

# A Nonconvex Approach for Phase Retrieval: Reshaped Wirtinger Flow and Incremental Algorithms

**Huishuai Zhang**

*Department of EECS  
Syracuse University  
Syracuse, NY 13244 USA*

HZHAN23@SYR.EDU

**Yi Zhou**

**Yingbin Liang**

**Yuejie Chi**

*Department of ECE  
The Ohio State University  
Columbus, OH 43210 USA*

ZHOU.1172@OSU.EDU

LIANG.889@OSU.EDU

CHI.97@OSU.EDU

**Editor:** Benjamin Recht

## Abstract

We study the problem of solving a quadratic system of equations, i.e., recovering a vector signal  $\mathbf{x} \in \mathbb{R}^n$  from its magnitude measurements  $y_i = |\langle \mathbf{a}_i, \mathbf{x} \rangle|$ ,  $i = 1, \dots, m$ . We develop a gradient descent algorithm (referred to as RWF for reshaped Wirtinger flow) by minimizing the quadratic loss of the magnitude measurements. Comparing with Wirtinger flow (WF) (Candès et al., 2015), the loss function of RWF is nonconvex and nonsmooth, but better resembles the least-squares loss when the phase information is also available. We show that for random Gaussian measurements, RWF enjoys linear convergence to the true signal as long as the number of measurements is  $\mathcal{O}(n)$ . This improves the sample complexity of WF ( $\mathcal{O}(n \log n)$ ), and achieves the same sample complexity as truncated Wirtinger flow (TWF) (Chen and Candès, 2015), but without any sophisticated truncation in the gradient loop. Furthermore, RWF costs less computationally than WF, and runs faster numerically than both WF and TWF. We further develop an incremental (stochastic) version of RWF (IRWF) and connect it with the randomized Kaczmarz method for phase retrieval. We demonstrate that IRWF outperforms existing incremental as well as batch algorithms with experiments.

**Keywords:** gradient descent, phase retrieval, nonconvex optimization, regularity condition, stochastic algorithms

## 1. Introduction

Many problems in machine learning and signal processing can be reduced to solving a quadratic system of equations. For instance, in phase retrieval applications, i.e., X-ray crystallography and coherent diffraction imaging (Drenth, 2007; Miao et al., 1999, 2008), the structure of an object is to be recovered from the measured far field diffracted intensity when the object is illuminated by a source light. Mathematically, such a problem amounts to recovering the signal from only the magnitudes of its linear measurements. Specifically, the problem is formulated as follows.

**Problem 1** Recover  $\mathbf{x} \in \mathbb{R}^n/\mathbb{C}^n$  from the measurements given as

$$y_i = |\langle \mathbf{a}_i, \mathbf{x} \rangle|, \quad \text{for } i = 1, \dots, m, \quad (1)$$

where  $\mathbf{a}_i \in \mathbb{R}^n/\mathbb{C}^n$  are known design vectors.

Various algorithms have been proposed to solve this problem since the 1970s. The error-reduction methods (Gerchberg, 1972; Fienup, 1982) work well empirically but lack theoretical guarantees. More recently, convex relaxation of the problem has been formulated, for example, via PhaseLift (Candès et al., 2013; Candès and Li, 2014; Chen et al., 2015) and PhaseCut (Waldspurger et al., 2015), and the correspondingly developed algorithms typically come with performance guarantees. The readers can refer to the review paper (Shechtman et al., 2015) to learn more about applications and algorithms of the phase retrieval problem.

Nonetheless, these convex approaches often suffer from high computational complexity particularly when the signal dimension is large. It is therefore desirable to develop more efficient nonconvex approaches that can provably recover the true signal. Netrapalli et al. (2013) proposed the *AltMinPhase* algorithm, which alternates between updates of the phase and the signal, and showed that it converges linearly and recovers the true signal with  $\mathcal{O}(n \log^3 n)$  Gaussian measurements, when  $\mathbf{a}_i$ 's are composed of *independent and identically distributed* (i.i.d.) standard Gaussian entries. More recently, Candès et al. (2015) introduced the *Wirtinger flow* (WF) algorithm, which guarantees signal recovery via a simple gradient descent algorithm with only  $\mathcal{O}(n \log n)$  Gaussian measurements and attains  $\epsilon$ -accuracy within  $\mathcal{O}(mn^2 \log 1/\epsilon)$  flops<sup>1</sup>. More specifically, WF obtains a good initialization by the spectral method, and then minimizes the following nonconvex loss function based on the quadratic loss of the squared magnitude measurements

$$\ell_{WF}(\mathbf{z}) := \frac{1}{4m} \sum_{i=1}^m (|\mathbf{a}_i^T \mathbf{z}|^2 - |y_i|^2)^2, \quad (2)$$

via the gradient descent scheme.

WF is further improved by *truncated Wirtinger flow* (TWF) (Chen and Candès, 2015), which adopts a Poisson loss function of  $|\mathbf{a}_i^T \mathbf{z}|^2$ , and keeps only well-behaved measurements based on carefully designed truncation thresholds for calculating the initialization and the gradient updates. Such truncation assists to achieve a linear convergence rate using a fixed step size; as a consequence, TWF reduces both the sample complexity to  $\mathcal{O}(n)$  and the computational complexity to  $\mathcal{O}(mn \log 1/\epsilon)$ .

Furthermore, incremental/stochastic methods have been proposed to solve Problem 1. Specifically, the randomized Kaczmarz method has been adopted to solve the phase retrieval problem (Wei, 2015; Li et al., 2015), which was shown to have excellent empirical performance, but no nonasymptotic global convergence guarantee was established. Incremental truncated Wirtinger flow (ITWF) (Kolte and Özgür, 2016) is another stochastic algorithm developed based on TWF.

---

1. WF has computational cost  $\mathcal{O}(mn^2 \log(1/\epsilon))$  (Candès et al., 2015, Theorem 3.1) because it needs  $\log(1/\epsilon) \cdot \frac{1}{-\log(1-c/n)} \approx \frac{n}{c} \log(1/\epsilon)$  iterations to achieve  $\epsilon$ -accuracy due to the convergence rate  $(1 - \frac{c}{n})^k$  and each iteration needs  $\mathcal{O}(mn)$  flops. Contrastingly, RWF has computational cost  $\mathcal{O}(mn \log 1/\epsilon)$  because it requires only  $c \log(1/\epsilon)$  iterations to achieve  $\epsilon$ -accuracy due to the convergence rate  $(1 - c)^k$ .

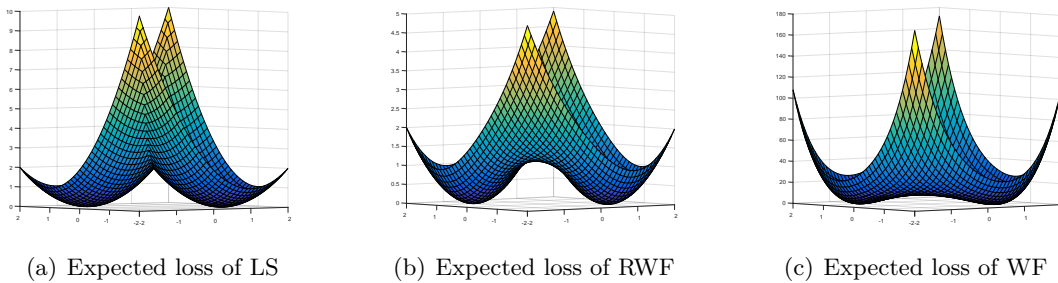


Figure 1: Surface of the expected loss function of (a) least-squares (mirrored symmetrically), (b) RWF, and (c) WF when  $\mathbf{x} = [1, -1]^T$ .

### 1.1 Our Contribution

This paper adopts the following loss function

$$\ell(\mathbf{z}) := \frac{1}{2m} \sum_{i=1}^m \left( |\mathbf{a}_i^T \mathbf{z}| - y_i \right)^2, \quad (3)$$

which is the quadratic loss of the magnitude measurements. Compared to the loss function (2) of WF that adopts the quadratic loss of  $|\mathbf{a}_i^T \mathbf{z}|^2$ , the above loss function adopts the quadratic loss of  $|\mathbf{a}_i^T \mathbf{z}|$  and hence has lower-order variables. While both loss functions are nonconvex in  $\mathbf{z}$ , the WF loss function (2) is a fourth-order smooth function of  $\mathbf{a}_i^T \mathbf{z}$  and our loss function (3) in contrast is nonsmooth.

To minimize such a nonconvex and nonsmooth loss function (3), we develop a gradient descent algorithm, which sets the “gradient” to zero corresponding to nonsmooth samples. We refer to such an algorithm together with an initialization using a new spectral method (different from that employed in TWF or WF) as *reshaped Wirtinger flow* (RWF). We show that the new loss function has great advantage in both statistical and computational efficiency, despite nonsmoothness. In fact, the curvature of such a loss function behaves similarly to that of a least-squares problem with phase information in the neighborhood of global optimizers, and hence yields faster convergence. To provide further insights, consider the standard problem of solving  $\mathbf{x}$  from *linear measurements*  $\langle \mathbf{a}_i, \mathbf{x} \rangle$ ,  $i = 1, \dots, m$ , where  $\mathbf{a}_i$ ’s are composed of i.i.d. standard Gaussian entries. In this case, it is natural to use the least-squares loss function

$$\ell_{LS}(\mathbf{z}) := \frac{1}{m} \sum_{i=1}^m \left( \mathbf{a}_i^T \mathbf{z} - \mathbf{a}_i^T \mathbf{x} \right)^2. \quad (4)$$

Examining the expected (with respect to  $\mathbf{a}_i$ ’s) loss surface of  $\min\{\ell_{LS}(\mathbf{z}), \ell_{LS}(-\mathbf{z})\}$  (to mimic sign ambiguity),  $\ell(\mathbf{z})$ , and  $\ell_{WF}(\mathbf{z})$  in Figure 1, whose expressions can be found in Appendix A, it can be seen that the loss of RWF, rather than the loss of WF, has a similar curvature to the quadratic least-squares loss around the global optimizers, which justifies its better performance than WF.

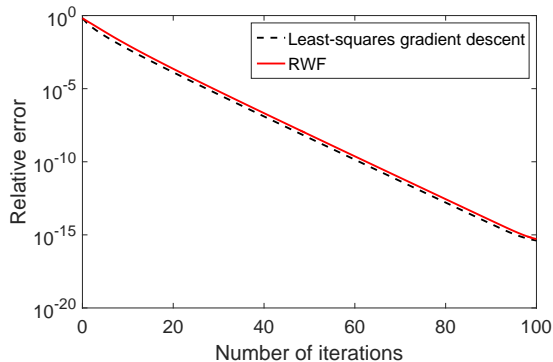


Figure 2: Comparison of convergence behavior between RWF and the least-squares gradient descent with the same initialization, the same parameters  $n = 1000$ ,  $m = 6n$ , and the same step size  $\mu = 0.8$ .

The nonsmoothness of the loss function (3) does not negatively impact the performance of RWF because only with negligible probability the algorithm encounters nonsmooth points for some samples, which furthermore are set not to contribute to the gradient direction by RWF. The gradient of the RWF loss (3) is given as

$$\nabla \ell(\mathbf{z}) := \frac{1}{m} \sum_{i=1}^m \left( \mathbf{a}_i^T \mathbf{z} - y_i \cdot \text{sgn}(\mathbf{a}_i^T \mathbf{z}) \right) \mathbf{a}_i, \quad (5)$$

where  $\text{sgn}(0) = 0$  by convention. Comparing this with the gradient of the least-squares loss

$$\nabla \ell_{LS}(\mathbf{z}) = \frac{1}{m} \sum_{i=1}^m \left( \mathbf{a}_i^T \mathbf{z} - \mathbf{a}_i^T \mathbf{x} \right) \mathbf{a}_i, \quad (6)$$

one can see that RWF uses estimated phase information  $\text{sgn}(\mathbf{a}_i^T \mathbf{z})$  to generate the gradient updates, and the convergence behavior of RWF is much similar to that of least-squares with phase information if initialized properly. Indeed, Figure 2 illustrates that RWF takes almost the *same* number of iterations for recovering a signal (with only the magnitude information) as the least-squares gradient descent method for recovering a signal (with both the magnitude and the sign information).

We further develop incremental/stochastic versions of RWF using mini-batches of measurements, called incremental RWF (IRWF), and show that IRWF also enjoys the advantageous local curvature of RWF, and achieves excellent statistical and computational performance. Along the way, we establish the performance guarantee of the Kaczmarz-PR update rule by interpreting it as a variant of IRWF. We conduct extensive numerical experiments to demonstrate that IRWF performs better than other competitive incremental algorithms (ITWF and Kaczmarz-PR) as well as batch algorithms (RWF, TWF, WF and AltMinPhase).

We summarize our main results as follows.

- Statistically, we show that RWF recovers the true signal with  $\mathcal{O}(n)$  Gaussian measurements, which is order-wise optimal. Thus, RWF improves the sample complexity

$\mathcal{O}(n \log n)$  of WF, and achieves the same sample complexity as TWF but without truncation in the gradient descent step.

- Computationally, RWF converges linearly to the true signal with a constant step size, requiring  $\mathcal{O}(mn \log(1/\epsilon))$  flops to reach  $\epsilon$ -accuracy. Again, without truncation in the gradient descent step, RWF improves the computational cost  $\mathcal{O}(mn^2 \log(1/\epsilon))$  of WF and achieves the same computational cost as TWF.
- We also show that RWF is robust to bounded additive noise. The estimation error is shown to diminish at a linear rate with respect to the noise level up to a scaling difference. Experiments on Poisson noise further corroborate the stability guarantee.
- The incremental versions of RWF (IRWF) is shown to shrink the estimation error in expectation for one update, which implies the linear convergence of the algorithm if the iteration path lies in the neighborhood of the true signal. Furthermore, we show that Kaczmarz-PR can be viewed as IRWF under a specific choice of step size, and establish a similar convergence result for Kaczmarz-PR under the same sample complexity.
- Numerically, RWF and its incremental versions require fewer parameters, e.g., truncation thresholds, than TWF in practice. RWF is generally two times faster than TWF and four to six times faster than WF in terms of both the number of iterations and time cost. IRWF also outperforms existing incremental as well as batch algorithms.

Compared to WF and TWF, the nonsmoothness in the loss function requires new bounding techniques in order to establish the theoretical guarantees. On the other hand, our technical proof is much simpler, because the lower-order loss function allows to bypass higher-order moments of variables as well as truncation of samples in the gradient descent steps.

## 1.2 Related Work

The quadratic loss function of magnitudes (3) was also used in the early literature of phase retrieval (Fienup, 1982) with Fourier magnitude measurements. However, no global convergence guarantee was available in (Fienup, 1982). Along the line of developing nonconvex algorithms with global performance guarantees for the phase retrieval problem, Netrapalli et al. (2013) and Waldspurger (2016) developed alternating minimization algorithms, Candès et al. (2015); Chen and Candès (2015); Zhang et al. (2016); Cai et al. (2016) developed first-order algorithms, and a recent study Sun et al. (2016) characterized the geometric structure of the nonconvex objective and designed a second-order trust-region algorithm. This paper is closely related to (Candès et al., 2015; Chen and Candès, 2015; Zhang et al., 2016), but develops a new gradient descent algorithm based on a lower-order nonsmooth (as well as nonconvex) loss function that yields advantageous statistical and computational efficiency.

Stochastic algorithms were also developed for the phase retrieval problem. Kolte and Özgür (2016) studied the incremental truncated Wirtinger flow (ITWF) and showed that ITWF needs much fewer passes of data than TWF to reach the same accuracy. Wei (2015) adapted the Kaczmarz method to solve the phase retrieval problem and demonstrated its

fast empirical convergence. We show that IRWF is closely related to Kaczmarz-PR, and empirically runs faster than ITWF thanks to the advantageous curvature of the loss function.

After our work was posted on arXiv, an independent work (Wang et al., 2016) was subsequently posted, which also adopts the same loss function but develops a slightly different algorithm called TAF (i.e., truncated amplitude flow). One major difference is that RWF does not require truncation in the gradient loops while TAF still employs truncation. Hence, RWF has fewer parameters to tune, and is easier to implement than TAF in practice. Furthermore, RWF demonstrates the performance advantage of adopting a lower-order loss function even without truncation, which cannot be observed from TAF. The two algorithms also employ different initialization strategies. Moreover, we analyze stochastic algorithms based on the new loss function while Wang et al. (2016) do not.

More generally, various high-dimensional signal estimation problems have been studied by minimizing nonconvex loss functions. For example, a partial list of these studies include matrix completion (Keshavan et al., 2010; Jain et al., 2013; Sun and Luo, 2016; Hardt, 2014; Sa et al., 2015; Zheng and Lafferty, 2016; Jin et al., 2016; Ge et al., 2016), low-rank matrix recovery (Bhojanapalli et al., 2016; Chen and Wainwright, 2015; Tu et al., 2015; Zheng and Lafferty, 2015; Park et al., 2016; Wei et al., 2016; Li et al., 2017), robust PCA (Netrapalli et al., 2014), robust tensor decomposition (Anandkumar et al., 2016), dictionary learning (Arora et al., 2015; Sun et al., 2015), community detection (Bandeira et al., 2016), phase synchronization (Boumal, 2016), blind deconvolution (Lee et al., 2016b; Li et al., 2016), etc.

For minimizing a general nonconvex nonsmooth objective, various algorithms have been proposed, such as gradient sampling (Burke et al., 2005; Kiwiel, 2007) and majorization-minimization (Ochs et al., 2015). These algorithms are often shown to converge to critical points which may be local minimizers or saddle points, without an explicit characterization of convergence rates. In contrast, our algorithm is specifically designed for the phase retrieval problem, and can be shown to converge linearly to the global optimum under an appropriate initialization.

The advantage of using nonsmooth loss functions in our study is analogous in spirit to that of the rectifier activation function (of the form  $\max\{0, \cdot\}$ ) in neural networks. It has been shown that rectified linear unit (*ReLU*) enjoys superb advantage in reducing the training time (Krizhevsky et al., 2012) and promoting sparsity (Glorot et al., 2011) over its counterparts of sigmoid and hyperbolic tangent functions, in spite of non-linearity and non-differentiability at zero. Our results in fact also demonstrate that a nonsmooth but simpler loss function yields improved performance.

### 1.3 Paper Organization and Notations

The rest of this paper is organized as follows. Section 2 describes the RWF algorithm in detail and establishes its performance guarantee. Section 3 introduces the IRWF algorithm, establishes its performance guarantee and compares it with existing stochastic algorithms. Section 4 compares RWF and IRWF with other competitive algorithms numerically. Finally, Section 5 concludes the paper with comments on future directions.

Throughout the paper, boldface lowercase letters such as  $\mathbf{a}_i, \mathbf{x}, \mathbf{z}$  denote vectors, and boldface capital letters such as  $\mathbf{A}, \mathbf{Y}$  denote matrices. For two matrices,  $\mathbf{A} \prec \mathbf{B}$  means that  $\mathbf{B} - \mathbf{A}$  is positive definite. For a complex matrix and a vector,  $\mathbf{A}^*$  and  $\mathbf{z}^*$  denote conjugate

transposes of  $\mathbf{A}$  and  $\mathbf{z}$ , respectively. For a real matrix and a vector,  $\mathbf{A}^T$  and  $\mathbf{z}^T$  denote transposes of  $\mathbf{A}$  and  $\mathbf{z}$ , respectively. We let  $\odot$  denotes element-wise product. The indicator function  $\mathbf{1}_A = 1$  if the event  $A$  occurs, and  $\mathbf{1}_A = 0$  otherwise. We let  $\|\mathbf{x}\|_1$  and  $\|\mathbf{x}\|$  denote the  $l_1$  norm and  $l_2$  norm of a vector  $\mathbf{x}$ , respectively. Moreover, let  $\|\mathbf{A}\|_F$  and  $\|\mathbf{A}\|$  denote the Frobenius norm and the spectral norm of a matrix  $\mathbf{A}$ , respectively. We note that the constants  $c, C, c_0, c_1, c_2$  may be different in different equations, for the sake of notational simplicity.

## 2. Reshaped Wirtinger Flow

Consider the Problem 1 for the complex case. It can be observed that if  $\mathbf{z}$  is a solution, i.e., satisfying Equation (1), then  $\mathbf{z}e^{-j\phi}$  is also the solution of the problem where  $\phi$  is an arbitrary phase constant. Therefore, the recovery is up to a phase difference. Thus, we define the Euclidean distance between two complex vectors up to a global phase difference (Candès et al., 2015) as,

$$\text{dist}(\mathbf{z}, \mathbf{x}) := \min_{\phi \in [0, 2\pi)} \|\mathbf{z}e^{-j\phi} - \mathbf{x}\|, \quad (7)$$

which is simply  $\min \|\mathbf{z} \pm \mathbf{x}\|$  for the real case.

In this paper, we focus on the real-valued case in our analysis, but the algorithm designed below is applicable to the complex-valued case and the case with other measurement vectors such as *coded diffraction pattern* (CDP) as we demonstrate via numerical experiments in Section 4.

We design RWF (see Algorithm 1) for solving the phase retrieval problem, which contains two stages: spectral initialization and gradient loop. The suggested values for the parameters<sup>2</sup> are given by  $\alpha_l = 1, \alpha_u = 5$  and  $\mu = 0.8$ . The rescaling coefficient in  $\lambda_0$  and the conjugate transpose  $\mathbf{a}_i^*$  allow the algorithm readily applicable to the complex and CDP cases. We next describe the two stages of RWF in details in Sections 2.1 and 2.2, respectively, and establish the convergence guarantee in Section 2.3. Finally, we provide the stability guarantee of RWF in Section 2.4.

### 2.1 Initialization via the Spectral Method

When solving nonconvex problems by iterative algorithms, the starting point is critical. The spectral method is a popular choice in the literature (Keshavan et al., 2010; Netrapalli et al., 2013; Candès et al., 2015; Chen and Candès, 2015), which often provides a good initialization. Different from the spectral initialization used in AltMinPhase (Netrapalli et al., 2013), WF (Candès et al., 2015) and TWF (Chen and Candès, 2015), which are based on the squared magnitudes as the weight of each rank-one matrix  $\mathbf{a}_i\mathbf{a}_i^*$ , we propose an alternative initialization in Algorithm 1 that uses the magnitudes instead, and truncates samples that are either too large or too small. We show that such an initialization achieves a smaller sample complexity than WF and the same sample complexity as TWF order-wise, and furthermore, performs better than both WF and TWF numerically.

Our initialization consists of estimation of both the norm and the direction of  $\mathbf{x}$ . The norm estimation of  $\mathbf{x}$  is given by  $\lambda_0$  in Algorithm 1. Intuitively, with real-valued Gaussian

---

2. For the complex Gaussian case, we suggest  $\mu = 1.2$ .

---

**Algorithm 1** Reshaped Wirtinger Flow
 

---

**Input:**  $\mathbf{y} = \{y_i\}_{i=1}^m$ ,  $\{\mathbf{a}_i\}_{i=1}^m$ ;

**Parameters:** Lower and upper thresholds  $\alpha_l, \alpha_u$  for truncation in initialization, step size  $\mu$ ;

**Initialization:** Let  $\mathbf{z}^{(0)} = \lambda_0 \tilde{\mathbf{z}}$ , where  $\lambda_0 = \frac{mn}{\sum_{i=1}^m \|\mathbf{a}_i\|_1} \cdot \left(\frac{1}{m} \sum_{i=1}^m y_i\right)$  and  $\tilde{\mathbf{z}}$  is the leading eigenvector of

$$\mathbf{Y} := \frac{1}{m} \sum_{i=1}^m y_i \mathbf{a}_i \mathbf{a}_i^* \mathbf{1}_{\{\alpha_l \lambda_0 < y_i < \alpha_u \lambda_0\}}. \quad (8)$$

**Gradient loop:** for  $t = 0 : T - 1$  do

$$\mathbf{z}^{(t+1)} = \mathbf{z}^{(t)} - \frac{\mu}{m} \sum_{i=1}^m \left( \mathbf{a}_i^* \mathbf{z}^{(t)} - y_i \cdot \frac{\mathbf{a}_i^* \mathbf{z}^{(t)}}{|\mathbf{a}_i^* \mathbf{z}^{(t)}|} \right) \mathbf{a}_i. \quad (9)$$

**Output**  $\mathbf{z}^{(T)}$ .
 

---

measurement vectors (i.e.,  $\mathbf{a}_i \sim \mathcal{N}(0, \mathbf{I}_{n \times n})$ ), the scaling coefficient  $\frac{mn}{\sum_{i=1}^m \|\mathbf{a}_i\|_1} \approx \sqrt{\frac{\pi}{2}}$ . Moreover, for  $i = 1, \dots, m$ ,  $y_i = |\mathbf{a}_i^T \mathbf{x}|$  are independent sub-Gaussian random variables with mean  $\sqrt{\frac{2}{\pi}} \|\mathbf{x}\|$ , and thus  $\frac{1}{m} \sum_{i=1}^m y_i \approx \sqrt{\frac{2}{\pi}} \|\mathbf{x}\|$ . Combining these two facts yields the desired argument.

The direction of  $\mathbf{x}$  is approximated by the leading eigenvector of  $\mathbf{Y}$ , because  $\mathbf{Y}$  approaches  $\mathbb{E}[\mathbf{Y}]$  by concentration of measure arguments, and the leading eigenvector of  $\mathbb{E}[\mathbf{Y}]$  takes the form  $c\mathbf{x}$  for some scalar  $c \in \mathbb{R}$ . We note that (8) involves truncation of samples from both sides, in contrast to truncation only by an upper threshold in TWF (Chen and Candès, 2015). The truncation parameters  $\alpha_l$  and  $\alpha_u$  are related to the eigenvalue gap between the top two eigenvalues of  $\mathbf{Y}$ , which determines the estimation accuracy of the eigenvector after running  $k$  iterations of the power method<sup>3</sup>. The larger the eigenvalue gap, the better the accuracy of the power method given a number of iterations. We note that  $\alpha_l = 1, \alpha_u = \infty$  yield the largest eigenvalue gap from the developments in Appendix B. We explicitly set a bounded  $\alpha_u$  for the development of the proof and for numerical stability.

We next provide the formal statement that with high probability, the proposed initialization lands in a small neighborhood around the true signal.

**Proposition 1** Fix  $\delta > 0, \alpha_l = 1$  and  $\alpha_u \gg 1$ . The initialization step in Algorithm 1 yields  $\mathbf{z}^{(0)}$  satisfying  $\text{dist}(\mathbf{z}^{(0)}, \mathbf{x}) \leq \delta \|\mathbf{x}\|$  with the probability at least  $1 - \exp(-c m \epsilon^2)$ , if  $m > C(\delta, \epsilon)n$ , where  $c$  is some positive constant and  $C$  is a positive constant affected by  $\delta$  and  $\epsilon$ .

**Proof.** See Appendix B. ■

Finally, we numerically compare different initialization methods in Figure 3. It is demonstrated that RWF achieves better initialization accuracy in terms of the relative error  $\text{dist}(\mathbf{z}^{(0)}, \mathbf{x}) / \|\mathbf{x}\|$  than WF, TWF, as well as the initialization method proposed in TAF (Wang et al., 2016).

---

3. Numerical linear algebra shows that after  $k$  iterations of the power method, the estimation accuracy of the eigenvector is given by  $\mathcal{O}((\lambda_2/\lambda_1)^k)$ , where  $\lambda_1$  and  $\lambda_2$  are the top two eigenvalues.



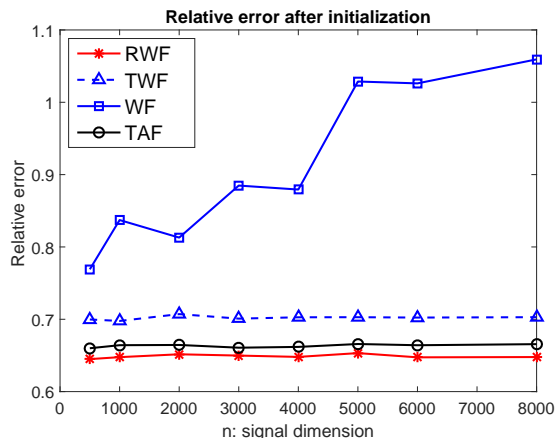


Figure 3: Comparison of different initialization methods with  $m = 6n$  and under 50 iterations of power method.

## 2.2 Gradient Loop

The gradient loop of Algorithm 1 is based on the loss function (3), and the update rule (9) makes itself suitable for the complex case. For the real case the update direction is given as follows:

$$\nabla \ell(\mathbf{z}) := \frac{1}{m} \sum_{i=1}^m \left( \mathbf{a}_i^T \mathbf{z} - y_i \cdot \text{sgn}(\mathbf{a}_i^T \mathbf{z}) \right) \mathbf{a}_i = \frac{1}{m} \sum_{i=1}^m \left( \mathbf{a}_i^T \mathbf{z} - y_i \cdot \frac{\mathbf{a}_i^T \mathbf{z}}{|\mathbf{a}_i^T \mathbf{z}|} \right) \mathbf{a}_i, \quad (10)$$

where  $\text{sgn}(\cdot)$  is the sign function for nonzero arguments. We further set  $\text{sgn}(0) = 0$  and  $\frac{0}{|0|} = 0$ . In fact,  $\nabla \ell(\mathbf{z})$  equals the gradient of the loss function (3) over the samples that satisfy  $\mathbf{a}_i^T \mathbf{z} \neq 0$ . For samples corresponding to the nonsmooth point, i.e.,  $\mathbf{a}_i^T \mathbf{z} = 0$ , we adopt the Fréchet superdifferential (Kruger, 2003) for nonconvex functions to set the gradient component to be zero (as zero is an element in the Fréchet superdifferential). With abuse of terminology, we still refer to  $\nabla \ell(\mathbf{z})$  in Equation (10) as the “gradient” for simplicity, which rather represents the update direction in the gradient loop of Algorithm 1.

## 2.3 Linear Convergence of RWF

We characterize the convergence guarantee of RWF in the following theorem.

**Theorem 2** *Consider the problem of solving any given  $\mathbf{x} \in \mathbb{R}^n$  from a system of equations (1) with Gaussian measurement vectors. There exist some universal constants  $\mu_0 > 0$  ( $\mu_0$  can be set as 0.8 in practice),  $0 < \rho, \nu < 1$  and  $c_0, c_1, c_2 > 0$  such that if  $m \geq c_0 n$  and  $\mu < \mu_0$ , then with probability at least  $1 - c_1 \exp(-c_2 m)$ , Algorithm 1 yields*

$$\text{dist}(\mathbf{z}^{(t)}, \mathbf{x}) \leq \nu(1 - \rho)^t \|\mathbf{x}\|, \quad \forall t \in \mathbb{N}. \quad (11)$$

**Proof** We outline the proof here with the details delegated to Appendix C. Compared to WF and TWF, our proof is much simpler due to the lower-order loss function that RWF

relies on. We first introduce a global phase notation for the real case as follows:

$$\Phi_{\mathbf{x}}(\mathbf{z}) := \begin{cases} 0, & \text{if } \|\mathbf{z} - \mathbf{x}\| \leq \|\mathbf{z} + \mathbf{x}\|, \\ \pi, & \text{otherwise.} \end{cases} \quad (12)$$

For the sake of simplicity, we let  $\mathbf{z}$  be  $e^{-j\Phi_{\mathbf{x}}(\mathbf{z})}\mathbf{z}$ , which indicates that  $\mathbf{z}$  is always in the neighborhood of  $\mathbf{x}$ .

Here, the central idea is to show that within the neighborhood of the global minimizer, RWF satisfies the *Regularity Condition*  $\text{RC}(\mu, \lambda, c)$  (Chen and Candès, 2015), i.e.,

$$\langle \nabla \ell(\mathbf{z}), \mathbf{z} - \mathbf{x} \rangle \geq \frac{\mu}{2} \|\nabla \ell(\mathbf{z})\|^2 + \frac{\lambda}{2} \|\mathbf{z} - \mathbf{x}\|^2 \quad (13)$$

for all  $\mathbf{z}$  obeying  $\|\mathbf{z} - \mathbf{x}\| \leq c\|\mathbf{x}\|$ , where  $0 < c < 1$  is some small constant. Then, as shown in (Candès et al., 2015; Chen and Candès, 2015), once the initialization lands into this neighborhood, linear convergence can be guaranteed with a proper choice of the constant step size, i.e.,

$$\text{dist}^2(\mathbf{z} - \mu \nabla \ell(\mathbf{z}), \mathbf{x}) \leq (1 - \mu\lambda) \text{dist}^2(\mathbf{z}, \mathbf{x}), \quad (14)$$

for any  $\mathbf{z}$  satisfying  $\|\mathbf{z} - \mathbf{x}\| \leq c\|\mathbf{x}\|$ .

Lemmas 6 and 7 in Appendix C yield that for all  $\mathbf{z}$  satisfying  $\|\mathbf{z} - \mathbf{x}\| \leq c\|\mathbf{x}\|$ ,

$$\langle \nabla \ell(\mathbf{z}), \mathbf{z} - \mathbf{x} \rangle \geq (1 - 0.26 - 3\epsilon) \|\mathbf{z} - \mathbf{x}\|^2 = (0.74 - 3\epsilon) \|\mathbf{z} - \mathbf{x}\|^2$$

with high probability. Moreover, Lemma 8 in Appendix C further yields that

$$\|\nabla \ell(\mathbf{z})\| \leq (1 + \delta) \cdot 2\|\mathbf{z} - \mathbf{x}\| \quad (15)$$

with high probability. Therefore, the above two bounds imply that the Regularity Condition (13) holds for  $\mu$  and  $\lambda$  satisfying

$$0.74 - 3\epsilon \geq \frac{\mu}{2} \cdot 4(1 + \delta)^2 + \frac{\lambda}{2} \quad (16)$$

for sufficiently small  $\epsilon$  and  $\delta$ . ■

We note that Equation (16) implies an upper bound for the step size  $\mu \leq \frac{0.74}{2} = 0.37$ , by taking  $\epsilon$  and  $\delta$  to be sufficiently small. However, in practice, the step size  $\mu$  can be set much larger than such a bound, say 0.8, while still keeping the algorithm convergent. This is because the coefficients in the proof are set for the convenience of the proof rather than being tightly chosen.

Theorem 2 indicates that RWF recovers the true signal with  $\mathcal{O}(n)$  samples, which is order-wise optimal. Such an algorithm improves the sample complexity  $\mathcal{O}(n \log n)$  of WF. Furthermore, RWF does not require truncation of samples in the gradient step to achieve the same sample complexity as TWF. This is mainly because RWF benefits from the lower-order loss function given in Equation (3), the curvature of which behaves similarly to the least-squares loss function locally as we explain in the introduction.

Theorem 2 also suggests that RWF converges linearly with a constant step size. To reach  $\epsilon$ -accuracy, it requires a computational cost of  $\mathcal{O}(mn \log(1/\epsilon))$  flops, which is better than WF ( $\mathcal{O}(mn^2 \log(1/\epsilon))$ ) and on par with TWF. Numerically, as we demonstrate in Section 4, RWF is two times faster than TWF and four to six times faster than WF in terms of both the iteration counts and the time cost in various examples.

## 2.4 Stability to Bounded Noise

We have established that RWF guarantees exact recovery at a linear convergence rate for noise-free measurements. We now study RWF in the presence of noise. Suppose the measurements are corrupted by bounded noise, and are given by

$$y_i = |\mathbf{a}_i^T \mathbf{x}| + w_i, \quad 1 \leq i \leq m, \quad (17)$$

where  $\mathbf{w} = \{w_i\}_{i=1}^m$  denote the additive noise. Then the following theorem shows that RWF is robust under bounded noise.

**Theorem 3** *Consider the model (17). Suppose that the measurement vectors are independently Gaussian, i.e.,  $\mathbf{a}_i \sim \mathcal{N}(0, \mathbf{I})$  for  $1 \leq i \leq m$ , and the noise is bounded, i.e.,  $\|\mathbf{w}\|/\sqrt{m} \leq c\|\mathbf{x}\|$  where  $c$  is a positive constant. Then there exist some universal constants  $\mu_0 > 0$  ( $\mu_0$  can be set as 0.8 in practice),  $0 < \rho < 1$  and  $c_0, c_1, c_2 > 0$  such that if  $m \geq c_0 n$  and  $\mu < \mu_0$ , then with probability at least  $1 - c_1 \exp(-c_2 m)$ , Algorithm 1 yields*

$$\text{dist}(\mathbf{z}^{(t)}, \mathbf{x}) \lesssim \frac{\|\mathbf{w}\|}{\sqrt{m}} + (1 - \rho)^t \|\mathbf{x}\|, \quad \forall t \in \mathbb{N}. \quad (18)$$

**Proof** See Appendix D. ■

Theorem 3 shows that under the same sample complexity, RWF converges at a linear rate to a neighborhood around the true signal, whose radius is on the level of the noise. The numerical result under the Poisson noise model in Section 4 further corroborates the stability of RWF.

## 3. Incremental Reshaped Wirtinger Flow

In large-sample and online scenarios, stochastic algorithms are preferred due to their potential advantage of faster convergence and lower memory requirement. Thus, in this section, we develop stochastic versions of RWF, referred to as incremental reshaped Wirtinger flow (IRWF). We show that IRWF guarantees exact recovery at a linear convergence rate under the same sample complexity. We further draw the connection between IRWF and the randomized Kaczmarz method recently developed for phase retrieval (Wei, 2015; Li et al., 2015; Chi and Lu, 2016), and establish its global convergence as a side product.

### 3.1 (Mini-batch) IRWF: Algorithm and Convergence

In order to fully exploit the processing throughput of CPU/GPU, we develop a mini-batch IRWF, described in Algorithm 2. The mini-batch IRWF applies the same initialization step as in RWF, and uses a mini-batch of measurements for each gradient update.

---

**Algorithm 2** Mini-batch Incremental Reshaped Wirtinger Flow (mini-batch IRWF)
 

---

**Input:**  $\mathbf{y} = \{y_i\}_{i=1}^m, \{\mathbf{a}_i\}_{i=1}^m$ , mini-batch size  $k$ ;

**Initialization:** Same as in RWF (Algorithm 1);

**Gradient loop:** for  $t = 0 : T - 1$  do

Choose  $\Gamma_t$  uniformly at random from the subsets of  $\{1, 2, \dots, m\}$  with the cardinality  $k$ , and let

$$\mathbf{z}^{(t+1)} = \mathbf{z}^{(t)} - \mu \cdot \mathbf{A}_{\Gamma_t}^* \left( \mathbf{A}_{\Gamma_t} \mathbf{z}^{(t)} - \mathbf{y}_{\Gamma_t} \odot \text{Ph}(\mathbf{A}_{\Gamma_t} \mathbf{z}^{(t)}) \right), \quad (19)$$

where  $\mathbf{A}_{\Gamma_t}$  is a matrix stacking  $\mathbf{a}_i^*$  for  $i \in \Gamma_t$  as its rows,  $\mathbf{y}_{\Gamma_t}$  is a vector stacking  $y_i$  for  $i \in \Gamma_t$  as its elements, and  $\text{Ph}(\mathbf{z})$  denotes the phase vector of  $\mathbf{z}$ .

**Output**  $\mathbf{z}^{(T)}$ .

---

If the gradient update uses only a single sample, i.e.,  $k = 1$ , we refer to Algorithm 2 as IRWF, where the step (19) becomes

$$\mathbf{z}^{(t+1)} = \mathbf{z}^{(t)} - \mu \left( \mathbf{a}_{i_t}^* \mathbf{z}^{(t)} - y_{i_t} \cdot \frac{\mathbf{a}_{i_t}^* \mathbf{z}^{(t)}}{|\mathbf{a}_{i_t}^* \mathbf{z}^{(t)}|} \right) \mathbf{a}_{i_t}. \quad (20)$$

We characterize the convergence of mini-batch IRWF in the following theorem.

**Theorem 4** Consider the problem of solving any given  $\mathbf{x} \in \mathbb{R}^n$  from a system of equations (1) with independent Gaussian measurement vectors  $\mathbf{a}_i \sim \mathcal{N}(0, \mathbf{I})$ . There exist some universal constants  $0 < \rho, \rho_0 < 1$  and  $c_0, c_1, c_2 > 0$  such that if  $m \geq c_0 n$  and  $\mu = \rho_0/n$  for the update rule (19), then with probability at least  $1 - c_1 \exp(-c_2 m)$ , we have that

$$\mathbb{E}_{\Gamma_t} \left[ \text{dist}^2(\mathbf{z}^{(t+1)}, \mathbf{x}) \right] \leq \left( 1 - \frac{k\rho}{n} \right) \cdot \text{dist}^2(\mathbf{z}^{(t)}, \mathbf{x}) \quad (21)$$

holds for all  $\mathbf{z}^{(t)}$  satisfying  $\frac{\text{dist}(\mathbf{z}^{(t)}, \mathbf{x})}{\|\mathbf{z}\|} \leq \frac{1}{10}$ .

**Proof** See Appendix E.1. ■

We suggest that  $\rho_0 = 1$  and hence the step size  $\mu = \frac{1}{n}$  in practice. Theorem 4 characterizes that the error decays exponentially fast in expectation if the estimate lands into the neighborhood of the global minimizer. For a generic optimization objective, it is not anticipated that incremental/stochastic first-order methods achieve linear convergence due to the variance of stochastic gradients. However, for our specific problem, the variance of stochastic gradients reduces as the estimate approaches the true signal, and hence a fixed step size can be employed and exponential decay can be established. A result similar in spirit was also established for the stochastic algorithm based on TWF (referred to as ITWF) (Kolte and Özgür, 2016). We provide further comparisons between IRWF and ITWF in Section 3.3. On the other hand, it was shown in (Moulines and Bach, 2011; Needell et al., 2016) that stochastic gradient methods yield linear convergence to the minimizer  $\mathbf{x}_*$  if the objective  $F(\mathbf{x}) = \sum_i f_i(\mathbf{x})$  is a smooth and strongly convex function and  $\mathbf{x}_*$  minimizes

all components  $f_i(\mathbf{x})$ . The summands of our objective (3) also share the same minimizer (although it is neither convex nor smooth), which also helps to explain the convergence property of the IRWF.

### 3.2 Connection to the Kaczmarz Method for Phase Retrieval

The *Kaczmarz method* was originally developed for solving systems of linear equations (Kaczmarz, 1937). In literature (Wei, 2015; Li et al., 2015), it was adapted to solve the phase retrieval problem, which we refer to as *Kaczmarz-PR*. It has been demonstrated in Wei (2015) that Kaczmarz-PR exhibits better empirical performance than error reduction (ER) (Gerchberg, 1972; Fienup, 1982) and WF (Candès et al., 2015). However, global convergence of Kaczmarz-PR has not been well established yet, although the randomized Kaczmarz method for the least-squares problem is known to converge at a linear rate (Strohmer and Vershynin, 2009; Zouzias and Freris, 2013). For instance, Wei (2015) obtained a bound on the estimation error which can be as large as the signal energy no matter how many iterations are taken. Li et al. (2015) established the asymptotic convergence in the regime when both  $m$  and  $n$  go to infinity but their ratio is fixed.

In this section, we draw connection between IRWF and Kaczmarz-PR, which enables us to establish the theoretical guarantee of Kaczmarz-PR by adapting that of IRWF. This is analogous to the connection made in Needell et al. (2016) between the Kaczmarz method and the stochastic gradient method for solving the least-squares problem. Here, the connection is made possible due to the lower-order loss function of RWF, which was not evident in previous studies of WF and TWF.

To be more specific, the Kaczmarz-PR (Wei, 2015, Algorithm 3) employs the following update rule

$$\mathbf{z}^{(t+1)} = \mathbf{z}^{(t)} - \frac{1}{\|\mathbf{a}_{i_t}\|^2} \left( \mathbf{a}_{i_t}^* \mathbf{z}^{(t)} - y_{i_t} \cdot \frac{\mathbf{a}_{i_t}^* \mathbf{z}^{(t)}}{|\mathbf{a}_{i_t}^* \mathbf{z}^{(t)}|} \right) \mathbf{a}_{i_t}, \quad (22)$$

where  $i_t$  is selected either in a deterministic manner or randomly. We focus on the randomized case where  $i_t$  is selected uniformly at random from  $\{1, \dots, m\}$ .

Comparing Equation (22) and Equation (20), the update rule of Kaczmarz-PR becomes equivalent to IRWF, if we replace the step size  $\mu$  by  $\frac{1}{\|\mathbf{a}_{i_t}\|^2}$ . Moreover, these two update rules are close if  $\mu$  is set as suggested, i.e.,  $\mu = \frac{1}{n}$ , because for Gaussian measurements,  $\|\mathbf{a}_{i_t}\|^2$  concentrates around  $n$  by the law of large numbers. As we demonstrate in the numerical experiments (see Table 1), Kaczmarz-PR and IRWF have similar performance as anticipated. Thus, following the convergence guarantee for IRWF in Theorem 4, we establish the convergence guarantee for the randomized Kaczmarz-PR as follows.

**Theorem 5** *Assume the measurement vectors are independent and each  $\mathbf{a}_i \sim \mathcal{N}(0, \mathbf{I})$ . There exist some universal constants  $0 < \rho < 1$  and  $c_0, c_1, c_2 > 0$  such that if  $m \geq c_0 n$ , then with probability at least  $1 - c_1 m \exp(-c_2 n)$ , the randomized Kaczmarz-PR update rule (22) yields*

$$\mathbb{E}_{i_t} \left[ \text{dist}^2(\mathbf{z}^{(t+1)}, \mathbf{x}) \right] \leq \left( 1 - \frac{\rho}{n} \right) \cdot \text{dist}^2(\mathbf{z}^{(t)}, \mathbf{x}) \quad (23)$$

for all  $\mathbf{z}^{(t)}$  satisfying  $\frac{\text{dist}(\mathbf{z}^{(t)}, \mathbf{x})}{\|\mathbf{z}^{(t)}\|} \leq \frac{1}{10}$ .

**Proof** See Appendix E.2. ■

Theorem 5 implies that as long as the iteration sequence  $\mathbf{z}^{(t)}$  lies in the neighborhood of the true signal, the error decays exponentially fast in expectation. However, Theorem 5 does not guarantee the full trajectory of the estimates to be in the neighborhood of the true signal. After submission of our work, more recent studies (Jeong and Gunturk, 2017; Tan and Vershynin, 2017) established the convergence of full trajectory by employing martingale theory.

Furthermore, Wei (2015) also provided a *block* Kaczmarz-PR (similar to the mini-batch version), whose update rule is given by

$$\mathbf{z}^{(t+1)} = \mathbf{z}^{(t)} - \mathbf{A}_{\Gamma_t}^\dagger \left( \mathbf{A}_{\Gamma_t} \mathbf{z}^{(t)} - \mathbf{y}_{\Gamma_t} \odot \text{Ph}(\mathbf{A}_{\Gamma_t} \mathbf{z}^{(t)}) \right), \quad (24)$$

where  $\Gamma_t$  is a selected block at iterate  $t$  containing row indices, and  $\dagger$  represents *Moore-Penrose pseudoinverse*, which is computed as follows:

$$\mathbf{A}^\dagger = \begin{cases} (\mathbf{A}^* \mathbf{A})^{-1} \mathbf{A}^*, & \text{if } \mathbf{A} \text{ has linearly independent columns;} \\ \mathbf{A}^* (\mathbf{A} \mathbf{A}^*)^{-1}, & \text{if } \mathbf{A} \text{ has linearly independent rows.} \end{cases} \quad (25)$$

Comparing Equation (24) and the mini-batch IRWF update in Equation (19), these two update rules are similar to each other if  $\mathbf{A}_{\Gamma_t} \mathbf{A}_{\Gamma_t}^*$  approaches  $\frac{n}{\rho_0} \mathbf{I}_{|\Gamma_t|}$ . For the case with Gaussian measurements,  $\mathbf{A}_{\Gamma_t}$  has linearly independent rows with high probability if  $|\Gamma_t| \leq n$  and hence  $\mathbf{A}_{\Gamma_t} \mathbf{A}_{\Gamma_t}^*$  is not far from  $n \mathbf{I}_{|\Gamma_t|}$ . Our numerical experiments (see Table 1) further suggest similar convergence rates for these two algorithms with the same block/mini-batch size.

Next, we argue that for the CDP setting, block Kaczmarz-PR is the same as the mini-batch IRWF with  $\mu = 1$ . The CDP measurements are collected in the following form

$$\mathbf{y}^{(l)} = |\mathbf{F} \mathbf{D}^{(l)} \mathbf{x}|, \quad 1 \leq l \leq L, \quad (26)$$

where  $\mathbf{F}$  represents the discrete Fourier transform (DFT) matrix,  $\mathbf{D}^{(l)}$  denotes a diagonal matrix (mask), and  $L$  denotes the number of masks. We choose the block size  $|\Gamma_t|$  to be the dimension  $n$  of the signal for the convenience of Fourier transform. Then  $\mathbf{A}_{\Gamma_t}$  becomes the Fourier transform composed with  $\mathbf{D}^{(l)}$  (mask effect) and  $\mathbf{A}_{\Gamma_t}^*$  becomes  $\mathbf{D}^{(l)*}$  multiplied by the inverse Fourier transform. Therefore,  $(\mathbf{A}_{\Gamma_t} \mathbf{A}_{\Gamma_t}^*) = \mathbf{I}$  if the diagonal elements of  $\mathbf{D}^{(l)}$  have unit magnitude. Taking the step size  $\mu = 1$ , the two algorithms are identical.

On the other hand, since the block Kaczmarz-PR needs to calculate the matrix inverse or to solve an inverse problem, the block size cannot be too large. However, mini-batch IRWF works well for a wide range of the mini-batch sizes, which can even grow with the signal dimension  $n$  as long as a batch of data is loadable into the memory.

### 3.3 Comparison with Incremental Truncated Wirtinger Flow (ITWF)

Recently, (Kolte and Özgür, 2016) designed and analyzed an incremental algorithm based on TWF, which is referred to as ITWF. More specifically, ITWF employs the same initialization procedure as TWF and randomly chooses one sample with the index  $i_t$  selected uniformly

at random from  $\{1, 2, \dots, m\}$  for the gradient update as follows:

$$\mathbf{z}^{(t+1)} = \mathbf{z}^{(t)} - \mu \cdot \frac{|\mathbf{a}_{i_t}^T \mathbf{z}|^2 - y_{i_t}^2}{\mathbf{a}_{i_t}^T \mathbf{z}} \mathbf{a}_{i_t} \mathbf{1}_{\mathcal{E}_{1,t}^{i_t} \cap \mathcal{E}_3^{i_t}}, \quad (27)$$

where  $\mathbf{1}_{\mathcal{E}_{1,t}^{i_t} \cap \mathcal{E}_3^{i_t}}$  represents the truncation rule determined by the events  $\mathcal{E}_{1,t}^{i_t}$  and  $\mathcal{E}_3^{i_t}$ . As a comparison, the update rule of IRWF is much simpler due to the use of the lower-order loss function and does not require any truncation in the gradient loop. Kolte and Özgür (2016) proved that the update (27) shrinks the estimate error as long as  $m/n$  is large enough and the estimate  $\mathbf{z}$ . Compared to ITWF, the IRWF update (19) also shrinks the estimate error, but runs faster than ITWF numerically as demonstrated in Section 4.

## 4. Numerical Experiments

In this section, we demonstrate the numerical efficiency of RWF and (mini-batch) IRWF by comparing their performance with other competitive algorithms. Our experiments are conducted not only for the real Gaussian case but also for the complex Gaussian and the CDP cases. All the experiments are implemented in Matlab 2015b and conducted on a computer equipped with Intel Core i7 3.4GHz CPU and 12GB RAM.

We first compare the sample complexity of RWF and IRWF with those of TWF, WF, Kaczmarz-PR and AltMinPhase via the empirical successful recovery rate versus the number of measurements. For RWF, we follow Algorithm 1 with the suggested parameters. For IRWF, we adopt a block size 64 for efficiency and set the step size  $\mu = 1/n$ . For WF, TWF, we use the code provided in the original papers with the suggested parameters. For ITWF, we also adopt a block size 64 and set the step size  $\mu = 0.6/n$  (optimal step size). We conduct the experiments for the real Gaussian, complex Gaussian and CDP cases respectively. For the real and complex cases, we set the signal dimension  $n$  to be 1000, and set the ratio  $m/n$  to take values from 2 to 6 with a step 0.1. For each  $m$ , we run 100 trials and count the number of successful trials. For each trial, we run a maximal number of iterations/passes  $T = 10000$  for all algorithms, and a trial is declared to be successful whenever the iterate satisfies  $\text{dist}(\mathbf{z}^{(T)}, \mathbf{x}) / \|\mathbf{x}\| \leq 10^{-5}$ . For the real Gaussian case, we generate the signal  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{n \times n})$ , and generate the measurement vectors  $\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{n \times n})$  i.i.d. for  $i = 1, \dots, m$ . For the complex Gaussian case, we generate the signal  $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_{n \times n}) + j\mathcal{N}(0, \mathbf{I}_{n \times n})$  and the measurement vectors  $\mathbf{a}_i \sim \frac{1}{2}\mathcal{N}(0, \mathbf{I}_{n \times n}) + j\frac{1}{2}\mathcal{N}(0, \mathbf{I}_{n \times n})$  i.i.d. for  $i = 1, \dots, m$ . For the CDP case (26), we set  $n = 1024$  for the convenience of FFT and  $m/n = L = 1, 2, \dots, 8$ . All other settings are the same as those for the real case.

We note that for the CDP case, the Kaczmarz-PR algorithm is identical to the IRWF with step size  $\mu = 1/n$  due to the argument in Section 3.2. Moreover under the CDP case, the AltMinPhase algorithm is identical to the RWF with step size  $\mu = 1$  because the inverse of the Fourier measurement matrix is nothing but its conjugate transpose. In the following experiments for the CDP case, we choose the step size  $\mu = 1/n$  for the IRWF and  $\mu = 1$  for the RWF, under which the Kaczmarz-PR algorithm coincides with the IRWF and the AltMinPhase algorithm coincides with the RWF.

Figure 4 plots the fraction of successful trials out of 100 trials for all algorithms, with respect to  $m/n$ . It can be seen that IRWF and Kaczmarz-PR exhibit a similar sample

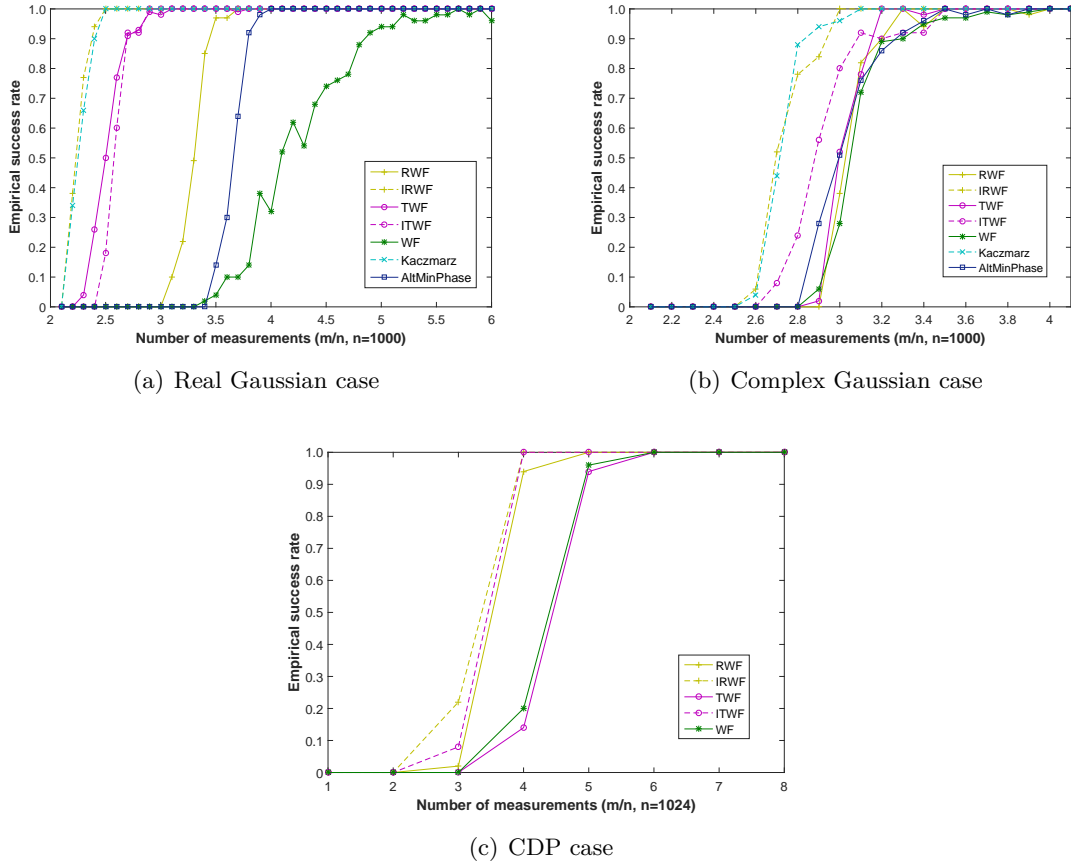


Figure 4: Comparison of sample complexity among RWF, IRWF, TWF, ITWF, WF, Kaczmarz-PR and AltMinPhase.

complexity, which is the best for all three cases, and is close to the theoretical limit (Bandeira et al., 2014). It can also be seen that the two incremental methods (IRWF and ITWF) outperform the batch methods (RWF, TWF, AltMinPhase and WF). This can be due to the inherent noise in incremental methods, which helps to escape bad local minima. This can be extremely helpful in the regime with a small number of samples, where local minima do exist near the global minima. Comparing among the batch methods (RWF, TWF, AltMinPhase and WF), it can be seen that although RWF outperforms only WF and AltMinPhase (not TWF) for the real Gaussian case, it has a comparable performance for the complex case and outperforms TWF and WF in the CDP case. An intuitive explanation for the real case is that a substantial number of samples with small  $|\mathbf{a}_i^T \mathbf{z}|$  can deviate the gradient direction so that truncation indeed helps to stabilize the algorithm if the number of measurements is not large.

We next compare the convergence rate of RWF, IRWF with those of TWF, ITWF, WF, Kaczmarz and AltMinPhase. We run all of the algorithms with the suggested parameters in the original code. We generate the signal and measurements in the same way as those



Table 1: Comparison of iteration count and time cost among algorithms ( $n = 5000, m = 8n$ ).

		Real Gaussian		Complex Gaussian	
		#passes	time(s)	# passes	time(s)
Batch methods	RWF	72	<b>12.66</b>	176	<b>122.4</b>
	TWF	186	32.36	487	395.1
	WF	319	54.83	932	887.8
	AltMinPhase	<b>6</b>	79.58	<b>159</b>	9637
Incremental methods	IRWF	9	44.77	21	233.2
	mini-batch IRWF (64)	9	<b>8.076</b>	<b>21</b>	<b>48.58</b>
	mini-batch ITWF (64)	16	37.38	29	149.5
	Kaczmarz-PR	9	50.68	21	248.4
	block Kaczmarz-PR (64)	<b>8</b>	28.50	22	89.31

in the first experiment with  $n = 5000, m = 8n$ . All algorithms are seeded with the RWF initialization. In Table 1, we list the number of passes and the time cost for all the algorithms to achieve a relative error of  $10^{-14}$  averaged over 10 trials. For the incremental methods, one update passes  $k$  samples and one pass amounts to  $m/k$  updates. Clearly, IRWF with mini-batch size 64 runs the fastest for both the real and complex cases. Moreover, among the batch (deterministic) algorithms, RWF takes much fewer passes as well as runs much faster than TWF and WF. Although RWF takes more iterations than AltMinPhase, it runs much faster than AltMinPhase due to the fact that each iteration of AltMinPhase needs to solve a least-squares problem that takes much longer than a simple gradient update in RWF.

We also compare the performance of the above algorithms on the recovery of a real image from the Fourier intensity measurements (the two dimensional CDP case). The image (see Figure 5) is the Milky Way Galaxy with resolution  $1920 \times 1080$ . Table 2 lists the number of passes and the time cost for the above six algorithms to achieve the relative error of  $10^{-15}$  for one R/G/B channel. All algorithms are seeded with the RWF initialization. To explore the advantage of FFT, we run the incremental/stochastic methods with the mini-batch size equal to the number of pixels for one R/G/B channel. We note that with such a mini-batch size, IRWF is equivalent to block Kaczmarz-PR from the discussion in Section 3.2. It can be seen that in general, the incremental/stochastic methods (IRWF and ITWF) run faster than the batch methods (RWF, TWF, WF). Moreover, among the batch methods, RWF outperforms the other three algorithms in both the number of passes and the computational time. In particular, RWF runs two times faster than TWF and six times faster than WF in terms of both the number of iterations and the computational time.

We next demonstrate the robustness of RWF to noise and compare it with TWF. We consider the phase retrieval problem in imaging applications, where Poisson noise is often used to model the sensor and electronic noise (Fogel et al., 2016). Specifically, the noisy measurements of intensity can be expressed as  $y_i = \sqrt{\alpha} \cdot \text{Poisson}(|\mathbf{a}_i^T \mathbf{x}|^2 / \alpha)$ , for



Figure 5: Milky way Galaxy.

Table 2: Comparison of iterations and time cost among algorithms on recovery of Galaxy image (shown in Figure 5), where  $L = m/n$  denotes the number of CDP masks.

	Algorithms	RWF	IRWF	TWF	ITWF	WF
$L = 6$	#passes	140	<b>24</b>	410	41	fail
	time cost(s)	110	<b>21.2</b>	406	43	fail
$L = 12$	#passes	70	<b>8</b>	190	12	315
	time cost(s)	107	<b>13.7</b>	363.6	25.9	426

$i = 1, 2, \dots, m$  where  $\alpha$  denotes the level of the input noise, and  $\text{Poisson}(\lambda)$  denotes a random sample generated by the Poisson distribution with mean  $\lambda$ . It can be observed from Figure 6 that RWF performs better than TWF in terms of the recovery accuracy under two different noise levels.

## 5. Conclusion

In this paper, we proposed RWF and its incremental version IRWF to recover a signal for the phase retrieval problem, based on a *nonconvex and nonsmooth* quadratic loss function of magnitude measurements. This loss function sacrifices the smoothness but enjoys advantages in statistical and computational efficiency. It has potential to be extended in various scenarios. One interesting direction is to extend such an algorithm to exploit signal structures (e.g., non-negativity, sparsity, etc) to assist the recovery. The lower-order loss function may offer great simplicity to prove the performance guarantee in such cases.

Another interesting direction is to study the convergence of algorithms from random initialization. In the regime with a large sample size ( $m \gg n$ ), the empirical loss surface

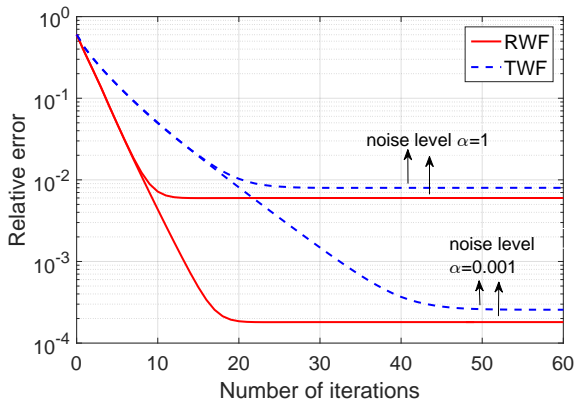


Figure 6: Comparison of relative error under Poisson noise between RWF and TWF.

approaches the asymptotic loss Figure 1(b) and hence has no spurious local minima. Due to the result (Lee et al., 2016a), it is conceivable that the gradient descent algorithm converges from any random starting point. Similar phenomena have been observed in (Sun et al., 2016; Ge et al., 2016). However, under a moderate number of measurements ( $m < 10n$ ), we observe that genuine local minima do exist and often locate not far from the global minima. In such a regime, the batch gradient method often fails with random initialization. As anticipated, stochastic algorithms are efficient in escaping bad local minima or saddle points in nonconvex optimization because of the inherent noise (Ge et al., 2015; Sa et al., 2015). We observe numerically that IRWF and block IRWF from a random starting point still converge to a global minimum even with a very small sample size which is close to the theoretical limit. It is of interest to analyze theoretically why stochastic methods escape these local minima (not just saddle points) efficiently.

## Acknowledgments

The work of H. Zhang, Y. Zhou and Y. Liang is supported in part by the grants NSF CCF-1704169, NSF ECCS-1609916 and AFOSR FA9550-16-1-0077. The work of Y. Chi is supported in part by the grants AFOSR under FA9550-15-1-0205, ONR N00014-15-1-2387, NSF CCF-1704245, and ECCS-1650449.

## Appendices

### Appendix A. Expected Loss Surfaces

It is easy to check that the expectation of the least-squares loss is given by  $\mathbb{E}[\ell_{LS}(\mathbf{z})] = \|\mathbf{z} - \mathbf{x}\|^2$ . The expectation of the WF loss function (2) is given by (Sun et al., 2016) as

$$\mathbb{E}[\ell_{WF}(\mathbf{z})] = \frac{3}{4}\|\mathbf{x}\|^4 + \frac{3}{4}\|\mathbf{z}\|^4 - \frac{1}{2}\|\mathbf{x}\|^2\|\mathbf{z}\|^2 - |\mathbf{z}^T \mathbf{x}|^2. \quad (28)$$

We next show that the expectation of the RWF loss function (3) has the following form:

$$\mathbb{E}[\ell(\mathbf{z})] = \frac{1}{2}\|\mathbf{x}\|^2 + \frac{1}{2}\|\mathbf{z}\|^2 - \|\mathbf{x}\|\|\mathbf{z}\| \cdot \mathbb{E}\left[\frac{|\mathbf{a}_i^T \mathbf{z}|}{\|\mathbf{z}\|} \cdot \frac{|\mathbf{a}_i^T \mathbf{x}|}{\|\mathbf{x}\|}\right], \quad (29)$$

where

$$\mathbb{E}\left[\frac{|\mathbf{a}_i^T \mathbf{z}|}{\|\mathbf{z}\|} \cdot \frac{|\mathbf{a}_i^T \mathbf{x}|}{\|\mathbf{x}\|}\right] = \begin{cases} \frac{(1-\rho^2)^{3/2}}{\pi} \int_0^\infty t(e^{\rho t} + e^{-\rho t})K_0(t)dt, & \text{if } |\rho| < 1; \\ 1, & \text{if } |\rho| = 1; \end{cases} \quad (30)$$

where  $\rho = \frac{\mathbf{z}^T \mathbf{x}}{\|\mathbf{x}\|\|\mathbf{z}\|}$  and  $K_0(\cdot)$  is the modified Bessel function of the second kind.

In order to derive (30), we first define

$$u := \frac{\mathbf{a}_i^T \mathbf{z}}{\|\mathbf{z}\|} \text{ and } v := \frac{\mathbf{a}_i^T \mathbf{x}}{\|\mathbf{x}\|},$$

and it suffices to drive  $\mathbb{E}[|uv|]$ . Note that  $(u, v) \sim \mathcal{N}(0, \Sigma)$ , where

$$\Sigma = \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}, \quad \text{and } \rho = \frac{\mathbf{z}^T \mathbf{x}}{\|\mathbf{x}\|\|\mathbf{z}\|}.$$

Following (Donahue, 1964), the density function of  $uv$  is given by

$$\phi_{uv}(x) = \frac{1}{\pi\sqrt{1-\rho^2}} \exp\left(\frac{\rho x}{1-\rho^2}\right) K_0\left(\frac{|x|}{1-\rho^2}\right), \quad x \neq 0.$$

Thus, the density function of  $|uv|$  is given by

$$\psi_{|uv|}(x) = \frac{1}{\pi\sqrt{1-\rho^2}} \left[ \exp\left(\frac{\rho x}{1-\rho^2}\right) + \exp\left(-\frac{\rho x}{1-\rho^2}\right) \right] K_0\left(\frac{|x|}{1-\rho^2}\right), \quad x > 0, \quad (31)$$

for  $|\rho| < 1$ . Therefore, if  $|\rho| < 1$ , we have

$$\begin{aligned} \mathbb{E}[|uv|] &= \int_0^\infty x \cdot \psi_\rho(x) dx \\ &= \int_0^\infty x \cdot \frac{1}{\pi\sqrt{1-\rho^2}} \left[ \exp\left(\frac{\rho x}{1-\rho^2}\right) + \exp\left(-\frac{\rho x}{1-\rho^2}\right) \right] K_0\left(\frac{|x|}{1-\rho^2}\right) dx \\ &= \frac{(1-\rho^2)^{3/2}}{\pi} \int_0^\infty t(e^{\rho t} + e^{-\rho t})K_0(t)dt, \end{aligned}$$

where the last step follows from change of variables.

If  $|\rho| = 1$ , then  $|uv|$  becomes a  $\chi_1^2$  random variable, with the density given by

$$\psi_{|uv|}(x) = \frac{1}{\sqrt{2\pi}} x^{-1/2} \exp(-x/2), \quad x > 0,$$

and hence  $\mathbb{E}[|uv|] = 1$ .

## Appendix B. Proof of Proposition 1

The idea of using truncation to bound some non-sub-Gaussian sequences has appeared in the literature (Candès et al., 2013, Lemma 2.3) and (Chen and Candès, 2015). Compared to the proof of the initialization for TWF (Chen and Candès, 2015), new technical developments are needed to address the magnitude measurements and truncation from both sides.

We first estimate the norm of  $\mathbf{x}$  as

$$\lambda_0 = \frac{mn}{\sum_{i=1}^m \|\mathbf{a}_i\|_1} \cdot \left( \frac{1}{m} \sum_{i=1}^m y_i \right). \quad (32)$$

Since  $\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{n \times n})$ , by Hoeffding-type inequality, it can be shown that

$$\left| \frac{\sum_{i=1}^m \|\mathbf{a}_i\|_1}{mn} - \sqrt{\frac{2}{\pi}} \right| < \frac{\epsilon}{3} \quad (33)$$

holds with probability at least  $1 - 2 \exp(-c_1 mn \epsilon^2)$  for some constant  $c_1 > 0$ .

Moreover, for a fixed  $\mathbf{x}$ ,  $y_i$ 's are independent sub-Gaussian random variables. Thus, by Hoeffding-type inequality, it can be shown that

$$\left| \sqrt{\frac{\pi}{2}} \left( \frac{1}{m} \sum_{i=1}^m y_i \right) - \|\mathbf{x}\| \right| < \frac{\epsilon}{3} \|\mathbf{x}\| \quad (34)$$

holds with probability at least  $1 - 2 \exp(-c_1 m \epsilon^2)$  for some constant  $c_1 > 0$ .

On the event  $E_1 = \{\text{both (33) and (34) hold}\}$ , it can be argued that

$$|\lambda_0 - \|\mathbf{x}\|| < \epsilon \|\mathbf{x}\|. \quad (35)$$

Without loss of generality, we let  $\|\mathbf{x}\| = 1$ . Then on the event  $E_1$ , the truncation function satisfies the following bounds

$$\mathbf{1}_{\{\alpha_l(1+\epsilon) < |\mathbf{a}_i^T \mathbf{x}| < \alpha_u(1-\epsilon)\}} \leq \mathbf{1}_{\{\alpha_l \lambda_0 < y_i < \alpha_u \lambda_0\}} \leq \mathbf{1}_{\{\alpha_l(1-\epsilon) < |\mathbf{a}_i^T \mathbf{x}| < \alpha_u(1+\epsilon)\}}.$$

Thus, define

$$\begin{aligned} \mathbf{Y}_1 &:= \frac{1}{m} \sum_{i=1}^m \mathbf{a}_i \mathbf{a}_i^T |\mathbf{a}_i^T \mathbf{x}| \mathbf{1}_{\{\alpha_l(1+\epsilon) < |\mathbf{a}_i^T \mathbf{x}| < \alpha_u(1-\epsilon)\}} \\ \mathbf{Y}_2 &:= \frac{1}{m} \sum_{i=1}^m \mathbf{a}_i \mathbf{a}_i^T |\mathbf{a}_i^T \mathbf{x}| \mathbf{1}_{\{\alpha_l(1-\epsilon) < |\mathbf{a}_i^T \mathbf{x}| < \alpha_u(1+\epsilon)\}}, \end{aligned}$$

and we have  $\mathbf{Y}_1 \prec \mathbf{Y} \prec \mathbf{Y}_2$ . We further compute the expectations of  $\mathbf{Y}_1$  and  $\mathbf{Y}_2$  and obtain

$$\mathbf{E}[\mathbf{Y}_1] = (\beta_1 \mathbf{x} \mathbf{x}^T + \beta_2 \mathbf{I}), \quad \mathbf{E}[\mathbf{Y}_2] = (\beta_3 \mathbf{x} \mathbf{x}^T + \beta_4 \mathbf{I}), \quad (36)$$

where

$$\begin{aligned} \beta_1 &:= \mathbf{E}[|\xi|^3 \mathbf{1}_{\{\alpha_l(1+\epsilon) < |\xi| < \alpha_u(1-\epsilon)\}}] - \mathbf{E}[|\xi| \mathbf{1}_{\{\alpha_l(1+\epsilon) < |\xi| < \alpha_u(1-\epsilon)\}}], \\ \beta_2 &:= \mathbf{E}[|\xi| \mathbf{1}_{\{\alpha_l(1+\epsilon) < |\xi| < \alpha_u(1-\epsilon)\}}] \\ \beta_3 &:= \mathbf{E}[|\xi|^3 \mathbf{1}_{\{\alpha_l(1-\epsilon) < |\xi| < \alpha_u(1+\epsilon)\}}] - \mathbf{E}[|\xi| \mathbf{1}_{\{\alpha_l(1-\epsilon) < |\xi| < \alpha_u(1+\epsilon)\}}], \\ \beta_4 &:= \mathbf{E}[|\xi| \mathbf{1}_{\{\alpha_l(1-\epsilon) < |\xi| < \alpha_u(1+\epsilon)\}}] \end{aligned}$$

where  $\xi \sim \mathcal{N}(0, 1)$ . For given  $\alpha_l$  and  $\alpha_u$ , a small value of  $\epsilon$  yields arbitrarily close  $\beta_1$  and  $\beta_3$ , as well as arbitrarily close  $\beta_2$  and  $\beta_4$ . For example, taking  $\alpha_l = 1, \alpha_u = 5$  and  $\epsilon = 0.01$ , we have  $\beta_1 = 0.9678, \beta_2 = 0.4791, \beta_3 = 0.9688, \beta_4 = 0.4888$ .

Now applying the standard results on random matrices with non-isotropic sub-Gaussian rows (Vershynin, 2012, equation (5.26)) and noticing that  $\mathbf{a}_i \mathbf{a}_i^T |\mathbf{a}_i^T \mathbf{x}| \mathbf{1}_{\{\alpha_l(1+\epsilon) < |\mathbf{a}_i^T \mathbf{x}| < \alpha_u(1-\epsilon)\}}$  can be rewritten as  $\mathbf{b}_i \mathbf{b}_i^T$  for a sub-Gaussian vector  $\mathbf{b}_i := \mathbf{a}_i \sqrt{|\mathbf{a}_i^T \mathbf{x}|} \mathbf{1}_{\{\alpha_l(1+\epsilon) < |\mathbf{a}_i^T \mathbf{x}| < \alpha_u(1-\epsilon)\}}$ , one can derive

$$\|\mathbf{Y}_1 - \mathbf{E}[\mathbf{Y}_1]\| \leq \delta, \quad \|\mathbf{Y}_2 - \mathbf{E}[\mathbf{Y}_2]\| \leq \delta \quad (37)$$

with probability at least  $1 - 4 \exp(-c_1(\delta)m)$  for some positive  $c_1$  which is only affected by  $\delta$ , provided that  $m/n$  exceeds a certain constant. Furthermore, when  $\epsilon$  is sufficiently small, one further has  $\|\mathbf{E}[\mathbf{Y}_1] - \mathbf{E}[\mathbf{Y}_2]\| \leq \delta$ . Combining the above facts together, one can show that

$$\|\mathbf{Y} - (\beta_1 \mathbf{x} \mathbf{x}^T + \beta_2 \mathbf{I})\| \leq 3\delta. \quad (38)$$

Let  $\tilde{\mathbf{z}}^{(0)}$  be the normalized leading eigenvector of  $\mathbf{Y}$ . Following the arguments in (Candès et al., 2015, Section 7.8) and taking  $\delta$  and  $\epsilon$  to be sufficiently small, one has

$$\text{dist}(\tilde{\mathbf{z}}^{(0)}, \mathbf{x}) \leq \tilde{\delta}, \quad (39)$$

for a given  $\tilde{\delta} > 0$ , as long as  $m/n$  exceeds a certain constant.

## Appendix C. Proof of Theorem 2

The general structure of the proof follows that for WF in (Candès et al., 2015) and TWF in (Chen and Candès, 2015). However, the proof requires the development of new bounds due to the nonsmoothness of the loss function and the magnitude measurements. On the other hand, the proof is much simpler due to the lower-order loss function adopted in RWF.

The idea of the proof is to show that within the neighborhood of the global optimizers, RWF satisfies the *Regularity Condition*  $\text{RC}(\mu, \lambda, c)$ , i.e.,

$$\langle \nabla \ell(\mathbf{z}), \mathbf{h} \rangle \geq \frac{\mu}{2} \|\nabla \ell(\mathbf{z})\|^2 + \frac{\lambda}{2} \|\mathbf{h}\|^2 \quad (40)$$

for all  $\mathbf{z}$  and  $\mathbf{h} = \mathbf{z} - \mathbf{x}$  obeying  $\|\mathbf{h}\| \leq c\|\mathbf{x}\|$ , where  $0 < c < 1$  is some constant. Then, as shown in (Chen and Candès, 2015), once the initialization lands into the neighborhood that satisfies  $\|\mathbf{z} - \mathbf{x}\| \leq c\|\mathbf{x}\|$ , linear convergence can be guaranteed, i.e.,

$$\text{dist}^2(\mathbf{z} - \mu\nabla\ell(\mathbf{z}), \mathbf{x}) \leq (1 - \mu\lambda)\text{dist}^2(\mathbf{z}, \mathbf{x}), \quad (41)$$

for any  $\mathbf{z}$  with  $\|\mathbf{z} - \mathbf{x}\| \leq c\|\mathbf{x}\|$ .

To show the regularity condition, we first define a set  $\mathcal{S} := \{i : 1 \leq i \leq m, (\mathbf{a}_i^T \mathbf{z})(\mathbf{a}_i^T \mathbf{x}) < 0\}$  (which depends on  $\mathbf{x}$  and  $\mathbf{z}$ ), and then derive the following bound:

$$\begin{aligned} \langle \nabla\ell(\mathbf{z}), \mathbf{h} \rangle &= \frac{1}{m} \sum_{i=1}^m \left( \mathbf{a}_i^T \mathbf{z} - |\mathbf{a}_i^T \mathbf{x}| \text{sgn}(\mathbf{a}_i^T \mathbf{z}) \right) (\mathbf{a}_i^T \mathbf{h}) \\ &= \frac{1}{m} \left[ \sum_{i=1}^m (\mathbf{a}_i^T \mathbf{h})^2 + 2 \sum_{i \in \mathcal{S}} (\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{h}) \right] \\ &\geq \frac{1}{m} \left[ \sum_{i=1}^m (\mathbf{a}_i^T \mathbf{h})^2 - 2 \left| \sum_{i \in \mathcal{S}} (\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{h}) \right| \right] \\ &\geq \frac{1}{m} \left[ \sum_{i=1}^m (\mathbf{a}_i^T \mathbf{h})^2 - \sum_{i \in \mathcal{S}} 2 \left| (\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{h}) \right| \right]. \end{aligned} \quad (42)$$

The first term in (42) can be bounded using Lemma 3.1 in (Candès et al., 2013), which we state below.

**Lemma 6** *For any  $0 < \epsilon < 1$ , there exist constants  $c_0, c_1 > 0$  such that if  $m > c_0 n \epsilon^{-2}$ , then with probability at least  $1 - 2 \exp(-c_1 \epsilon^2 m)$ ,*

$$(1 - \epsilon)\|\mathbf{h}\|^2 \leq \frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^T \mathbf{h})^2 \leq (1 + \epsilon)\|\mathbf{h}\|^2 \quad (43)$$

holds for all non-zero vectors  $\mathbf{h} \in \mathbb{R}^n$ .

For the second term in (42), we derive

$$\begin{aligned} \sum_{i \in \mathcal{S}} 2 \left| \mathbf{a}_i^T \mathbf{x} \right| \left| \mathbf{a}_i^T \mathbf{h} \right| &\leq \sum_{i \in \mathcal{S}} \left[ (\mathbf{a}_i^T \mathbf{x})^2 + (\mathbf{a}_i^T \mathbf{h})^2 \right] \\ &= \sum_{i=1}^m \left[ (\mathbf{a}_i^T \mathbf{x})^2 + (\mathbf{a}_i^T \mathbf{h})^2 \right] \cdot \mathbf{1}_{\{(\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{z}) < 0\}} \\ &= \sum_{i=1}^m \left[ (\mathbf{a}_i^T \mathbf{x})^2 + (\mathbf{a}_i^T \mathbf{h})^2 \right] \cdot \mathbf{1}_{\{(\mathbf{a}_i^T \mathbf{x})^2 + (\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{h}) < 0\}} \\ &\leq \sum_{i=1}^m \left[ (\mathbf{a}_i^T \mathbf{x})^2 + (\mathbf{a}_i^T \mathbf{h})^2 \right] \cdot \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{x}| < |\mathbf{a}_i^T \mathbf{h}|\}} \\ &\leq 2 \sum_{i=1}^m (\mathbf{a}_i^T \mathbf{h})^2 \cdot \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{x}| < |\mathbf{a}_i^T \mathbf{h}|\}}. \end{aligned} \quad (44)$$

The right hand side of the above equation can be further upper bounded by the following lemma.

**Lemma 7** For any  $\epsilon > 0$ , there exist constants  $c_0, c_1, C > 0$  such that if  $m > c_0 n \epsilon^{-2} \log \epsilon^{-1}$ , then with probability at least  $1 - C \exp(-c_1 \epsilon^2 m)$ ,

$$\frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^T \mathbf{h})^2 \cdot \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{x}| < |\mathbf{a}_i^T \mathbf{h}|\}} \leq (0.13 + \epsilon) \|\mathbf{h}\|^2 \quad (45)$$

holds for all non-zero vectors  $\mathbf{h} \in \mathbb{R}^n$  satisfying  $\|\mathbf{h}\| \leq \frac{1}{10} \|\mathbf{x}\|$ .

**Proof** See Section C.1. ■

Therefore, combining Lemmas 6 and 7 with (42) yields

$$\langle \nabla \ell(\mathbf{z}), \mathbf{h} \rangle \geq (1 - 0.26 - 3\epsilon) \|\mathbf{h}\|^2 = (0.74 - 3\epsilon) \|\mathbf{h}\|^2. \quad (46)$$

We further provide an upper bound on  $\|\nabla \ell(\mathbf{z})\|$  in the following lemma.

**Lemma 8** Fix  $\delta > 0$ . There exist constants  $c_0, c, C > 0$  such that if  $m \geq c_0 n$  then with probability at least  $1 - C \exp(-cm)$ ,

$$\|\nabla \ell(\mathbf{z})\| \leq (1 + \delta) \cdot 2 \|\mathbf{h}\| \quad (47)$$

holds for all non-zero vectors  $\mathbf{h}, \mathbf{z} \in \mathbb{R}^n$  satisfying  $\mathbf{z} = \mathbf{x} + \mathbf{h}$  and  $\|\mathbf{h}\| \leq \frac{1}{10} \|\mathbf{x}\|$ .

**Proof** See Section C.2. ■

Thus, applying Lemma 8 and (46) to (40), we conclude that *Regularity Condition* holds for  $\mu$  and  $\lambda$  satisfying

$$0.74 - 3\epsilon \geq \frac{\mu}{2} \cdot 4(1 + \delta)^2 + \frac{\lambda}{2}, \quad (48)$$

which concludes the proof.

### C.1 Proof of Lemma 7

We first prove bounds for any fixed  $\mathbf{h} \leq \frac{1}{10} \|\mathbf{x}\|$ , and then develop a uniform bound later on. We introduce a series of auxiliary random Lipschitz functions to approximate the indicator functions. For  $i = 1, \dots, m$ , define

$$\chi_i(t) := \begin{cases} t, & \text{if } t > (\mathbf{a}_i^T \mathbf{x})^2; \\ \frac{1}{\delta}(t - (\mathbf{a}_i^T \mathbf{x})^2) + (\mathbf{a}_i^T \mathbf{x})^2, & \text{if } (1 - \delta)(\mathbf{a}_i^T \mathbf{x})^2 \leq t \leq (\mathbf{a}_i^T \mathbf{x})^2; \\ 0, & \text{else;} \end{cases} \quad (49)$$

and then  $\chi_i(t)$ 's are random Lipschitz functions with Lipschitz constant  $\frac{1}{\delta}$ . We further have

$$|\mathbf{a}_i^T \mathbf{h}|^2 \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{x}| < |\mathbf{a}_i^T \mathbf{h}|\}} \leq \chi_i(|\mathbf{a}_i^T \mathbf{h}|^2) \leq |\mathbf{a}_i^T \mathbf{h}|^2 \mathbf{1}_{\{(1-\delta)|\mathbf{a}_i^T \mathbf{x}|^2 < |\mathbf{a}_i^T \mathbf{h}|^2\}}. \quad (50)$$



For convenience, we denote  $\gamma_i := \frac{|\mathbf{a}_i^T \mathbf{h}|^2}{\|\mathbf{h}\|^2} \mathbf{1}_{\{(1-\delta)|\mathbf{a}_i^T \mathbf{x}|^2 < |\mathbf{a}_i^T \mathbf{h}|^2\}}$  and  $\theta := \|\mathbf{h}\|/\|\mathbf{x}\|$ . We next estimate the expectation of  $\gamma_i$ , by the following conditional expectation,

$$\mathbb{E}[\gamma_i] = \int_{\Omega} \gamma_i d\mathbb{P} = \iint_{-\infty}^{\infty} \mathbb{E} \left[ \gamma_i | \mathbf{a}_i^T \mathbf{x} = \tau_1 \|\mathbf{x}\|, \mathbf{a}_i^T \mathbf{h} = \tau_2 \|\mathbf{h}\| \right] \cdot f(\tau_1, \tau_2) d\tau_1 d\tau_2, \quad (51)$$

where  $f(\tau_1, \tau_2)$  is the density of two joint Gaussian random variables with the correlation  $\rho = \frac{\mathbf{h}^T \mathbf{x}}{\|\mathbf{h}\| \|\mathbf{x}\|} \neq \pm 1$ . We then continue to derive

$$\begin{aligned} \mathbb{E}[\gamma_i] &= \iint_{-\infty}^{\infty} \tau_2^2 \cdot \mathbf{1}_{\{\sqrt{1-\delta}|\tau_1| < |\tau_2|\theta\}} \cdot f(\tau_1, \tau_2) d\tau_1 d\tau_2 \\ &= \frac{1}{2\pi\sqrt{1-\rho^2}} \int_{-\infty}^{\infty} \tau_2^2 \exp\left(-\frac{\tau_2^2}{2}\right) \cdot \int_{\frac{-|\tau_2|\theta}{\sqrt{1-\delta}}}^{\frac{|\tau_2|\theta}{\sqrt{1-\delta}}} \exp\left(-\frac{(\tau_1 - \rho\tau_2)^2}{2(1-\rho^2)}\right) d\tau_1 d\tau_2 \\ &= \frac{1}{2\pi} \int_{-\infty}^{\infty} \tau_2^2 \exp\left(-\frac{\tau_2^2}{2}\right) \cdot \int_{\frac{-\frac{|\tau_2|\theta}{\sqrt{1-\delta}} - \rho\tau_2}{\sqrt{1-\rho^2}}}^{\frac{\frac{|\tau_2|\theta}{\sqrt{1-\delta}} - \rho\tau_2}{\sqrt{1-\rho^2}}} \exp\left(-\frac{\tau^2}{2}\right) d\tau d\tau_2 \quad \text{by change of variables} \\ &= \frac{1}{2\pi} \int_{-\infty}^{\infty} \tau_2^2 \exp\left(-\frac{\tau_2^2}{2}\right) \cdot \sqrt{\frac{\pi}{2}} \left( \operatorname{erf}\left(\frac{\frac{|\tau_2|\theta}{\sqrt{1-\delta}} - \rho\tau_2}{\sqrt{1-\rho^2}}\right) - \operatorname{erf}\left(\frac{-\frac{|\tau_2|\theta}{\sqrt{1-\delta}} - \rho\tau_2}{\sqrt{1-\rho^2}}\right) \right) d\tau_2 \\ &= \frac{1}{\sqrt{2\pi}} \int_0^{\infty} \tau_2^2 \exp\left(-\frac{\tau_2^2}{2}\right) \cdot \left( \operatorname{erf}\left(\frac{(\frac{\theta}{\sqrt{1-\delta}} - \rho)\tau_2}{\sqrt{1-\rho^2}}\right) + \operatorname{erf}\left(\frac{(\frac{\theta}{\sqrt{1-\delta}} + \rho)\tau_2}{\sqrt{1-\rho^2}}\right) \right) d\tau_2. \quad (53) \end{aligned}$$

For  $|\rho| < 1$ ,  $\mathbb{E}[\gamma_i]$  is a continuous function of  $\rho$ . For  $|\rho| = 1$ ,  $\mathbb{E}[\gamma_i] = 0$ . The last integral (53) can be calculated numerically. Figure 7 plots  $\mathbb{E}[\gamma_i]$  for  $\theta = 0.1$  and  $\delta = 0.01$  over  $\rho \in [-1, 1]$ . Furthermore, (52) indicates that  $\mathbb{E}[\gamma_i]$  is monotonically increasing with both  $\theta$  and  $\delta$ . Thus, we obtain a universal bound

$$\mathbb{E}[\gamma_i] \leq 0.13 \quad \text{for } \theta < 0.1 \text{ and } \delta = 0.01, \quad (54)$$

which implies  $\mathbb{E}[\chi_i(|\mathbf{a}_i^T \mathbf{h}|^2)] \leq 0.13\|\mathbf{h}\|^2$  for  $\theta < 0.1$  and  $\delta = 0.01$ .

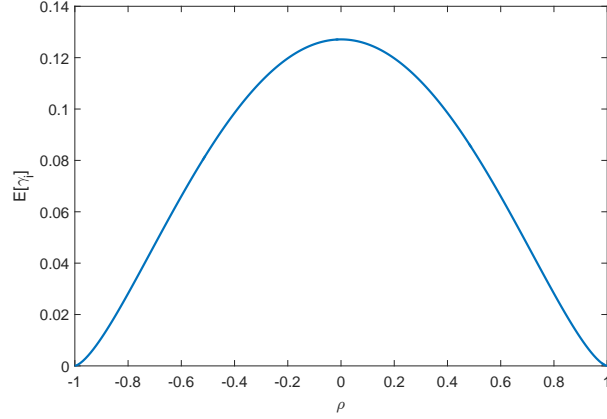
Furthermore,  $\chi_i(|\mathbf{a}_i^T \mathbf{h}|^2)$ 's are sub-exponential with the sub-exponential norm  $\mathcal{O}(\|\mathbf{h}\|^2)$ . By the sub-exponential tail bound (Bernstein type) (Vershynin, 2012), we have

$$\mathbb{P} \left[ \frac{1}{m} \sum_{i=1}^m \frac{\chi_i(|\mathbf{a}_i^T \mathbf{h}|^2)}{\|\mathbf{h}\|^2} > (0.13 + \epsilon) \right] < \exp(-cm\epsilon^2), \quad (55)$$

for some universal constant  $c$ , as long as  $\|\mathbf{h}\| \leq \frac{1}{10}\|\mathbf{x}\|$ .

We have proved so far that the claim holds for a fixed  $\mathbf{h}$ . We next obtain a uniform bound over all  $\mathbf{h}$  satisfying  $\|\mathbf{h}\| \leq \frac{1}{10}\|\mathbf{x}\|$ . We first show the claim holds for all  $\mathbf{h}$  with  $\|\mathbf{h}\| = \frac{1}{10}\|\mathbf{x}\|$  and then argue that the claim holds when  $\|\mathbf{h}\| < \frac{1}{10}\|\mathbf{x}\|$  towards the end of the proof. Let  $\epsilon' = \frac{\epsilon\|\mathbf{x}\|}{10}$  and we construct an  $\epsilon'$ -net  $\mathcal{N}_{\epsilon'}$  covering the sphere with the radius  $\frac{1}{10}\|\mathbf{x}\|$  in  $\mathbb{R}^n$  with the cardinality  $|\mathcal{N}_{\epsilon'}| \leq (1 + \frac{2}{\epsilon'})^n$ . Then for any  $\|\mathbf{h}\| = \frac{1}{10}\|\mathbf{x}\|$ , there exists an  $\mathbf{h}_0 \in \mathcal{N}_{\epsilon'}$  such that  $\|\mathbf{h} - \mathbf{h}_0\| \leq \epsilon\|\mathbf{h}\|$ . Taking the union bound for all the points on the net, we claim that

$$\frac{1}{m} \sum_{i=1}^m \chi_i \left( |\mathbf{a}_i^T \mathbf{h}_0|^2 \right) \leq (0.13 + \epsilon) \|\mathbf{h}_0\|^2, \quad \forall \mathbf{h}_0 \in \mathcal{N}_{\epsilon'} \quad (56)$$


 Figure 7:  $E[\gamma_i]$  with respect to  $\rho$ .

holds with probability at least  $1 - (1 + 2/\epsilon)^n \exp(-cm\epsilon^2)$ .

Since  $\chi_i(t)$ 's are Lipschitz functions with the constant  $1/\delta$ , we have the following bound

$$\left| \chi_i(|\mathbf{a}_i^T \mathbf{h}|^2) - \chi_i(|\mathbf{a}_i^T \mathbf{h}_0|^2) \right| \leq \frac{1}{\delta} \left| |\mathbf{a}_i^T \mathbf{h}|^2 - |\mathbf{a}_i^T \mathbf{h}_0|^2 \right|. \quad (57)$$

Moreover, by (Chen and Candès, 2015, Lemma 1), we have

$$\frac{1}{m} \sum_{i=1}^m |\mathbf{a}_i^T \mathbf{M} \mathbf{a}_i| \leq c_2 \|\mathbf{M}\|_F, \quad \text{for all symmetric rank-2 matrices } \mathbf{M} \in \mathbb{R}^{n \times n}, \quad (58)$$

holds with probability at least  $1 - C \exp(-c_1 m)$  as long as  $m > c_0 n$  for some constants  $C, c_0, c_1, c_2 > 0$ . Consequently, on the event that (58) holds, we have

$$\begin{aligned} & \left| \frac{1}{m} \sum_{i=1}^m \chi_i(|\mathbf{a}_i^T \mathbf{h}|^2) - \frac{1}{m} \sum_{i=1}^m \chi_i(|\mathbf{a}_i^T \mathbf{h}_0|^2) \right| \\ & \leq \frac{1}{m} \sum_{i=1}^m \left| \chi_i(|\mathbf{a}_i^T \mathbf{h}|^2) - \chi_i(|\mathbf{a}_i^T \mathbf{h}_0|^2) \right| \\ & \leq \frac{1}{\delta} \cdot \frac{1}{m} \sum_{i=1}^m \left| \mathbf{a}_i^T (\mathbf{h} \mathbf{h}^T - \mathbf{h}_0 \mathbf{h}_0^T) \mathbf{a}_i \right| \quad \text{because of (57)} \\ & \leq \frac{1}{\delta} \cdot c_2 \|\mathbf{h} \mathbf{h}^T - \mathbf{h}_0 \mathbf{h}_0^T\|_F \quad \text{because of (58)} \\ & \leq \frac{1}{\delta} \cdot 3c_2 \|\mathbf{h} - \mathbf{h}_0\| \cdot \|\mathbf{h}\| \\ & \leq \frac{3c_2 \epsilon}{\delta} \|\mathbf{h}\|^2, \end{aligned}$$

where the second to last inequality is due to (Chen and Candès, 2015, Lemma 2).

On the intersection of the events over which (56) and (58) hold respectively, we have

$$\frac{1}{m} \sum_{i=1}^m \chi_i(|\mathbf{a}_i^T \mathbf{h}|^2) \leq (0.13 + \epsilon + 3c_2 \epsilon / \delta) \|\mathbf{h}\|^2, \quad (59)$$

for all  $\mathbf{h}$  with  $\|\mathbf{h}\| = \frac{1}{10}\|\mathbf{x}\|$ .

For the case when  $\|\mathbf{h}'\| < \frac{1}{10}\|\mathbf{x}\|$ ,  $\mathbf{h}' = \omega\mathbf{h}$  for some  $\mathbf{h}$  satisfying  $\|\mathbf{h}\| = \frac{1}{10}\|\mathbf{x}\|$  and  $0 < \omega < 1$ . By the definition of  $\chi_i(\cdot)$ , it can be verified that

$$\chi_i(|\mathbf{a}_i^T \mathbf{h}'|^2) = \chi_i(|\mathbf{a}_i^T (\omega\mathbf{h})|^2) \leq \omega^2 \chi_i(|\mathbf{a}_i^T \mathbf{h}|^2). \quad (60)$$

Applying (59), on the same event over which (56) and (58) hold, we have

$$\frac{1}{m} \sum_{i=1}^m \chi_i(|\mathbf{a}_i^T \mathbf{h}'|^2) \leq (0.13 + \epsilon + 3c_2\epsilon/\delta) \|\mathbf{h}'\|^2, \quad (61)$$

for all  $\|\mathbf{h}'\| < \frac{1}{10}\|\mathbf{x}\|$ . Since  $\epsilon$  can be arbitrarily small, the proof is completed.

## C.2 Proof of Lemma 8

Denote  $v_i := \mathbf{a}_i^T \mathbf{z} - |\mathbf{a}_i^T \mathbf{x}| \text{sgn}(\mathbf{a}_i^T \mathbf{z})$  for  $i = 1, \dots, m$ . Then

$$\nabla \ell(\mathbf{z}) = \frac{1}{m} \mathbf{A}^T \mathbf{v}, \quad (62)$$

where  $\mathbf{A}$  is a matrix with each row being  $\mathbf{a}_i^T$  and  $\mathbf{v} = [v_1, \dots, v_m]^T$ . Thus,

$$\|\nabla \ell(\mathbf{z})\| = \left\| \frac{1}{m} \mathbf{A}^T \mathbf{v} \right\| \leq \frac{1}{m} \|\mathbf{A}\| \cdot \|\mathbf{v}\| \leq (1 + \delta) \frac{\|\mathbf{v}\|}{\sqrt{m}} \quad (63)$$

as long as  $m \geq c_1 n$  for some sufficiently large  $c_1 > 0$ , where the spectral norm bound  $\|\mathbf{A}\| \leq \sqrt{m}(1 + \delta)$  follows from (Vershynin, 2012, Theorem 5.32).

We next bound  $\|\mathbf{v}\|$ . Let  $\mathbf{v} = \mathbf{v}^{(1)} + \mathbf{v}^{(2)}$ , where  $v_i^{(1)} = \mathbf{a}_i^T \mathbf{h}$  and  $v_i^{(2)} = 2\mathbf{a}_i^T \mathbf{x} \mathbf{1}_{\{(\mathbf{a}_i^T \mathbf{z})(\mathbf{a}_i^T \mathbf{x}) < 0\}}$ . By the triangle inequality, we have  $\|\mathbf{v}\| \leq \|\mathbf{v}^{(1)}\| + \|\mathbf{v}^{(2)}\|$ . Furthermore, given  $m > c_0 n$ , following from (Candès et al., 2013, Lemma 3.1), with probability at least  $1 - \exp(-cm)$ , we have

$$\frac{1}{m} \|\mathbf{v}^{(1)}\|^2 = \frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^T \mathbf{h})^2 \leq (1 + \delta) \|\mathbf{h}\|^2. \quad (64)$$

By Lemma 7, we have with probability at least  $1 - C \exp(-c_1 m)$

$$\frac{1}{m} \|\mathbf{v}^{(2)}\|^2 = \frac{1}{m} \sum_{i=1}^m 4(\mathbf{a}_i^T \mathbf{x})^2 \cdot \mathbf{1}_{\{(\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{z}) < 0\}} \leq 4(0.13 + \epsilon) \|\mathbf{h}\|^2. \quad (65)$$

Hence,

$$\frac{\|\mathbf{v}\|}{\sqrt{m}} \leq \left( \sqrt{1 + \delta} + 2\sqrt{0.13 + \epsilon} \right) \|\mathbf{h}\|. \quad (66)$$

This concludes the proof.

### Appendix D. Proof of Theorem 3

The initialization analysis is similar to Appendix B and is hence omitted. To analyze the gradient loop, we consider the following two regimes.

• **Regime 1:**  $c_4\|\mathbf{z}\| \geq \|\mathbf{h}\| \geq c_3\frac{\|\mathbf{w}\|}{\sqrt{m}}$ . In this regime, the error contraction by each gradient descent step is given by

$$\text{dist}(\mathbf{z} - \mu\nabla\ell(\mathbf{z}), \mathbf{x}) \leq (1 - \rho)\text{dist}(\mathbf{z}, \mathbf{x}). \quad (67)$$

It suffices to justify that  $\nabla\ell(\mathbf{z})$  satisfies the RC. We first have

$$\begin{aligned} \nabla\ell(\mathbf{z}) &= \frac{1}{m} \sum_{i=1}^m \left( \mathbf{a}_i^T \mathbf{z} - y_i \cdot \frac{\mathbf{a}_i^T \mathbf{z}}{|\mathbf{a}_i^T \mathbf{z}|} \right) \mathbf{a}_i \\ &= \underbrace{\frac{1}{m} \sum_{i=1}^m \left( \mathbf{a}_i^T \mathbf{z} - |\mathbf{a}_i^T \mathbf{x}| \cdot \frac{\mathbf{a}_i^T \mathbf{z}}{|\mathbf{a}_i^T \mathbf{z}|} \right) \mathbf{a}_i}_{\nabla^{\text{clean}}\ell(\mathbf{z})} - \underbrace{\frac{1}{m} \sum_{i=1}^m \left( w_i \cdot \frac{\mathbf{a}_i^T \mathbf{z}}{|\mathbf{a}_i^T \mathbf{z}|} \right) \mathbf{a}_i}_{\nabla^{\text{noise}}\ell(\mathbf{z})}. \end{aligned} \quad (68)$$

All the proofs for Lemmas 6, 7 and 8 are still valid for  $\nabla^{\text{clean}}\ell(\mathbf{z})$ , and thus we have

$$\langle \nabla^{\text{clean}}\ell(\mathbf{z}), \mathbf{h} \rangle \geq 0.74\|\mathbf{h}\|^2, \quad (69)$$

$$\|\nabla^{\text{clean}}\ell(\mathbf{z})\| \leq 2(1 + \delta)\|\mathbf{h}\|. \quad (70)$$

Next, we analyze the contribution of the noise. Let  $\tilde{w}_i = w_i \frac{\mathbf{a}_i^T \mathbf{z}}{|\mathbf{a}_i^T \mathbf{z}|}$ , and then for sufficiently large  $m/n$ , we have

$$\|\nabla^{\text{noise}}\ell(\mathbf{z})\| = \left\| \frac{1}{m} \mathbf{A}^T \tilde{\mathbf{w}} \right\| \leq \left\| \frac{1}{\sqrt{m}} \mathbf{A}^T \right\| \left\| \frac{\tilde{\mathbf{w}}}{\sqrt{m}} \right\| \leq (1 + \delta) \frac{\|\tilde{\mathbf{w}}\|}{\sqrt{m}} \leq (1 + \delta) \frac{\|\mathbf{w}\|}{\sqrt{m}}, \quad (71)$$

where the second inequality is due to the spectral norm bound  $\|\mathbf{A}\| \leq \sqrt{m}(1 + \delta)$  (Vershynin, 2012, Theorem 5.32). Given the regime condition  $\|\mathbf{h}\| \geq c_3\frac{\|\mathbf{w}\|}{\sqrt{m}}$ , we further have

$$\|\nabla^{\text{noise}}\ell(\mathbf{z})\| \leq \frac{(1 + \delta)}{c_3} \|\mathbf{h}\|, \quad (72)$$

$$\left| \langle \nabla^{\text{noise}}\ell(\mathbf{z}), \mathbf{h} \rangle \right| \leq \|\nabla^{\text{noise}}\ell(\mathbf{z})\| \cdot \|\mathbf{h}\| \leq \frac{(1 + \delta)}{c_3} \|\mathbf{h}\|^2. \quad (73)$$

Combining these together, one has

$$\langle \nabla\ell(\mathbf{z}), \mathbf{h} \rangle \geq \langle \nabla^{\text{clean}}\ell(\mathbf{z}), \mathbf{h} \rangle - \left| \langle \nabla^{\text{noise}}\ell(\mathbf{z}), \mathbf{h} \rangle \right| \geq \left( 0.74 - \frac{(1 + \delta)}{c_3} \right) \|\mathbf{h}\|^2, \quad (74)$$

and

$$\|\nabla\ell(\mathbf{z})\| \leq \|\nabla^{\text{clean}}\ell(\mathbf{z})\| + \|\nabla^{\text{noise}}\ell(\mathbf{z})\| \leq (1 + \delta) \left( 2 + \frac{1}{c_3} \right) \|\mathbf{h}\|. \quad (75)$$

The RC is guaranteed if  $\mu, \lambda$  are chosen properly,  $c_3$  is sufficiently large, and  $\delta$  is sufficiently small.

• **Regime 2:** Once the iterate enters the regime with  $\|\mathbf{h}\| \leq \frac{c_3\|\mathbf{w}\|}{\sqrt{m}}$ , the gradient descent update may not reduce the estimation error. However, in this regime, the size  $\mu\nabla\ell(\mathbf{z})$  of each move is at most  $\mathcal{O}(\|\mathbf{w}\|/\sqrt{m})$ . Then the estimation error cannot increase by more than  $\|\mathbf{w}\|/\sqrt{m}$  with a constant factor. Thus, one has

$$\text{dist}(\mathbf{z} + \mu\nabla\ell(\mathbf{z}), \mathbf{x}) \leq c_5 \frac{\|\mathbf{w}\|}{\sqrt{m}} \quad (76)$$

for some constant  $c_5$ . As long as  $\|\mathbf{w}\|/\sqrt{m}$  is sufficiently small, it is guaranteed that  $c_5 \frac{\|\mathbf{w}\|}{\sqrt{m}} \leq c_4\|\mathbf{x}\|$ . If the iterate jumps out of *Regime 2*, it falls back into *Regime 1*.

## Appendix E. Convergence of Mini-batch IRWF and Kaczmarz-PR

### E.1 Proof of Theorem 4

Without loss of generality, we assume  $\mathbf{z}^{(t)}$  is in the neighborhood of  $\mathbf{x}$  (otherwise it is in the neighborhood of  $-\mathbf{x}$ ). Let  $\mathbf{h} = \mathbf{z}^{(t)} - \mathbf{x}$ . We follow the notations in Appendix C and let  $\mathcal{S} = \{i : (\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{z}^{(t)}) < 0\}$ . Then we have

$$\begin{aligned} & \mathbb{E}_{\Gamma_t} \left[ \text{dist}^2(\mathbf{z}^{(t+1)}, \mathbf{x}) \right] \\ &= \mathbb{E}_{\Gamma_t} \left[ \left\| \mathbf{z}^{(t)} - \mu \mathbf{A}_{\Gamma_t}^T \left( \mathbf{A}_{\Gamma_t} \mathbf{z}^{(t)} - \mathbf{y}_{\Gamma_t} \odot \text{sgn}(\mathbf{A}_{\Gamma_t} \mathbf{z}^{(t)}) \right) - \mathbf{x} \right\|^2 \right] \\ &= \|\mathbf{h}\|^2 - 2\mu \mathbb{E}_{\Gamma_t} \left[ \left( \mathbf{A}_{\Gamma_t} \mathbf{z}^{(t)} - \mathbf{y}_{\Gamma_t} \odot \text{sgn}(\mathbf{A}_{\Gamma_t} \mathbf{z}^{(t)}) \right)^T \mathbf{A}_{\Gamma_t} \mathbf{h} \right] \\ &\quad + \mu^2 \mathbb{E}_{\Gamma_t} \left[ \left( \mathbf{A}_{\Gamma_t} \mathbf{z}^{(t)} - \mathbf{y}_{\Gamma_t} \odot \text{sgn}(\mathbf{A}_{\Gamma_t} \mathbf{z}^{(t)}) \right)^T \mathbf{A}_{\Gamma_t} \mathbf{A}_{\Gamma_t}^T \left( \mathbf{A}_{\Gamma_t} \mathbf{z}^{(t)} - \mathbf{y}_{\Gamma_t} \odot \text{sgn}(\mathbf{A}_{\Gamma_t} \mathbf{z}^{(t)}) \right) \right] \\ &\stackrel{(a)}{=} \|\mathbf{h}\|^2 - \frac{2\mu k}{m} \sum_{i=1}^m \mathbf{a}_i^T \mathbf{h} \left( \mathbf{a}_i^T \mathbf{z}^{(t)} - y_i \cdot \frac{\mathbf{a}_i^T \mathbf{z}^{(t)}}{|\mathbf{a}_i^T \mathbf{z}^{(t)}|} \right) + \frac{\mu^2 k}{m} \sum_{i=1}^m \|\mathbf{a}_i\|^2 \left( \mathbf{a}_i^T \mathbf{z}^{(t)} - y_i \cdot \frac{\mathbf{a}_i^T \mathbf{z}^{(t)}}{|\mathbf{a}_i^T \mathbf{z}^{(t)}|} \right)^2 \\ &= \|\mathbf{h}\|^2 - \frac{2\mu k}{m} \left( \sum_{i=1}^m (\mathbf{a}_i^T \mathbf{h})^2 + \sum_{i \in \mathcal{S}} 2(\mathbf{a}_i^T \mathbf{h})(\mathbf{a}_i^T \mathbf{x}) \right) \\ &\quad + \frac{\mu^2 k}{m} \left( \sum_{i=1}^m \|\mathbf{a}_i\|^2 (\mathbf{a}_i^T \mathbf{h})^2 + 4 \sum_{i \in \mathcal{S}} \|\mathbf{a}_i\|^2 (\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{z}^{(t)}) \right) \\ &\leq \|\mathbf{h}\|^2 - \frac{2\mu k}{m} \sum_{i=1}^m (\mathbf{a}_i^T \mathbf{h})^2 + \frac{4\mu k}{m} \sum_{i \in \mathcal{S}} |(\mathbf{a}_i^T \mathbf{h})(\mathbf{a}_i^T \mathbf{x})| + \frac{\mu^2 k}{m} \sum_{i=1}^m \|\mathbf{a}_i\|^2 (\mathbf{a}_i^T \mathbf{h})^2, \end{aligned} \quad (77)$$

where (a) is due to the fact that  $\Gamma_t$  is uniformly chosen from all subsets of  $\{1, 2, \dots, m\}$  with the cardinality  $k$ .

By Lemma 6, we have that if  $m \geq c_0 \epsilon^{-2} n$  then with probability  $1 - 2 \exp(-c_1 m \epsilon^2)$

$$(1 - \epsilon) \|\mathbf{h}\|^2 \leq \frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^T \mathbf{h})^2 \leq (1 + \epsilon) \|\mathbf{h}\|^2.$$

holds for all vectors  $\mathbf{h}$ . By Lemma 7, we have that with probability  $1 - C \exp(-c_1 m \epsilon^2)$

$$\frac{1}{m} \sum_{i \in \mathcal{S}} |(\mathbf{a}_i^T \mathbf{h})(\mathbf{a}_i^T \mathbf{x})| \leq (0.13 + \epsilon) \|\mathbf{h}\|^2$$

holds for all  $\mathbf{h}$  satisfying  $\|\mathbf{h}\|/\|\mathbf{x}\| \leq \frac{1}{10}$ .

Define an event  $E_1 := \{\max_{1 \leq i \leq m} \|\mathbf{a}_i\|^2 \leq 6n\}$ . It can be shown that  $E_1$  holds with probability  $1 - m \exp(-1.5n)$ . Then on the event  $E_1$ , Equation (77) is further upper bounded by

$$\begin{aligned} \mathbb{E}_{\Gamma_t} \left[ \text{dist}^2 \left( \mathbf{z}^{(t+1)}, \mathbf{x} \right) \right] &\leq \left( 1 - 2\mu k(1 - \epsilon) + 4\mu k(0.13 + \epsilon) + \mu^2 k \cdot 6n(1 + \epsilon) \right) \|\mathbf{h}\|^2 \\ &\leq (1 - 2\mu k(0.74 - 3\epsilon - 3n\mu(1 + \epsilon))) \|\mathbf{h}\|^2. \end{aligned} \quad (78)$$

By choosing the step size  $\mu \leq \frac{0.24}{n}$ , the proposition is proved.

## E.2 Proof of Theorem 5

Without loss of generality, we assume the  $\mathbf{z}^{(t)}$  is in the neighborhood of  $\mathbf{x}$  (otherwise it is in the neighborhood of  $-\mathbf{x}$ ). Let  $\mathbf{h} = \mathbf{z}^{(t)} - \mathbf{x}$ . We follow the notations in Appendix C and let  $\mathcal{S} = \{i : (\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{z}^{(t)}) < 0\}$ . Then we have

$$\begin{aligned} &\mathbb{E}_{i_t} \text{dist}^2 \left( \mathbf{z}^{(t+1)}, \mathbf{x} \right) \\ &= \mathbb{E}_{i_t} \left\| \mathbf{z}^{(t)} - \frac{1}{\|\mathbf{a}_{i_t}\|^2} \left( \mathbf{a}_{i_t}^T \mathbf{z}^{(t)} - y_{i_t} \cdot \frac{\mathbf{a}_{i_t}^T \mathbf{z}^{(t)}}{|\mathbf{a}_{i_t}^T \mathbf{z}^{(t)}|} \right) \mathbf{a}_{i_t} - \mathbf{x} \right\|^2 \\ &= \|\mathbf{h}\|^2 - 2\mathbb{E}_{i_t} \frac{(\mathbf{a}_{i_t}^T \mathbf{h})}{\|\mathbf{a}_{i_t}\|^2} \left( \mathbf{a}_{i_t}^T \mathbf{z}^{(t)} - y_{i_t} \cdot \frac{\mathbf{a}_{i_t}^T \mathbf{z}^{(t)}}{|\mathbf{a}_{i_t}^T \mathbf{z}^{(t)}|} \right) + \mathbb{E}_{i_t} \frac{1}{\|\mathbf{a}_{i_t}\|^2} \left( \mathbf{a}_{i_t}^T \mathbf{z}^{(t)} - y_{i_t} \cdot \frac{\mathbf{a}_{i_t}^T \mathbf{z}^{(t)}}{|\mathbf{a}_{i_t}^T \mathbf{z}^{(t)}|} \right)^2 \\ &\stackrel{(a)}{=} \|\mathbf{h}\|^2 - \frac{2}{m} \sum_{i=1}^m \frac{(\mathbf{a}_i^T \mathbf{h})}{\|\mathbf{a}_i\|^2} \left( \mathbf{a}_i^T \mathbf{z}^{(t)} - y_i \cdot \frac{\mathbf{a}_i^T \mathbf{z}^{(t)}}{|\mathbf{a}_i^T \mathbf{z}^{(t)}|} \right) + \frac{1}{m} \sum_{i=1}^m \frac{1}{\|\mathbf{a}_i\|^2} \left( \mathbf{a}_i^T \mathbf{z}^{(t)} - y_i \cdot \frac{\mathbf{a}_i^T \mathbf{z}^{(t)}}{|\mathbf{a}_i^T \mathbf{z}^{(t)}|} \right)^2 \\ &= \|\mathbf{h}\|^2 - \frac{2}{m} \left( \sum_{i=1}^m \frac{(\mathbf{a}_i^T \mathbf{h})^2}{\|\mathbf{a}_i\|^2} + \sum_{i \in \mathcal{S}} \frac{2(\mathbf{a}_i^T \mathbf{h})(\mathbf{a}_i^T \mathbf{x})}{\|\mathbf{a}_i\|^2} \right) + \frac{1}{m} \left( \sum_{i=1}^m \frac{(\mathbf{a}_i^T \mathbf{h})^2}{\|\mathbf{a}_i\|^2} + 4 \sum_{i \in \mathcal{S}} \frac{(\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{z}^{(t)})}{\|\mathbf{a}_i\|^2} \right) \\ &= \|\mathbf{h}\|^2 - \frac{1}{m} \sum_{i=1}^m \frac{(\mathbf{a}_i^T \mathbf{h})^2}{\|\mathbf{a}_i\|^2} + \frac{4}{m} \sum_{i \in \mathcal{S}} \frac{(\mathbf{a}_i^T \mathbf{x})^2}{\|\mathbf{a}_i\|^2} \end{aligned} \quad (79)$$

where (a) is due to the fact that  $i_t$  is sampled uniformly at random from  $\{1, 2, \dots, m\}$ . By the special case of Lemma 5.20 in (Vershynin, 2012),  $\{\sqrt{n} \frac{\mathbf{a}_i}{\|\mathbf{a}_i\|}\}_{i=1}^m$  are independent isotropic random vectors in  $\mathbb{R}^n$  and hence

$$\mathbb{E} \left[ n \frac{(\mathbf{a}_i^T \mathbf{h})^2}{\|\mathbf{a}_i\|^2} \right] = \|\mathbf{h}\|^2.$$

Moreover,  $\{\sqrt{n} \frac{\mathbf{a}_i}{\|\mathbf{a}_i\|}\}_{i=1}^m$  are sub-Gaussian and the sub-Gaussian norm is bounded by an absolute constant. Thus, we have that if  $m \geq c_0 \epsilon^{-2} n$ , then with probability  $1 - 2 \exp(-c_1 m \epsilon^2)$ ,

$$\frac{1}{m} \sum_{i=1}^m \frac{(\mathbf{a}_i^T \mathbf{h})^2}{\|\mathbf{a}_i\|^2} \geq \frac{(1 - \epsilon)}{n} \|\mathbf{h}\|^2.$$

holds for all vectors  $\mathbf{h}$ . By Lemma 7, we have that with probability  $1 - C \exp(-c_1 m \epsilon^2)$

$$\frac{1}{m} \sum_{i \in \mathcal{S}} \left| \mathbf{a}_i^T \mathbf{x} \right|^2 \leq \frac{1}{m} \sum_{i=1}^m \left| \mathbf{a}_i^T \mathbf{h} \right|^2 \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{x}| < |\mathbf{a}_i^T \mathbf{h}|\}} \leq (0.13 + \epsilon) \|\mathbf{h}\|^2$$

holds for all  $\mathbf{h}$  satisfying  $\|\mathbf{h}\|/\|\mathbf{x}\| \leq \frac{1}{10}$ .

Define an event  $E_2 := \{\min_{1 \leq i \leq m} \|\mathbf{a}_i\|^2 \geq \frac{2}{3}n\}$ . It can be shown that  $\mathbb{P}\{E_2\} \geq 1 - m \exp(-n/12)$ . Then on the event  $E_2$ , (79) is further upper bounded by

$$E_{i_t} \left[ \text{dist}^2 \left( \mathbf{z}^{(t+1)}, \mathbf{x} \right) \right] \leq \left( 1 - \frac{1 - \epsilon}{n} + \frac{6(0.13 + \epsilon)}{n} \right) \|\mathbf{h}\|^2 \leq \left( 1 - \frac{0.22 - 7\epsilon}{n} \right) \|\mathbf{h}\|^2, \quad (80)$$

which concludes the proof.

## References

- Anima Anandkumar, Prateek Jain, Yang Shi, and Uma Naresh Niranjan. Tensor vs. matrix methods: Robust tensor decomposition under block sparse perturbations. In *Artificial Intelligence and Statistics*, pages 268–276, 2016.
- Sanjeev Arora, Rong Ge, Tengyu Ma, and Ankur Moitra. Simple, efficient, and neural algorithms for sparse coding. In *COLT*, pages 113–149, 2015.
- Afonso S Bandeira, Jameson Cahill, Dustin G Mixon, and Aaron A Nelson. Saving phase: Injectivity and stability for phase retrieval. *Applied and Computational Harmonic Analysis*, 37(1):106–125, 2014.
- Afonso S Bandeira, Nicolas Boumal, and Vladislav Voroninski. On the low-rank approach for semidefinite programs arising in synchronization and community detection. In *29th Annual Conference on Learning Theory*, 2016.
- Srinadh Bhojanapalli, Behnam Neyshabur, and Nati Srebro. Global optimality of local search for low rank matrix recovery. In *Advances in Neural Information Processing Systems*, pages 3873–3881, 2016.
- Nicolas Boumal. Nonconvex phase synchronization. *SIAM Journal on Optimization*, 26(4):2355–2377, 2016.
- James V Burke, Adrian S Lewis, and Michael L Overton. A robust gradient sampling algorithm for nonsmooth, nonconvex optimization. *SIAM Journal on Optimization*, 15(3):751–779, 2005.
- Tony T Cai, Xiaodong Li, and Zongming Ma. Optimal rates of convergence for noisy sparse phase retrieval via thresholded Wirtinger flow. *The Annals of Statistics*, 44(5):2221–2251, 2016.
- E. J. Candès and X. Li. Solving quadratic equations via phaselift when there are about as many equations as unknowns. *Foundations of Computational Mathematics*, 14(5):1017–1026, 2014.
- Emmanuel J Candès, Thomas Strohmer, and Vladislav Voroninski. Phaselift: Exact and stable signal recovery from magnitude measurements via convex programming. *Communications on Pure and Applied Mathematics*, 66(8):1241–1274, 2013.

- Emmanuel J Candès, Xiaodong Li, and Mahdi Soltanolkotabi. Phase retrieval via Wirtinger flow: Theory and algorithms. *IEEE Transactions on Information Theory*, 61(4):1985–2007, 2015.
- Yudong Chen and Martin J Wainwright. Fast low-rank estimation by projected gradient descent: General statistical and algorithmic guarantees. *arXiv preprint arXiv:1509.03025*, 2015.
- Yuxin Chen and Emmanuel Candès. Solving random quadratic systems of equations is nearly as easy as solving linear systems. In *Advances in Neural Information Processing Systems (NIPS)*. 2015.
- Yuxin Chen, Yuejie Chi, and Andrea J Goldsmith. Exact and stable covariance estimation from quadratic sampling via convex programming. *IEEE Transactions on Information Theory*, 61(7):4034–4059, 2015.
- Yuejie Chi and Yue M Lu. Kaczmarz method for solving quadratic equations. *IEEE Signal Processing Letters*, 23(9):1183–1187, 2016.
- James D Donahue. Products and quotients of random variables and their applications. Technical report, DTIC Document, 1964.
- Jan Drenth. *Principles of protein X-ray crystallography*. Springer Science & Business Media, 2007.
- James R Fienup. Phase retrieval algorithms: a comparison. *Applied Optics*, 21(15):2758–2769, 1982.
- Fajwel Fogel, Irène Waldspurger, and Alexandre dAspremont. Phase retrieval for imaging problems. *Mathematical programming computation*, 8(3):311–335, 2016.
- Rong Ge, Furong Huang, Chi Jin, and Yang Yuan. Escaping from saddle points-online stochastic gradient for tensor decomposition. In *Conference on Learning Theory (COLT)*, pages 797–842, 2015.
- Rong Ge, Jason D Lee, and Tengyu Ma. Matrix completion has no spurious local minimum. In *Advances in Neural Information Processing Systems*, pages 2973–2981, 2016.
- Ralph W Gerchberg. A practical algorithm for the determination of phase from image and diffraction plane pictures. *Optik*, 35:237, 1972.
- Xavier Glorot, Antoine Bordes, and Yoshua Bengio. Deep sparse rectifier neural networks. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2011.
- Moritz Hardt. Understanding alternating minimization for matrix completion. In *Foundations of Computer Science (FOCS), 2014 IEEE 55th Annual Symposium on*, pages 651–660. IEEE, 2014.
- Prateek Jain, Praneeth Netrapalli, and Sujay Sanghavi. Low-rank matrix completion using alternating minimization. In *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*, 2013.



- Halyun Jeong and C. Sinan Gunturk. Convergence of the randomized Kaczmarz method for phase retrieval. *arXiv preprint arXiv:1706.10291*, 2017.
- Chi Jin, Sham M Kakade, and Praneeth Netrapalli. Provable efficient online matrix completion via non-convex stochastic gradient descent. In *Advances in Neural Information Processing Systems*, pages 4520–4528, 2016.
- Stefan Kaczmarz. Angenäherte auflösung von systemen linearer gleichungen. *Bulletin International de l'Academie Polonaise des Sciences et des Lettres*, 35:355–357, 1937.
- Raghuveer H Keshavan, Andrea Montanari, and Sewoong Oh. Matrix completion from a few entries. *IEEE Transactions on Information Theory*, 56(6):2980–2998, June 2010.
- Krzysztof C Kiwiel. Convergence of the gradient sampling algorithm for nonsmooth nonconvex optimization. *SIAM Journal on Optimization*, 18(2):379–388, 2007.
- Ritesh Kolte and Ayfer Özgür. Phase retrieval via incremental truncated Wirtinger flow. *arXiv preprint arXiv:1606.03196*, 2016.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems (NIPS)*, 2012.
- A Ya Kruger. On Fréchet subdifferentials. *Journal of Mathematical Sciences*, 116(3):3325–3358, 2003.
- Jason D Lee, Max Simchowitz, Michael I Jordan, and Benjamin Recht. Gradient descent converges to minimizers. *arXiv preprint arXiv:1602.04915*, 2016a.
- Kiryung Lee, Yanjun Li, Marius Junge, and Yoram Bresler. Blind recovery of sparse signals from subsampled convolution. *IEEE Transactions on Information Theory*, 2016b.
- Gen Li, Yuantao Gu, and Yue M Lu. Phase retrieval using iterative projections: Dynamics in the large systems limit. In *The 53rd Annual Allerton Conference on Communication, Control, and Computing*, 2015.
- Xiaodong Li, Shuyang Ling, Thomas Strohmer, and Ke Wei. Rapid, robust, and reliable blind deconvolution via nonconvex optimization. *arXiv preprint arXiv:1606.04933*, 2016.
- Yuanxin Li, Yuejie Chi, Huishuai Zhang, and Yingbin Liang. Non-convex low-rank matrix recovery from corrupted random linear measurements. In *Sampling Theory and Applications (SampTA), International Conference on*, 2017.
- Jianwei Miao, Pambos Charalambous, Janos Kirz, and David Sayre. Extending the methodology of X-ray crystallography to allow imaging of micrometre-sized non-crystalline specimens. *Nature*, 400(6742):342–344, 1999.
- Jianwei Miao, Tetsuya Ishikawa, Qun Shen, and Thomas Earnest. Extending X-ray crystallography to allow the imaging of noncrystalline materials, cells, and single protein complexes. *Annu. Rev. Phys. Chem.*, 59:387–410, 2008.

- Eric Moulines and Francis R Bach. Non-asymptotic analysis of stochastic approximation algorithms for machine learning. In *Advances in Neural Information Processing Systems (NIPS)*, 2011.
- Deanna Needell, Nathan Srebro, and Rachel Ward. Stochastic gradient descent, weighted sampling, and the randomized Kaczmarz algorithm. *Mathematical Programming*, 155(1-2): 549–573, 2016.
- Praneeth Netrapalli, Prateek Jain, and Sujay Sanghavi. Phase retrieval using alternating minimization. *Advances in Neural Information Processing Systems (NIPS)*, 2013.
- Praneeth Netrapalli, UN Niranjan, Sujay Sanghavi, Animashree Anandkumar, and Prateek Jain. Non-convex robust PCA. In *Advances in Neural Information Processing Systems (NIPS)*, 2014.
- Peter Ochs, Alexey Dosovitskiy, Thomas Brox, and Thomas Pock. On iteratively reweighted algorithms for nonsmooth nonconvex optimization in computer vision. *SIAM Journal on Imaging Sciences*, 8(1):331–372, 2015.
- Dohyung Park, Anastasios Kyrillidis, Srinadh Bhojanapalli, Constantine Caramanis, and Sujay Sanghavi. Provable non-convex projected gradient descent for a class of constrained matrix optimization problems. *arXiv preprint arXiv:1606.01316*, 2016.
- Christopher De Sa, Christopher Re, and Kunle Olukotun. Global convergence of stochastic gradient descent for some non-convex matrix problems. In *Proceedings of International Conference on Machine Learning (ICML)*, pages 2332–2341, 2015.
- Yoav Shechtman, Yonina C Eldar, Oren Cohen, Henry Nicholas Chapman, Jianwei Miao, and Mordechai Segev. Phase retrieval with application to optical imaging: a contemporary overview. *IEEE Signal Processing Magazine*, 32(3):87–109, 2015.
- Thomas Strohmer and Roman Vershynin. A randomized kaczmarz algorithm with exponential convergence. *Journal of Fourier Analysis and Applications*, 15(2):262–278, 2009.
- Ju Sun, Qing Qu, and John Wright. Complete dictionary recovery using nonconvex optimization. In *Proceedings of International Conference on Machine Learning (ICML)*, 2015.
- Ju Sun, Qing Qu, and John Wright. A geometric analysis of phase retrieval. *arXiv preprint arXiv:1602.06664*, 2016.
- Ruoyu Sun and Zhi-Quan Luo. Guaranteed matrix completion via non-convex factorization. *IEEE Transactions on Information Theory*, 62(11):6535–6579, 2016.
- Yan Shuo Tan and Roman Vershynin. Phase retrieval via randomized Kaczmarz: Theoretical guarantees. *arXiv preprint arXiv:1706.09993*, 2017.
- Stephen Tu, Ross Boczar, Max Simchowitz, Mahdi Soltanolkotabi, and Benjamin Recht. Low-rank solutions of linear matrix equations via Procrustes flow. *arXiv preprint arXiv:1507.03566*, 2015.

- Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. *Compressed Sensing, Theory and Applications*, pages 210 – 268, 2012.
- Irène Waldspurger. Phase retrieval with random gaussian sensing vectors by alternating projections. *arXiv preprint arXiv:1609.03088*, 2016.
- Irène Waldspurger, Alexandre d’Aspremont, and Stéphane Mallat. Phase recovery, maxcut and complex semidefinite programming. *Mathematical Programming*, 149(1-2):47–81, 2015.
- Gang Wang, Georgios B Giannakis, and Yonina C Eldar. Solving systems of random quadratic equations via truncated amplitude flow. *arXiv preprint arXiv:1605.08285*, 2016.
- Ke Wei. Solving systems of phaseless equations via Kaczmarz methods: a proof of concept study. *Inverse Problems*, 31(12):125008, 2015.
- Ke Wei, Jian-Feng Cai, Tony F Chan, and Shingyu Leung. Guarantees of Riemannian optimization for low rank matrix recovery. *SIAM Journal on Matrix Analysis and Applications*, 37(3):1198–1222, 2016.
- Huishuai Zhang, Yuejie Chi, and Yingbin Liang. Median-truncated nonconvex approach for phase retrieval with outliers. In *Proceedings of International Conference on Machine Learning (ICML)*, 2016.
- Qinqing Zheng and John Lafferty. A convergent gradient descent algorithm for rank minimization and semidefinite programming from random linear measurements. In *Advances in Neural Information Processing Systems (NIPS)*, 2015.
- Qinqing Zheng and John Lafferty. Convergence analysis for rectangular matrix completion using burer-monteiro factorization and gradient descent. *arXiv preprint arXiv:1605.07051*, 2016.
- Anastasios Zouzias and Nikolaos M Freris. Randomized extended Kaczmarz for solving least squares. *SIAM Journal on Matrix Analysis and Applications*, 34(2):773–793, 2013.