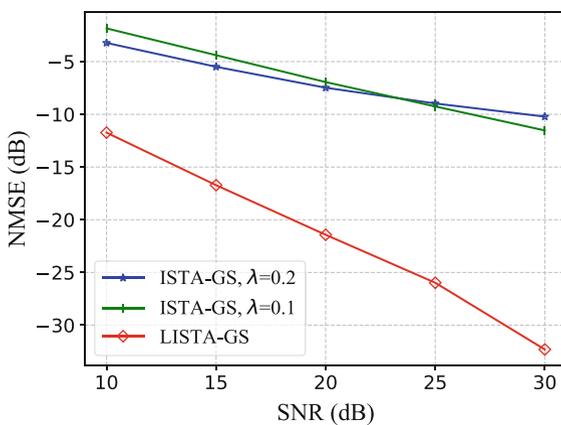
**Fig. 6.2** Performance comparison between the proposed LISTA-GS and baseline ISTA-GS in terms of NMSE



(a)



(b)



(c)

**Table 6.1** Summary of learning augmented techniques, methods, application, and corresponding references

Techniques	Methods	Application/Reference
Structured signal processing under a generative prior	Learning a generator $G(z) : z \in \mathbb{R}^k \rightarrow G(z) \in \mathbb{R}^n$	Compressed sensing [3], sparse linear model (2.3), Blind deconvolution [1]
Joint design of measurement matrix and sparse support recovery	Sample scheduling	Compressed sensing [10]
	Sensing matrix optimization	Compressed sensing [6, 7, 14, 15, 24]
	Learning augmented methods: learning an auto-encoder or generative adversarial networks	Linear compression for real signals [18, 22], nonlinear compression for real signals was considered in [17, 23], linear compression for complex signals [13]
Deep-learning-based AMP	Learned AMP	Compressed sensing [4], sparse linear model (2.3), (2.8)
	Learned Vector-AMP	
	Learned ISTA for group sparsity	

## References

1. Asim, M., Shamshad, F., Ahmed, A.: Blind image deconvolution using deep generative priors (2018). Preprint. arXiv: 1802.04073
2. Bonnefoy, A., Emiya, V., Ralaivola, L., Gribonval, R.: Dynamic screening: accelerating first-order algorithms for the lasso and group-lasso. *IEEE Trans. Signal Process.* **63**(19), 5121–5132 (2015)
3. Bora, A., Jalal, A., Price, E., Dimakis, A.G.: Compressed sensing using generative models. In: *Proceedings of the International Conference on Machine Learning (ICML)*, pp. 537–546 (2017). *JMLR. org*
4. Borgerding, M., Schniter, P., Rangan, S.: AMP-inspired deep networks for sparse linear inverse problems. *IEEE Trans. Signal Process.* **65**(16), 4293–4308 (2017)
5. Chen, X., Liu, J., Wang, Z., Yin, W.: Theoretical linear convergence of unfolded ISTA and its practical weights and thresholds. In: *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 9061–9071 (2018)
6. Duarte-Carvajalino, J.M., Sapiro, G.: Learning to sense sparse signals: simultaneous sensing matrix and sparsifying dictionary optimization. *IEEE Trans. Image Process.* **18**(7), 1395–1408 (2009)
7. Elad, M.: Optimized projections for compressed sensing. *IEEE Trans. Signal Process.* **55**(12), 5695–5702 (2007)
8. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. In: *Proceedings of the Neural Information Processing Systems (NeurIPS)*, pp. 2672–2680 (2014)

9. Gregor, K., LeCun, Y.: Learning fast approximations of sparse coding. In: Proceedings of the International Conference on Machine Learning (ICML), pp. 399–406. Omnipress, Madison (2010)
10. Hao, J., Zhang, B., Jiao, Z., Mao, S.: Adaptive compressive sensing based sample scheduling mechanism for wireless sensor networks. *Pervasive Mob. Comput.* **22**, 113–125 (2015)
11. Kingma, D.P., Ba, J.: ADAM: a method for stochastic optimization. In: Proceedings of the International Conference on Learning Representations (ICLR) (2015)
12. Kingma, D.P., Welling, M.: Auto-encoding variational Bayes (2013). Preprint. arXiv: 1312.6114
13. Li, S., Zhang, W., Cui, Y., Cheng, H.V., Yu, W.: Joint design of measurement matrix and sparse support recovery method via deep auto-encoder. Preprint. arXiv: 1910.04330 (2019)
14. Mangia, M., Rovatti, R., Setti, G.: Rakeness in the design of analog-to-information conversion of sparse and localized signals. *IEEE Trans. Circuits Syst. I, Reg. Papers* **59**(5), 1001–1014 (2012)
15. Mangia, M., Pareschi, F., Rovatti, R., Setti, G.: Adaptive matrix design for boosting compressed sensing. *IEEE Trans. Circuits Syst. I, Reg. Papers* **65**(3), 1016–1027 (2017)
16. Mousavi, A., Baraniuk, R.G.: Learning to invert: signal recovery via deep convolutional networks. In: Proceedings of the IEEE International Conference on Acoustics Speech Signal Processing (ICASSP), pp. 2272–2276. IEEE, Piscataway (2017)
17. Mousavi, A., Dasarathy, G., Baraniuk, R.G.: DeepCodec: adaptive sensing and recovery via deep convolutional neural networks. In: Annual Allerton Conference on Communication, Control, and Computing (Allerton), pp. 744–744. IEEE, Piscataway (2017)
18. Mousavi, A., Dasarathy, G., Baraniuk, R.G.: A data-driven and distributed approach to sparse signal representation and recovery. In: Proceedings of the International Conference on Learning Representations (ICLR) (2019)
19. Rangan, S., Schniter, P., Fletcher, A.K.: Vector approximate message passing. *IEEE Trans. Inf. Theory* **65**, 6664–6684 (2019)
20. Taha, A., Alrabeiah, M., Alkhateeb, A.: Enabling large intelligent surfaces with compressive sensing and deep learning (2019). Preprint. arXiv: 1904.10136
21. Tibshirani, R.: Regression shrinkage and selection via the lasso. *J. R. Stat. Soc.* **58**(1), 267–288 (1996)
22. Wu, S., Dimakis, A., Sanghavi, S., Yu, F., Holtmann-Rice, D., Storchus, D., Rostamizadeh, A., Kumar, S.: Learning a compressed sensing measurement matrix via gradient unrolling. In: Proceedings of the International Conference on Machine Learning (ICML), pp. 6828–6839 (2019)
23. Wu, Y., Rosca, M., Lillcrap, T.: Deep compressed sensing. In: Proceedings of the International Conference on Machine Learning (ICML), pp. 6850–6860 (2019)
24. Xu, J., Pi, Y., Cao, Z.: Optimized projection matrix for compressive sensing. *EURASIP J. Adv. Signal Process.* **2010**(1), 560349 (2010)
25. Yao, S., Zhao, Y., Zhang, A., Su, L., Abdelzaher, T.: DeepIoT: compressing deep neural network structures for sensing systems with a compressor-critic framework. In: Proceedings of the 15th ACM Conference on Embedded Network Sensor Systems, p. 4. ACM, New York (2017)

# Chapter 7

## Conclusions and Discussions



**Abstract** This chapter concludes the monograph. A summary is first provided for the main results of each chapter, and two reference tables are provided that contain the main analytical results and algorithms. Furthermore, we provide discussions on the future research directions of low-overhead communications and the corresponding structured signal processing approaches.

### 7.1 Summary

This monograph investigated different structured signal processing approaches to support low-overhead communications in IoT networks. Chapter 1 provided some background for low-overhead communications in IoT networks, introducing three main techniques that exclude particular parts of the metadata: grant-free random access, pilot-free communications, and identification-free communications. Furthermore, four general structured signal processing models, i.e., a sparse linear model, blind demixing, and sparse blind demixing, a shuffled linear regression, were introduced. Chapters 2–5 formed the core of this monograph, where four general structured signal processing models with corresponding applications in low-overhead communications were presented. For each chapter, the corresponding analysis results and algorithms were provided, which are summarized in Tables 7.1 and 7.2. More details on the proof of analysis results can be referred to Chap. 8. In Chap. 6, some cutting-edge learning augmented based techniques for structured signal processing were introduced, which represent an interesting direction for future research.

**Table 7.1** Analytical results in Chaps. 2–6

Result	Description
Theorem 2.1	Approximate kinematic formula captures a phase transition on whether the two randomly rotated cones share a ray
Proposition 2.1	Statistical dimension bound for the smoothed regularizer $\tilde{\mathcal{H}}_G$ (2.22)
Theorem 3.1	The least value of sample size required for exact recovery of problem (3.21)
Theorem 3.2	The convergence analysis of the regularized Wirtinger flow algorithm with spectral initialization for solving the blind demixing problem (3.22)
Theorem 3.3	The convergence analysis of the regularization-free Wirtinger flow algorithm with spectral initialization for solving the blind demixing problem (3.28)
Theorem 3.4	The convergence analysis of the Riemannian gradient with spectral initialization for solving the blind demixing problem (3.41)
Theorem 5.4	The recovery guarantees of Algorithm 5.3 when the shuffled data come from the Gaussian measurement model (5.8)
Theorem 5.5	Demonstrate that the system of polynomial equation (5.27) for exact data has a finite number of solutions, providing theoretical guarantees for developing algebraical method to solving the shuffled linear regression problem
Theorem 5.6	Demonstrate that the system of polynomial equation (5.28) for corrupted data has a finite number of solutions, providing theoretical guarantees for developing algebraical method to solving the shuffled linear regression

**Table 7.2** Algorithms in Chaps. 2–6

Algorithm	Description
Algorithm 2.1	Lan, Lu, and Monteiro's
Algorithm 3.1	Initialization via spectral method and projection
Algorithm 3.2	Riemannian optimization on product manifolds
Algorithm 3.3	Riemannian gradient descent with spectral initialization
Algorithm 4.1	DC algorithm for the sparse blind demixing problem (4.18)
Algorithm 5.1	Exact algorithm for calculating the maximum likelihood estimate of the permutation, i.e., (5.12)
Algorithm 5.2	“Row Sampling” algorithm that is exploited as the initialization of Algorithm 5.3
Algorithm 5.3	Approximation algorithm for computing the maximum likelihood estimator (5.9) for shuffled linear regression
Algorithm 5.4	Algebraically-initialized expectation-maximization for shuffled linear regression
Algorithm 6.1	Tied LAMP- $\ell_1$ parameter learning for solving lasso (6.1)
Algorithm 6.2	Untied LAMP- $\ell_1$ parameter learning for solving lasso (6.1)
Algorithm 6.3	Vector AMP for solving lasso (6.1)
Algorithm 6.4	Learned ISTA for the sparse linear model endowed with group row sparsity, i.e., (6.19)

## 7.2 Discussions

From Fig. 1.1 which illustrates an exemplary packet structure, we observe other opportunities to further reduce the overhead of the packet. For instance, for device activity detection or blind demixing, a smaller pilot length is preferred. According to recent studies [2, 3], generative models can yield a much more precise representation of the sparse signals. That is, much fewer measurements are required for the recovery of structural signal processing under a generative prior, compared with traditional analytical models. These methods are known as learning augmented methods, some of which were introduced in Chap. 6. They provide a promising direction for future study. Moreover, from the structured signal processing point of view, more sophisticated models are expected to be proposed via exploiting the sporadic activity pattern in massive connectivity networks, by exploiting spatial and temporal correlation of device activities. Furthermore, both convex methods and nonconvex methods with theoretical guarantees have been evoking researcher's interests. For instance, the rigorous statistical analysis for sparse blind demixing is called for further investigation. This is more challenging than the state-of-the-art model [1, 5] which assumes that  $\{\mathbf{x}_i\}$  in

$$y_j = \sum_{i=1}^s \mathbf{b}_j^H \mathbf{h}_i \mathbf{x}_i^H \mathbf{a}_{ij}, \quad 1 \leq j \leq L,$$

are sparse. From the algorithmic perspective, more efficient and robust algorithms are expected to be developed. For example, the learning argument methods have been exploited to solve the sparse linear model, e.g., Learned Vector-AMP, Learned AMP [4], and Learned ISTA for group sparsity, as introduced in Sect. 6.3.3.

Overall, given the promising results in applying structured signal processing for low-overhead communications reported in this monograph, we expect these methods will find abundant applications in practical IoT networks. We also hope the methods introduced in the monograph will lead to more effective algorithms, and inspire innovative approaches to exploit various structures in IoT systems, thus enable more application scenarios.

## References

1. Ahmed, A., Demanet, L.: Leveraging diversity and sparsity in blind deconvolution. *IEEE Trans. Inf. Theory* **64**(6), 3975–4000 (2018)
2. Asim, M., Shamshad, F., Ahmed, A.: Blind image deconvolution using deep generative priors. *arXiv preprint:1802.04073* (2018)
3. Bora, A., Jalal, A., Price, E., Dimakis, A.G.: Compressed sensing using generative models. In: *Proceedings of the 34th International Conference on Machine Learning (ICML)*, vol. 70, pp. 537–546 (2017). [JMLR.org](http://jmlr.org)

4. Borgerding, M., Schniter, P., Rangan, S.: AMP-inspired deep networks for sparse linear inverse problems. *IEEE Trans. Signal Process.* **65**(16), 4293–4308 (2017)
5. Flinth, A.: Sparse blind deconvolution and demixing through  $\ell_{1,2}$ -minimization. *Adv. Comput. Math.* **44**(1), 1–21 (2018)



### 8.1 Conic Integral Geometry

In this section, several basic concepts of conic integral geometry theory are introduced. We begin with the kinematic formula for cones which is the probability that a randomly rotated convex cone shares a ray with a fixed convex cone. This formula plays a vital role in characterizing the success or failure probability of an estimation problem. The following introduction is based on [3].

#### 8.1.1 The Kinematic Formula for Convex Cones

In the area of conic integral geometry, it is critical to identify the probability that a randomly rotated convex cone shares a ray with a fixed convex cone. Considering convex cones  $C$  and  $S$  in  $\mathbb{R}^d$ , and a random orthogonal basis  $A \in \mathbb{R}^{d \times d}$ , we aim to find an effective expression for the probability

$$P\{C \cap AS \neq \{\mathbf{0}\}\}. \tag{8.1}$$

Studying this probability enables to understand the phase transition phenomena in convex optimization problems with random data.

We start with the simple case of two dimensions where the solution to the problem can be quickly computed. Consider two convex cones  $C$  and  $S$  in  $\mathbb{R}^2$ , and assume that neither cone is a linear subspace. Then

$$P\{C \cap AS \neq \{\mathbf{0}\}\} = \min\{v_2(C) + v_2(S), 1\}, \tag{8.2}$$

where  $v_2(\cdot)$  returns the portion of the unit circle united by a convex cone in  $\mathbb{R}^2$ . If one of the cones is a subspace, a similar formula can be derived. In spaces with

higher dimensions, the representation of convex cones becomes more complicated. In three dimensions, it might be troublesome to find a reasonable solution in general. To address this issue, an extraordinary tool called the *conic kinematic formula* [19, Thm. 6.5.6] has been developed. It shows that there exists an *exact* formula to identify the probability that a randomly rotated convex cone shares a ray with a fixed convex cone. Moreover, only  $d + 1$  numbers are needed to summarize each cone in  $d$  dimensions.

**Fact 8.1 (The Kinematic Formula for Cones)** *Let  $C$  and  $S$  be closed convex cones in  $\mathbb{R}^n$ , one of which is not a subspace. Assuming a random orthogonal basis  $A \in \mathbb{R}^{n \times n}$ , then*

$$P\{C \cap AS \neq \{\mathbf{0}\}\} = \sum_{i=0}^n (1 + (-1)^{i+1}) \sum_{j=i}^n v_i(C) \cdot v_{n+i-j}(S). \quad (8.3)$$

For each  $k = 0, 1, 2, \dots, n$ , the operation  $v_k$  maps a closed convex cone to a nonnegative number, called the  $k$ -th intrinsic volume of the cone.

Even though the conic kinematic formula is beneficial for studying random instances of convex optimization problems [2, 15], this approach suffers a strenuous computation of expressions for the intrinsic volumes of a cone, except in the simplest cases. To address this issue, the paper [3] provided a novel method that makes the kinematic formula effective, which is elaborated in the following.

### 8.1.2 Intrinsic Volumes and the Statistical Dimension

The conic intrinsic volumes, illustrated in Fact 8.1, are the elemental geometric invariants of a closed convex cone. That is, the conic intrinsic volumes do not depend on the orientation of the cone within the space in which the cone is embedded, nor on the dimension of that space. This quantity is similar to some quantity defined for compact convex sets in Euclidean geometry, such as the usual volume, the surface area, the mean width, and the Euler characteristic [18].

The intrinsic volume of a closed convex cone  $C$  in  $\mathbb{R}^n$  consists of a sequel of probability distributions on  $\{0, 1, 2, \dots, n\}$ , represented as

$$\sum_{i=0}^n v_i(C) = 1 \quad \text{and} \quad v_i(C) \geq 0 \quad \text{for } i = 0, 1, 2, \dots, n. \quad (8.4)$$

The work [3] established an extraordinary fact about conic geometry: for each closed convex cone, the distribution of conic intrinsic volumes sharply concentrates around its mean value. The precise statement on the concentration of intrinsic is illustrated in Theorem 8.1. To begin with, we introduce several definitions that contribute to Theorem 8.1.

**Definition 8.1 (Tail Functionals)** Let  $C$  be a closed convex cone in  $\mathbb{R}^n$ . For every  $k = 0, 1, 2, \dots, n$ , the  $k$ -th *tail functional* is defined as

$$t_k(C) := v_k(C) + v_{k+1}(C) + \dots = \sum_{j=k}^n v_j(C). \quad (8.5)$$

The  $k$ -th *half-tail functional* is given by

$$h_k(C) := v_k(C) + v_{k+2}(C) + \dots = \sum_{\substack{j=k \\ j-k \text{ even}}}^n v_j(C). \quad (8.6)$$

**Definition 8.2 (Statistical Dimension)** Let  $C$  be a closed convex cone in  $\mathbb{R}^n$ . The *statistical dimension*  $\delta(C)$  of the cone is given by

$$\delta(C) := \sum_{k=0}^n k v_k(C). \quad (8.7)$$

As Definition 8.1 shows, the statistical dimension indicates the dimensionality of a convex cone. In particular, the statistical dimension is a *canonical extension* of the dimension of a linear subspace to the class of convex cones.

Based on the aforementioned definition, we arrive at the theorem that demonstrates the concentration of intrinsic volumes.

**Theorem 8.1 (Concentration of Intrinsic Volumes [3])** *Assuming that  $C$  is a closed convex cone, the transition width is given by*

$$\rho(C) := \sqrt{\delta(C^\circ) \wedge \delta(C)}.$$

Define a function

$$p_C(\gamma) := 4 \exp\left(\frac{-\gamma^2/8}{\rho^2(C) + \gamma}\right) \quad \text{for } \gamma \geq 0. \quad (8.8)$$

Then

$$k_- \leq \delta(C) - \gamma + 1 \quad \implies \quad t_{k_-}(C) \geq 1 - p_C(\gamma); \quad (8.9)$$

$$k_+ \geq \delta(C) + \gamma \quad \implies \quad t_{k_+}(C) \leq p_C(\gamma), \quad (8.10)$$

where  $t_k$  (8.5) is the tail functional, and  $\wedge$  is the operator that returns the minimum of two numbers.

### 8.1.3 The Approximate Kinematic Formula

Based on the concentration of intrinsic volumes provided in Theorem 8.1 and the conic kinematic formula illustrated in Fact 8.1, we can arrive at the following approximate kinematic formula.

**Theorem 8.2 (Approximate Kinematic Formula [3])** *Define a fix parameter  $\alpha \in (0, 1)$ . Let  $C$  and  $S$  be convex cones in  $\mathbb{R}^n$ , and assume a random orthogonal basis  $A \in \mathbb{R}^{n \times n}$ . Then*

$$\begin{aligned} \delta(C) + \delta(S) \leq n - \sqrt{n8 \log(4/\alpha)} &\implies P\{C \cap AS \neq \{\mathbf{0}\}\} \leq \alpha; \\ \delta(C) + \delta(S) \geq n + \sqrt{n8 \log(4/\alpha)} &\implies P\{C \cap AS \neq \{\mathbf{0}\}\} \geq 1 - \alpha. \end{aligned}$$

Theorem 8.2 demonstrates that two rotated cones are prone to share a ray in the case that the total statistical dimension of the two cones exceeds the ambient dimension. For problems in conic integral geometry, the cone is analogous to a subspace with approximate dimension  $\delta(C)$ . In the paper [3], a large class of random convex optimization problems have been proved to exhibit a phase transition, and the statistical dimension corresponding to each convex optimization problem characterizes the location of the phase transition.

### 8.1.4 Computing the Statistical Dimension

The statistical dimension plays a vital role in conic integral geometry, which can be used to identify that phase transitions occur in random convex optimization problems. To efficiently compute the statistical dimension, the method proposed in [3] is presented in the follows. We begin with several basic definitions. For a closed convex cone  $C$ , the projection  $\text{Proj}_C(\mathbf{x})$  that maps a point  $\mathbf{x}$  onto a point on the cone  $C$  which is nearest to  $\mathbf{x}$ :

$$\Pi_C(\mathbf{x}) := \operatorname{argmin} \{\|\mathbf{x} - \mathbf{y}\| : \mathbf{y} \in C\}. \quad (8.11)$$

For a general cone  $C \subset \mathbb{R}^n$ , the *polar cone*  $C^\circ$  is defined as the set of outward normals of  $C$ :

$$C^\circ := \{\mathbf{y} \in \mathbb{R}^n : \langle \mathbf{y}, \mathbf{x} \rangle \leq 0 \text{ for all } \mathbf{x} \in C\}. \quad (8.12)$$

**Proposition 8.1 (Statistical Dimension (Recall Definition 2.2))** *The statistical dimension  $\delta(C)$  of a closed convex cone  $C$  in  $\mathbb{R}^n$  satisfies*

$$\delta(C) = \mathbb{E}[\|\Pi_C(\mathbf{g})\|_2], \quad (8.13)$$

where  $\mathbf{g} \in \mathbb{R}^d$  is a standard Gaussian vector, and  $\Pi_C$  is defined in (8.11).

The metric characterization of the statistical dimension illustrated in Proposition 8.1 enables to connect the approach based on integral geometry and to the approach based on Gaussian process theory. The results can be obtained by a classic argument called the spherical Steiner formula [19, Thm. 6.5.1]. Furthermore, the formula (8.13) is related to another definition of parameter for convex cones called the *Gaussian width*, i.e., for a convex cone  $C \subset \mathbb{R}^n$ , the width is defined as

$$w(C) := \mathbb{E}[\sup_{\mathbf{y} \in C \cap \mathbb{S}^{n-1}} \langle \mathbf{y}, \mathbf{g} \rangle],$$

where  $\mathbf{g} \in \mathbb{R}^d$  is a standard Gaussian vector. This relation enables us to compute the statistical dimension by exploiting methods [4, 17] developed for the Gaussian width.

## 8.2 Proof of Proposition 2.1

Without loss of generality, define that

$$\boldsymbol{\theta}_0 = \left[ \left( \boldsymbol{\theta}_0^1 \right)^T, \dots, \left( \boldsymbol{\theta}_0^S \right)^T, \mathbf{0}_{M \times (N-S)} \right]^T \in \mathbb{C}^{N \times M},$$

where  $\boldsymbol{\theta}_0^i$  are nonzero. Hence, (2.21) is reformulated as

$$\delta \left( \mathcal{D} \left( \tilde{\mathcal{R}}_G; \tilde{\boldsymbol{\theta}}_0 \right) \right) \leq \inf_{\eta \geq 0} \mathbb{E} \left[ \text{dist}^2 \left( \mathbf{G}, \eta \cdot \partial \tilde{\mathcal{R}}_G \left( \tilde{\boldsymbol{\theta}}_0 \right) \right) \right], \quad (8.14)$$

where  $\mathbf{G} \in \mathbb{R}^{2N \times M}$  is a standard Gaussian matrix. Since  $\partial \tilde{\mathcal{R}}_G(\tilde{\boldsymbol{\theta}}_0) = \partial \mathcal{R}_G(\tilde{\boldsymbol{\theta}}_0) + \frac{\mu}{2} \partial \|\tilde{\boldsymbol{\theta}}_0\|_F^2$ , we have

$$\begin{aligned} & U \in \partial \mathcal{R}_G \left( \tilde{\boldsymbol{\theta}}_0 \right) \\ \Leftrightarrow & \begin{cases} U \gamma_j = \left( \tilde{\boldsymbol{\theta}}_0 \right)_{\gamma_j} / \left\| \left( \tilde{\boldsymbol{\theta}}_0 \right)_{\gamma_j} \right\|_F + \mu \left( \tilde{\boldsymbol{\theta}}_0 \right)_{\gamma_j} & \text{if } j = 1, \dots, S, \\ \|U \gamma_j\|_F \leq 1 & \text{if } j = S+1, \dots, N, \end{cases} \end{aligned} \quad (8.15)$$

where  $\tilde{\boldsymbol{\theta}}_{\gamma_j} = \mathbf{0}$  for  $j \neq i$  for some  $\tilde{\boldsymbol{\theta}} \in \mathbb{R}^{2N \times M}$ , defined in (2.23). Hence,

$$\begin{aligned} \text{dist}^2 \left( \mathbf{G}, \eta \cdot \partial \tilde{\mathcal{R}}_G \left( \tilde{\boldsymbol{\theta}}_0 \right) \right) &= \sum_{i=1}^S \left\| \mathbf{G} \gamma_i - \eta \left( \left( \tilde{\boldsymbol{\theta}}_0 \right)_{\gamma_i} / \left\| \left( \tilde{\boldsymbol{\theta}}_0 \right)_{\gamma_i} \right\|_F + \mu \left( \tilde{\boldsymbol{\theta}}_0 \right)_{\gamma_i} \right) \right\|_F^2 \\ &\quad + \sum_{i=S+1}^N \max \{ \|\mathbf{G} \gamma_i\|_2 - \eta, 0 \}^2. \end{aligned} \quad (8.16)$$

Taking the expectation over the Gaussian matrix  $\mathbf{G}$ , it arrives

$$\begin{aligned} \mathbb{E} \left[ \text{dist}^2 \left( \mathbf{G}, \eta \cdot \partial \mathcal{R}_G \left( \tilde{\Theta}_0 \right) \right) \right] &= S \left( 2M + \eta^2 \left( 1 + 2\mu\bar{a} + \mu^2\bar{b} \right) \right) \\ &\quad + (N - S) \frac{2^{1-M}}{\Gamma(M)} \int_{\eta}^{\infty} (u - \eta)^2 u^{2M-1} e^{-\frac{u^2}{2}} du, \end{aligned}$$

where  $\bar{a} = \frac{1}{S} \sum_{i=1}^S \|(\tilde{\Theta}_0)_{\mathcal{Y}_i}\|_F$  and  $\bar{b} = \frac{1}{S} \sum_{i=1}^S \|(\tilde{\Theta}_0)_{\mathcal{Y}_i}\|_F^2$ . Letting  $\rho = S/N$  and taking the infimum over  $\eta \geq 0$  completes the proof of (2.24).

### 8.3 Proof of Theorem 3.3

Theorem 3.3 can be justified via trajectory analysis for blind demixing via the Wirtinger flow algorithm. This is achieved by proving that iterates of Wirtinger flow sustain in the region of incoherence and contraction by exploiting the local geometry of blind demixing. The steps of proving Theorem 3.3 are summarized as follows.

- **Identifying local geometry in the region of incoherence and contraction (RIC).** First identify a region  $\mathcal{R}$ , i.e., RIC, where the objective function enjoys restricted strong convexity and smoothness near the ground truth  $\mathbf{z}^\natural$ . Furthermore, any point  $\mathbf{z} \in \mathcal{R}$  obeys the  $\ell_2$  error contraction and the incoherence conditions. Please refer to Lemma 8.1 for details. Hence, the convergence rate of the algorithm can be established according to Lemma 8.2, if and only if all the iterates of Wirtinger flow with spectral initialization are in the region  $\mathcal{R}$ .
- **Establishing the auxiliary sequences via the leave-one-out approach.** To justify that the Wirtinger Flow algorithm enforces the iterates to stay within the RIC, we introduce the leave-one-out sequences. Specifically, the leave-one-out sequences are denoted by  $\{\mathbf{h}_i^{t,(l)}, \mathbf{x}_i^{t,(l)}\}_{t \geq 0}$  for each  $1 \leq i \leq s$ ,  $1 \leq l \leq m$  obtained by removing the  $l$ -th measurement from the objective function  $f(\mathbf{h}, \mathbf{x})$ . Hence,  $\{\mathbf{h}_i^{t,(l)}\}$  and  $\{\mathbf{x}_i^{t,(l)}\}$  are independent with  $\{\mathbf{b}_j\}$  and  $\{\mathbf{a}_{ij}\}$ , respectively.
- **Establishing the incoherence condition via induction.** In this step, we employ the auxiliary sequences to establish the incoherence condition via induction. For brief, with  $\tilde{\mathbf{z}}_i^t = [\tilde{\mathbf{z}}_1^{t*}, \dots, \tilde{\mathbf{z}}_s^{t*}]^*$  where  $\tilde{\mathbf{z}}_i^t = [\tilde{\mathbf{h}}_i^{t*} \tilde{\mathbf{x}}_i^{t*}]^*$ , the set of induction hypotheses of local geometry is listed as follows:

$$\text{dist} \left( \mathbf{z}^t, \mathbf{z}^\natural \right) \leq C_1 \frac{1}{\log^2 m}, \quad (8.17a)$$

$$\text{dist} \left( \mathbf{z}^{t,(l)}, \tilde{\mathbf{z}}^t \right) \leq C_2 \frac{s\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 K \log^9 m}{m}}, \quad (8.17b)$$

$$\max_{1 \leq i \leq s, 1 \leq j \leq m} \left| \mathbf{a}_{ij}^* \left( \tilde{\mathbf{x}}_i^t - \mathbf{x}_i^{\natural} \right) \right| \leq C_3 \frac{1}{\sqrt{s} \log^{3/2} m}, \quad (8.17c)$$

$$\max_{1 \leq i \leq s, 1 \leq j \leq m} \left| \mathbf{b}_i^* \tilde{\mathbf{h}}_i^t \right| \leq C_4 \frac{\mu}{\sqrt{m}} \log^2 m, \quad (8.17d)$$

where  $C_1, C_3$  are some sufficiently small constants, while  $C_2, C_4$  are some sufficiently large constants. That is, as long as the current iterate stays within the RIC, the next iterate remains in the RIC.

- **Concentration between original and auxiliary sequences.** The gap between  $\{\mathbf{z}^t\}$  and  $\{\mathbf{z}^{t,(l)}\}$  can be established via employing the restricted strong convexity of the objective function in RIC.
- **Incoherence condition of auxiliary sequences.** Based on the fact that  $\{\mathbf{z}^t\}$  and  $\{\mathbf{z}^{t,(l)}\}$  are sufficiently close, we can instead bound the incoherence of  $\mathbf{h}_i^{t,(l)}$  (resp.  $\mathbf{x}_i^{t,(l)}$ ) in terms of  $\{\mathbf{b}_j\}$  (resp.  $\{\mathbf{a}_{ij}\}$ ), which turns out to be much easier due to the statistical independence between  $\{\mathbf{h}_i^{t,(l)}\}$  (resp.  $\{\mathbf{x}_i^{t,(l)}\}$ ) and  $\{\mathbf{b}_j\}$  (resp.  $\{\mathbf{a}_{ij}\}$ ).
- **Establishing iterates in RIC.** By combining the above bounds together, we arrive at  $|\mathbf{a}_{ij}^* (\mathbf{x}_i^t - \mathbf{x}_i^{\natural})| \leq \|\mathbf{a}_{ij}\|_2 \cdot \|\mathbf{x}_i^t - \mathbf{x}_i^{t,(l)}\|_2 + \|\mathbf{a}_{ij}^* (\mathbf{x}_i^{t,(l)} - \mathbf{x}_i^{\natural})\|$  via the triangle inequality. Based on the similar arguments, the other incoherence condition can be established in Lemma 8.3.
- **Establishing initial point in RIC.** Lemmas 8.6–8.8 are integrated to justify that the spectral initialization point is in RIC.

**Lemma 8.1 (Restricted Strong Convexity and Smoothness for Blind Demixing Problem  $\mathcal{P}$ )** *Let  $\delta > 0$  be a sufficiently small constant. If the number of measurements satisfies  $m \gg \mu^2 s^2 \kappa^2 K \log^5 m$ , then with probability at least  $1 - O(m^{-10})$ , the Wirtinger Hessian  $\nabla^2 f_{\text{clean}}(\mathbf{z})$  obeys*

$$\mathbf{u}^* \left[ \mathbf{D} \nabla^2 f_{\text{clean}}(\mathbf{z}) + \nabla^2 f_{\text{clean}}(\mathbf{z}) \mathbf{D} \right] \mathbf{u} \geq \frac{1}{4\kappa} \|\mathbf{u}\|_2^2 \quad \text{and} \\ \left\| \nabla^2 f_{\text{clean}}(\mathbf{z}) \right\| \leq 2 + s \quad (8.18)$$

simultaneously for all

$$\mathbf{u} = \begin{bmatrix} \mathbf{u}_1 \\ \vdots \\ \mathbf{u}_s \end{bmatrix} \quad \text{with } \mathbf{u}_i = \begin{bmatrix} \mathbf{h}_i - \mathbf{h}'_i \\ \mathbf{x}_i - \mathbf{x}'_i \\ \mathbf{h}_i - \mathbf{h}'_i \\ \mathbf{x}_i - \mathbf{x}'_i \end{bmatrix},$$

and  $\mathbf{D} = \text{diag}(\{\mathbf{W}_i\}_{i=1}^s)$

with  $\mathbf{W}_i = \text{diag}([\bar{\beta}_{i1} \mathbf{I}_K \bar{\beta}_{i2} \mathbf{I}_K \bar{\beta}_{i1} \mathbf{I}_K \bar{\beta}_{i2} \mathbf{I}_K]^*)$ .

Here  $\mathbf{z}$  satisfies

$$\max_{1 \leq i \leq s} \max \left\{ \left\| \mathbf{h}_i - \mathbf{h}_i^{\natural} \right\|_2, \left\| \mathbf{x}_i - \mathbf{x}_i^{\natural} \right\|_2 \right\} \leq \frac{\delta}{\kappa \sqrt{s}}, \quad (8.19a)$$

$$\max_{1 \leq i \leq s, 1 \leq j \leq m} \left| \mathbf{a}_{ij}^* \left( \mathbf{x}_i - \mathbf{x}_i^{\natural} \right) \right| \cdot \left\| \mathbf{x}_i^{\natural} \right\|_2^{-1} \leq \frac{2C_3}{\sqrt{s} \log^{3/2} m}, \quad (8.19b)$$

$$\max_{1 \leq i \leq s, 1 \leq j \leq m} \left| \mathbf{b}_j^* \mathbf{h}_i \right| \cdot \left\| \mathbf{h}_i^{\natural} \right\|_2^{-1} \leq \frac{2C_4 \mu}{\sqrt{m}} \log^2 m, \quad (8.19c)$$

where  $(\mathbf{h}_i, \mathbf{x}_i)$  is aligned with  $(\mathbf{h}'_i, \mathbf{x}'_i)$ , and one has  $\max\{\|\mathbf{h}_i - \mathbf{h}_i^{\natural}\|_2, \|\mathbf{h}'_i - \mathbf{h}_i^{\natural}\|_2, \|\mathbf{x}_i - \mathbf{x}_i^{\natural}\|_2, \|\mathbf{x}'_i - \mathbf{x}_i^{\natural}\|_2\} \leq \delta/(\kappa\sqrt{s})$ , for  $i = 1, \dots, s$  and  $\mathbf{W}_i$ 's satisfy that for  $\beta_{i1}, \beta_{i2} \in \mathbb{R}$ , for  $i = 1, \dots, s$   $\max_{1 \leq i \leq s} \max \left\{ \left| \beta_{i1} - \frac{1}{\kappa} \right|, \left| \beta_{i2} - \frac{1}{\kappa} \right| \right\} \leq \frac{\delta}{\kappa\sqrt{s}}$ . Therein,  $C_3, C_4 \geq 0$  are numerical constants.

Based on the local geometry in the region of incoherence and contraction, we further establish contraction of the error measured by the distance function.

**Lemma 8.2** *Suppose the number of measurements satisfies  $m \gg \mu^2 s^2 \kappa^2 K \log^5 m$  and the step size obeys  $\eta > 0$  and  $\eta \asymp s^{-1}$ . Then with probability at least  $1 - O(m^{-10})$ , we have*

$$\text{dist} \left( \mathbf{z}^{t+1}, \mathbf{z}^{\natural} \right) \leq (1 - \eta/(16\kappa)) \text{dist} \left( \mathbf{z}^t, \mathbf{z}^{\natural} \right) + 3\kappa\sqrt{s} \max_{1 \leq k \leq s} \|\mathcal{A}_k(\mathbf{e})\|,$$

provided that

$$\text{dist} \left( \mathbf{z}^t, \mathbf{z}^{\natural} \right) \leq \xi, \quad (8.20a)$$

$$\max_{1 \leq i \leq s, 1 \leq j \leq m} \left| \mathbf{a}_{ij}^* \left( \tilde{\mathbf{x}}_i^t - \mathbf{x}_i^{\natural} \right) \right| \cdot \left\| \mathbf{x}_i^{\natural} \right\|_2^{-1} \leq \frac{2C_3}{\sqrt{s} \log^{3/2} m}, \quad (8.20b)$$

$$\max_{1 \leq i \leq s, 1 \leq j \leq m} \left| \mathbf{b}_j^* \tilde{\mathbf{h}}_i^t \right| \cdot \left\| \mathbf{h}_i^{\natural} \right\|_2^{-1} \leq \frac{2C_4 \mu}{\sqrt{m}} \log^2 m, \quad (8.20c)$$

for some constants  $C_3, C_4 > 0$  and a sufficiently small constant  $\xi > 0$ . Here,  $\tilde{\mathbf{h}}_i^t$  and  $\tilde{\mathbf{x}}_i^t$  are defined as  $\tilde{\mathbf{h}}_i^t = \frac{1}{\alpha_i^t} \mathbf{h}_i^t$  and  $\tilde{\mathbf{x}}_i^t = \alpha_i^t \mathbf{x}_i^t$  for  $i = 1, \dots, s$ .

**Proof** From the definition of  $\alpha_k^{t+1}$ ,  $k = 1, \dots, s$ , one has

$$\begin{aligned} \text{dist}(\mathbf{z}^{t+1}, \mathbf{z}^{\natural})^2 &\leq \sum_{k=1}^s \text{dist}(\mathbf{z}_k^{t+1}, \mathbf{z}_k^{\natural})^2 \\ &\stackrel{\text{(i)}}{\leq} s \left\| \frac{1}{\alpha_k^{t+1}} \mathbf{h}_k^{t+1} - \mathbf{h}_k^{\natural} \right\|_2^2 + s \left\| \alpha_k^{t+1} \mathbf{x}_k^{t+1} - \mathbf{x}_k^{\natural} \right\|_2^2 \\ &\leq s \left\| \frac{1}{\alpha_k^t} \mathbf{h}_k^{t+1} - \mathbf{h}_k^{\natural} \right\|_2^2 + s \left\| \alpha_k^t \mathbf{x}_k^{t+1} - \mathbf{x}_k^{\natural} \right\|_2^2, \end{aligned} \quad (8.21)$$

where  $k$  in the step (i) satisfies that  $k = \arg \max_{1 \leq i \leq s} \text{dist}(\mathbf{z}_i^{t+1}, \mathbf{z}_i^{\natural})^2$ .

By denoting  $\tilde{\mathbf{h}}_k^t = \frac{1}{\alpha_k^t} \mathbf{h}_k^t$ ,  $\tilde{\mathbf{x}}_k^t = \alpha_k^t \mathbf{x}_k^t$ ,  $\hat{\mathbf{h}}_k^{t+1} = \frac{1}{\alpha_k^t} \mathbf{h}_k^{t+1}$  and  $\hat{\mathbf{x}}_k^{t+1} = \alpha_k^t \mathbf{x}_k^{t+1}$ , we have

$$\begin{bmatrix} \hat{\mathbf{h}}_k^{t+1} - \mathbf{h}_k^{\natural} \\ \hat{\mathbf{x}}_k^{t+1} - \mathbf{x}_k^{\natural} \\ \hat{\mathbf{h}}_k^{t+1} - \mathbf{h}_k^{\natural} \\ \hat{\mathbf{x}}_k^{t+1} - \mathbf{x}_k^{\natural} \end{bmatrix} = \begin{bmatrix} \tilde{\mathbf{h}}_k^t - \mathbf{h}_k^{\natural} \\ \tilde{\mathbf{x}}_k^t - \mathbf{x}_k^{\natural} \\ \tilde{\mathbf{h}}_k^t - \mathbf{h}_k^{\natural} \\ \tilde{\mathbf{x}}_k^t - \mathbf{x}_k^{\natural} \end{bmatrix} - \eta \mathbf{W}_k \begin{bmatrix} \nabla_{\mathbf{h}_k} f(\tilde{\mathbf{z}}^t) \\ \nabla_{\mathbf{x}_k} f(\tilde{\mathbf{z}}^t) \\ \nabla_{\mathbf{h}_k} f(\tilde{\mathbf{z}}^t) \\ \nabla_{\mathbf{x}_k} f(\tilde{\mathbf{z}}^t) \end{bmatrix}, \quad (8.22)$$

where

$$\mathbf{W}_k = \text{diag} \left( \left[ \|\tilde{\mathbf{x}}_k^t\|_2^{-2} \mathbf{I}_K \quad \|\tilde{\mathbf{h}}_k^t\|_2^{-2} \mathbf{I}_K \quad \|\tilde{\mathbf{x}}_k^t\|_2^{-2} \mathbf{I}_K \quad \|\tilde{\mathbf{h}}_k^t\|_2^{-2} \mathbf{I}_K \right] \right). \quad (8.23)$$

The Wirtinger Hessian without noise of  $f_{\text{clean}}(\mathbf{z})$  in terms of  $\mathbf{z}_i$  can be written as

$$\nabla_{\mathbf{z}_i}^2 f_{\text{clean}} := \begin{bmatrix} \mathbf{C} & \mathbf{E} \\ \mathbf{E}^* & \bar{\mathbf{C}} \end{bmatrix}, \quad (8.24)$$

where  $\mathbf{C} := \frac{\partial}{\partial \mathbf{z}_i} \left( \frac{\partial f_{\text{clean}}}{\partial \mathbf{z}_i} \right)^*$  and  $\mathbf{E} := \frac{\partial}{\partial \mathbf{z}_i} \left( \frac{\partial f_{\text{clean}}}{\partial \mathbf{z}_i} \right)^*$ . The Wirtinger Hessian of  $f_{\text{clean}}(\mathbf{z})$  in terms of  $\mathbf{z}$  is thus represented as

$$\nabla^2 f_{\text{clean}}(\mathbf{z}) := \text{diag} \left( \left\{ \nabla_{\mathbf{z}_i}^2 f_{\text{clean}} \right\}_{i=1}^s \right), \quad (8.25)$$

where the operation  $\text{diag}(\{\mathbf{A}_i\}_{i=1}^s)$  generates a block diagonal matrix with the diagonal elements being matrices  $\mathbf{A}_1, \dots, \mathbf{A}_s$ . According to the fundamental theorem of calculus provided in [13], together with the definition of the noiseless objective

function  $f_{\text{clean}}$  and the noiseless Wirtinger Hessian  $\nabla_{z_k}^2 f_{\text{clean}}$ , we get  $\nabla_{z_k}^2 f_{\text{clean}}$ ,

$$\begin{aligned} \begin{bmatrix} \nabla_{\mathbf{h}_k} f(\tilde{\mathbf{z}}^t) \\ \nabla_{\mathbf{x}_k} f(\tilde{\mathbf{z}}^t) \\ \nabla_{\mathbf{h}_k} f(\tilde{\mathbf{z}}^t) \\ \nabla_{\mathbf{x}_k} f(\tilde{\mathbf{z}}^t) \end{bmatrix} &= \begin{bmatrix} \nabla_{\mathbf{h}_k} f_{\text{clean}}(\tilde{\mathbf{z}}^t) \\ \nabla_{\mathbf{x}_k} f_{\text{clean}}(\tilde{\mathbf{z}}^t) \\ \nabla_{\mathbf{h}_k} f_{\text{clean}}(\tilde{\mathbf{z}}^t) \\ \nabla_{\mathbf{x}_k} f_{\text{clean}}(\tilde{\mathbf{z}}^t) \end{bmatrix} + \begin{bmatrix} \mathcal{A}_k(\mathbf{e})\mathbf{x}_k^t \\ \mathcal{A}_k^*(\mathbf{e})\mathbf{h}_k^t \\ \mathcal{A}_k(\mathbf{e})\mathbf{x}_k^t \\ \mathcal{A}_k^*(\mathbf{e})\mathbf{h}_k^t \end{bmatrix} \\ &= \mathbf{H}_k \begin{bmatrix} \tilde{\mathbf{h}}_k^t - \mathbf{h}_k^\natural \\ \tilde{\mathbf{x}}_k^t - \mathbf{x}_k^\natural \\ \tilde{\mathbf{h}}_k^t - \mathbf{h}_k^\natural \\ \tilde{\mathbf{x}}_k^t - \mathbf{x}_k^\natural \end{bmatrix} + \begin{bmatrix} \mathcal{A}_k(\mathbf{e})\mathbf{x}_k^t \\ \mathcal{A}_k^*(\mathbf{e})\mathbf{h}_k^t \\ \mathcal{A}_k(\mathbf{e})\mathbf{x}_k^t \\ \mathcal{A}_k^*(\mathbf{e})\mathbf{h}_k^t \end{bmatrix}, \end{aligned} \quad (8.26)$$

where  $\mathbf{H}_k = \int_0^1 \nabla_{z_k}^2 f_{\text{clean}}(z(\tau)) d\tau$  with  $z(\tau) := z^\natural + \tau(\tilde{\mathbf{z}}^t - z^\natural)$  and  $\mathcal{A}_k(\mathbf{e}) = \sum_{j=1}^m e_j \mathbf{b}_j \mathbf{a}_{kj}^*$  and  $\mathcal{A}_k^*(\mathbf{e}) = \sum_{j=1}^m \bar{e}_j \mathbf{a}_{kj} \mathbf{b}_j^*$ . Since  $z(\tau)$  lies between  $\tilde{\mathbf{z}}^t$  and  $z^\natural$ , we derive from the assumption (8.20) that for all  $\tau \in [0, 1]$ ,

$$\text{dist}(z(\tau), z^\natural) \leq \xi \leq \delta,$$

$$\begin{aligned} \max_{1 \leq i \leq s, 1 \leq j \leq m} \left| \mathbf{a}_{ij}^* \left( \mathbf{x}_i(\tau) - \mathbf{x}_i^\natural \right) \right| &\leq \frac{C_3}{\sqrt{s} \log^{3/2} m}, \\ \max_{1 \leq i \leq s, 1 \leq j \leq m} \left| \mathbf{b}_j^* \mathbf{h}_i(\tau) \right| &\leq \frac{C_4 \mu}{\sqrt{m}} \log^2 m, \end{aligned}$$

for some constants  $C_3, C_4 > 0$  and the constant  $\xi > 0$  being sufficiently small.

For simplicity, denote  $\tilde{\mathbf{z}}_k^{t+1} = [\tilde{\mathbf{h}}_k^{t+1*} \tilde{\mathbf{x}}_k^{t+1*}]^*$ . Substituting (8.26) to (8.22), one has

$$\begin{bmatrix} \tilde{\mathbf{z}}_k^{t+1} - \mathbf{z}_k^\natural \\ \tilde{\mathbf{z}}_k^{t+1} - \mathbf{z}_k^\natural \end{bmatrix} = \boldsymbol{\varphi}_k^t + \boldsymbol{\psi}_k^t, \quad (8.27)$$

where

$$\boldsymbol{\varphi}_k^t = (\mathbf{I} - \eta \mathbf{W}_k \mathbf{H}_k) \begin{bmatrix} \tilde{\mathbf{z}}_k^t - \mathbf{z}_k^\natural \\ \tilde{\mathbf{z}}_k^t - \mathbf{z}_k^\natural \end{bmatrix}, \quad \boldsymbol{\psi}_k^t = \begin{bmatrix} \mathcal{A}_k(\mathbf{e})\mathbf{x}_k^t \\ \mathcal{A}_k^*(\mathbf{e})\mathbf{h}_k^t \\ \mathcal{A}_k(\mathbf{e})\mathbf{x}_k^t \\ \mathcal{A}_k^*(\mathbf{e})\mathbf{h}_k^t \end{bmatrix}.$$

Take the Euclidean norm of both sides of (8.27) to arrive

$$\|\boldsymbol{\varphi}_k^t + \boldsymbol{\psi}_k^t\|_2 \leq \|\boldsymbol{\varphi}_k^t\|_2 + \|\boldsymbol{\psi}_k^t\|_2. \quad (8.28)$$

We first control the second Euclidean norm at the right-hand side of Eq. (8.28):

$$\|\boldsymbol{\psi}_k^t\|_2^2 = 2 \left( \|\mathcal{A}_k(\mathbf{e})\|^2 \|\mathbf{x}_k^t\|_2^2 + \|\mathcal{A}_k^*(\mathbf{e})\|^2 \|\mathbf{h}_k^t\|_2^2 \right) \leq 16 \|\mathcal{A}_k(\mathbf{e})\|^2, \quad (8.29)$$

where we use the fact that  $\max\{\|\mathbf{x}_k\|_2, \|\mathbf{h}_k\|_2\} \leq 2$  for  $1 \leq k \leq s$ . Based on the paper [13, Section C.2], the squared Euclidean norm of  $\boldsymbol{\varphi}_k^t$  is bounded by

$$\|\boldsymbol{\varphi}_k^t\|_2^2 \leq 2(1 - \eta/8) \left\| \tilde{\mathbf{z}}_k^t - \mathbf{z}_k^{\natural} \right\|_2^2, \quad (8.30)$$

under the assumption (8.20). We thus conclude that

$$\|\boldsymbol{\varphi}_k^t + \boldsymbol{\psi}_k^t\|_2 \leq \sqrt{2}(1 - \eta/8)^{1/2} \left\| \tilde{\mathbf{z}}_k^t - \mathbf{z}_k^{\natural} \right\|_2 + 4 \|\mathcal{A}_k(\mathbf{e})\|, \quad (8.31)$$

and hence

$$\begin{aligned} \left\| \tilde{\mathbf{z}}_k^{t+1} - \mathbf{z}_k^{\natural} \right\|_2 &\leq \left\| \tilde{\mathbf{z}}_k^{t+1} - \mathbf{z}_k^{\natural} \right\|_2 \leq \sqrt{2}/2 \|\boldsymbol{\varphi}_k^t + \boldsymbol{\psi}_k^t\|_2 \\ &\leq (1 - \eta/16) \left\| \tilde{\mathbf{z}}_k^t - \mathbf{z}_k^{\natural} \right\|_2 + 3 \|\mathcal{A}_k(\mathbf{e})\|. \end{aligned} \quad (8.32)$$

Integrate the above inequality (8.32) for  $i = 1, \dots, s$ , we further get

$$\text{dist} \left( \mathbf{z}^{t+1}, \mathbf{z}^{\natural} \right) \leq (1 - \eta/16) \text{dist} \left( \mathbf{z}^t, \mathbf{z}^{\natural} \right) + 3\sqrt{s} \max_{1 \leq k \leq s} \|\mathcal{A}_k(\mathbf{e})\|. \quad (8.33)$$

**Lemma 8.3** *Suppose the induction hypotheses hold true for  $t$ -th iteration and the number of measurements obeys  $m \gg (\mu^2 + \sigma^2)s^2 K \log^8 m$ . Then with probability at least  $1 - O(m^{-9})$ ,*

$$\max_{1 \leq i \leq s, 1 \leq j \leq m} \left| \mathbf{b}_l^* \tilde{\mathbf{h}}_i^{t+1} \right| \cdot \|\mathbf{h}_i^{\natural}\|_2^{-1} \leq C_4 \frac{\mu}{\sqrt{m}} \log^2 m, \quad (8.34)$$

provided that  $C_4$  is sufficiently large and the step size obeys  $\eta > 0$  and  $\eta \asymp s^{-1}$ .

**Proof** Similar to the strategy used in [13, Section C.4], it suffices to control  $|\mathbf{b}_l^* \frac{1}{\alpha_i^t} \mathbf{h}_i^{t+1}|$  to finish the proof, since

$$\begin{aligned} \max_{1 \leq i \leq s, 1 \leq l \leq m} \left| \mathbf{b}_l^* \frac{1}{\alpha_i^{t+1}} \mathbf{h}_i^{t+1} \right| &\leq \left| \frac{\alpha_i^t}{\alpha_i^{t+1}} \right| \max_{1 \leq i \leq s, 1 \leq l \leq m} \left| \mathbf{b}_l^* \frac{1}{\alpha_i^t} \mathbf{h}_i^{t+1} \right| \\ &\leq (1 + \delta) \left| \mathbf{b}_l^* \frac{1}{\alpha_i^t} \mathbf{h}_i^{t+1} \right| \end{aligned} \quad (8.35)$$

for some small  $\delta \asymp 1/\log^2 m$ , where the last step bases on

$$\left| \frac{\alpha_i^{t+1}}{\alpha_i^t} - 1 \right| \lesssim \frac{1}{\log^2 m} \leq \delta. \quad (8.36)$$

The gradient update rule for  $\mathbf{h}_i^{t+1}$  is written as

$$\frac{1}{\alpha_i^t} \mathbf{h}_i^{t+1} = \tilde{\mathbf{h}}_i^t - \eta \xi_i \sum_{j=1}^m \sum_{k=1}^s \mathbf{b}_j \mathbf{b}_j^* \left( \tilde{\mathbf{h}}_k^t \tilde{\mathbf{x}}_k^{t*} - \mathbf{h}_k^{\natural} \mathbf{h}_k^{\natural*} \right) \mathbf{a}_{kj} \mathbf{a}_{ij}^* \tilde{\mathbf{x}}_i^t + \eta \xi_i \sum_{j=1}^m e_j \mathbf{b}_j \mathbf{a}_{ij}^* \tilde{\mathbf{x}}_i^t, \quad (8.37)$$

where  $\xi_i = \frac{1}{\|\tilde{\mathbf{x}}_i^t\|_2}$  and  $\tilde{\mathbf{h}}_i^t = \frac{1}{\alpha_i^t} \mathbf{h}_i^t$  and  $\tilde{\mathbf{x}}_i^t = \alpha_i^t \mathbf{x}_i^t$  for  $i = 1, \dots, s$ . The formula (8.37) can be further decomposed into the following terms:

$$\begin{aligned} \frac{1}{\alpha_i^t} \mathbf{h}_i^{t+1} &= \tilde{\mathbf{h}}_i^t - \eta \xi_i \sum_{j=1}^m \sum_{k=1}^s \mathbf{b}_j \mathbf{b}_j^* \tilde{\mathbf{h}}_k^t \tilde{\mathbf{x}}_k^{t*} \mathbf{a}_{kj} \mathbf{a}_{ij}^* \tilde{\mathbf{x}}_i^t + \eta \xi_i \sum_{j=1}^m \sum_{k=1}^s \mathbf{b}_j \mathbf{b}_j^* \mathbf{h}_k^{\natural} \mathbf{x}_k^{\natural*} \mathbf{a}_{kj} \mathbf{a}_{ij}^* \tilde{\mathbf{x}}_i^t \\ &\quad + \eta \xi_i \sum_{j=1}^m e_j \mathbf{b}_j \mathbf{a}_{ij}^* \tilde{\mathbf{x}}_i^t \\ &= \tilde{\mathbf{h}}_i^t - \eta \xi_i \sum_{k=1}^s \tilde{\mathbf{h}}_k^t \left\| \mathbf{x}_k^{\natural} \right\|_2^2 - \eta \xi_i \mathbf{v}_{i1} - \eta \xi_i \mathbf{v}_{i2} + \eta \xi_i \mathbf{v}_{i3} + \eta \xi_i \mathbf{v}_{i4}, \end{aligned} \quad (8.38)$$

where

$$\begin{aligned} \mathbf{v}_{i1} &= \sum_{j=1}^m \sum_{k=1}^s \mathbf{b}_j \mathbf{b}_j^* \tilde{\mathbf{h}}_k^t \left( \tilde{\mathbf{x}}_k^{t*} \mathbf{a}_{kj} \mathbf{a}_{ij}^* \tilde{\mathbf{x}}_i^t - \mathbf{x}_k^{\natural*} \mathbf{a}_{kj} \mathbf{a}_{ij}^* \mathbf{x}_i^{\natural} \right) \\ \mathbf{v}_{i2} &= \sum_{j=1}^m \sum_{k=1}^s \mathbf{b}_j \mathbf{b}_j^* \tilde{\mathbf{h}}_k^t \left( \mathbf{x}_k^{\natural*} \mathbf{a}_{kj} \mathbf{a}_{ij}^* \mathbf{x}_i^{\natural} - \|\mathbf{x}_k^{\natural}\|_2^2 \right) \\ \mathbf{v}_{i3} &= \sum_{j=1}^m \sum_{k=1}^s \mathbf{b}_j \mathbf{b}_j^* \mathbf{h}_k^{\natural} \mathbf{x}_k^{\natural*} \mathbf{a}_{kj} \mathbf{a}_{ij}^* \tilde{\mathbf{x}}_i^t \\ \mathbf{v}_{i4} &= \sum_{j=1}^m e_j \mathbf{b}_j \mathbf{a}_{ij}^* \tilde{\mathbf{x}}_i^t, \end{aligned}$$

which bases on the fact that  $\sum_{j=1}^m \mathbf{b}_j \mathbf{b}_j^* = \mathbf{I}_K$ . In what follows, we bound the above four terms, respectively.

1. We start with  $|\mathbf{b}_l^* \mathbf{v}_{i1}|$  via the operation that

$$\begin{aligned} |\mathbf{b}_l^* \mathbf{v}_{i1}| &= \left| \sum_{j=1}^m \mathbf{b}_l^* \mathbf{b}_j \mathbf{b}_j^* \left[ \sum_{k=1}^s \tilde{\mathbf{h}}_k^t \left( \mathbf{a}_{ij}^* \left( \tilde{\mathbf{x}}_i^t - \mathbf{x}_i^{\natural} \right) \left( \mathbf{a}_{kj}^* \tilde{\mathbf{x}}_k^t \right)^* \right. \right. \right. \\ &\quad \left. \left. \left. + \mathbf{a}_{ij}^* \mathbf{x}_i^{\natural} \left( \mathbf{a}_{kj}^* \left( \tilde{\mathbf{x}}_k^t - \mathbf{x}_k^{\natural} \right) \right)^* \right) \right] \right| \leq s \sum_{j=1}^m |\mathbf{b}_l^* \mathbf{b}_j| \left\{ \max_{1 \leq k \leq s, 1 \leq j \leq m} |\mathbf{b}_j^* \tilde{\mathbf{h}}_k^t| \right\} \cdot \\ &\quad \left\{ \max_{1 \leq k \leq s, 1 \leq j \leq m} |\mathbf{a}_{kj}^* \left( \tilde{\mathbf{x}}_k^t - \mathbf{x}_k^{\natural} \right)| \left( |\mathbf{a}_{kj}^* \tilde{\mathbf{x}}_k^t| + |\mathbf{a}_{kj}^* \mathbf{x}_k^{\natural}| \right) \right\}. \end{aligned} \quad (8.39)$$

Based on the inductive hypothesis (8.17c) and the concentration inequality [13]

$$\max_{1 \leq i \leq s, 1 \leq j \leq m} |\mathbf{a}_{ij}^* \mathbf{x}_i^{\natural}| \leq 5\sqrt{\log m}, \quad (8.40)$$

with probability at least  $1 - O(m^{-10})$ , it yields

$$\begin{aligned} \max_{1 \leq k \leq s, 1 \leq j \leq m} |\mathbf{a}_{kj}^* \tilde{\mathbf{x}}_k^t| &\leq \max_{1 \leq k \leq s, 1 \leq j \leq m} |\mathbf{a}_{kj}^* \left( \tilde{\mathbf{x}}_k^t - \mathbf{x}_k^{\natural} \right)| \\ &\quad + \max_{1 \leq k \leq s, 1 \leq j \leq m} |\mathbf{a}_{kj}^* \mathbf{x}_k^{\natural}| \leq 6\sqrt{\log m}, \end{aligned} \quad (8.41)$$

as long as  $m$  is sufficiently large. We further derive that

$$\begin{aligned} \max_{1 \leq k \leq s, 1 \leq j \leq m} |\mathbf{a}_{kj}^* \left( \tilde{\mathbf{x}}_k^t - \mathbf{x}_k^{\natural} \right)| \left( |\mathbf{a}_{kj}^* \tilde{\mathbf{x}}_k^t| + |\mathbf{a}_{kj}^* \mathbf{x}_k^{\natural}| \right) \\ \leq \frac{1}{\sqrt{s} \log^{3/2} m} \cdot 11\sqrt{\log m} \leq 11C_3 \frac{1}{\log m}. \end{aligned} \quad (8.42)$$

Substituting (8.42) into (8.39) and combining lemma [13, Lemma 48] such that

$$\sum_{j=1}^m |\mathbf{b}_l^* \mathbf{b}_j| \leq 4 \log m, \quad (8.43)$$

we get

$$\begin{aligned} |\mathbf{b}_l^* \mathbf{v}_{i1}| &\lesssim s \log m \cdot \left\{ \max_{1 \leq k \leq s, 1 \leq j \leq m} |\mathbf{b}_j^* \tilde{\mathbf{h}}_k^t| \right\} \cdot C_3 \frac{1}{\log m} \\ &\lesssim s C_3 \max_{1 \leq k \leq s, 1 \leq j \leq m} |\mathbf{b}_j^* \tilde{\mathbf{h}}_k^t| \\ &\leq 0.1s \max_{1 \leq k \leq s, 1 \leq j \leq m} |\mathbf{b}_j^* \tilde{\mathbf{h}}_k^t|, \end{aligned} \quad (8.44)$$

as long as  $C_3$  is sufficiently small.

2. Regarding to  $|\mathbf{b}_l^* \mathbf{v}_{i3}|$ , one has

$$\begin{aligned} |\mathbf{b}_l^* \mathbf{v}_{i3}| &\leq \left| \sum_{j=1}^m \mathbf{b}_l^* \mathbf{b}_j \mathbf{b}_j^* \left( \sum_{k=1}^s \mathbf{h}_k^\natural \mathbf{x}_k^{\natural*} \mathbf{a}_{kj} \right) \mathbf{a}_{ij}^* \mathbf{x}_i^\natural \right| \\ &\quad + \left| \sum_{j=1}^m \mathbf{b}_l^* \mathbf{b}_j \mathbf{b}_j^* \left( \sum_{k=1}^s \mathbf{h}_k^\natural \mathbf{x}_k^{\natural*} \mathbf{a}_{kj} \right) \mathbf{a}_{ij}^* (\tilde{\mathbf{x}}_i^t - \mathbf{x}_i^\natural) \right|. \end{aligned} \quad (8.45)$$

**Lemma 8.4** Suppose  $m \gg s^2 K \log m$  for some sufficiently large constant  $C > 0$ . Then with probability at least  $1 - O(m^{-10})$ , there is

$$\left| \sum_{j=1}^m \mathbf{b}_l^* \mathbf{b}_j \mathbf{b}_j^* \left( \sum_{k=1}^s \mathbf{h}_k^\natural \mathbf{x}_k^{\natural*} \mathbf{a}_{kj} \right) \mathbf{a}_{ij}^* \mathbf{x}_i^\natural - \mathbf{b}_l^* \mathbf{h}_i^\natural \right| \lesssim \frac{\mu}{\sqrt{m}}. \quad (8.46)$$

*Proof* See Appendix 8.3.1.

Regarding to the second term in (8.45), we exploit the same technical method as in controlling  $|\mathbf{b}_l^* \mathbf{v}_{i1}|$ , which yields

$$\begin{aligned} &\left| \sum_{j=1}^m \mathbf{b}_l^* \mathbf{b}_j \mathbf{b}_j^* \left( \sum_{k=1}^s \mathbf{h}_k^\natural \mathbf{x}_k^{\natural*} \mathbf{a}_{kj} \right) \mathbf{a}_{ij}^* (\tilde{\mathbf{x}}_i^t - \mathbf{x}_i^\natural) \right| \\ &\leq s \sum_{j=1}^m |\mathbf{b}_l^* \mathbf{b}_j| \left\{ \max_{1 \leq k \leq s, 1 \leq j \leq m} |\mathbf{b}_j^* \mathbf{h}_k^\natural| \right\} \cdot \left\{ \max_{1 \leq k \leq s, 1 \leq j \leq m} |\mathbf{a}_{kj}^* (\tilde{\mathbf{x}}_k^t - \mathbf{x}_k^\natural)| \right\} \\ &\quad \left\{ \max_{1 \leq k \leq s, 1 \leq j \leq m} |\mathbf{a}_{kj}^* \mathbf{x}_k^\natural| \right\} \\ &\leq 4s \log m \cdot \frac{\mu}{\sqrt{m}} \cdot C_3 \frac{1}{\sqrt{s} \log^{3/2} m} \cdot 5\sqrt{\log m} \\ &\lesssim C_3 \frac{\sqrt{s}\mu}{\sqrt{m}}, \end{aligned} \quad (8.47)$$

where the second step arises from the incoherence, the induction hypothesis (8.17c) and the condition (8.40) and [13, Lemma 48]. Combining the above inequalities and the incoherence, one achieves

$$|\mathbf{b}_l^* \mathbf{v}_{i3}| \lesssim |\mathbf{b}_l^* \mathbf{h}_i^\natural| + \frac{\mu}{\sqrt{m}} + C_3 \frac{\sqrt{s}\mu}{\sqrt{m}} \lesssim (1 + C_3 \sqrt{s}) \frac{\mu}{\sqrt{m}}, \quad (8.48)$$

as long as picking up sufficiently small  $C_3 > 0$ .

3. We further move to control  $|\mathbf{b}_l^* \mathbf{v}_{i2}|$ . The idea of proof is based on the strategy in [13, Section C.4], which groups  $\{\mathbf{b}_j\}_{1 \leq j \leq m}$  into bins each containing  $\tau$  adjacent vectors. Similarly to the paper [13], we assume  $m/\tau$  to be an integer. For  $0 \leq l \leq m - \tau$ , one has

$$\begin{aligned} \mathbf{b}_1^* \sum_{j=1}^{\tau} \mathbf{b}_{l+j} \mathbf{b}_{l+j}^* \left( \sum_{k=1}^s \tilde{\mathbf{h}}_k^t z_{ijkl} \right) &= \mathbf{b}_1^* \sum_{j=1}^{\tau} \mathbf{b}_{l+1} \mathbf{b}_{l+1}^* \left( \sum_{k=1}^s \tilde{\mathbf{h}}_k^t z_{ijkl} \right) \\ &\quad + \mathbf{b}_1^* \sum_{j=1}^{\tau} \left( \mathbf{b}_{l+j} \mathbf{b}_{l+j}^* - \mathbf{b}_{l+1} \mathbf{b}_{l+1}^* \right) \left( \sum_{k=1}^s \tilde{\mathbf{h}}_k^t z_{ijkl} \right) \\ &= p_{i\tau 1} + p_{i\tau 2} + p_{i\tau 3}, \end{aligned} \quad (8.49)$$

where

$$\begin{aligned} z_{ijkl} &= \mathbf{x}_k^{\natural*} \mathbf{a}_{k,j+l} \mathbf{a}_{i,j+l}^* \mathbf{x}_i^{\natural} - \left\| \mathbf{x}_i^{\natural} \right\|_2^2, \\ p_{i\tau 1} &= \sum_{k=1}^s \left( \sum_{j=1}^{\tau} z_{ijkl} \right) \mathbf{b}_1^* \mathbf{b}_{l+1} \mathbf{b}_{l+1}^* \tilde{\mathbf{h}}_k^t, \\ p_{i\tau 2} &= \mathbf{b}_1^* \sum_{j=1}^{\tau} (\mathbf{b}_{l+j} - \mathbf{b}_{l+1}) \mathbf{b}_{l+j}^* \sum_{k=1}^s \tilde{\mathbf{h}}_k^t z_{ijkl}, \\ p_{i\tau 3} &= \mathbf{b}_1^* \sum_{j=1}^{\tau} \mathbf{b}_{l+1} (\mathbf{b}_{l+j} - \mathbf{b}_{l+1})^* \sum_{k=1}^s \tilde{\mathbf{h}}_k^t z_{ijkl}. \end{aligned}$$

We will control three terms in (8.49), respectively.

- (a) According to [13, Section C.4], with probability at least  $1 - O(m^{-10})$ ,

$$\left| \sum_{j=1}^{\tau} z_{ijkl} \right| \leq \left| \sum_{j=1}^{\tau} \left( \max \left\{ \left| \mathbf{a}_{k,l+j}^* \mathbf{x}_k^{\natural} \right|^2, \left| \mathbf{a}_{i,l+j}^* \mathbf{x}_i^{\natural} \right|^2 \right\} - \left\| \mathbf{x}_i^{\natural} \right\|_2^2 \right) \right| \lesssim \sqrt{\tau \log m}. \quad (8.50)$$

Combining above bound, we control the first term in (8.49) as

$$|p_{i\tau 1}| \lesssim s \sqrt{\tau \log m} |\mathbf{b}_1^* \mathbf{b}_{l+1}| \max_{1 \leq k \leq s, 1 \leq j \leq m} \left| \mathbf{b}_l^* \tilde{\mathbf{h}}_k^t \right|. \quad (8.51)$$

The summation over all bins is given as

$$\begin{aligned} & \left| \sum_{d=0}^{\frac{m}{\tau}-1} \sum_{j=1}^{\tau} z_{ijk(d\tau)} \mathbf{b}_1^* \mathbf{b}_{d\tau+1} \mathbf{b}_{d\tau+1}^* \tilde{\mathbf{h}}_k^t \right| \\ & \lesssim s \sqrt{\tau \log m} \sum_{d=0}^{\frac{m}{\tau}-1} \left| \mathbf{b}_1^* \mathbf{b}_{d\tau+1} \right| \max_{1 \leq k \leq s, 1 \leq j \leq m} \left| \mathbf{b}_l^* \tilde{\mathbf{h}}_k^t \right|. \end{aligned} \quad (8.52)$$

Substituting the bound

$$\sum_{d=0}^{\frac{m}{\tau}-1} \left| \mathbf{b}_1^* \mathbf{b}_{d\tau+1} \right| \leq \frac{K}{m} + O\left(\frac{\log m}{\tau}\right), \quad (8.53)$$

provided in [13, Section C.4], into the inequality (8.52) yields

$$\begin{aligned} & \left| \sum_{d=0}^{\frac{m}{\tau}-1} \sum_{j=1}^{\tau} z_{ijk(d\tau)} \mathbf{b}_1^* \mathbf{b}_{d\tau+1} \mathbf{b}_{d\tau+1}^* \tilde{\mathbf{h}}_k^t \right| \\ & \lesssim \left( \frac{sK\sqrt{\tau \log m}}{m} + \sqrt{\frac{s^2 \log^3 m}{\tau}} \right) \max_{1 \leq k \leq s, 1 \leq j \leq m} \left| \mathbf{b}_l^* \tilde{\mathbf{h}}_k^t \right| \\ & \leq 0.1 \max_{1 \leq k \leq s, 1 \leq j \leq m} \left| \mathbf{b}_l^* \tilde{\mathbf{h}}_k^t \right|, \end{aligned} \quad (8.54)$$

as long as  $m \gg Ks\sqrt{\tau \log m}$  and  $\tau \gg s^2 \log^3 m$ .

(b) The second term of (8.49),  $p_{i\tau 2}$ , is controlled by

$$\begin{aligned} |p_{i\tau 2}| & \leq \max_{1 \leq k \leq s, 1 \leq l \leq m} \left| \mathbf{b}_l^* \tilde{\mathbf{h}}_k^t \right| \sqrt{\sum_{j=1}^{\tau} \left| \mathbf{b}_1^* (\mathbf{b}_{l+j} - \mathbf{b}_{l+1}) \right|^2}. \\ & \sqrt{\sum_{j=1}^{\tau} \sum_{k=1}^s \left( \left| \max \left\{ \left| \mathbf{a}_{k,l+j}^* \mathbf{x}_k^{\natural} \right|^2, \left| \mathbf{a}_{i,l+j}^* \mathbf{x}_i^{\natural} \right|^2 \right\} - \left\| \mathbf{x}_k^{\natural} \right\|_2^2 \right)^2 \right)} \\ & \lesssim \sqrt{s\tau} \max_{1 \leq k \leq s, 1 \leq l \leq m} \left| \mathbf{b}_l^* \tilde{\mathbf{h}}_k^t \right| \sqrt{\sum_{i=1}^{\tau} \left| \mathbf{b}_1^* (\mathbf{b}_{l+j} - \mathbf{b}_{l+1}) \right|^2}, \end{aligned} \quad (8.55)$$

where the first inequality is due to Cauchy–Schwarz and the second step holds because of the following lemma.

**Lemma 8.5** Suppose  $\tau \geq C \log^4 m$  for some sufficiently large constant  $C > 0$ , with probability at least  $1 - O(m^{-10})$ , one has

$$\sum_{j=1}^{\tau} \sum_{k=1}^s \left( \max \left\{ \left| \mathbf{a}_{k,l+j}^* \mathbf{x}_k^{\flat} \right|^2, \left| \mathbf{a}_{i,l+j}^* \mathbf{x}_i^{\flat} \right|^2 \right\} - \left\| \mathbf{x}_k^{\flat} \right\|_2^2 \right)^2 \lesssim s\tau. \quad (8.56)$$

**Proof** This claim can be identified easily from [13, Appendix D.3.1].

We further sum over all bins of size  $\tau$  to obtain

$$\begin{aligned} & \left| \mathbf{b}_1^* \sum_{d=0}^{\frac{m}{\tau}-1} \sum_{j=1}^{\tau} (\mathbf{b}_{d\tau+j} - \mathbf{b}_{d\tau+1}) \mathbf{b}_{d\tau+j}^* \sum_{k=1}^s \tilde{\mathbf{h}}_k^t z_{ijk}(d\tau) \right| \\ & \leq \left\{ \sqrt{s\tau} \sum_{d=0}^{\frac{m}{\tau}-1} \sqrt{\sum_{i=1}^{\tau} |\mathbf{b}_1^* (\mathbf{b}_{d\tau+j} - \mathbf{b}_{d\tau+1})|^2} \right\} \cdot \max_{1 \leq k \leq s, 1 \leq l \leq m} \left| \mathbf{b}_l^* \tilde{\mathbf{h}}_k^t \right| \\ & \leq 0.1 \sqrt{s} \max_{1 \leq k \leq s, 1 \leq l \leq m} \left| \mathbf{b}_l^* \tilde{\mathbf{h}}_k^t \right|. \end{aligned} \quad (8.57)$$

Here, the last line arises from [13, Lemma 51] such that for any small constant  $c > 0$ ,

$$\sum_{d=0}^{\frac{m}{\tau}-1} \sqrt{\sum_{j=1}^{\tau} |\mathbf{b}_1^* (\mathbf{b}_{d\tau+j} - \mathbf{b}_{d\tau+1})|^2} \leq c \frac{1}{\sqrt{\tau}}, \quad (8.58)$$

as long as  $m \gg \tau K \log m$ .

(c) The third term of (8.49),  $p_{i\tau 3}$ , obeys that

$$\begin{aligned} |p_{i\tau 3}| & \leq |\mathbf{b}_1^* \mathbf{b}_{l+1}| \left\{ \sum_{j=1}^{\tau} \sum_{k=1}^s \left( \max \left\{ \left| \mathbf{a}_{k,l+j}^* \mathbf{x}_k^{\flat} \right|^2, \left| \mathbf{a}_{i,l+j}^* \mathbf{x}_i^{\flat} \right|^2 \right\} - \left\| \mathbf{x}_k^{\flat} \right\|_2^2 \right)^2 \right\} \\ & \quad \max_{1 \leq k \leq s, 1 \leq l \leq m-\tau, 1 \leq j \leq \tau} \left| (\mathbf{b}_{l+j} - \mathbf{b}_{l+1})^* \tilde{\mathbf{h}}_k^t \right| \\ & \lesssim s\tau |\mathbf{b}_1^* \mathbf{b}_{l+1}| \max_{1 \leq k \leq s, 1 \leq l \leq m-\tau, 1 \leq j \leq \tau} \left| (\mathbf{b}_{l+j} - \mathbf{b}_{l+1})^* \tilde{\mathbf{h}}_k^t \right|, \end{aligned} \quad (8.59)$$

where the last line relies on the inequality (8.56) and the Cauchy–Schwarz inequality.

The summation over all bins is given as

$$\begin{aligned}
& \sum_{d=0}^{\frac{m}{\tau}-1} \left| \mathbf{b}_1^* \sum_{j=1}^{\tau} \mathbf{b}_{d\tau+1} (\mathbf{b}_{d\tau+j} - \mathbf{b}_{d\tau+1})^* \sum_{k=1}^s \tilde{\mathbf{h}}_k^t z_{ijk(d\tau)} \right| \\
& \lesssim \tau s \sum_{d=0}^{\frac{m}{\tau}-1} |\mathbf{b}_1^* \mathbf{b}_{d\tau+1}| \cdot \max_{1 \leq k \leq s, 1 \leq l \leq m-\tau, 1 \leq j \leq \tau} |(\mathbf{b}_{l+j} - \mathbf{b}_{l+1})^* \tilde{\mathbf{h}}_k^t| \\
& \lesssim s \log m \max_{1 \leq k \leq s, 1 \leq l \leq m-\tau, 1 \leq j \leq \tau} |(\mathbf{b}_{l+j} - \mathbf{b}_{l+1})^* \tilde{\mathbf{h}}_k^t| \\
& \lesssim cC_4 \frac{s\mu}{\sqrt{m}} \log^2 m, \tag{8.60}
\end{aligned}$$

where the last relation makes use of (8.53) and the claim

$$\max_{1 \leq j \leq \tau} |(\mathbf{b}_j - \mathbf{b}_1)^* \tilde{\mathbf{h}}_k^t| \leq cC_4 \frac{\mu}{\sqrt{m}} \log m, \tag{8.61}$$

for some sufficiently small constant  $c > 0$ , provided that  $m \gg \tau K \log^4 m$ .

(d) Combining the above results together, we get

$$|\mathbf{b}_1^* \mathbf{v}_{i2}| \leq (0.1 + 0.1\sqrt{s}) \max_{1 \leq k \leq s, 1 \leq l \leq m} |\mathbf{b}_l^* \tilde{\mathbf{h}}_k^\top| + O\left(cC_4 \frac{s\mu}{\sqrt{m}} \log^2 m\right). \tag{8.62}$$

4. We end the proof with controlling  $|\mathbf{b}_l^* \mathbf{v}_{i4}|$ :

$$\begin{aligned}
|\mathbf{b}_l^* \mathbf{v}_{i4}| &= \left| \sum_{j=1}^m \mathbf{b}_l^* \mathbf{b}_j e_j \mathbf{a}_{ij}^* \tilde{\mathbf{x}}_i^t \right| \leq \sum_{j=1}^m |\mathbf{b}_l^* \mathbf{b}_j| \cdot \left\{ \max_{1 \leq k \leq s, 1 \leq j \leq m} |\mathbf{a}_{kj}^* \tilde{\mathbf{x}}_k^t| \right\} \cdot |e_j| \\
&\leq 4 \log m \cdot 6\sqrt{\log m} \cdot \frac{\sigma^2}{m}, \tag{8.63}
\end{aligned}$$

where the last step arises from the inequality (8.41), (8.43) and the assumption  $|e_j| \lesssim \sigma^2/m \ll 1$ . It thus yields

$$|\mathbf{b}_l^* \mathbf{v}_{i4}| \lesssim \sigma^2 \frac{\log^{3/2} m}{m} \leq \log m, \tag{8.64}$$

as long as  $m \gg \sigma^2 \sqrt{\log m}$ .

5. Putting the above results together, there exists some constant  $C_8 > 0$  such that

$$\begin{aligned}
\left| \mathbf{b}_l^* \tilde{\mathbf{h}}_i^{t+1} \right| &\leq (1 + \delta) \left\{ \left| \mathbf{b}_l^* \tilde{\mathbf{h}}_i^t \right| - \eta \xi_i \sum_{k=1}^s \left| \mathbf{b}_l^* \tilde{\mathbf{h}}_k^t \right| + (1 + 0.1\sqrt{s} + 0.1s) \right. \\
&\quad \left. \max_{1 \leq k \leq s, 1 \leq j \leq m} \left| \mathbf{b}_j^* \tilde{\mathbf{h}}_k^t \right| + C_8(1 + C_3\sqrt{s})\eta \xi_i \frac{\mu}{\sqrt{m}} \right. \\
&\quad \left. + C_8 c C_4 \eta \xi_i \frac{s\mu}{\sqrt{m}} \log^2 m + C_8 \eta \xi_i \log m \right\} \\
&\stackrel{(i)}{\leq} \left( 1 + \mathcal{O}\left(\frac{1}{\log^2 m}\right) \right) \left\{ (1 - 0.7s\eta \xi_i) C_4 \frac{\mu}{\sqrt{m}} \log^2 m \right. \\
&\quad \left. + C_8(1 + C_3\sqrt{s})\eta \xi_i \frac{\mu}{\sqrt{m}} + C_8 c C_4 \eta \xi_i \frac{s\mu}{\sqrt{m}} \log^2 m + C_8 \eta \xi_i \log m \right\} \\
&\stackrel{(ii)}{\leq} C_4 \frac{\mu}{\sqrt{m}} \log^2 m. \tag{8.65}
\end{aligned}$$

Here, step (i) arises from the induction hypothesis (8.17d), step (ii) holds as  $\log$  as  $c > 0$  is sufficiently small, i.e.,  $(1 + \delta)C_8\eta\xi_i c \gg 1$ , and  $\eta > 0$  is some sufficiently small constant, i.e.,  $\eta \asymp s^{-1}$ . In order for the proof to go through, we need to pick the sample size that

$$m \gg (\mu^2 + \sigma^2)\tau K \log^4 m, \tag{8.66}$$

where  $\tau = c_{10}s^2 \log^4 m$  with some sufficiently large constant  $c_{10} > 0$ .

### 8.3.1 Proof of Lemma 8.4

Denote

$$w_{ij} = \mathbf{b}_l^* \mathbf{b}_j \mathbf{b}_j^* \left( \sum_{k=1}^s \mathbf{h}_k^{\natural} \mathbf{x}_k^{\natural*} \mathbf{a}_{kj} \right) \mathbf{a}_{ij}^* \mathbf{x}_i^{\natural}. \tag{8.67}$$

Combining the fact that  $\mathbb{E}[\mathbf{a}_{ij} \mathbf{a}_{ij}^*] = \mathbf{I}_K$ ,  $\mathbb{E}[\mathbf{a}_{kj} \mathbf{a}_{ij}^*] = 0$  for  $k \neq i$  and  $\sum_{j=1}^m \mathbf{b}_j \mathbf{b}_j^* = \mathbf{I}_K$ , we can represent the objective quantity as the sum of independent random variables,

$$\sum_{j=1}^m \mathbf{b}_l^* \mathbf{b}_j \mathbf{b}_j^* \left( \sum_{k=1}^s \mathbf{h}_k^{\natural} \mathbf{x}_k^{\natural*} \mathbf{a}_{kj} \right) \mathbf{a}_{ij}^* \mathbf{x}_i^{\natural} - \mathbf{b}_l^* \mathbf{h}_i^{\natural} = \sum_{j=1}^m (w_{ij} - \mathbb{E}(w_{ij})). \tag{8.68}$$

Based on the definition of sub-exponential norm [21], i.e., denoted as  $\|\cdot\|_{\psi_1}$ , we get

$$\begin{aligned} \|w_{ij} - \mathbb{E}[w_{ij}]\|_{\psi_1} &\stackrel{(i)}{\leq} 2\|w_{ij}\|_{\psi_1} \\ &\stackrel{(ii)}{\leq} 4 \sum_{k=1}^s |\mathbf{b}_l^* \mathbf{b}_j| \left| \mathbf{b}_j^* \mathbf{h}_k^{\natural} \right| \max_{1 \leq q \leq s} \left\| \mathbf{a}_{qj}^* \mathbf{x}_q^{\natural} \right\|_{\psi_2}^2 \end{aligned} \quad (8.69)$$

$$\stackrel{(iii)}{\lesssim} |\mathbf{b}_l^* \mathbf{b}_j| \frac{s\mu}{\sqrt{m}}, \quad (8.70)$$

where (i) uses the centering property of the sub-exponential norm [21, Remark 5.18], (ii) arises from the relationship between the sub-exponential norm and the sub-Gaussian norm [21, Lemma 5.14], and (iii) occurs since the incoherence condition and the fact that  $\|\mathbf{a}_{kj}^* \mathbf{x}_k^{\natural}\|_{\psi_2} \lesssim 1$ . According to [13, Section C.4.1], one has

$$W := \sqrt{\frac{1}{m} \sum_{j=1}^m \|w_{ij} - \mathbb{E}[w_{ij}]\|_{\psi_1}^2} \asymp \frac{\mu}{\sqrt{m}} \frac{s\sqrt{K}}{m}. \quad (8.71)$$

It can further invoke [10, Corollary 1] to obtain

$$\mathbb{P} \left( \left| \frac{1}{m} \sum_{j=1}^m (w_{ij} - \mathbb{E}[w_{ij}]) \right| \geq t \right) \leq \exp \left( 1 - \frac{m}{8} \min \left\{ \frac{t}{2W}, \left( \frac{t}{2W} \right)^2 \right\} \right). \quad (8.72)$$

By taking  $t = 2\varepsilon W$  for  $\varepsilon \in (0, 1)$ , we obtain with probability at least  $1 - \exp(1 - m\varepsilon^2/8)$ ,

$$\sum_{j=1}^m (w_{ij} - \mathbb{E}[w_{ij}]) \leq 2\varepsilon W m \lesssim \varepsilon s \sqrt{K} \frac{\mu}{\sqrt{m}}. \quad (8.73)$$

Thus, choosing  $\varepsilon \asymp 1/s\sqrt{K}$ , we conclude that with probability at least  $1 - \exp(1 - cm/s^2K)$  for some constant  $c > 0$ ,

$$\left| \sum_{j=1}^m (w_{ij} - \mathbb{E}[w_{ij}]) \right| \lesssim \frac{\mu}{\sqrt{m}}. \quad (8.74)$$

We finished the proof by observing that  $m \gg s^2K \log m$  as claimed in the assumption.

**Lemma 8.6** *With probability at least  $1 - O(m^{-9})$ , there exists some constant  $C > 0$  such that*

$$\min_{\alpha_i \in \mathbb{C}, |\alpha_i|=1} \left\{ \left\| \alpha_i \mathbf{h}_i^0 - \mathbf{h}_i^{\natural} \right\| + \left\| \alpha_i \mathbf{x}_i^0 - \mathbf{x}_i^{\natural} \right\| \right\} \leq \frac{\xi}{\kappa \sqrt{s}} \text{ and} \quad (8.75)$$

$$\min_{\alpha_i \in \mathbb{C}, |\alpha_i|=1} \left\{ \left\| \alpha_i \mathbf{h}_i^{0,(l)} - \mathbf{h}_i^{\natural} \right\| + \left\| \alpha_i \mathbf{x}_i^{0,(l)} - \mathbf{x}_i^{\natural} \right\| \right\} \leq \frac{\xi}{\kappa \sqrt{s}}, \quad (8.76)$$

and  $|\alpha_i^0| - 1 < 1/4$ , for each  $1 \leq i \leq s$ ,  $1 \leq l \leq m$ , provided that

$$m \geq C(\mu^2 + \sigma^2)s\kappa^2 K \log m / \xi^2.$$

**Lemma 8.7** *Suppose that  $m \gg (\mu^2 + \sigma^2)s^2\kappa^2 K \log^3 m$ . Then with probability at least  $1 - O(m^{-9})$ ,*

$$\text{dist}(\mathbf{z}^{0,(l)}, \tilde{\mathbf{z}}^0) \leq C_2 \frac{s\kappa\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 s K \log^5 m}{m}} \text{ and} \quad (8.77)$$

$$\max_{1 \leq i \leq m} \left| \mathbf{b}_i^* \tilde{\mathbf{h}}_i^0 \right| \cdot \left\| \mathbf{h}_i^{\natural} \right\|_2^{-1} \leq C_4 \frac{\mu \log^2 m}{\sqrt{m}}. \quad (8.78)$$

**Proof** With a similar strategy as in [13, Section C.6], we first show that the normalized singular vectors of  $\mathbf{M}_i$  and  $\mathbf{M}_i^{(l)}$ ,  $i = 1, \dots, s$  are close enough. We further extend this inequality to the scaled singular vectors, thereby converting the  $\ell_2$  metric to the distance function defined in (3.30). We finally prove the incoherence in terms of  $\{\mathbf{b}_j\}_{j=1}^m$ .

Recall that  $\check{\mathbf{h}}_i^0$  and  $\check{\mathbf{x}}_i^0$  are the leading left and right singular vectors of  $\mathbf{M}_i$ ,  $i = 1, \dots, s$ , and  $\check{\mathbf{h}}_i^{0,(l)}$  and  $\check{\mathbf{x}}_i^{0,(l)}$  are the leading left and right singular vectors of  $\mathbf{M}_i^{(l)}$ ,  $i = 1, \dots, s$ . By exploiting a variant of Wedin's  $\sin\Theta$  theorem [8, Theorem 2.1], we derive that

$$\leq \frac{c_1 \left\| \left( \mathbf{M}_i - \mathbf{M}_i^{(l)} \right) \check{\mathbf{x}}_i^{0,(l)} \right\|_2 + c_1 \left\| \check{\mathbf{h}}_i^{0,(l)*} \left( \mathbf{M}_i - \mathbf{M}_i^{(l)} \right) \right\|_2}{\sigma_1 \left( \mathbf{M}_i^{(l)} \right) - \sigma_2 \left( \mathbf{M}_i \right)}, \quad (8.79)$$

for  $i = 1, \dots, s$  with some constant  $c_1 > 0$ . According to [13, Section C.6], for  $i = 1, \dots, s$ , we have

$$\sigma_1 \left( \mathbf{M}_i^{(l)} \right) - \sigma_2 \left( \mathbf{M}_i \right) \geq 3/4 - \left\| \mathbf{M}_i^{(l)} - \mathbb{E}[\mathbf{M}_i^{(l)}] \right\| - \left\| \mathbf{M}_i - \mathbb{E}[\mathbf{M}_i] \right\| \geq 1/2, \quad (8.80)$$

where the last step comes from [12, Lemma 6.16] provided that  $m \gg (\mu^2 + \sigma^2)sK \log m$ . As a result, we obtain that for  $i = 1, \dots, s$

$$\begin{aligned} & \left\| \beta_i^{0,(l)} \check{\mathbf{h}}_i^0 - \check{\mathbf{h}}_i^{0,(l)} \right\|_2 + \left\| \beta_i^{0,(l)} \check{\mathbf{x}}_i^0 - \check{\mathbf{x}}_i^{0,(l)} \right\|_2 \\ & \leq 2c_1 \left\{ \left\| (\mathbf{M}_i - \mathbf{M}_i^{(l)}) \check{\mathbf{x}}_i^{0,(l)} \right\|_2 + \left\| \check{\mathbf{h}}_i^{0,(l)*} (\mathbf{M}_i - \mathbf{M}_i^{(l)}) \right\|_2 \right\}, \end{aligned} \quad (8.81)$$

where

$$\beta_i^{0,(l)} := \operatorname{argmin}_{\alpha \in \mathbb{C}, |\alpha|=1} \left\| \alpha \check{\mathbf{h}}_i^0 - \check{\mathbf{h}}_i^{0,(l)} \right\|_2 + \left\| \alpha \check{\mathbf{x}}_i^0 - \check{\mathbf{x}}_i^{0,(l)} \right\|_2. \quad (8.82)$$

It thus suffices to control the two terms on the right-hand side of (8.81). Therein,

$$\mathbf{M}_i - \mathbf{M}_i^{(l)} = \mathbf{b}_l \mathbf{b}_l^* \sum_{k=1}^s \mathbf{h}_k^\natural \mathbf{x}_k^{\natural*} \mathbf{a}_{kl} \mathbf{a}_{il}^* + e_l \mathbf{b}_l \mathbf{a}_{il}^*. \quad (8.83)$$

1. To bound the first term, we observe that

$$\begin{aligned} & \left\| (\mathbf{M}_i - \mathbf{M}_i^{(l)}) \check{\mathbf{x}}_i^{0,(l)} \right\|_2 \\ & = \left\| \mathbf{b}_l \mathbf{b}_l^* \sum_{k=1}^s \mathbf{h}_k^\natural \mathbf{x}_k^{\natural*} \mathbf{a}_{kl} \mathbf{a}_{il}^* \check{\mathbf{x}}_i^{0,(l)} + e_l \mathbf{b}_l \mathbf{a}_{il}^* \check{\mathbf{x}}_i^{0,(l)} \right\|_2 \\ & \leq \max_{1 \leq k \leq s} s \|\mathbf{b}_l\|_2 \cdot \left| \mathbf{b}_l^* \mathbf{h}_k^\natural \right| \cdot \left| \mathbf{a}_{kl}^* \mathbf{x}_i^\natural \right| \cdot \left| \mathbf{a}_{il}^* \check{\mathbf{x}}_i^{0,(l)} \right| + \|\mathbf{b}_l\|_2 \cdot |e_l| \cdot \left| \mathbf{a}_{il}^* \check{\mathbf{x}}_i^{0,(l)} \right| \\ & \leq 30 \frac{s\mu}{\sqrt{m}} \cdot \sqrt{\frac{K \log^2 m}{m}} + \frac{5\sigma^2}{m} \sqrt{\frac{K \log m}{m}}, \end{aligned} \quad (8.84)$$

where we use the fact that  $\|\mathbf{b}_l\|_2 = \sqrt{K/m}$ , the incoherence condition, the bound (8.40), the assumption  $|e_j| \leq \frac{\sigma^2}{m}$ , and the condition that with probability exceeding  $1 - O(m^{-10})$ ,

$$\max_{1 \leq l \leq m} \left| \mathbf{a}_{il}^* \check{\mathbf{x}}_i^{0,(l)} \right| \leq 5\sqrt{\log m}, \quad (8.85)$$

due to the independence between  $\check{\mathbf{x}}_i^{0,(l)}$  and  $\mathbf{a}_{il}$  [13, Section C.6].

2. To control the second term, we observe that

$$\begin{aligned}
& \left\| \check{\mathbf{h}}_i^{0,(l)*} \left( \mathbf{M}_i - \mathbf{M}_i^{(l)} \right) \right\|_2 \\
&= \left\| \check{\mathbf{h}}_i^{0,(l)*} \mathbf{b}_l \mathbf{b}_l^* \sum_{k=1}^s \mathbf{h}_k^\natural \mathbf{x}_k^{\natural*} \mathbf{a}_{kl} \mathbf{a}_{il}^* + e_l \check{\mathbf{h}}_i^{0,(l)*} \mathbf{b}_l \mathbf{a}_{il}^* \right\|_2 \\
&\leq s \max_{1 \leq k \leq s} \left\| \mathbf{a}_{il}^* \right\|_2 \cdot \left| \mathbf{b}_l^* \mathbf{h}_k^\natural \right| \cdot \left| \mathbf{a}_{kl}^* \mathbf{x}_k^\natural \right| \cdot \left| \mathbf{b}_l^* \check{\mathbf{h}}_i^{0,(l)} \right| + \left\| \mathbf{a}_{il}^* \right\|_2 \cdot |e_l| \cdot \left| \mathbf{b}_l^* \check{\mathbf{h}}_i^{0,(l)} \right| \\
&\leq \left( 15 \sqrt{\frac{\mu^2 s^2 K \log m}{m}} + 3\sqrt{K} \frac{\sigma^2}{m} \right) \cdot \left( \left| \mathbf{b}_l^* \check{\mathbf{h}}_i^{0,(l)} \right| + \sqrt{\frac{K}{m}} \left\| \tilde{\alpha}_i \check{\mathbf{h}}_i^0 - \check{\mathbf{h}}_i^{0,(l)} \right\|_2 \right), \tag{8.86}
\end{aligned}$$

where  $|\tilde{\alpha}_i| = 1$ . Here, the last step easily arises from the similar strategy used in (8.84) and [13, Section C.6]. Substitution of the bounds (8.84) and (8.86) into (8.81) yields

$$\begin{aligned}
& \left\| \beta_i^{0,(l)} \check{\mathbf{h}}_i^0 - \check{\mathbf{h}}_i^{0,(l)} \right\|_2 + \left\| \beta_i^{0,(l)} \check{\mathbf{x}}_i^0 - \check{\mathbf{x}}_i^{0,(l)} \right\|_2 \\
&\leq 2c_1 \left\{ 30 \frac{\mu}{\sqrt{m}} \cdot \sqrt{\frac{s^2 K \log^2 m}{m}} + \frac{5\sigma^2}{m} \sqrt{\frac{K \log m}{m}} \right. \\
&\quad \times \left. \left( 15 \sqrt{\frac{\mu^2 s^2 K \log m}{m}} + 3\sqrt{K} \frac{\sigma^2}{m} \right) \right. \\
&\quad \left. \left| \mathbf{b}_l^* \check{\mathbf{h}}_i^0 \right| + \left( 15 \sqrt{\frac{\mu^2 s^2 K \log m}{m}} \sqrt{\frac{K}{m}} + 3\sqrt{K} \frac{\sigma^2}{m} \right) \left\| \tilde{\alpha}_i \check{\mathbf{h}}_i^0 - \check{\mathbf{h}}_i^{0,(l)} \right\|_2 \right\}. \tag{8.87}
\end{aligned}$$

Since the inequality (8.87) holds for any  $|\tilde{\alpha}_i| = 1$ , we can pick up  $\tilde{\alpha}_i = \beta_i^{0,(l)}$  and reformulate (8.87) as

$$\begin{aligned}
& \left( 1 - 30c_1 \sqrt{\frac{\mu^2 s^2 K \log m}{m}} \cdot \sqrt{\frac{K}{m}} - 6\sqrt{K} \frac{\sigma^2}{m} \right) \cdot \left\| \beta_i^{0,(l)} \check{\mathbf{h}}_i^0 - \check{\mathbf{h}}_i^{0,(l)} \right\|_2 \\
&\quad + \left\| \beta_i^{0,(l)} \check{\mathbf{x}}_i^0 - \check{\mathbf{x}}_i^{0,(l)} \right\|_2 \\
&\leq 60c_1 \frac{\mu}{\sqrt{m}} \cdot \sqrt{\frac{s^2 K \log^2 m}{m}} + \frac{10c_1 \sigma^2}{m} \sqrt{\frac{K \log m}{m}} \\
&\quad + \left( 30c_1 \sqrt{\frac{\mu^2 s^2 K \log m}{m}} + 6c_1 \sqrt{K} \frac{\sigma^2}{m} \right) \left| \mathbf{b}_l^* \check{\mathbf{h}}_i^0 \right|. \tag{8.88}
\end{aligned}$$

With the assumption that  $m \gg (\mu + \sigma^2)sK \log^{1/2} m$ , it yields  $1 - 30c_1 \sqrt{\frac{\mu^2 s^2 K \log m}{m}} \cdot \sqrt{\frac{K}{m}} - 6\sqrt{K} \frac{\sigma^2}{m} \leq \frac{1}{2}$ . Hence,

$$\begin{aligned} & \max_{1 \leq i \leq s, 1 \leq j \leq m} \left\| \beta_i^{0,(l)} \check{\mathbf{h}}_i^0 - \check{\mathbf{h}}_i^{0,(l)} \right\|_2 + \left\| \beta_i^{0,(l)} \check{\mathbf{x}}_i^0 - \check{\mathbf{x}}_i^{0,(l)} \right\|_2 \\ & \leq 120c_1 \frac{\mu}{\sqrt{m}} \cdot \sqrt{\frac{s^2 K \log^2 m}{m}} + \frac{20c_1 \sigma^2}{m} \sqrt{\frac{K \log m}{m}} \\ & \quad + \left( 60c_1 \sqrt{\frac{\mu^2 s^2 K \log m}{m}} + 12c_1 \sqrt{K} \frac{\sigma^2}{m} \right) \cdot \max_{1 \leq i \leq s, 1 \leq j \leq m} \left| \mathbf{b}_i^* \check{\mathbf{h}}_i^0 \right|. \end{aligned} \quad (8.89)$$

It thus suffices to control  $\max_{1 \leq i \leq s, 1 \leq j \leq m} |\mathbf{b}_i^* \check{\mathbf{h}}_i^0|$ . Denote  $\mathbf{M}_i \check{\mathbf{x}}^0 = \sigma_1(\mathbf{M}_i) \check{\mathbf{h}}_i^0$  and

$$\mathbf{W}_i = \sum_{j=1}^m \mathbf{b}_j \left( \sum_{k \neq i} \mathbf{b}_j^* \mathbf{h}_k^{\natural} \mathbf{x}_k^{\natural*} \mathbf{a}_{kj} + e_j \right) \mathbf{a}_{ij}^*, \quad (8.90)$$

which further leads to

$$\begin{aligned} \left| \mathbf{b}_i^* \check{\mathbf{h}}_i^0 \right| &= \frac{1}{\sigma_1(\mathbf{M}_i)} \left| \mathbf{b}_i^* \mathbf{M}_i \check{\mathbf{x}}_i^0 \right| \\ &\stackrel{(i)}{\leq} 2 \left| \sum_{j=1}^m (\mathbf{b}_i^* \mathbf{b}_j) \mathbf{b}_j^* \mathbf{h}_i^{\natural} \mathbf{x}_i^{\natural*} \mathbf{a}_{ij} \mathbf{a}_{ij}^* \check{\mathbf{x}}_i^0 \right| + 2 \left| \mathbf{b}_i^* \mathbf{W}_i \check{\mathbf{x}}_i^0 \right| \\ &\stackrel{(ii)}{\leq} 2 \left( \sum_{j=1}^m |\mathbf{b}_i^* \mathbf{b}_j| \right) \max_{1 \leq j \leq m} \left\{ \left| \mathbf{b}_j^* \mathbf{h}_i^{\natural} \right| \cdot \left| \mathbf{a}_{ij}^* \mathbf{x}_i^{\natural} \right| \cdot \left| \mathbf{a}_{ij}^* \check{\mathbf{x}}_i^0 \right| \right\} \\ &\quad + 2 \|\mathbf{b}_i\|_2 \cdot \|\mathbf{W}_i\| \cdot \|\check{\mathbf{x}}_i^0\|_2 \\ &\stackrel{(iii)}{\leq} 2 \sqrt{\frac{K}{m}} \cdot \|\mathbf{W}_i\| + 8 \log m \cdot \frac{\mu}{\sqrt{m}} \cdot (5\sqrt{\log m}) \cdot \\ &\quad \max_{1 \leq j \leq m} \left\{ \left| \mathbf{a}_j^* \check{\mathbf{x}}_i^{0,(j)} \right| + \|\mathbf{a}_{ij}\|_2 \left\| \beta_i^{0,(j)} \check{\mathbf{x}}_i^0 - \check{\mathbf{x}}_i^{0,(j)} \right\|_2 \right\} \\ &\stackrel{(iv)}{\leq} \sqrt{\frac{K}{m \log m}} + 200 \frac{\mu \log^2 m}{\sqrt{m}} + 120 \sqrt{\frac{\mu^2 K \log^3 m}{m}} \\ &\quad \cdot \max_{1 \leq j \leq m} \left\| \beta_i^{0,(j)} \check{\mathbf{x}}_i^0 - \check{\mathbf{x}}_i^{0,(j)} \right\|_2, \end{aligned} \quad (8.91)$$

where  $\beta_i^{0,(j)}$  is defined in (8.82). Here, (i) arises from the low bound  $\sigma_1(\mathbf{M}_i) \geq \frac{1}{2}$  and the triangle inequality. (ii) uses the Cauchy–Schwarz inequality. The step (iii) comes from combining the incoherence condition, the bound (8.40), the triangle inequality, the estimate:  $\sum_{j=1}^m |\mathbf{b}_i^* \mathbf{b}_j| \leq 4 \log m$  [13, Lemma 48],  $\|\mathbf{b}_i\| = \sqrt{K/m}$  and  $\|\check{\mathbf{x}}_i^0\|_2 = 1$ . The last step (iv) exploits the inequality (8.85) to yield that with probability  $1 - O(m^{-9})$  [12],

$$\|\mathbf{W}_i\| \leq \frac{1}{2\sqrt{\log m}}, \quad (8.92)$$

if  $m \gg (\mu^2 + \sigma^2)sK \log^2 m$ . The bound (8.91) further leads to

$$\begin{aligned} \max_{1 \leq i \leq s} |\mathbf{b}_i^* \check{\mathbf{h}}_i^0| &\leq \sqrt{\frac{K}{m \log m}} + 200 \frac{\mu \log^2 m}{\sqrt{m}} + 120 \sqrt{\frac{\mu^2 K \log^3 m}{m}} \\ &\quad \max_{1 \leq i \leq s, 1 \leq j \leq m} \left\| \beta_i^{0,(j)} \check{\mathbf{x}}_i^0 - \check{\mathbf{x}}_i^{0,(j)} \right\|_2. \end{aligned} \quad (8.93)$$

Combining the bound (8.89) and (8.93) gives

$$\begin{aligned} &\max_{1 \leq i \leq s, 1 \leq l \leq m} \left\| \beta_i^{0,(l)} \check{\mathbf{h}}_i^0 - \check{\mathbf{h}}_i^{0,(l)} \right\|_2 + \left\| \beta_i^{0,(l)} \check{\mathbf{x}}_i^0 - \check{\mathbf{x}}_i^{0,(l)} \right\|_2 \\ &\leq 120c_1 \frac{\mu}{\sqrt{m}} \cdot \sqrt{\frac{s^2 K \log^2 m}{m}} + \frac{20c_1 \sigma^2}{m} \sqrt{\frac{K \log m}{m}} \\ &\quad + \left( 60c_1 \sqrt{\frac{\mu^2 s^2 K \log m}{m}} + 12c_1 \sqrt{K \frac{\sigma^2}{m}} \right) \\ &\quad \left( \sqrt{\frac{K}{m \log m}} + 200 \frac{\mu \log^2 m}{\sqrt{m}} + 120 \sqrt{\frac{\mu^2 K \log^3 m}{m}} \right) \\ &\quad \max_{1 \leq i \leq s, 1 \leq j \leq m} \left\| \beta_i^{0,(j)} \check{\mathbf{x}}_i^0 - \check{\mathbf{x}}_i^{0,(j)} \right\|_2. \end{aligned} \quad (8.94)$$

As long as  $m \gg (\mu^2 + \sigma^2)s^2 K \log^2 m$ , we have

$$\left( 60c_1 \sqrt{\frac{\mu^2 s^2 K \log m}{m}} + 12c_1 \sqrt{K \frac{\sigma^2}{m}} \right) \cdot 120 \sqrt{\frac{\mu^2 s^2 K \log^3 m}{m}} \leq 1/2. \quad (8.95)$$

Reformulating the inequality (8.94), we have

$$\begin{aligned} & \max_{1 \leq i \leq s, 1 \leq l \leq m} \left\| \beta_i^{0,(l)} \check{\mathbf{h}}_i^0 - \check{\mathbf{h}}_i^{0,(l)} \right\|_2 + \left\| \beta_i^{0,(l)} \check{\mathbf{x}}_i^0 - \check{\mathbf{x}}_i^{0,(l)} \right\|_2 \\ & \leq C_4 \frac{\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 s^2 K \log^5 m}{m}}, \end{aligned} \quad (8.96)$$

for some constant  $C_4 > 0$ . Taking the bound (8.96) together with (8.93), it yields

$$\begin{aligned} \max_{1 \leq i \leq s, 1 \leq l \leq m} \left| \mathbf{b}_i^* \check{\mathbf{h}}_i^0 \right| & \leq \sqrt{\frac{K}{m \log m}} + 200 \frac{\mu \log^2 m}{\sqrt{m}} + 120 \sqrt{\frac{\mu^2 K \log^3 m}{m}} \\ & \max_{1 \leq i \leq s, 1 \leq j \leq m} C_4 \frac{\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 s^2 K \log^5 m}{m}} \\ & \leq c_2 \frac{\mu \log^2 m}{\sqrt{m}}, \end{aligned} \quad (8.97)$$

for some constant  $c_2 > 0$ , as long as  $m \gg (\mu^2 + \sigma^2) s K \log^2 m$ .

We further scale the preceding bounds to the final version. Based on [13, Section C.6], one has

$$\begin{aligned} & \left\| \alpha \mathbf{h}^0 - \mathbf{h}^{0,(l)} \right\|_2 + \left\| \alpha \mathbf{x}^0 - \mathbf{x}^{0,(l)} \right\|_2 \\ & \leq \left\| \left( \mathbf{M}_i - \mathbf{M}_i^{(l)} \right) \check{\mathbf{x}}_i^{0,(l)} \right\|_2 + 6 \left\{ \left\| \alpha \check{\mathbf{h}}_i^0 - \check{\mathbf{h}}_i^{0,(l)} \right\|_2 + \left\| \alpha \check{\mathbf{x}}_i^0 - \check{\mathbf{x}}_i^{0,(l)} \right\|_2 \right\}. \end{aligned} \quad (8.98)$$

Taking the bounds (8.84), (8.96), and (8.98) collectively yields

$$\min_{\alpha_i \in \mathbb{C}, |\alpha_i|=1} \left\| \alpha_i \mathbf{h}_i^0 - \mathbf{h}_i^{0,(l)} \right\|_2 + \left\| \alpha_i \mathbf{x}_i^0 - \mathbf{x}_i^{0,(l)} \right\|_2 \leq c_5 \frac{\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 s^2 K \log^5 m}{m}}, \quad (8.99)$$

for some constant  $c_5 > 0$ , as long as  $m \gg (\mu^2 + \sigma^2) s^2 K \log^2 m$ .

Furthermore, by exploiting the technical methods provided in [13, Section C.6], we have

$$\begin{aligned}
\text{dist}\left(\mathbf{z}^{0,(l)}, \tilde{\mathbf{z}}^0\right) &= \min_{\alpha_i \in \mathbb{C}} \sqrt{\sum_{i=1}^s \left\| \frac{1}{\alpha_i} \mathbf{h}^{0,(l)} - \frac{1}{\alpha_i^0} \mathbf{h}^0 \right\|_2^2 + \|\alpha_i \mathbf{x}^{0,(l)} - \alpha_i^0 \mathbf{x}^0\|_2^2} \\
&\stackrel{(i)}{\leq} \min_{\alpha_i \in \mathbb{C}, |\alpha_i|=1} \sqrt{\sum_{i=1}^s \left\| \frac{1}{\alpha_i^0} \mathbf{h}^0 - \frac{\alpha_i}{\alpha_i^0} \mathbf{h}^{0,(l)} \right\|_2^2 + \|\alpha_i^0 \mathbf{x}^0 - \alpha_i \alpha_i^0 \mathbf{x}^{0,(l)}\|_2^2} \\
&\stackrel{(ii)}{\leq} 2\sqrt{s} \min_{\alpha_i \in \mathbb{C}, |\alpha_i|=1} \left\{ \left\| \mathbf{h}_i^0 - \alpha_i \mathbf{h}_i^{0,(l)} \right\|_2 + \left\| \mathbf{x}_i^0 - \alpha_i \mathbf{x}_i^{0,(l)} \right\|_2 \right\} \\
&\leq 2c_5 \frac{s\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 s K \log^5 m}{m}}, \tag{8.100}
\end{aligned}$$

where  $\alpha_i^0$  is defined in (3.30) and satisfies

$$\frac{1}{2} \leq |\alpha_i^0| \leq 2. \tag{8.101}$$

Here, the step (i) occurs since the feasible set for the latter optimization problem is smaller, and (ii) follows directly from [13, Lemma 19], [13, Lemma 52]. This accomplishes the proof for the claim (8.77). We further move to the proof for the claim (8.78).

In terms of  $|\mathbf{b}_i^* \tilde{\mathbf{h}}_i^0|$ , one has

$$\left| \mathbf{b}_i^* \tilde{\mathbf{h}}_i^0 \right| \leq \left| \mathbf{b}_i^* \frac{1}{\alpha_i^0} \mathbf{h}_i^0 \right| \leq \left| \frac{1}{\alpha_i^0} \right| \left| \mathbf{b}_i^* \mathbf{h}_i^0 \right| \leq 2 \left| \sqrt{\sigma_1(\mathbf{M}_i)} \mathbf{b}_i^* \check{\mathbf{h}}_i^0 \right| \leq 2\sqrt{2}c_2 \frac{\mu \log^2 m}{\sqrt{m}}, \tag{8.102}$$

based on fact that

$$\frac{1}{2} \leq \sigma_1(\mathbf{M}_i) \leq 2. \tag{8.103}$$

**Lemma 8.8** *Suppose the sample complexity  $m \gg (\mu^2 + \sigma^2)s^{3/2}K \log^5 m$ . Then with probability at least  $1 - O(m^{-9})$ ,*

$$\max_{1 \leq i \leq s, 1 \leq j \leq m} \left| \mathbf{a}_{ij}^* \left( \tilde{\mathbf{x}}_i^0 - \mathbf{x}_i^{\natural} \right) \right| \cdot \left\| \mathbf{x}_i^{\natural} \right\|_2^{-1} \leq C_3 \frac{1}{\sqrt{s} \log^{3/2} m}. \tag{8.104}$$

**Proof** Recall several alignment parameters defined before:

$$\begin{aligned}\alpha_i^0 &:= \operatorname{argmin}_{\alpha \in \mathbb{C}} \left\| \frac{1}{\alpha} \mathbf{h}_i^0 - \mathbf{h}_i^\natural \right\|_2^2 + \left\| \alpha \mathbf{x}_i^0 - \mathbf{x}_i^\natural \right\|_2^2, \\ \alpha_i^{0,(l)} &:= \operatorname{argmin}_{\alpha \in \mathbb{C}} \left\| \frac{1}{\alpha} \mathbf{h}_i^{0,(l)} - \mathbf{h}_i^\natural \right\|_2^2 + \left\| \alpha \mathbf{x}_i^{0,(l)} - \mathbf{x}_i^\natural \right\|_2^2, \\ \alpha_{i,\text{mutual}}^{0,(l)} &:= \operatorname{argmin}_{\alpha \in \mathbb{C}} \left\| \frac{1}{\alpha} \mathbf{h}_i^{0,(l)} - \frac{1}{\alpha_i^0} \mathbf{h}_i^0 \right\|_2^2 + \left\| \alpha \mathbf{x}_i^{0,(l)} - \alpha_i^0 \mathbf{x}_i^0 \right\|_2^2.\end{aligned}$$

Combining (8.17a) and (8.100) with the triangle inequality yields that

$$\begin{aligned}&\leq \sqrt{\left\| \frac{1}{\alpha_{i,\text{mutual}}^{0,(l)}} \mathbf{h}_i^{0,(l)} - \frac{1}{\alpha_i^0} \mathbf{h}_i^0 \right\|_2^2 + \left\| \alpha_{i,\text{mutual}}^{0,(l)} \mathbf{x}_i^{0,(l)} - \alpha_i^0 \mathbf{x}_i^0 \right\|_2^2} \\ &\leq 2c_5 \frac{s\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 K \log^5 m}{m}} + C_1 \frac{1}{\sqrt{s} \log^2 m} \\ &\leq 2C_1 \frac{1}{\sqrt{s} \log^2 m},\end{aligned}\tag{8.105}$$

where the last step holds as long as  $m \gg (\mu^2 + \sigma^2) s \sqrt{sK} \log^{9/2} m$ .

According to [13, Section C.7], [13, Lemma 55], and the bound (8.100), it implies that

$$\begin{aligned}&\sqrt{\left\| \frac{1}{\alpha_i^{0,(l)}} \mathbf{h}_i^{0,(l)} - \frac{1}{\alpha_i^0} \mathbf{h}_i^0 \right\|_2^2 + \left\| \alpha_i^{0,(l)} \mathbf{x}_i^{0,(l)} - \alpha_i^0 \mathbf{x}_i^0 \right\|_2^2} \\ &\lesssim \sqrt{\left\| \frac{1}{\alpha_i^0} \mathbf{h}_i^0 - \frac{1}{\alpha_{i,\text{mutual}}^{0,(l)}} \mathbf{h}_i^{0,(l)} \right\|_2^2 + \left\| \alpha_i^0 \mathbf{x}_i^0 - \alpha_{i,\text{mutual}}^{0,(l)} \mathbf{x}_i^{0,(l)} \right\|_2^2} \\ &\lesssim \frac{s\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 K \log^5 m}{m}}.\end{aligned}\tag{8.106}$$

Based on the above estimate, we can show that with high probability,

$$\begin{aligned}&\left| \mathbf{a}_{il}^* \left( \alpha_i^0 \mathbf{x}_i^0 - \mathbf{x}_i^\natural \right) \right| \\ &\stackrel{(i)}{\leq} \left| \mathbf{a}_{il}^* \left( \alpha_i^{0,(l)} \mathbf{x}_i^{0,(l)} - \mathbf{x}_i^\natural \right) \right| + \left| \mathbf{a}_{il}^* \left( \alpha_i^0 \mathbf{x}_i^0 - \alpha_i^{0,(l)} \mathbf{x}_i^{0,(l)} \right) \right|\end{aligned}$$

$$\begin{aligned}
& \stackrel{\text{(ii)}}{\leq} 5\sqrt{\log m} \left\| \mathbf{a}_{il}^* \left( \alpha^{0,(l)} \mathbf{x}_i^{0,(l)} - \mathbf{x}_i^{\natural} \right) \right\|_2 + \|\mathbf{a}_{il}\|_2 \left\| \mathbf{a}_{il}^* \left( \alpha_i^0 \mathbf{x}_i^0 - \alpha_i^{0,(l)} \mathbf{x}_i^{0,(l)} \right) \right\|_2 \\
& \stackrel{\text{(iii)}}{\lesssim} \sqrt{\log m} \cdot \frac{1}{\sqrt{s \log^2 m}} + \sqrt{K} \frac{s\mu}{\sqrt{m}} \sqrt{\frac{\mu^2 K \log^5 m}{m}} \\
& \stackrel{\text{(iv)}}{\lesssim} \frac{1}{\sqrt{s \log^{3/2} m}}, \tag{8.107}
\end{aligned}$$

where (i) arises from the triangle inequality, (ii) uses Cauchy–Schwarz inequality and the independence between  $\mathbf{x}_i^{0,(l)}$  and  $\mathbf{a}_{il}$ , (iii) holds since (8.106), and (iv) occurs as long as  $m \gg (\mu^2 + \sigma^2)s^{3/2}K \log^4 m$ .

## 8.4 Theoretical Analysis of Wirtinger Flow with Random Initialization for Blind Demixing

Based on the notations for blind demixing introduced in Sect. 3.4.2, we present the theoretical analysis of Wirtinger flow with random initialization in this section, which is based on [6]. It demonstrates that random initialization which enjoys a model-agnostic and natural initialization implementation for practitioners is good enough to guarantee Wirtinger flow to linearly converge to the optimal solution.

To present the theorem, we begin with introducing some notations. Let  $\tilde{\mathbf{h}}_i^t$  and  $\tilde{\mathbf{x}}_i^t$  be

$$\tilde{\mathbf{h}}_i^t = \frac{1}{\omega_i^t} \mathbf{h}_i^t \quad \text{and} \quad \tilde{\mathbf{x}}_i^t = \omega_i^t \mathbf{x}_i^t, \tag{8.108}$$

for  $i = 1, \dots, s$ , respectively, where alignment parameters are denoted as  $\omega_i$ . Without loss of the generality, we assume that the ground truth  $\mathbf{x}_i^{\natural} = q_i \mathbf{e}_1$  for  $i = 1, \dots, s$ , where  $0 < q_i \leq 1$ ,  $i = 1, \dots, s$  are some constants and define a parameter  $\kappa = \frac{\max_i q_i}{\min_i q_i}$ . For simplification, for  $i = 1, \dots, s$ , the first entry and the rest entries of  $\mathbf{x}_i^t$  are denoted as

$$x_{i1}^t \quad \text{and} \quad \mathbf{x}_{i\perp}^t := \left[ x_{ij}^t \right]_{2 \leq j \leq N}, \tag{8.109}$$

respectively. Hence, (8.113) and (8.114) can be reformulated as

$$\alpha_{x_i} := \tilde{x}_{i1}^t \quad \text{and} \quad \beta_{x_i} := \|\tilde{\mathbf{x}}_{i\perp}^t\|_2. \tag{8.110}$$

Define the norm of signal component and the perpendicular component in terms of  $\mathbf{h}_i$  for  $i = 1, \dots, s$ , as

$$\alpha_{h_i} := \langle \mathbf{h}_i^{\natural}, \tilde{\mathbf{h}}_i^t \rangle / \|\mathbf{h}_i^{\natural}\|_2, \quad (8.111)$$

$$\beta_{h_i} := \left\| \tilde{\mathbf{h}}_i^t - \frac{\langle \mathbf{h}_i^{\natural}, \tilde{\mathbf{h}}_i^t \rangle}{\|\mathbf{h}_i^{\natural}\|_2^2} \mathbf{h}_i^{\natural} \right\|_2, \quad (8.112)$$

respectively. Likewise, the norms of the signal component and the perpendicular component in terms of  $\mathbf{x}_i$  for  $i = 1, \dots, s$  are given by

$$\alpha_{x_i} := \langle \mathbf{x}_i^{\natural}, \tilde{\mathbf{x}}_i^t \rangle / \|\mathbf{x}_i^{\natural}\|_2, \quad (8.113)$$

$$\beta_{x_i} := \left\| \tilde{\mathbf{x}}_i^t - \frac{\langle \mathbf{x}_i^{\natural}, \tilde{\mathbf{x}}_i^t \rangle}{\|\mathbf{x}_i^{\natural}\|_2^2} \mathbf{x}_i^{\natural} \right\|_2, \quad (8.114)$$

respectively.

**Theorem 8.3 ([6])** *Assuming that the initial points are randomly generated as (3.33), and the stepsize  $\eta > 0$  obeys  $\eta \asymp s^{-1}$ . 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  $1 - c_1 m^{-\nu} - c_1 m e^{-c_2 N}$  with some constants  $\nu, c_1, c_2 > 0$ , for a sufficiently small constant  $0 \leq \gamma \leq 1$  and  $T_\gamma \lesssim s \log(\max\{K, N\})$ , it holds that

1. *The randomly initialized Wirtinger flow 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 in terms of  $\mathbf{h}_i^t$  and  $\mathbf{x}_i^t$  obey*

$$\max_{1 \leq i \leq s} \frac{\alpha_{h_i^t}}{\beta_{h_i^t}} \gtrsim \frac{1}{\sqrt{K \log K}} (1 + c_3 \eta)^t, \quad (8.115a)$$

$$\max_{1 \leq i \leq s} \frac{\alpha_{x_i^t}}{\beta_{x_i^t}} \gtrsim \frac{1}{\sqrt{N \log N}} (1 + c_4 \eta)^t, \quad (8.115b)$$

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

The precise statistical analysis on the computational efficiency of Wirtinger flow with random initialization is illustrated in Theorem 8.3. In Stage I, it takes randomly initialized Wirtinger flow  $T_\gamma = \mathcal{O}(s \log(\max\{K, N\}))$  iterations to reach a local region near the ground truth that enjoys strong convexity and strong smoothness. In Stage II, it takes  $\mathcal{O}(s \log(1/\varepsilon))$  iterations to linearly converge  $\varepsilon$ -accurate point. Hence, the randomly initialized Wirtinger flow is guaranteed to converge to the ground truth with the iteration complexity  $\mathcal{O}(s \log(\max\{K, N\}) + s \log(1/\varepsilon))$  where the sample size is  $m \gtrsim s^2 \max\{K, N\} \text{poly log } m$ .

The proof of Theorem 8.3 is briefly summarized in the following. The key idea is to investigate the dynamics of the iterates of Wirtinger flow with random initialization.

### 1. Stage I:

- **Dynamics of population-level state evolution.** Establish the population-level state evolution of  $\alpha_{x_i}$  (8.116a) and  $\beta_{x_i}$  (8.116b),  $\alpha_{h_i}$  (8.117a),  $\beta_{h_i}$  (8.117b), respectively:

$$\alpha_{x_i^{t+1}} = (1 - \eta)\alpha_{x_i^t} + \eta \frac{q_i \alpha_{h_i^t}}{\alpha_{h_i^t}^2 + \beta_{h_i^t}^2}, \quad (8.116a)$$

$$\beta_{x_i^{t+1}} = (1 - \eta)\beta_{x_i^t}. \quad (8.116b)$$

Similarly, the population-level state evolution for both  $\alpha_{h_i^t}$  and  $\beta_{h_i^t}$ :

$$\alpha_{h_i^{t+1}} = (1 - \eta)\alpha_{h_i^t} + \eta \frac{q_i \alpha_{x_i^t}}{\alpha_{x_i^t}^2 + \beta_{x_i^t}^2}, \quad (8.117a)$$

$$\beta_{h_i^{t+1}} = (1 - \eta)\beta_{h_i^t}, \quad (8.117b)$$

where the sample size approaches infinity. The approximate state evolution (8.118) is then established, which is significantly close to the population-level state evolution:

$$\alpha_{h_i^{t+1}} = \left( 1 - \eta + \frac{\eta q_i \gamma_{h_i^t}}{\alpha_{x_i^t}^2 + \beta_{x_i^t}^2} \right) \alpha_{h_i^t} + \eta (1 - \eta) \frac{q_i \alpha_{x_i^t}}{\alpha_{x_i^t}^2 + \beta_{x_i^t}^2}, \quad (8.118a)$$

$$\beta_{h_i^{t+1}} = \left( 1 - \eta + \frac{\eta q_i \varphi_{h_i^t}}{\alpha_{x_i^t}^2 + \beta_{x_i^t}^2} \right) \beta_{h_i^t}, \quad (8.118b)$$

$$\alpha_{x_i^{t+1}} = \left( 1 - \eta + \frac{\eta q_i \gamma_{x_i^t}}{\alpha_{\nu} h_i^{t2} + \beta_{h_i^2}^2} \right) \alpha_{x_i^t} + \eta \left( 1 - \nu_{x_i^t} \right) \frac{q_i \alpha_{h_i^t}}{\alpha_{h_i^2}^2 + \beta_{h_i^2}^2}, \quad (8.118c)$$

$$\beta_{x_i^{t+1}} = \left( 1 - \eta + \frac{\eta q_i \varphi_{x_i^t}}{\alpha_{h_i^2}^2 + \beta_{h_i^2}^2} \right) \beta_{x_i^t}, \quad (8.118d)$$

where the perturbation terms are denoted as  $\{\gamma_{h_i^t}\}, \{\gamma_{x_i^t}\}, \{\varphi_{h_i^t}\}, \{\varphi_{x_i^t}\}, \{\nu_{h_i^t}\},$  and  $\{\nu_{x_i^t}\}.$

- **Dynamics of approximate state evolution.** Show that if  $\alpha_{h_i}$  (8.111),  $\beta_{h_i}$  (8.112),  $\alpha_{x_i}$  (8.113), and  $\beta_{x_i}$  (8.114) obey the approximate state evolution (8.118), it has some  $T_\gamma = \mathcal{O}(s \log(\max\{K, N\}))$  such that  $\text{dist}(z^{T_\gamma}, z^{\natural}) \leq \gamma.$  Furthermore, the ratio  $\alpha_{h_i}/\beta_{h_i}$  and  $\alpha_{x_i}/\beta_{x_i}$  enjoys exponential growth.
- **Leave-one-out arguments.** Identify the conditions where  $\alpha_{h_i}, \beta_{h_i}, \alpha_{x_i},$  and  $\beta_{x_i}$  obey the approximate state evolution (8.118) with high probability, followed by demonstrating the iterates of randomly initialized Wirtinger flow that solve the blind demixing problem satisfy the conditions.

## 2. Stage II: Local geometry in the region of incoherence and contraction.

Invoke the prior theory provided in [5] to show local convergence of the random initialized Wirtinger flow in Stage II.

## 8.5 The Basic Concepts on Riemannian Optimization

As a supplementary of Sect. 3.4.3, we introduce some basic concepts on Riemannian optimization via some examples. More details can be referred to the book [1].

**Embedded Submanifolds** Denote  $\mathcal{N}$  as a subset of a manifold  $\mathcal{M},$  which  $\mathcal{N}$  admits at most one differentiable structure [1, Section 3.3]. Some examples of embedded submanifold are provided in the following.

*Example 8.1 (The Stiefel Manifold)* The Stiefel manifold is an embedded submanifold of  $\mathbb{R}^{m \times n}.$  For  $n \leq m,$  a Stiefel manifold can be denoted as

$$\text{St}(n, m) := \left\{ X \in \mathbb{R}^{m \times n} : X^\top X = I_n \right\}, \quad (8.119)$$

which is the set of all  $m \times n$  orthonormal matrices, where  $I_n$  denotes the  $n \times n$  identity matrix. For  $n = 1,$  the Stiefel manifold  $\text{St}(n, m)$  reduces to the unit sphere  $\mathbb{S}^{m-1}.$

**Tangent Vectors and Tangent Spaces** Denote  $\mathfrak{F}_x(\mathcal{M})$  as the set of smooth real-valued functions. A tangent vector  $\xi_x$  to a manifold  $\mathcal{M}$  at a point  $x$  is a mapping from  $\mathfrak{F}_x(\mathcal{M})$  to  $\mathbb{R}$ . For  $f \in \mathfrak{F}_x(\mathcal{M})$ , there exists a curve  $\phi$  on  $\mathcal{M}$  with  $\phi(0) = x$ , such that

$$\xi_x f = \left. \frac{d(f(\phi(t)))}{dt} \right|_{t=0}. \quad (8.120)$$

Based on the curve  $\phi$ , it yields a straightforward identification of the tangent space  $T_x\mathcal{M}$  with the set

$$\{\phi'(0) : \phi \text{ curve in } \mathcal{M}, \phi(0) = x\}. \quad (8.121)$$

An example of tangent space for sphere is presented in the following.

*Example 8.2 (Tangent Space to a Sphere)* Define a curve in the unit sphere  $S^{n-1}$  as  $t \mapsto \gamma(t)$ , and there is  $\gamma_0$  at  $t = 0$ . Since  $\gamma(t) \in S^{n-1}$  for all  $t$ , it holds that

$$\gamma^\top(t)\gamma(t) = 1 \quad (8.122)$$

for all  $t$ . Equation (8.122) is differentiated in terms of  $t$ , yielding

$$\dot{\gamma}^\top(t)\gamma(t) + \gamma^\top(t)\dot{\gamma}(t) = 0. \quad (8.123)$$

Thus  $\dot{\gamma}(0)$  is an entry of the set

$$\{\mathbf{x} \in \mathbb{R}^n : \gamma_0^\top \mathbf{x} = 0\}. \quad (8.124)$$

Furthermore, let  $\mathbf{x}$  belong to the set (8.124). Then the curve

$$t \mapsto \gamma(t) := \frac{\gamma_0 + t\mathbf{x}}{\|\gamma_0 + t\mathbf{x}\|}$$

is on  $S^{n-1}$  and it holds  $\dot{\gamma}(0) = \mathbf{x}$ . Hence (8.124) is a subset of  $T_{\gamma_0}S^{n-1}$ , and

$$T_\gamma S^{n-1} = \{\mathbf{x} \in \mathbb{R}^n : \gamma^\top \mathbf{x} = 0\} \quad (8.125)$$

is the set of vectors orthogonal to the curve  $\gamma$  in  $\mathbb{R}^n$ .

**Riemannian Metric** As mentioned above, the notion of a directional derivative can be generalized by tangent vectors. To further identify which direction of act from  $\mathbf{x}$  yields the steepest decrease in  $f$ , a notion of length with respect to tangent vectors is required. This can be achieved by assigning an inner product  $\langle \cdot, \cdot \rangle_x$ , i.e., a symmetric positive-definite or bilinear operator, to each tangent space  $T_x\mathcal{M}$ . The inner product  $\langle \cdot, \cdot \rangle_x$  for the point  $\mathbf{x} \in \mathcal{M}$  is called the Riemannian metric, which can be represented as  $g_x$ .

**Product Manifolds** The differentiable structure defined by two compact manifold  $\mathcal{M}_1$  and  $\mathcal{M}_2$ ,  $\mathcal{M}_1 \times \mathcal{M}_2$  is called the product of the manifolds  $\mathcal{M}_1$  and  $\mathcal{M}_2$ . Its manifold topology is equivalent to the product topology [1, Section 3.1.6], which means that the geometry concepts on the product manifolds can be represented by the set of elementwise geometry concepts on individual manifold. An example of product manifolds in the blind demixing problem (3.41) is introduced in the following.

*Example 8.3 (Product Manifolds in the Blind Demixing Problem (3.41))* Taking the individual manifold  $\mathcal{M}$  as an example, a smoothly varying inner product  $g_X(\zeta_X, \eta_X)$ , where  $\zeta_X, \eta_X \in T_X \mathcal{M}$ , characterizes the notion of length that applies to each tangent space  $T_X \mathcal{M}$ . With a smoothly varying inner product  $g_X$ , the manifold  $\mathcal{M}$  is called the *Riemannian manifold*, and the inner product is called the *Riemannian metric*. Denote  $\mathcal{M}$  as the Riemannian manifold endowed with the Riemannian metric  $g_{X_k}$ , where  $k \in [s]$  with  $[s] = \{1, 2, \dots, s\}$ . The set of matrices  $(X_1, \dots, X_s)$  where  $X_k \in \mathcal{M}, k = 1, 2, \dots, s$  is denoted as  $\mathcal{M}^s = \underbrace{\mathcal{M} \times \mathcal{M} \times \dots \times \mathcal{M}}_s$ , and is called product manifold.

Based on the Riemannian geometry concepts, the notion of length on the product manifold can be characterized via endowing tangent space  $T_V \mathcal{M}^s$  with the smoothly varying inner product, given by

$$g_V(\zeta_V, \eta_V) := \sum_{k=1}^s g_{X_k}(\zeta_{X_k}, \eta_{X_k}), \quad (8.126)$$

where  $\zeta_V, \eta_V \in T_V \mathcal{M}^s$  and  $\zeta_{X_k}, \eta_{X_k} \in T_{X_k} \mathcal{M}$ . Since  $\mathcal{M}$  is the Riemannian manifold endowed with the Riemannian metric  $g_{X_k}$  for  $\forall k \in [s]$ , the product manifold  $\mathcal{M}^s$  is also a Riemannian manifold, endowed with the Riemannian metric  $g_V$ .

**Quotient Manifolds** Computations related to subspaces are generally operated via representing the subspace by the span of corresponding matrices' columns. For a given subspace, to represent the subspace with a unique matrix, it is beneficial to divide the set of matrices into classes of "equivalent" elements that serve as the same object. This operation yields the geometry concept of quotient spaces, which is called *quotient manifolds* when concerning the Riemannian manifold optimization. We first present the general theory of quotient manifolds, then we introduce the corresponding representations of the blind demixing problem.

Denote a manifold endowed with an *equivalence relation* as  $\sim \mathcal{M}$ . Then the equivalence class containing  $\mathbf{x}$  can be represented by the set

$$[\mathbf{x}] := \{\mathbf{y} \in \mathcal{M} : \mathbf{y} \sim \mathbf{x}\}, \quad (8.127)$$

which contains all elements that are equivalent to a point  $\mathbf{x}$ . The quotient of  $\mathcal{M}$  by  $\sim$  is defined as

$$\mathcal{M} / \sim := \{[\mathbf{x}] : \mathbf{x} \in \mathcal{M}\}, \quad (8.128)$$

which contains all equivalence classes of  $\sim$  in  $\mathcal{M}$ . Thus, the points of  $\mathcal{M}/\sim$  are subsets of  $\mathcal{M}$ , and the set  $\mathcal{M}/\sim$  is called the total space of the quotient  $\mathcal{M}/\sim$ .

Furthermore, a natural projection that maps the elements in the manifold  $\mathcal{M}$  to the quotient manifold  $\mathcal{M}/\sim$  is defined as  $\pi : \mathcal{M} \rightarrow \mathcal{M}/\sim$ . If and only if  $\mathbf{x} \sim \mathbf{y}$ , it holds that  $\pi(\mathbf{x}) = \pi(\mathbf{y})$  such that  $[\mathbf{x}] = \pi^{-1}(\pi(\mathbf{x}))$ .

An example of quotient manifold in the blind demixing problem (3.41) is provided in the following.

*Example 8.4 (Quotient Manifold in the Blind Demixing Problem)* According to the blind demixing problem given by (3.20):

$$\begin{aligned} & \text{find } \text{rank}(\mathbf{W}_i) = 1, \text{ for } \mathbf{W}_1, \dots, \mathbf{W}_s \\ & \text{subject to } \left\| \sum_{i=1}^s \mathcal{A}_i(\mathbf{W}_i) - \mathbf{y} \right\|_2 \leq \varepsilon, \end{aligned}$$

this problem is a rank-constrained optimization problem. The key idea of Riemannian optimization for rank-constrained problem is based on matrix factorization [16, 23]. Specifically, the factorization  $\mathbf{M}_k = \mathbf{w}_k \mathbf{w}_k^H$  in problem (3.41) is established to identify rank-one Hermitian positive semidefinite matrices [22, 23]. Nevertheless, the factorization  $\mathbf{M}_k = \mathbf{w}_k \mathbf{w}_k^H$  is not unique since the transformation  $\mathbf{w}_k \mapsto a_k \mathbf{w}_k$  with  $a_k \in \{a_k \in \mathbb{C} : a_k a_k^* = a_k^* a_k = 1\}$  makes the matrix  $\mathbf{w}_k \mathbf{w}_k^H$  unchanged. To address this indeterminacy, the transformation  $\mathbf{w}_k \mapsto a_k \mathbf{w}_k$  where  $k = 1, 2, \dots, s$ , is embedded in an abstract search space, which constructs the equivalence class:

$$[\mathbf{M}_k] = \{a_k \mathbf{w}_k : a_k a_k^* = a_k^* a_k = 1, a_k \in \mathbb{C}\}. \quad (8.129)$$

The product of  $[\mathbf{M}_k]$ 's yields the equivalence class

$$[\mathbf{V}] = \{[\mathbf{M}_k]\}_{k=1}^s, \quad (8.130)$$

which is denoted as  $\mathcal{M}^s/\sim$ , called the *quotient space*. Since the quotient manifold  $\mathcal{M}^s/\sim$  is an abstract space, the matrix representations defined in the computational space are needed to represent corresponding abstract geometric objects in the abstract space [1], thereby implementing the optimization algorithms. Denote an element of the quotient space  $\mathcal{M}^s/\sim$  as  $\tilde{\mathbf{V}}$  and its matrix representation in the computational space  $\mathcal{M}^s$  as  $\mathbf{V}$ . Hence, there exists  $\tilde{\mathbf{V}} = \pi(\mathbf{V})$  and  $[\mathbf{V}] = \pi^{-1}(\pi(\mathbf{V}))$ , where the mapping  $\pi : \mathcal{M}^s \rightarrow \mathcal{M}^s/\sim$  is the natural projection.

**Riemannian Submanifolds** Denote an embedded submanifold of a Riemannian manifold  $\overline{\mathcal{M}}$  as  $\mathcal{M}$ . Note that each tangent space  $T_{\mathbf{x}}\mathcal{M}$  can be termed as a subspace of  $T_{\mathbf{x}}\overline{\mathcal{M}}$ . Hence, the Riemannian metric  $g$  on  $\mathcal{M}$  can be derived by a Riemannian metric  $\bar{g}$  of  $\overline{\mathcal{M}}$  given by

$$g_{\mathbf{x}}(\boldsymbol{\eta}, \boldsymbol{\zeta}) = \bar{g}_{\mathbf{x}}(\boldsymbol{\eta}, \boldsymbol{\zeta}), \quad \boldsymbol{\eta}, \boldsymbol{\zeta} \in T_{\mathbf{x}}\mathcal{M}. \quad (8.131)$$

This implies that  $\mathcal{M}$  is a Riemannian manifold. The manifold  $\mathcal{M}$  with the Riemannian metric  $g_x$  is called a Riemannian submanifold of  $\overline{\mathcal{M}}$ . The orthogonal complement of the tangent space  $T_x\mathcal{M}$  in  $T_x\overline{\mathcal{M}}$  is called the normal space to  $\mathcal{M}$  at  $x$ , which is denoted by  $(T_x\mathcal{M})^\perp$ :

$$(T_x\mathcal{M})^\perp = \left\{ \eta \in T_x\overline{\mathcal{M}} : g_x(\eta, \zeta) = 0 \text{ for all } \zeta \in T_x\mathcal{M} \right\}. \quad (8.132)$$

Thus, the sum of an element of  $T_x\mathcal{M}$  and an element of  $(T_x\mathcal{M})^\perp$  can yield an element  $\eta \in T_x\overline{\mathcal{M}}$ :

$$\eta = P_x\eta + P_x^\perp\eta, \quad (8.133)$$

where  $P_x^\perp$  is the orthogonal projection onto  $(T_x\mathcal{M})^\perp$  and  $P_x$  is the orthogonal projection onto  $T_x\mathcal{M}$ . We first present a simple example, i.e., sphere  $S^{n-1}$  which is a Riemannian submanifold of  $\mathbb{R}^n$ . Based on the Riemannian submanifold of product manifolds, we further introduce the decomposition of the tangent space  $T_V\mathcal{M}^S$  of the product manifold  $\mathcal{M}^S$  in the blind demixing problem.

*Example 8.5 (Sphere)* On the unit sphere  $S^{n-1}$  which is a Riemannian submanifold of  $\mathbb{R}^n$ , the inner product derived from the Euclidean inner product on  $\mathbb{R}^n$  is represented as

$$\langle \xi, \zeta \rangle_x := \xi^\top \zeta. \quad (8.134)$$

The normal space is

$$(T_x S^{n-1})^\perp = \{x\theta : \theta \in \mathbb{R}\}, \quad (8.135)$$

and the projections are given by  $P_x\xi = (I - xx^\top)\xi$ ,  $P_x^\perp\xi = xx^\top\xi$  for  $x \in S^{n-1}$ .

*Example 8.6 (Product Manifolds in Problem (3.41))* Considering the product manifold  $\mathcal{M}^S$  endowed with the Riemannian metric (8.126), the tangent space  $T_V\mathcal{M}^S$  can be decomposed into two complementary vector spaces, given by Absil et al. [1]

$$T_V\mathcal{M}^S = \mathcal{V}_V\mathcal{M}^S \oplus \mathcal{H}_V\mathcal{M}^S, \quad (8.136)$$

where  $\oplus$  is the direct sum operator. Particularly, the set of vectors which are tangent to the set of equivalence class (8.130) is denoted as the *vertical space*, i.e.,  $\mathcal{V}_V\mathcal{M}^S$ . While the set of vectors which are orthogonal to the equivalence class (8.130) is denoted as the *horizontal space*, i.e.,  $\mathcal{H}_V\mathcal{M}^S$ . Hence the tangent space  $T_{\tilde{V}}(\mathcal{M}^S/\sim)$  at the point  $\tilde{V} \in \mathcal{M}^S/\sim$  can be represented by the horizontal space  $\mathcal{H}_{\tilde{V}}\mathcal{M}^S$  at point  $V \in \mathcal{M}^S$ . Hence, the matrix representation of  $\eta_{\tilde{V}} \in T_{\tilde{V}}(\mathcal{M}^S/\sim)$  [1, Section 3.5.8] can be represented by a unique element  $\eta_V \in \mathcal{H}_V\mathcal{M}^S$ . Additionally, for each

$\xi_V, \eta_V \in T_V \mathcal{M}^S$ , by exploiting the Riemannian metric  $g_V(\xi_V, \eta_V)$  (8.126),

$$g_{\tilde{V}}(\xi_{\tilde{V}}, \eta_{\tilde{V}}) := g_V(\xi_V, \eta_V) \quad (8.137)$$

defines a Riemannian metric on the quotient space  $\mathcal{M}^S / \sim$  [1, Section 3.6.2], where  $\xi_{\tilde{V}}, \eta_{\tilde{V}} \in T_{\tilde{V}} \mathcal{M}^S$ . With the Riemannian metric (8.137), the natural projection  $\pi : \mathcal{M}^S \rightarrow \mathcal{M}^S / \sim$  is a mapping from the quotient manifold  $\mathcal{M}^S / \sim$  to the computational space, which is also called *Riemannian submersion*  $\mathcal{M}^S$  [1, Section 3.6.2]. According to the Riemannian submersion theory, the objects on the quotient manifold can be represented by corresponding objects in the computational space, which facilitates to develop Riemannian optimization algorithm on the Riemannian manifold.

## 8.6 Proof of Theorem 3.4

Based on the notions mentioned in Sect. 3.4.3.2, the Euclidean gradient of  $f(\mathbf{v})$  in problem (3.41) with respect to  $\mathbf{w}_k$  is given by

$$\nabla_{\mathbf{w}_k} f(\mathbf{v}) = 2 \cdot \sum_{i=1}^L \left( c_i \mathbf{J}_{ki} + c_i^* \mathbf{J}_{ki}^H \right) \cdot \mathbf{w}_k, \quad (8.138)$$

where  $c_i = \sum_{k=1}^S [\mathcal{J}_k(\mathbf{w}_k \mathbf{w}_k^H)]_i - y_i$ . According to the definition of the horizontal space given in Table 3.4, it yields that  $\nabla_{\mathbf{w}_k} f(\mathbf{v})$  is in the horizontal space due to  $\nabla_{\mathbf{w}_k} f(\mathbf{v})^H \mathbf{w}_k = \mathbf{w}_k^H \nabla_{\mathbf{w}_k} f(\mathbf{v})$ . Thus, the update rule in the Riemannian gradient descent algorithm, i.e., Algorithm 3.3, can be reformulated as

$$\mathbf{w}_k^{[t+1]} = \mathbf{w}_k^{[t]} - \frac{\alpha_t}{2 \left\| \mathbf{w}_k^{[t]} \right\|_2} \nabla_{\mathbf{w}_k} f(\mathbf{v})|_{\mathbf{w}_k^{[t]}}, \quad (8.139)$$

according to the definition of the Riemannian metric  $g_{\mathbf{w}_k}$  and the retraction  $\mathcal{R}_{\mathbf{w}_k}$  in Table 3.4. The update rule can be reformulated as

$$\begin{bmatrix} \mathbf{w}_k^{[t+1]} \\ \mathbf{w}_k^{[t+1]} \end{bmatrix} = \begin{bmatrix} \mathbf{w}_k^{[t]} \\ \mathbf{w}_k^{[t]} \end{bmatrix} - \frac{\alpha_t}{\left\| \mathbf{w}_k^{[t]} \right\|_2} \begin{bmatrix} \frac{\partial f}{\partial \mathbf{w}_k^H} |_{\mathbf{w}_k^{[t]}} \\ \frac{\partial f}{\partial \mathbf{w}_k^H} |_{\mathbf{w}_k^{[t]}} \end{bmatrix}, \quad (8.140)$$

according to the fact that  $\nabla_{\mathbf{w}_k} f(\mathbf{v}) = 2 \frac{\partial f(\mathbf{v})}{\partial \mathbf{w}_k^H}$ .

The proof of Theorem 3.4 is summarized as follows.

- Lemma 8.9 characterizes the local geometry in the region of incoherence and contraction (RIC) illustrated in Definition 8.3, where the objective function  $f(\mathbf{v})$  (3.41) enjoys restricted strong convexity and smoothness near the ground truth  $\mathbf{v}^\natural$ .

- Based on the property of the local geometry, Lemma 8.10 establishes the error contraction, i.e., convergence analysis.
- Lemma 8.11 demonstrates that the iterates of Algorithm 3.3, including the spectral initialization point, stay within the RIC. This is achieved by exploiting the induction arguments.

**Definition 8.3** ( $(\phi, \beta, \gamma, \mathbf{v}^\natural) - \mathcal{R}$  **the Region of Incoherence and Contraction**) Define  $\mathbf{v}_i = [\mathbf{x}_i^H \mathbf{h}_i^H]^H \in \mathbb{C}^{N+K}$  and  $\mathbf{v} = [\mathbf{v}_1^H \dots \mathbf{v}_s^H]^H \in \mathbb{C}^{s(N+K)}$ . For  $\mathbf{v} \in (\phi, \theta, \gamma, \mathbf{v}^\natural) - \mathcal{R}$ , there exists

$$\text{dist}(\mathbf{v}^t, \mathbf{v}^\natural) \leq \phi, \quad (8.141a)$$

$$\max_{1 \leq i \leq s, 1 \leq j \leq m} \left| \mathbf{c}_{ij}^H (\tilde{\mathbf{x}}_i^t - \mathbf{x}_i^\natural) \right| \cdot \|\mathbf{x}_i^\natural\|_2^{-1} \leq C_2 \theta, \quad (8.141b)$$

$$\max_{1 \leq i \leq s, 1 \leq j \leq m} \left| \mathbf{b}_j^H \tilde{\mathbf{h}}_i^t \right| \cdot \|\mathbf{h}_i^\natural\|_2^{-1} \leq C_3 \gamma, \quad (8.141c)$$

where some constants  $C_2, C_3 > 0$  and some sufficiently small constants  $\phi, \theta, \gamma > 0$ . Additionally,  $\tilde{\mathbf{h}}_i^t$  and  $\tilde{\mathbf{x}}_i^t$  are defined as  $\tilde{\mathbf{h}}_i^t = \frac{1}{\psi_i^t} \mathbf{h}_i^t$  and  $\tilde{\mathbf{x}}_i^t = \psi_i^t \mathbf{x}_i^t$  for  $i = 1, \dots, s$ , with the alignment parameter  $\psi_i^t$ .

The Riemannian Hessian is denoted as  $\text{Hess} f(\mathbf{v}) := \text{diag}(\{\text{Hess}_{\mathbf{w}_i} f\}_{i=1}^s)$ .

**Lemma 8.9** ([7]) *Assuming a sufficiently small constant  $\delta > 0$ . If the number of measurements satisfies  $m \gg \mu^2 s^2 \kappa^2 \max\{N, K\} \log^5 m$ , then with probability exceeding  $1 - \mathcal{O}(m^{-10})$ ,  $\text{Hess} f(\mathbf{v})$  satisfies*

$$\begin{aligned} \mathbf{z}^H [\mathbf{D} \text{Hess} f(\mathbf{v}) + \text{Hess} f(\mathbf{v}) \mathbf{D}] \mathbf{z} &\geq \frac{1}{4\kappa} \|\mathbf{z}\|_2^2 \\ \text{and } \|\text{Hess} f(\mathbf{v})\| &\leq 2 + s \end{aligned} \quad (8.142)$$

simultaneously for all

$$\mathbf{z} = \left[ \mathbf{z}_1^H \dots \mathbf{z}_s^H \right]^H \text{ with } \mathbf{z}_i = \left[ (\mathbf{x}_i - \mathbf{x}_i')^H \ (\mathbf{h}_i - \mathbf{h}_i')^H \ (\mathbf{x}_i - \mathbf{x}_i')^\top \ (\mathbf{h}_i - \mathbf{h}_i')^\top \right]^H,$$

and  $\mathbf{D} = \text{diag}(\{\mathbf{W}_i\}_{i=1}^s)$  with

$$\mathbf{W}_i = \text{diag} \left( [\bar{\beta}_{i1} \mathbf{I}_K \ \bar{\beta}_{i2} \mathbf{I}_N \ \bar{\beta}_{i1} \mathbf{I}_K \ \bar{\beta}_{i2} \mathbf{I}_N]^* \right).$$

Here  $\mathbf{v}$  is in the region  $(\delta, \frac{1}{\sqrt{s} \log^{3/2} m}, \frac{\mu}{\sqrt{m}} \log^2 m, \mathbf{v}^\natural) - \mathcal{R}$ , and one has

$$\max \left\{ \|\mathbf{h}_i - \mathbf{h}_i^\natural\|_2, \|\mathbf{h}_i' - \mathbf{h}_i^\natural\|_2, \|\mathbf{x}_i - \mathbf{x}_i^\natural\|_2, \|\mathbf{x}_i' - \mathbf{x}_i^\natural\|_2 \right\} \leq \delta / (\kappa \sqrt{s}),$$

for  $i = 1, \dots, s$  and  $\mathbf{W}_i$ 's satisfy that for  $\beta_{i1}, \beta_{i2} \in \mathbb{R}$ , for  $i = 1, \dots, s$

$$\max_{1 \leq i \leq s} \max \left\{ \left| \beta_{i1} - \frac{1}{\kappa} \right|, \left| \beta_{i2} - \frac{1}{\kappa} \right| \right\} \leq \frac{\delta}{\kappa \sqrt{s}}.$$

Therein,  $C_2, C_3 \geq 0$  are numerical constants.

**Lemma 8.10 ([7])** Assuming that the step size satisfies  $\alpha_t > 0$  and  $\alpha_t \equiv \alpha \asymp s^{-1}$ , then with probability exceeding  $1 - \mathcal{O}(m^{-10})$ ,

$$\text{dist} \left( \mathbf{v}^{t+1}, \mathbf{v}^\natural \right) \leq \left( 1 - \frac{\alpha}{16\kappa} \right) \text{dist} \left( \mathbf{v}^t, \mathbf{v}^\natural \right), \quad (8.143)$$

provided that the number of measurements follows  $m \gg \mu^2 s^2 \kappa^4 \max \{N, K\} \log^5 m$  and  $\mathbf{v}$  is in the region  $(\delta, \frac{1}{\sqrt{s} \log^{3/2} m}, \frac{\mu}{\sqrt{m}} \log^2 m, \mathbf{v}^\natural) - \mathcal{R}$ , which is denoted as  $\mathcal{R}_{bd}$ .

**Lemma 8.11 ([7])** Assuming the number of measurements

$$m \gg \mu^2 s^2 \kappa^2 \max \{K, N\} \log^6 m,$$

then the spectral initialization point  $\mathbf{v}^0$  is in the region  $\mathcal{R}_{bd}$  with probability exceeding  $1 - \mathcal{O}(m^{-9})$ .

Assuming that  $t$ -th iteration  $\mathbf{v}^t$  is in the region  $\mathcal{R}_{bd}$  and the number of measurements satisfy

$$m \gg \mu^2 s^2 \kappa^2 \max \{K, N\} \log^8 m,$$

then with probability exceeding  $1 - \mathcal{O}(m^{-9})$ , the  $(t+1)$ -th iteration  $\mathbf{v}^{t+1}$  is also in the region  $\mathcal{R}_{bd}$ , which the step size satisfies  $\alpha_t > 0$  and  $\alpha_t \equiv \alpha \asymp s^{-1}$ .

## 8.7 Basic Concepts in Algebraic–Geometric Theory

In this section, we will introduce some basic concepts in algebraic–geometry theory, which contribute to the proof of Theorems 5.5 and 5.6. The content in this section is based on the paper [20].

### 8.7.1 Geometric Characterization of Dimension

Denote polynomials in  $\mathbb{R}[\mathbf{x}] = \mathbb{R}[x_1, \dots, x_n]$  as  $f_1, \dots, f_s$ , their common root position  $\mathcal{V}_{\mathbb{R}^n}(f_1, \dots, f_s)$ , called an *algebraic variety*, is defined as

$$\mathcal{V}_{\mathbb{R}^n}(f_1, \dots, f_s) := \{ \boldsymbol{\zeta} \in \mathbb{R}^n : f_i(\boldsymbol{\zeta}) = 0, \forall i \in [s] \}. \quad (8.144)$$

Considering the dimension of  $\mathcal{V}_{\mathbb{R}^n}(f_1, \dots, f_s)$ , for the case of a single equation, i.e.,  $s = 1$ ,  $\mathcal{V}_{\mathbb{R}^n}(f_1)$  is a hypersurface of  $\mathbb{R}^n$  endowed with dimension  $n - 1$ ; this is similar to the situation where a single linear equation identifies a linear subspace of dimension one less than the ambient dimension. There are other more complicated cases where  $\mathcal{V}_{\mathbb{R}^n}(f_1)$  consists of a single point or no points, which is similar to algebra where a linear subspace has zero dimension only if it contains a single point or the origin point. In these cases, the common root position of the polynomials in the algebraic closure  $\mathbb{C}$  of  $\mathbb{R}$  is considered:

$$\mathcal{V}_{\mathbb{C}^n}(f_1, \dots, f_s) := \{\zeta \in \mathbb{C}^n : f_i(\zeta) = 0, \forall i \in [s]\}. \quad (8.145)$$

The dimension of (8.145)  $\mathcal{V}_{\mathbb{C}^n}(f_1, \dots, f_s) \subset \mathbb{C}^n$  can be characterized by a well-developed theory [9, 11, 14]. The geometric characterization of  $\dim \mathcal{V}_{\mathbb{C}^n}(f_1, \dots, f_s)$  is presented by the following definition.

**Definition 8.4** If  $\mathcal{Y} = \mathcal{V}(g_1, \dots, g_r)$  for some polynomials  $g_1, \dots, g_r \in \mathbb{C}[\mathbf{x}]$ ,  $\mathcal{Y} \subset \mathbb{C}^n$  is defined to be closed. If  $\mathcal{Y} = \mathcal{V}(g_1, \dots, g_r)$  is not the union of two proper closed subsets,  $\mathcal{Y}$  is defined to be irreducible. The dimension of geometric object  $\mathcal{V}_{\mathbb{C}^n}(f_1, \dots, f_s)$  is defined to be the largest non-negative integer  $d$  such that there is a chain of the form

$$\mathcal{V}_{\mathbb{C}^n}(f_1, \dots, f_s) \supset \mathcal{Y}_0 \supsetneq \mathcal{Y}_1 \supsetneq \mathcal{Y}_2 \supsetneq \dots \supsetneq \mathcal{Y}_d, \quad (8.146)$$

where  $\mathcal{Y}_i$  for any  $i \in \{1, \dots, d\}$  is a closed irreducible subset of  $\mathcal{V}_{\mathbb{C}^n}(f_1, \dots, f_s)$ .

Definition 8.4 can be termed as a generalization of the notion of dimension in linear algebra. For instance, considering that  $\mathcal{Y}$  is a linear subspace of  $\mathbb{C}^n$ ,  $\dim \mathcal{Y}$  is the same as the maximal length of a descending chain of linear subspaces. The descending chain can be derived by removing a single basis vector of  $\mathcal{Y}$  at each step. Please refer to Example 8.7 for details.

*Example 8.7* Define a unit vector  $\mathbf{e}_i$  with the value of 1 at position  $i$  zeros and zeros at the rest positions. Then  $\mathcal{Y}_i = \text{Span}(\mathbf{e}_1, \dots, \mathbf{e}_{n-i})$ ,  $\mathbb{C}^n$  admits a chain

$$\mathbb{C}^n = \mathcal{Y}_0 \supsetneq \mathcal{Y}_1 \supsetneq \mathcal{Y}_2 \supsetneq \dots \supsetneq \mathcal{Y}_{n-1} \supsetneq \mathcal{Y}_n := \{0\}. \quad (8.147)$$

Furthermore, the following propositions present several structural property about algebraic varieties.

**Proposition 8.2** Define  $\mathcal{Y} = \mathcal{V}_{\mathbb{C}^n}(f_1, \dots, f_s)$  for some  $f_i \in \mathbb{C}[\mathbf{x}]$ . With irreducible closed sets of  $\mathbb{C}^n$  defined in Definition 8.4, i.e.,  $\mathcal{Y}_i$ ,  $\mathcal{Y}$  can be represented as  $\mathcal{Y} = \mathcal{Y}_1 \cup \dots \cup \mathcal{Y}_d$  for some positive integer  $d$ . The set  $\mathcal{Y}$  is minimal, that is, removing one of the  $\mathcal{Y}_i$  would yield a union that is a strictly smaller set than  $\mathcal{Y}$ . The  $\mathcal{Y}_i$  for  $i \in \{1, \dots, d\}$  are called the irreducible components of  $\mathcal{Y}$ .

Definition 8.4 along with Proposition 8.2 demonstrate that the dimension of  $\mathcal{V}_{\mathbb{C}^n}(f_1, \dots, f_s)$  is zero if and only if the algebraic varieties  $\mathcal{V}_{\mathbb{C}^n}(f_1, \dots, f_s)$  consist of a finite number of points. These varieties are concerned in the paper [20].

**Proposition 8.3** *Define  $\mathcal{Y} = \mathcal{V}_{\mathbb{C}^n}(f_1, \dots, f_s)$ . Then the dimension of  $\mathcal{Y}$  is 0 if and only if  $\mathcal{Y}$  consists of a finite number of points of  $\mathbb{C}^n$ .*

## References

1. Absil, P.A., Mahony, R., Sepulchre, R.: Optimization Algorithms on Matrix Manifolds. Princeton University Press, Princeton (2009)
2. Amelunxen, D., Bürgisser, P.: Intrinsic volumes of symmetric cones and applications in convex programming. *Math. Program.* **149**(1–2), 105–130 (2015)
3. Amelunxen, D., Lotz, M., McCoy, M.B., Tropp, J.A.: Living on the edge: phase transitions in convex programs with random data. *Inf. Inference* **3**(3), 224–294 (2014)
4. Chandrasekaran, V., Recht, B., Parrilo, P.A., Willsky, A.S.: The convex geometry of linear inverse problems. *Found. Comput. Math.* **12**(6), 805–849 (2012)
5. Dong, J., Shi, Y.: Nonconvex demixing from bilinear measurements. *IEEE Trans. Signal Process.* **66**(19), 5152–5166 (2018)
6. Dong, J., Shi, Y.: Blind demixing via Wirtinger flow with random initialization. In: The 22nd International Conference on Artificial Intelligence and Statistics (AISTATS), vol. 89, pp. 362–370 (2019)
7. Dong, J., Yang, K., Shi, Y.: Blind demixing for low-latency communication. *IEEE Trans. Wireless Commun.* **18**(2), 897–911 (2019)
8. Dopico, F.M.: A note on  $\sin \theta$  theorems for singular subspace variations. *BIT Numer. Math.* **40**(2), 395–403 (2000)
9. Eisenbud, D.: Commutative Algebra: With a View Toward Algebraic Geometry, vol. 150. Springer, Berlin (2013)
10. Fan, J., Wang, D., Wang, K., Zhu, Z., et al.: Distributed estimation of principal eigenspaces. *Ann. Stat.* **47**(6), 3009–3031 (2019)
11. Hartshorne, R.: Algebraic Geometry, vol. 52. Springer, Berlin (2013)
12. Ling, S., Strohmer, T.: Regularized gradient descent: a nonconvex recipe for fast joint blind deconvolution and demixing. *Inf. Inference J. IMA* **8**(1), 1–49 (2018)
13. Ma, C., Wang, K., Chi, Y., Chen, Y.: Implicit regularization in nonconvex statistical estimation: gradient descent converges linearly for phase retrieval, matrix completion and blind deconvolution. arXiv preprint:1711.10467 (2017)
14. Matsumura, H.: Commutative Ring Theory, vol. 8. Cambridge University, Cambridge (1989)
15. McCoy, M., Tropp, J.: Sharp recovery bounds for convex demixing, with applications. *Found. Comput. Math.* **14**(3), 503–567 (2014)
16. Mishra, B., Meyer, G., Bonnabel, S., Sepulchre, R.: Fixed-rank matrix factorizations and Riemannian low-rank optimization. *Comput. Stat.* **29**(3–4), 591–621 (2014)
17. Rudelson, M., Vershynin, R.: On sparse reconstruction from Fourier and Gaussian measurements. *Commun. Pure Appl. Math.* **61**(8), 1025–1045 (2008)
18. Schneider, R.: Convex Bodies: The Brunn-Minkowski Theory 151. Cambridge University, Cambridge (2014)
19. Schneider, R., Weil, W.: Stochastic and Integral Geometry. Springer, Berlin (2008)
20. Tsakiris, M.C., Peng, L., Conca, A., Kneip, L., Shi, Y., Choi, H.: An algebraic-geometric approach to shuffled linear regression. arXiv preprint:1810.05440 (2018)
21. Vershynin, R.: Introduction to the non-asymptotic analysis of random matrices. In: Compressed Sensing, Theory and Applications, pp. 210–268 (2010)

22. Yatawatta, S.: On the interpolation of calibration solutions obtained in radio interferometry. *Mon. Not. R. Astron. Soc.* **428**(1), 828–833 (2013)
23. Yatawatta, S.: Radio interferometric calibration using a Riemannian manifold. In: 2013 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 3866–3870. IEEE, Piscataway (2013)