

Phase Retrieval via Reweighted Amplitude Flow

Gang Wang , *Student Member, IEEE*, Georgios B. Giannakis , *Fellow, IEEE*, Yousef Saad,
and Jie Chen , *Senior Member, IEEE*

Abstract—This paper deals with finding an n -dimensional solution x to a system of quadratic equations of the form $y_i = |\langle a_i, x \rangle|^2$ for $1 \leq i \leq m$, which is also known as the generalized phase retrieval problem. For this NP-hard problem, a novel approach is developed for minimizing the amplitude-based least-squares empirical loss, which starts with a weighted maximal correlation initialization obtainable through a few power or Lanczos iterations, followed by successive refinements based on a sequence of iteratively reweighted gradient iterations. The two stages (initialization and gradient flow) distinguish themselves from prior contributions by the inclusion of a fresh (re)weighting regularization procedure. For certain random measurement models, the novel scheme is shown to be able to recover the true solution x in time proportional to reading the data $\{(a_i; y_i)\}_{1 \leq i \leq m}$. This holds with high probability and without extra assumption on the signal vector x to be recovered, provided that the number m of equations is some constant $c > 0$ times the number n of unknowns in the signal vector, namely $m > cn$. Empirically, the upshots of this contribution are: first, (almost) 100% perfect signal recovery in the high-dimensional (say $n \geq 2000$) regime given only an information-theoretic limit number of noiseless equations, namely $m = 2n - 1$, in the real Gaussian case; and second, (nearly) optimal statistical accuracy in the presence of additive noise of bounded support. Finally, substantial numerical tests using both synthetic data and real images corroborate markedly improved recovery performance and computational efficiency of the novel scheme relative to the state-of-the-art approaches.

Index Terms—Non-convex non-smooth optimization, regularization, iteratively reweighted gradient flow, convergence to the global optimum.

Manuscript received August 15, 2017; revised January 25, 2018 and March 10, 2018; accepted March 13, 2018. Date of publication March 23, 2018; date of current version April 19, 2018. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Marco Moretti. The work of G. Wang and G. B. Giannakis was partially supported by the National Science Foundation (NSF) under Grant 1500713 and Grant 1514056. The work of Y. Saad was supported by the NSF under Grant 1505970. The work of J. Chen was supported in part by the National Natural Science Foundation of China under Grant U1509215 and Grant 61621063, and in part by the Program for Changjiang Scholars and Innovative Research Team in University (IRT1208). This paper was presented in part at the Thirty-second Annual Conference on Neural Information Processing Systems, Long Beach, CA, USA, Dec. 4–9, 2017 [1]. (Corresponding author: Georgios B. Giannakis.)

G. Wang and G. B. Giannakis are with the Digital Technology Center and the Department of Electrical and Computer Engineering, University of Minnesota, Minneapolis, MN 55455 USA (e-mail: gangwang@umn.edu; georgios@umn.edu).

Y. Saad is with the Department of Computer Science and Engineering, University of Minnesota, Minneapolis, MN 55455 USA (e-mail: saad@umn.edu).

J. Chen is with the State Key Lab of Intelligent Control and Decision of Complex Systems and the School of Automation, Beijing Institute of Technology, Beijing 100081, China (e-mail: chenjie@bit.edu.cn).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TSP.2018.2818077

I. INTRODUCTION

ONE is often faced with solving quadratic equations of the form $y_i = |\langle a_i, x \rangle|^2$, or equivalently,

$$\psi_i = |\langle a_i, x \rangle|, \quad 1 \leq i \leq m \quad (1)$$

where $x \in \mathbb{R}^n$ is the wanted unknown $n \times 1$ signal vector, given observations ψ_i and feature/sensing vectors $a_i \in \mathbb{R}^n$ that are collectively stacked in the data vector $\psi := [\psi_i]_{1 \leq i \leq m}$, and the $m \times n$ sensing matrix $A := [a_i]_{1 \leq i \leq m}$, respectively. Phrased differently, when information about the (squared) modulus of the inner products of x and several known measurement vectors a_i is provided, can one reconstruct exactly (up to a global sign) x , or alternatively, the missing signs of $\langle a_i, x \rangle$? In fact, much effort has recently been devoted to determining the number of such equations necessary and/or sufficient to ensure uniqueness of the solution x ; see, for instance, [2], [3]. It has been proved that a number $m \geq 2n - 1$ of generic measurement vectors a_i ¹ (which includes the case of random vectors) are sufficient for uniquely determining an n -dimensional real vector x (up to a global sign), while $m = 2n - 1$ has also been shown necessary [2]. In this sense, the number $m = 2n - 1$ of equations as in (1) can be thought of as the information-theoretic limit for such a quadratic system to be uniquely solvable. Nevertheless, even for random measurement vectors, despite the existence of a unique solution given the minimal number $2n - 1$ of quadratic equations, it is unclear so far whether there is a numerical polynomial-time algorithm that is able to stably find the true solution (say with probability $\geq 99\%$).

In diverse physical sciences and engineering fields, it is impossible or very difficult to record phase measurements. Recovering the signal or phase from magnitude measurements only, also commonly known as the phase retrieval problem, emerges naturally [4]–[6]. Relevant application domains include e.g., X-ray crystallography, ptychography, astronomy, and coherent diffraction imaging [6]. In such setups however, optical measurement and detection systems record only the photon flux, which is proportional to the (squared) magnitude of the field, but not the phase. A related task of this kind is that of estimating a mixture of linear regressions, where the latent membership indicators can be converted into the missing phases [7]. Although of simple form and practical relevance across different fields, solving systems of nonlinear equations is arguably the most difficult task numerically [8, Page 355].

¹It is out of the scope of the present paper to explain the meaning of generic vectors, whereas interested readers are referred to [2].

Regarding notation used in this paper, lower-(upper-) case boldface letters denote vectors (matrices). Calligraphic letters are reserved for sets, e.g., \mathcal{S} . Fractions are denoted by A/B or $\frac{A}{B}$, but with a slight abuse of notation, we also use y_i/ψ_i , to denote either y_i or ψ_i . The floor operation $\lfloor c \rfloor$ denotes the largest integer no greater than the given number $c > 0$, $|\mathcal{S}|$ the number of entries in set \mathcal{S} , and $\|\mathbf{x}\|$ is the Euclidean norm. Since $\mathbf{x} \in \mathbb{R}^n$ and $-\mathbf{x}$ are indistinguishable given $\{\psi_i\}$ in (1), let $\text{dist}(\mathbf{z}, \mathbf{x}) = \min\{\|\mathbf{z} + \mathbf{x}\|, \|\mathbf{z} - \mathbf{x}\|\}$ be the Euclidean distance of any estimate $\mathbf{z} \in \mathbb{R}^n$ to the solution set $\{\pm \mathbf{x}\}$ of (1).

A. Prior Contributions

Following the least-squares criterion (which coincides with the maximum likelihood one when assuming additive white Gaussian noise), the problem of solving systems of quadratic equations can be recast as the ensuing empirical loss minimization

$$\underset{\mathbf{z} \in \mathbb{R}^n}{\text{minimize}} \quad L(\mathbf{z}) := \frac{1}{m} \sum_{i=1}^m \ell(\mathbf{z}; \psi_i/y_i) \quad (2)$$

where one can choose to work with the *amplitude-based* loss function $\ell(\mathbf{z}; \psi_i) := (\psi_i - |\langle \mathbf{a}_i, \mathbf{z} \rangle|)^2/2$ [9], or the *intensity-based* ones $\ell(\mathbf{z}; y_i) := (y_i - |\langle \mathbf{a}_i, \mathbf{z} \rangle|^2)^2/2$ [10], [11], and its related Poisson likelihood $\ell(\mathbf{z}; y_i) := -y_i \log(|\langle \mathbf{a}_i, \mathbf{z} \rangle|^2) + |\langle \mathbf{a}_i, \mathbf{z} \rangle|^2$ [12]. Either way, $L(\mathbf{z})$ is non-convex; hence, it is in general NP-hard, and computationally intractable to compute the least-squares or the maximum likelihood estimate [13].

Minimizing the squared amplitude-based least-squares loss in (2), several numerical polynomial-time algorithms have been devised based on convex programming for certain choices of design vectors \mathbf{a}_i [14]–[18], [19], [20]. Relying upon the so-called matrix-lifting technique semidefinite programming (SDP) based convex approaches first express all intensity data into linear terms in a new rank one matrix variable, followed by solving a convex SDP after dropping the rank constraint (a.k.a. semidefinite relaxation). It has been established that perfect recovery and (near-)optimal statistical accuracy can be achieved in noiseless and noisy settings, respectively, with an optimal-order number of measurements [18]. Another line of convex relaxation [21], [22], [23] reformulated the problem of phase retrieval as that of sparse signal recovery, and solved a linear program in the natural parameter vector domain. Although exact signal recovery can be established assuming an accurate enough anchor vector, its empirical performance is not competitive with state-of-the-art non-convex phase retrieval approaches.

Instead of convex relaxation, recent proposals also advocate judiciously initialized iterative procedures for coping with certain non-convex formulations directly, which include solvers based on e.g., alternating minimization [24], Wirtinger flows [10], [12], [25]–[32], amplitude flows [1], [9], [33]–[36], as well as a prox-linear procedure via composite optimization [37], [38], [39]. These non-convex approaches operate directly upon vector optimization variables, therefore leading to significant computational advantages over matrix-lifting based convex counterparts. With random features, they can be interpreted as performing stochastic optimization over acquired data samples $\{(\mathbf{a}_i; \psi_i/y_i)\}_{1 \leq i \leq m}$ to approximately minimize the population

risk functional $\bar{L}(\mathbf{z}) := \mathbb{E}_{(\mathbf{a}_i, \psi_i/y_i)}[\ell(\mathbf{z}; \psi_i/y_i)]$. It is well documented that minimizing non-convex functionals is computationally intractable in general due to existence of many stationary points [13]. Assuming random Gaussian sampling vectors however, such non-convex paradigms can provably locate the global optimum under suitable conditions, some of which also achieve optimal (statistical) guarantees. Specifically, starting with a judiciously designed initial guess, successive improvement is effected through a sequence of (truncated) (generalized) gradient iterations given by

$$\mathbf{z}^{t+1} := \mathbf{z}^t - \frac{\mu^t}{m} \sum_{i \in \mathcal{T}^{t+1}} \nabla \ell(\mathbf{z}^t; \psi_i/y_i), \quad t = 0, 1, \dots \quad (3)$$

where \mathbf{z}^t denotes the estimate returned by the algorithm at the t -th iteration, $\mu^t > 0$ the learning rate, and $\nabla \ell(\mathbf{z}^t; \psi_i/y_i)$ is the (generalized) gradient of the modulus- or squared modulus-based least-squares loss evaluated at \mathbf{z}^t [40]. Here, $\mathcal{T}^{t+1} \subseteq \{1, 2, \dots, m\}$ represents some time-varying index set signifying the truncation.

Although they achieve optimal statistical guarantees in both noiseless and noisy settings, state-of-the-art (convex and non-convex) approaches studied under random Gaussian designs, empirically require stable recovery of a number of equations (several) times larger than the aforementioned information-theoretic limit [10], [12], [27]. As a matter of fact, when there are numerous enough measurements (on the order of the signal dimension n up to some polylog factors), the amplitude-square based least-squares loss functional admits benign geometric structure in the sense that [41]: with high probability, i) all local minimizers are global; and, ii) there always exists a negative directional curvature at every saddle point. In a nutshell, the grand challenge of solving systems of random quadratic equations remains to develop numerical polynomial-time algorithms capable of achieving perfect recovery and optimal statistical accuracy when the number of measurements approaches the information-theoretic limit.

B. This Contribution

Building upon but going well beyond the scope of the aforementioned non-convex paradigms, the present paper puts forth a novel iterative linear-time procedure, meaning proportional to that required by the processor to scan the entire data $\{(\mathbf{a}_i; \psi_i)\}_{1 \leq i \leq m}$, which we term *reweighted amplitude flow* and abbreviate as RAF. Our methodology is capable of solving noiseless random quadratic equations exactly, and constructing an estimate of (near-)optimal statistical accuracy from noisy modulus observations. Exactness and accuracy hold with high probability and without any extra assumption on the signal \mathbf{x} to be recovered, provided that the ratio m/n of the number of measurements to that of the unknowns exceeds some large constant. Empirically, our procedure is demonstrated to be able to achieve perfect recovery of arbitrary high-dimensional signals given a minimal number of equations, which in the real case is $m = 2n - 1$. The new twist here is to leverage judiciously designed yet conceptually simple (iterative) (re)weighting regularization techniques to enhance existing initializations and also gradient refinements. An informal depiction of our RAF

methodology is given in two stages below, with rigorous algorithmic details deferred to Section III.

- S1) Weighted maximal correlation initialization:** Obtain an initialization \mathbf{z}^0 maximally correlated with a carefully selected subset $\mathcal{S} \subsetneq \mathcal{M} := \{1, 2, \dots, m\}$ of feature vectors \mathbf{a}_i , whose contributions toward constructing \mathbf{z}^0 are judiciously weighted by suitable parameters $\{w_i^0 > 0\}_{i \in \mathcal{S}}$; and
- S2) Iteratively reweighted “gradient-like” iterations:** Loop over $0 \leq t \leq T$

$$\mathbf{z}^{t+1} = \mathbf{z}^t - \frac{\mu^t}{m} \sum_{i=1}^m w_i^t \nabla \ell(\mathbf{z}^t; \psi_i) \quad (4)$$

for some time-varying weights $w_i^t \geq 0$ that are adapted in time, each depending on the current iterate \mathbf{z}_t and the datum $(\mathbf{a}_i; \psi_i)$.

Two attributes of our novel methodology are worth highlighting. First, albeit being a variant of the orthogonality-promoting initialization [9], the initialization here [cf. S1)] is distinct in the sense that different importance is attached to each selected datum $(\mathbf{a}_i; \psi_i)$, or more precisely, to each selected directional vector \mathbf{a}_i . Likewise, the gradient flow [cf. S2)] weighs judiciously the search direction suggested by each datum $(\mathbf{a}_i; \psi_i)$. In this manner, more accurate and robust initializations as well as more stable overall search directions in the gradient flow stage can be obtained even based only on a relatively limited number of data samples. Moreover, with particular choices of weights w_i^t 's (for example, when they take 0/1 values), our methodology subsumes as special cases the recently proposed truncated amplitude flow (TAF) [9], and the reshaped Wirtinger flow (RWF) [27].

II. ALGORITHM: REWEIGHTED AMPLITUDE FLOW

This section explains the intuition and the basic principles behind each stage of RAF in detail. For concreteness, we focus on the real Gaussian model with a real signal vector \mathbf{x} , and independent Gaussian random measurement vectors $\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, $1 \leq i \leq m$. Nevertheless, RAF can be applied without algorithmic changes for the complex Gaussian model with $\mathbf{x} \in \mathbb{C}^n$ and independent $\mathbf{a}_i \sim \mathcal{CN}(\mathbf{0}, \mathbf{I}_n) := \mathcal{N}(\mathbf{0}, \mathbf{I}_n/2) + j\mathcal{N}(\mathbf{0}, \mathbf{I}_n/2)$, and also when coded diffraction pattern (CDP) models [40] are considered.

A. Weighted Maximal Correlation Initialization

For general non-convex iterative heuristics to succeed in finding the global optimum is to seed them with an excellent starting point [43]. In fact, several smart initialization strategies have been advocated for iterative phase retrieval algorithms; see e.g., the spectral [24], [10], truncated spectral [12], [27], and orthogonality-promoting [9] initializations. One promising approach among them is the one proposed in [9], which is robust to outliers [37], and also enjoys better phase transitions than the spectral procedures [44]. To hopefully achieve perfect signal recovery at the information-theoretic limit however, its numerical performance may still need further enhancement. On the other hand, it is intuitive that improving the initialization perfor-

mance (over state-of-the-art procedures) becomes increasingly challenging as the number of acquired data samples approaches the information-theoretic limit of $m = 2n - 1$.

In this context, we develop below a more flexible initialization scheme based on the correlation property (as opposed to orthogonality), in which the added benefit relative to the initialization procedure in [9] is the inclusion of a flexible weighting regularization technique to better balance the useful information exploited in all selected data. In words, we introduce carefully designed weights to the initialization procedure developed in [9]. Similar to related approaches, our strategy entails estimating both the norm $\|\mathbf{x}\|$ and the unit direction $\mathbf{x}/\|\mathbf{x}\|$. Leveraging the strong law of large numbers and the rotational invariance of Gaussian \mathbf{a}_i sampling vectors (the latter suffices to assume $\mathbf{x} = \|\mathbf{x}\|\mathbf{e}_1$, with \mathbf{e}_1 being the first canonical vector in \mathbb{R}^n), it is clear that

$$\sum_{i=1}^m \psi_i^2 = \sum_{i=1}^m |\langle \mathbf{a}_i, \|\mathbf{x}\|\mathbf{e}_1 \rangle|^2 = \sum_{i=1}^m a_{i,1}^2 \|\mathbf{x}\|^2 \approx m \|\mathbf{x}\|^2 \quad (5)$$

whereby $\|\mathbf{x}\|$ can be estimated as $\sum_{i=1}^m \psi_i^2 / m$. This estimate proves very accurate even with an information-theoretic limit number of data samples, because it is unbiased and tightly concentrated.

The challenge thus lies in accurately estimating the direction of \mathbf{x} , or seeking a unit vector maximally aligned with \mathbf{x} , which is a bit tricky. To gain intuition for our initialization strategy, let us first present a variant of the initialization in [9], whose generalizations have been discussed in [37], [45]. Note that the larger the amplitude ψ_i of the inner-product between \mathbf{a}_i and \mathbf{x} is, the known design vector \mathbf{a}_i is deemed *more correlated* to the unknown solution \mathbf{x} , hence bearing useful directional information of \mathbf{x} . Inspired by this fact and based on available data $\{(\mathbf{a}_i; \psi_i)\}_{1 \leq i \leq m}$, one can sort all (absolute) correlation coefficients $\{\psi_i\}_{1 \leq i \leq m}$ in an ascending order, to yield ordered coefficients denoted by $0 < \psi_{[m]} \leq \dots \leq \psi_{[2]} \leq \psi_{[1]}$. Sorting m records takes time proportional to $\mathcal{O}(m \log m)$.² Let $\mathcal{S} \subsetneq \mathcal{M}$ represent the set of selected feature vectors \mathbf{a}_i to be used for computing the initialization, which is to be designed next. Fix *a priori* the cardinality $|\mathcal{S}|$ to some integer on the order of m , say $|\mathcal{S}| := \lfloor 3m/13 \rfloor$. It is then natural to *define* \mathcal{S} to collect the \mathbf{a}_i vectors that correspond to one of the largest $|\mathcal{S}|$ correlation coefficients $\{\psi_{[i]}\}_{1 \leq i \leq |\mathcal{S}|}$, each of which can be thought of as pointing to (roughly) the direction of \mathbf{x} . Approximating the direction of \mathbf{x} thus boils down to finding a vector to maximize its correlation with the subset \mathcal{S} of selected directional vectors \mathbf{a}_i . Succinctly stated, the wanted approximation vector can be efficiently found as the solution of

$$\underset{\|\mathbf{z}\|=1}{\text{maximize}} \quad \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} |\langle \mathbf{a}_i, \mathbf{z} \rangle|^2 = \mathbf{z}^* \left(\frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \mathbf{a}_i \mathbf{a}_i^* \right) \mathbf{z} \quad (6)$$

where the superscript $*$ represents transposition. Upon scaling the solution of (6) by the norm estimate $\sum_{i=1}^m \psi_i^2 / m$ in (5) to match the size of \mathbf{x} , we obtain what we will henceforth refer to as maximal correlation initialization.

²For a given function $g(n)$ of integer $n > 0$, $\mathcal{O}(g(n))$ denotes the set of functions $\mathcal{O}(g(n)) = \{f(n) : \text{there exist positive constants } C \text{ and } n_0 \text{ such that } 0 \leq f(n) \leq Cg(n) \text{ for all } n \geq n_0\}$.

As long as $|\mathcal{S}|$ is chosen on the order of m , the maximal correlation method outperforms the spectral ones in [10], [12], [24], and has comparable performance to the orthogonality-promoting method [9]. Its empirical performance around the information-theoretic limit however, is still not the best that we can hope for. Observe that all directional vectors $\{\mathbf{a}_i\}_{i \in \mathcal{S}}$ selected for forming the matrix $\bar{\mathbf{Y}} := (1/|\mathcal{S}|) \sum_{i \in \mathcal{S}} \mathbf{a}_i \mathbf{a}_i^*$ in (6) are treated *the same* in terms of their contributions to constructing the (direction of the) initialization. Nevertheless, according to our starting principle, this ordering information carried by the selected \mathbf{a}_i vectors has *not* been exploited by the initialization scheme in (6) (see also [9], [37]). In words, if for selected data $i, j \in \mathcal{S}$, the correlation coefficient of ψ_i with \mathbf{a}_i is larger than that of ψ_j with \mathbf{a}_j , then \mathbf{a}_i is deemed more correlated (with \mathbf{x}) than \mathbf{a}_j is, hence bearing more useful information about the wanted direction of \mathbf{x} . This prompts one to weight more (i.e., attach more importance to) the selected \mathbf{a}_i vectors corresponding to larger ψ_i values. Given the ordering information $\psi_{[|\mathcal{S}|]} \leq \dots \leq \psi_{[2]} \leq \psi_{[1]}$ available from the sorting procedure, a natural way to achieve this goal is by weighting each \mathbf{a}_i vector with simple functions of ψ_i , say e.g., taking the weights $w_i^0 := \psi_i^\gamma, \forall i \in \mathcal{S}$, with the parameter $\gamma \geq 0$ chosen to maintain the wanted ordering $w_{[|\mathcal{S}|]}^0 \leq \dots \leq w_{[2]}^0 \leq w_{[1]}^0$. In a nutshell, a more flexible initialization scheme, that we refer to as *weighted maximal correlation*, can be summarized as follows

$$\tilde{\mathbf{z}}_0 := \arg \max_{\|\mathbf{z}\|=1} \mathbf{z}^* \left(\frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \psi_i^\gamma \mathbf{a}_i \mathbf{a}_i^* \right) \mathbf{z}. \quad (7)$$

The upshot of (7) is that the objective can be efficiently minimized in time proportional to $\mathcal{O}(n|\mathcal{S}|)$ by means of the power method or the Lanczos algorithm [46]. The proposed initialization can be obtained after scaling $\tilde{\mathbf{z}}^0$ from (7) with the estimate of its norm, to obtain $\mathbf{z}^0 := (\sum_{i=1}^m \psi_i^2/m) \tilde{\mathbf{z}}^0$. By default, we take $\gamma := 1/2$ in all reported numerical implementations, yielding weights $w_i^0 := \sqrt{|\langle \mathbf{a}_i, \mathbf{x} \rangle|}$ for all $i \in \mathcal{S}$.

Regarding the initialization procedure in (7), we next highlight two features, while details and theoretical performance guarantees are provided in Section III:

- F1)** The weights $\{w_i^0\}$ in the maximal correlation scheme enable leveraging useful information that each feature vector \mathbf{a}_i may bear regarding the direction of \mathbf{x} .
- F2)** Taking $w_i^0 := \psi_i^\gamma$ for all $i \in \mathcal{S}$ and 0 otherwise, (7) can be equivalently rewritten as

$$\tilde{\mathbf{z}}^0 := \arg \max_{\|\mathbf{z}\|=1} \mathbf{z}^* \left(\frac{1}{m} \sum_{i=1}^m w_i^0 \mathbf{a}_i \mathbf{a}_i^* \right) \mathbf{z} \quad (8)$$

which subsumes existing initialization schemes with particular weight selections; e.g., the “plain-vanilla” spectral initialization in [10], [24] is recovered by choosing $\mathcal{S} := \mathcal{M}$, and $w_i^0 := \psi_i^2, \forall i = 1, \dots, m$.

For numerical comparison, define the Relative error $:= \text{dist}(\mathbf{z}, \mathbf{x})/\|\mathbf{x}\|$. All simulated tests reported here were averaged over 100 Monte Carlo realizations. Fig. 1 depicts the performance of the proposed initialization relative to several state-of-the-art strategies, and also with the information limit number benchmarking the minimal number of samples required. It is clear that our initialization is: i) consistently better than the

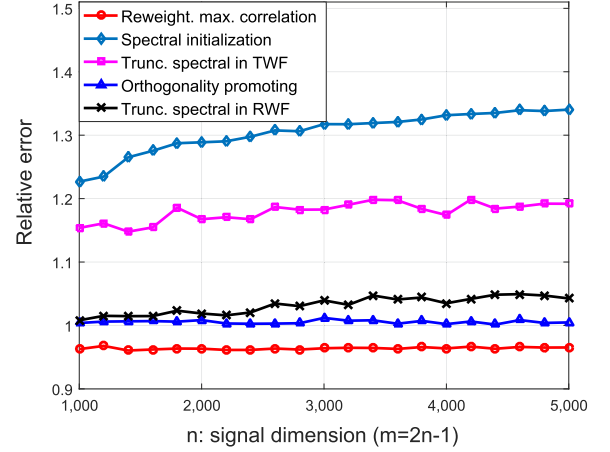


Fig. 1. Relative initialization error for the real Gaussian model with $n = 1,000$ and $m = 2n - 1 = 1,999$.

state-of-the-art; and, ii) stable as the signal dimension n grows, which is in sharp contrast to the instability encountered by the spectral ones [10], [12], [24], [27]. It is also worth stressing that about 5% empirical advantage is shown over the best in [9] at the challenging information-theoretic benchmark, which is indeed nontrivial, and constitutes one of the main advantages of RAF. This numerical advantage becomes increasingly pronounced as the ratio m/n grows. This suggests that our proposed initialization procedure may be combined with other iterative phase retrieval approaches to improve their numerical performance.

B. Adaptively Reweighted Gradient Flow

For independent data adhering to the real Gaussian model, the direction that TAF moves along in stage S2) presented earlier is given by the following (generalized) gradient [9], [41]

$$\frac{1}{m} \sum_{i \in \mathcal{T}} \nabla \ell(\mathbf{z}; \psi_i) = \frac{1}{m} \sum_{i \in \mathcal{T}} \left(\mathbf{a}_i^* \mathbf{z} - \psi_i \frac{\mathbf{a}_i^* \mathbf{z}}{|\mathbf{a}_i^* \mathbf{z}|} \right) \mathbf{a}_i \quad (9)$$

where the dependence on the iterate count t is neglected for notational brevity, and the convention $\frac{\mathbf{a}_i^* \mathbf{z}}{|\mathbf{a}_i^* \mathbf{z}|} := 0$ is adopted if $\mathbf{a}_i^* \mathbf{z} = 0$.

Unfortunately, the (negative) gradient of the average in (9) may not point towards the true \mathbf{x} , unless the current iterate \mathbf{z} is already very close to \mathbf{x} . As a consequence, moving along such a descent direction may not drag \mathbf{z} closer to \mathbf{x} . To see this, consider an initial guess \mathbf{z}_0 that has already been in a basin of attraction (i.e., a region within which there is only a unique stationary point) of \mathbf{x} . Certainly, there are summands $(\mathbf{a}_i^* \mathbf{z} - \psi_i \frac{\mathbf{a}_i^* \mathbf{z}}{|\mathbf{a}_i^* \mathbf{z}|}) \mathbf{a}_i$ in (9), that could give rise to “bad/misleading” search directions due to the erroneously estimated signs $\frac{\mathbf{a}_i^* \mathbf{z}}{|\mathbf{a}_i^* \mathbf{z}|} \neq \frac{\mathbf{a}_i^* \mathbf{x}}{|\mathbf{a}_i^* \mathbf{x}|}$ in (9) [9]. Those gradients as a whole may drag \mathbf{z} away from \mathbf{x} , and hence out of the basin of attraction. Such an effect becomes increasingly severe as the number m of acquired examples approaches the information-theoretic limit of $2n - 1$, thus rendering past approaches less effective in this case. Although this issue is somewhat remedied by TAF with a truncation procedure, its efficacy is limited due to misses of bad gradients and mis-rejections of meaningful ones at the information-theoretic limit.

To address this challenge, our reweighted gradient flow effecting suitable search directions from *almost all* acquired data samples $\{(\mathbf{a}_i; \psi_i)\}_{1 \leq i \leq m}$ will be adopted in a (timely) adaptive fashion; that is,

$$\mathbf{z}^{t+1} = \mathbf{z}^t - \mu^t \nabla \ell_{\text{rw}}(\mathbf{z}^t; \psi_i), \quad t = 0, 1, \dots \quad (10)$$

The *reweighted gradient* $\nabla \ell_{\text{rw}}(\mathbf{z}^t)$ evaluated at the current point \mathbf{z}^t is given as

$$\nabla \ell_{\text{rw}}(\mathbf{z}) := \frac{1}{m} \sum_{i=1}^m w_i \nabla \ell(\mathbf{z}; \psi_i) \quad (11)$$

for suitable weights $\{w_i\}_{1 \leq i \leq m}$ to be designed shortly.

Toward that end, we observe that the truncation criterion $\mathcal{T} := \{1 \leq i \leq m : |\mathbf{a}_i^* \mathbf{z}| \geq \alpha |\mathbf{a}_i^* \mathbf{x}|\}$ with some given parameter $\alpha > 0$ suggests to include only gradients associated with $|\mathbf{a}_i^* \mathbf{z}|$ of relatively large sizes. This is because gradients of sizable $\frac{|\mathbf{a}_i^* \mathbf{z}|}{|\mathbf{a}_i^* \mathbf{x}|}$ offer reliable and meaningful directions pointing to the true \mathbf{x} with large probability [9]. As such, the ratio $\frac{|\mathbf{a}_i^* \mathbf{z}|}{|\mathbf{a}_i^* \mathbf{x}|}$ can be viewed as a confidence score on the reliability or meaningfulness of the corresponding gradient $\nabla \ell(\mathbf{z}; \psi_i)$. Recognizing that confidence can vary, it is natural to distinguish the contributions that different gradients make to the overall search direction. An easy way is to attach large weights to the reliable gradients, and small weights to the spurious ones. Assume without loss of generality that $0 \leq w_i \leq 1$ for all $1 \leq i \leq m$; otherwise, lump the normalization factor achieving this into the learning rate μ^t . Building upon this observation and leveraging the gradient reliability confidence score $\frac{|\mathbf{a}_i^* \mathbf{z}|}{|\mathbf{a}_i^* \mathbf{x}|}$, the weight per gradient $\nabla \ell(\mathbf{z}; \psi_i)$ in our proposed RAF algorithm is

$$w_i := \frac{1}{1 + \beta_i / (|\mathbf{a}_i^* \mathbf{z}| / |\mathbf{a}_i^* \mathbf{x}|)}, \quad 1 \leq i \leq m \quad (12)$$

where $\{\beta_i > 0\}_{1 \leq i \leq m}$ are some pre-selected parameters.

Regarding the weighting criterion in (28), three remarks are in order.

Remark 1: The weights $\{w_i\}_{1 \leq i \leq m}$ are time adapted to the iterate \mathbf{z}^t . One can also interpret the reweighted gradient flow \mathbf{z}^{t+1} in (10) as performing a single gradient step to minimize the *smooth reweighted loss* $(1/m) \sum_{i=1}^m w_i \ell(\mathbf{z}; \psi_i)$ with starting point \mathbf{z}^t ; see also [47] for related ideas successfully exploited in the *iteratively reweighted least-squares* approach to compressive sampling.

Remark 2: The larger the confidence score $\frac{|\mathbf{a}_i^* \mathbf{z}|}{|\mathbf{a}_i^* \mathbf{x}|}$ is, the larger the corresponding weight w_i will be. More importance will be then attached to reliable gradients than to spurious ones. Gradients from *almost all* data are accounted for, which is in contrast to [9], where withdrawn gradients do not contribute the information they carry.

Remark 3: At the points $\{\mathbf{z}\}$ where $\mathbf{a}_i^* \mathbf{z} = 0$ for some datum $i \in \mathcal{M}$, the i -th weight will be $w_i = 0$. In other words, the squared losses $\ell(\mathbf{z}; \psi_i)$ in (2) that are non-smooth at points \mathbf{z} will be eliminated, to prevent their contribution to the reweighted gradient update in (10). This simplifies the convergence analysis of RAF considerably because it does not have to cope with the non-smoothness of the objective function in (2).

Having elaborated on the two stages, RAF can be readily summarized in Algorithm 1.

Algorithm 1: Reweighted Amplitude Flow (RAF).

- 1: **Input:** Data $\{(\mathbf{a}_i; \psi_i)\}_{1 \leq i \leq m}$; maximum number of iterations T ; step sizes $\mu^t = 2/6$ and weighting parameters $\beta_i = 10/5$ for real and complex Gaussian models; subset cardinality $|\mathcal{S}| = \lfloor 3m/13 \rfloor$, and exponent $\gamma = 0.5$.
- 2: **Construct** \mathcal{S} to include indices associated with the $|\mathcal{S}|$ largest entries among $\{\psi_i\}_{1 \leq i \leq m}$.
- 3: **Initialize** $\mathbf{z}^0 := \sqrt{\sum_{i=1}^m \psi_i^2 / m} \tilde{\mathbf{z}}^0$ with $\tilde{\mathbf{z}}^0$ being the unit-norm principal eigenvector of

$$\frac{1}{m} \sum_{i=1}^m w_i^0 \mathbf{a}_i \mathbf{a}_i^*, \quad \text{where } w_i^0 := \begin{cases} \psi_i^\gamma, & i \in \mathcal{S} \subseteq \mathcal{M} \\ 0, & \text{otherwise.} \end{cases}$$

- 4: **Loop:** for $t = 0$ to $T - 1$

$$\mathbf{z}^{t+1} = \mathbf{z}^t - \frac{\mu^t}{m} \sum_{i=1}^m w_i^t \left(\mathbf{a}_i^* \mathbf{z}^t - \psi_i \frac{\mathbf{a}_i^* \mathbf{z}^t}{|\mathbf{a}_i^* \mathbf{z}^t|} \right) \mathbf{a}_i \quad (13)$$

$$\text{where } w_i^t := \frac{|\mathbf{a}_i^* \mathbf{z}^t| / \psi_i}{|\mathbf{a}_i^* \mathbf{z}^t| / \psi_i + \beta_i} \text{ for all } 1 \leq i \leq m.$$

- 5: **Output:** \mathbf{z}^T .
-

C. Parameters of the Algorithm

To optimize the empirical performance and facilitate numerical implementations, the choice of pertinent RAF parameters is outlined here. For the four RAF parameters, our theory and experiments are based on: i) $|\mathcal{S}|/m \leq 0.25$; ii) $0 \leq \beta_i \leq 10$ for all $1 \leq i \leq m$; and, iii) $0 \leq \gamma \leq 1$. For convenience, a constant step size $\mu^t \equiv \mu > 0$ is suggested, but other step size rules such as backtracking line search with the reweighted objective would work as well. As will be formalized in Section III, RAF converges if the constant μ is not too large, with the upper bound depending in part on the selection of $\{\beta_i\}_{1 \leq i \leq m}$.

In the numerical tests presented in Sections II and IV, we take $|\mathcal{S}| := \lfloor 3m/13 \rfloor$, $\beta_i \equiv \beta := 10$, $\gamma := 0.5$, and $\mu := 2$ (larger step sizes can be afforded for larger m/n values).

III. MAIN RESULTS

Our main results summarized below establish exact recovery under the real Gaussian model, whose proof is postponed to Section V for readability. Our RAF methodology however, can be generalized to the complex Gaussian as well as the CDP models.

Theorem 1 (Exact recovery): Consider m noiseless measurements $\psi = |\mathbf{A}\mathbf{x}|$ for an arbitrary signal $\mathbf{x} \in \mathbb{R}^n$. If $m \geq c_0 |\mathcal{S}| \geq c_1 n$ with $|\mathcal{S}|$ being the pre-selected subset cardinality in the initialization step and the learning rate $\mu \leq \mu_0$, then with probability at least $1 - c_3 e^{-c_2 m}$, the RAF estimates \mathbf{z}^t in Algorithm 1 obey

$$\text{dist}(\mathbf{z}^t, \mathbf{x}) \leq \frac{1}{10} (1 - \nu)^t \|\mathbf{x}\|, \quad t = 0, 1, \dots \quad (14)$$

where $c_0, c_1, c_2, c_3 > 0$, $0 < \nu < 1$, and $\mu_0 > 0$ are certain numerical constants depending on the choice of algorithmic parameters $|\mathcal{S}|$, β , γ , and μ .

According to Theorem 1, a few interesting properties of our RAF algorithm are worth highlighting. To start, RAF recovers

the true solution exactly with high probability whenever the ratio m/n of the number of equations to the unknowns exceeds some numerical constant. Expressed differently, RAF achieves the information-theoretic optimal order of sample complexity, which is consistent with the state-of-the-art including truncated Wirtinger flow (TWF) [12], TAF [9], and RWF [27]. Notice that the error contraction in (14) also holds at $t = 0$, namely $\text{dist}(\mathbf{z}^0, \mathbf{x}) \leq \|\mathbf{x}\|/10$, therefore providing theoretical performance guarantees for the proposed initialization strategy (cf. Step 1 of Algorithm 1). Moreover, starting from this initial estimate, RAF converges exponentially fast to the true solution \mathbf{x} . In other words, to reach any ϵ -relative solution accuracy (i.e., $\text{dist}(\mathbf{z}^T, \mathbf{x}) \leq \epsilon \|\mathbf{x}\|$), it suffices to run at most $T = \mathcal{O}(\log 1/\epsilon)$ RAF iterations in Step 1 of Algorithm 1. This in conjunction with the per-iteration complexity $\mathcal{O}(mn)$ (namely, the complexity of one reweighted gradient update in (66)) confirms that RAF solves exactly a quadratic system in time $\mathcal{O}(mn \log 1/\epsilon)$, which is linear in $\mathcal{O}(mn)$, the time required by the processor to read the entire data $\{(\mathbf{a}_i; \psi_i)\}_{1 \leq i \leq m}$. Given the fact that the initialization stage can be performed in time $\mathcal{O}(n|S|)$ and $|S| < m$, the overall linear-time complexity of RAF is order-optimal.

IV. SIMULATED TESTS

Our theoretical findings about RAF have been corroborated with comprehensive numerical experiments, a sample of which are presented next. Performance of RAF is evaluated relative to the state-of-the-art (T)WF [10], [12], RWF [27], and TAF [9] in terms of the empirical success rate among 100 MC realizations, where a success will be declared for an independent trial if the returned estimate incurs error $\|\psi - \mathbf{A}\mathbf{z}^T\|/\|\mathbf{x}\| \leq 10^{-5}$. Both the real Gaussian and the physically realizable CDP models were simulated. For fairness, all procedures were implemented with their suggested parameter values. We generated the true $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, and i.i.d. measurement vectors $\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, $1 \leq i \leq m$. Each iterative scheme obtained its initial guess based on 200 power or Lanczos iterations, followed by a sequence of $T = 2,000$ (which can be set smaller as the ratio m/n grows away from the limit of 2) gradient-type iterations. All the numerical experiments in this paper were implemented with MATLAB R2016a on an Intel CPU @ 3.4 GHz (32 GB RAM) computer. For reproducibility, the Matlab code of our RAF algorithm is publicly available at <https://gangwg.github.io/RAF>.

To examine how the parameter value of γ in (7) influences our initialization performance, the relative error versus the parameter value ranging from 0 to 1 is presented in Fig. 2, where the real Gaussian model is simulated with n varying from 1,000 to 5,000 and $m = 2n - 1$ fixed. Evidently, the plots clearly validate our choice of the default parameter value $\gamma = 0.5$.

To show the power of RAF in the high-dimensional regime, the function value $L(\mathbf{z})$ in (2) evaluated at the returned estimate \mathbf{z}^T (cf. Step 1 of Algorithm 1) after 200 MC realizations is plotted (in negative logarithmic scale) in Fig. 3, where the number of simulated noiseless measurements was set to be the information-theoretic limit, namely $m = 2n - 1 = 3,999$ for $n = 2,000$. It is evident that our proposed RAF approach returns a solution of function value $L(\mathbf{z}^T)$ smaller than 10^{-25} in all 200 independent realizations even at this challenging

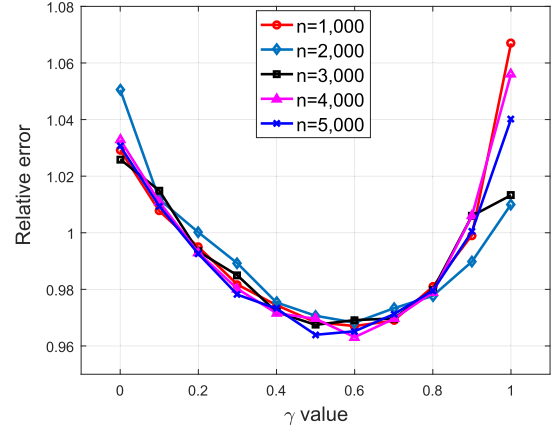


Fig. 2. Relative error versus γ for the proposed initialization scheme with n varying from 1,000 to 5,000 and $m = 2n - 1$ fixed under the real Gaussian model.

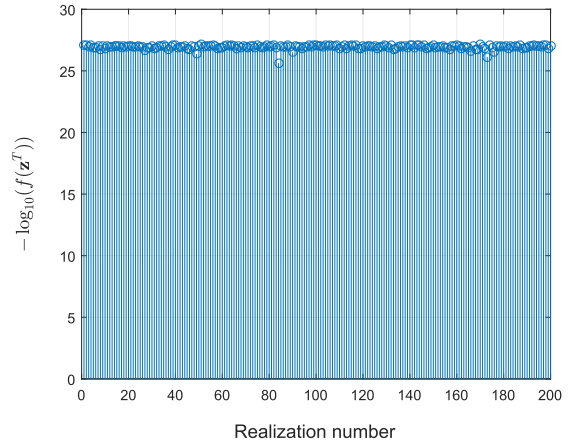


Fig. 3. Function value $L(\mathbf{z}^T)$ evaluated at the returned RAF estimate \mathbf{z}^T for 200 trials with $n = 2,000$ and $m = 2n - 1 = 3,999$.

information-theoretic limit condition. To the best of our knowledge, RAF is the first algorithm that empirically reconstructs any high-dimensional (say e.g., $n \geq 1,500$) signals exactly from a minimal number of random quadratic equations, which also provides a positive answer to the question posed earlier in the Introduction.

Fig. 4 compares the empirical success rate of the five schemes with the signal dimension being fixed at $n = 1,000$ while the ratio m/n increasing by 0.1 from 1 to 5. Specifically, in the top panel, each scheme uses its own initialization, while in the bottom panel, all schemes start with the same maximally reweighted correlation initialization. As clearly depicted by the plots, our RAF approach (color coded red) outperforms its competing alternatives in both cases. Moreover, it also achieves 100% signal recovery as soon as m is about $2n$, where the others do not show perfect recovery. Through comparing the two figures, it is clear that the performance of TAF and TWF can benefit from using the proposed initialization.

Fig. 5 further compares the convergence speed of various schemes in terms of the number of iterations to produce solutions of a given accuracy. Evidently, RAF converges faster than WF and TWF, and it has comparable efficiency as TAF and RWF when using the real Gaussian model with $\mathbf{x} \in \mathbb{R}^{1,000}$ and $m =$

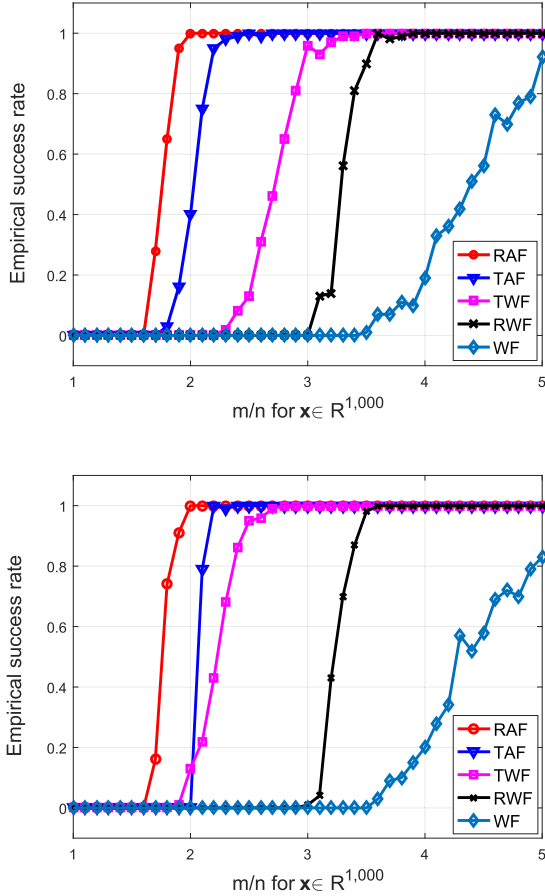


Fig. 4. Empirical success rate under the real Gaussian model using: different initializations (top); and, the same reweighted maximal correlation initialization (bottom).

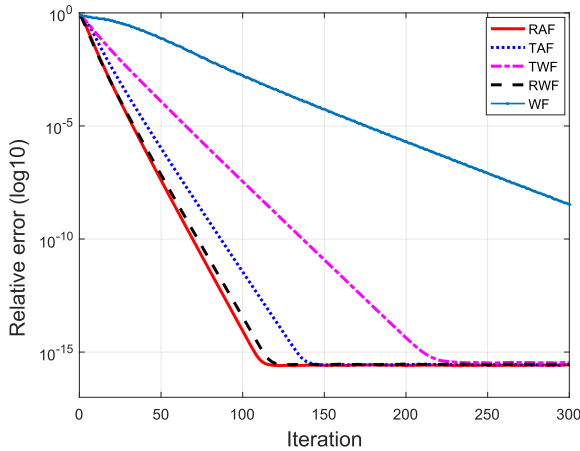


Fig. 5. Relative error versus iterations using: i) RAF; ii) TAF; iii) TWF; iv) RWF; and v) WF with $n = 1,000$ and $m/n = 5$ under the real Gaussian model.

5,000. Regarding running times, to reach solution accuracy of relative error 10^{-15} or a maximum of 500 iterations, the computational times for RAF, TAF, TWF, RWF, and WF are 0.63 s, 1.12 s, 1.49 s, 0.94 s, and 19.16 s, respectively.

To numerically demonstrate the stability and robustness of RAF in the presence of additive noise, Fig. 6 examines the normalized mean-square error $\text{NMSE} := \text{dist}^2(\mathbf{z}^T, \mathbf{x})/\|\mathbf{x}\|^2$ as a function of the signal-to-noise ratio (SNR) for m/n tak-

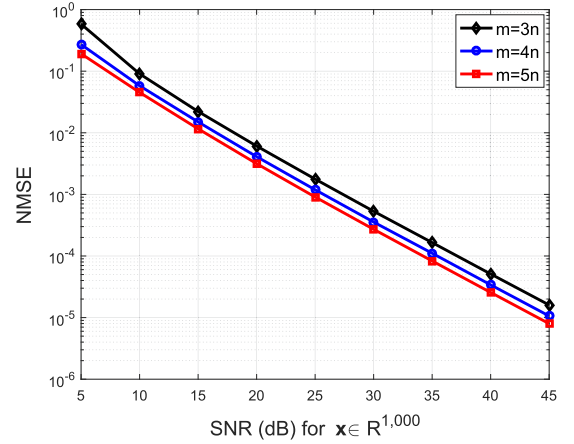


Fig. 6. NMSE vs. SNR for RAF under the real Gaussian model.

ing values $\{3, 4, 5\}$. The noise model $\psi_i = |\langle \mathbf{a}_i, \mathbf{x} \rangle| + \eta_i$ with $\boldsymbol{\eta} := [\eta_i]_{1 \leq i \leq m} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_m)$ was simulated, where σ^2 was set such that certain $\text{SNR} := 10 \log_{10}(\|\mathbf{A}\mathbf{x}\|^2/m\sigma^2)$ values were achieved. For all choices of m (as small as $3n$ which is nearly minimal), the numerical experiments illustrate that the NMSE scales inversely proportional to the SNR, which corroborates the stability of our RAF approach.

To demonstrate the efficacy and scalability of RAF in real-world conditions, the last experiment entails the Galaxy image³ depicted by a three-way array $\mathbf{X} \in \mathbb{R}^{1,080 \times 1,920 \times 3}$, whose first two coordinates encode the pixel locations, and the third the RGB color bands. Consider the physically realizable CDP model with random masks [10]. Letting $\mathbf{x} \in \mathbb{R}^n$ ($n \approx 2 \times 10^6$) be a vectorization of a certain band of \mathbf{X} , the CDP model with K masks is

$$\psi^{(k)} = |\mathbf{F}\mathbf{D}^{(k)}\mathbf{x}|, \quad 1 \leq k \leq K \quad (15)$$

where $\mathbf{F} \in \mathbb{C}^{n \times n}$ is a discrete Fourier transform matrix, and diagonal matrices $\mathbf{D}^{(k)}$ have their diagonal entries sampled uniformly at random from $\{1, -1, j, -j\}$ with $j := \sqrt{-1}$. Implementing $K = 4$ masks, each algorithm performs independently over each band 100 power iterations to obtain the initial guess, which was refined by 100 gradient iterations. Recovered images of TAF (top) and RAF (bottom) are displayed in Fig. 7, whose relative errors were 1.0347 and 1.0715×10^{-3} , respectively. WF and TWF returned images of corresponding relative error 1.6870 and 1.4211 , which are far away from the ground truth.

It is worth pointing out that RAF converges faster both in time and in the number of iterations required to achieve certain solution accuracy than TWF and WF in all our simulated experiments, and it has comparable computational efficiency as TAF and RWF.

V. PROOFS

To prove Theorem 1, this section establishes a few lemmas and the main ideas, whereas technical details are postponed to the Appendix to facilitate readability. It is clear from Algorithm 1 that the weighted maximal correlation initialization (cf. Step 3) and the reweighted gradient flow (cf. Step 4)

³Downloaded from <http://pics-about-space.com/milky-way-galaxy>.

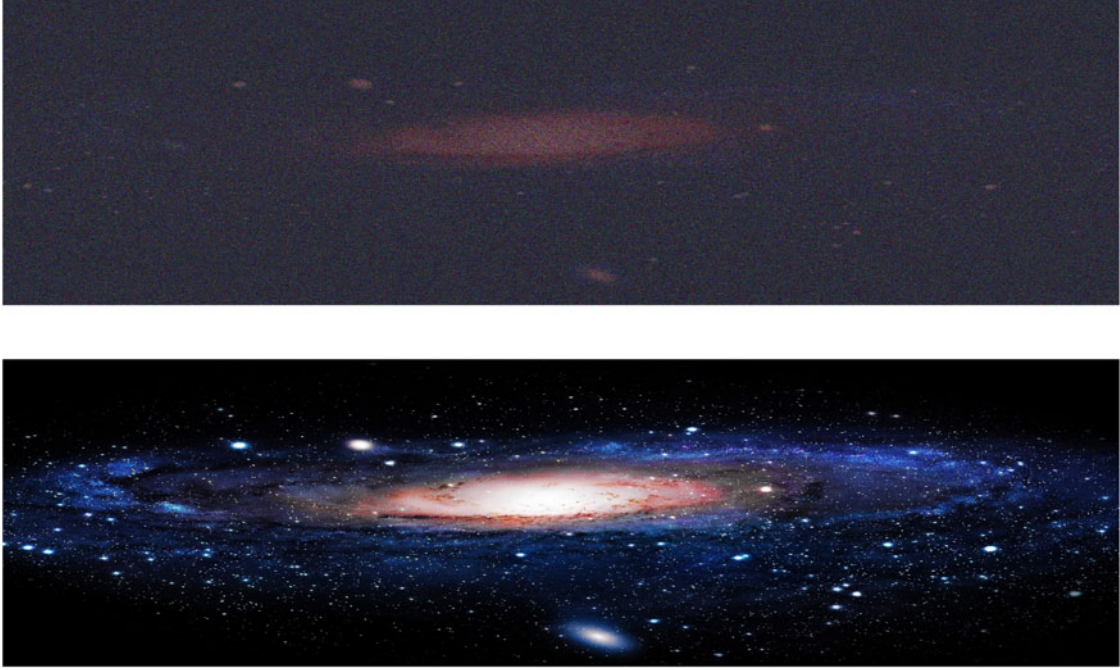


Fig. 7. The recovered Galaxy images after 100 truncated gradient iterations of TAF (top); and after 100 reweighted gradient iterations of RAF (bottom).

distinguish themselves from those procedures in (T)WF [10], [12], TAF [9], and RWF [27]. Hence, new proof techniques to cope with the weighting in both the initialization and the gradient flow, as well as the non-smoothness and non-convexity of the amplitude-based least-squares functional are required. Nevertheless, part of the proof is built upon those in [9], [10], [27], [48].

The proof of Theorem 1 consists of two parts: Section V-A below asserts guaranteed theoretical performance of the proposed initialization, which essentially achieves any given constant relative error as soon as the number of equations is on the order of the number of unknowns; that is, $m \geq c_1 n$ for some constant $c_1 > 0$. It is worth mentioning that we reserve c and its subscripted versions for absolute constants, even though their values may vary with the context. Under the sample complexity of order $\mathcal{O}(n)$, Section V-B further shows that RAF converges to the true signal \mathbf{x} exponentially fast whenever the initial estimate lands within a relatively small-size neighborhood of \mathbf{x} defined by $\text{dist}(\mathbf{z}^0, \mathbf{x}) \leq (1/10)\|\mathbf{x}\|$.

A. Weighted Maximal Correlation Initialization

This section is devoted to establishing analytical guarantees for the novel initialization procedure, which is summarized in the following proposition.

Proposition 1: For an arbitrary $\mathbf{x} \in \mathbb{R}^n$, consider the noiseless measurements $\psi_i = |\mathbf{a}_i^* \mathbf{x}|$, $1 \leq i \leq m$. If $m \geq c_0 |\mathcal{S}| \geq c_1 n$, then with probability exceeding $1 - c_3 e^{-c_2 m}$, the initial guess \mathbf{z}^0 obtained by the weighted maximal correlation method in Step 3 of Algorithm 1 satisfies

$$\text{dist}(\mathbf{z}^0, \mathbf{x}) \leq \rho \|\mathbf{x}\| \quad (16)$$

for $\rho = 1/10$ (or any sufficiently small positive number). Here, $c_0, c_1, c_2, c_3 > 0$ are some absolute constants.

Due to the homogeneity, it suffices to prove the result when $\|\mathbf{x}\| = 1$. Assume first that the norm $\|\mathbf{x}\| = 1$ is also perfectly known, and \mathbf{z}^0 has already been scaled such that $\|\mathbf{z}^0\| = 1$. At the end of this proof, this approximation error between the actually employed norm estimate $\sqrt{\sum_{i=1}^m y_i/m}$ found based on the strong law of large numbers and the unknown norm $\|\mathbf{x}\| = 1$, will be taken care of. Consider independent Gaussian random measurement vectors $\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ and an arbitrary unit-norm vector \mathbf{x} . Since Gaussian distributions are rotationally invariant, it suffices to prove the results for $\mathbf{x} = \mathbf{e}_1$, where \mathbf{e}_1 is the first canonical vector in \mathbb{R}^n .

Since the norm $\|\mathbf{x}\| = 1$ is assumed known, the weighted maximal correlation initialization in Step 3 finds the initial estimate $\mathbf{z}^0 = \tilde{\mathbf{z}}^0$ (the scaling factor is the exactly known norm 1 in this case) as the principal eigenvector of

$$\frac{1}{|\mathcal{S}|} \mathbf{B}^* \mathbf{B} = \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \psi_i^\gamma \mathbf{a}_i \mathbf{a}_i^* \quad (17)$$

where $\mathbf{B} := [\psi_i^{\gamma/2} \mathbf{a}_i]_{i \in \mathcal{S}}$ is an $|\mathcal{S}| \times n$ matrix, and $\mathcal{S} \subsetneq \{1, 2, \dots, m\}$ includes the indices of the $|\mathcal{S}|$ largest entities among all modulus data $\{\psi_i\}_{1 \leq i \leq m}$. The following result is a modification of [9, Lemma 1], which is key to proving Proposition 1.

Lemma 1: Consider m noiseless measurements $\psi_i = |\mathbf{a}_i^* \mathbf{x}|$, $1 \leq i \leq m$. For an arbitrary $\mathbf{x} \in \mathbb{R}^n$ of unit norm, the next result holds for all unit-norm vectors $\mathbf{u} \in \mathbb{R}^n$ perpendicular to \mathbf{x} ; that is, for all $\mathbf{u} \in \mathbb{R}^n$ satisfying $\mathbf{u}^* \mathbf{x} = 0$ and $\|\mathbf{u}\| = 1$, we have

$$\frac{1}{2} \|\mathbf{x} \mathbf{x}^* - \tilde{\mathbf{z}}^0 (\tilde{\mathbf{z}}^0)^*\|_F^2 \leq \frac{\|\mathbf{B} \mathbf{u}\|^2}{\|\mathbf{B} \mathbf{x}\|^2} \quad (18)$$

where $\tilde{\mathbf{z}}^0$ is given by

$$\tilde{\mathbf{z}}^0 := \arg \max_{\|\mathbf{z}\|=1} \frac{1}{|\mathcal{S}|} \mathbf{z}^* \mathbf{B}^* \mathbf{B} \mathbf{z}. \quad (19)$$

Let us start with the proof of Proposition 1. The first step consists in upper-bounding the quantity on the right-hand-side of (18). This involves upper bounding its numerator, and lower bounding its denominator, tasks summarized in Lemmas 2 and 3, whose proofs are deferred to Appendix B and Appendix C, accordingly.

Lemma 2: In the setting of Lemma 1, if $|\mathcal{S}|/n \geq c_4$, then the inequality

$$\|\mathbf{B}\mathbf{u}\|^2 \leq 1.01\sqrt{2^\gamma/\pi}\Gamma(\gamma+1/2)|\mathcal{S}| \quad (20)$$

holds with probability at least $1 - 2e^{-c_5 n}$, where $\Gamma(\cdot)$ is the Gamma function, and c_4, c_5 are certain universal constants.

Lemma 3: In the setting of Lemma 1, the following holds with probability exceeding $1 - e^{-c_6 m}$:

$$\|\mathbf{B}\mathbf{x}\|^2 \geq 0.99 \times 1.14^\gamma |\mathcal{S}| [1 + \log(m/|\mathcal{S}|)] \quad (21)$$

provided that $m \geq c_0 |\mathcal{S}| \geq c_1 n$ for some absolute constants $c_0, c_1, c_6 > 0$.

Taking together, the upper bound in (20) and the lower bound in (21), one arrives at

$$\frac{\|\mathbf{B}\mathbf{u}\|^2}{\|\mathbf{B}\mathbf{x}\|^2} \leq \frac{C}{1 + \log(m/|\mathcal{S}|)} \triangleq \kappa \quad (22)$$

where $C := 1.02 \times 1.14^{-\gamma} \sqrt{2^\gamma/\pi} \Gamma(\gamma+1/2)$, and (22) holds with probability at least $1 - 2e^{-c_5 n} - e^{-c_6 m}$, with the proviso that $m \geq c_0 |\mathcal{S}| \geq c_1 n$. Since $m = \mathcal{O}(n)$, one can rewrite the probability as $1 - c_3 e^{-c_2 m}$ for certain constants $c_2, c_3 > 0$. To have a sense of the size of C , taking our default value $\gamma = 0.5$ for instance gives rise to $C = 0.7854$.

It is clear that the bound κ in (22) can be rendered arbitrarily small by taking sufficiently large $m/|\mathcal{S}|$ values (while maintaining $|\mathcal{S}|/n$ to be some constant based on Lemma 3). With no loss of generality, let us work with $\kappa := 0.001$ in the following.

The wanted upper bound on the distance between the initialization \mathbf{z}^0 and the true \mathbf{x} can be obtained based upon similar arguments found in [10, Section 7.8], which are delineated next. For unit-norm \mathbf{x} and $\mathbf{z}^0 = \tilde{\mathbf{z}}^0$, if $0 \leq \theta \leq \pi/2$ denotes the angle between the spaces spanned by \mathbf{z}^0 and \mathbf{x} , using (18) and (22) yields

$$\begin{aligned} |\mathbf{x}^* \mathbf{z}^0|^2 &= \cos^2 \theta = 1 - \sin^2 \theta \\ &= 1 - \frac{\|\mathbf{B}\mathbf{u}\|^2}{\|\mathbf{B}\mathbf{x}\|^2} \\ &\geq 1 - \kappa \end{aligned} \quad (23)$$

thus giving rise to

$$\begin{aligned} \text{dist}^2(\mathbf{z}^0, \mathbf{x}) &\leq \|\mathbf{z}^0\|^2 + \|\mathbf{x}\|^2 - 2|\mathbf{x}^* \mathbf{z}^0| \\ &\leq (2 - 2\sqrt{1 - \kappa}) \|\mathbf{x}\|^2 \\ &\approx \kappa \|\mathbf{x}\|^2. \end{aligned} \quad (24)$$

As discussed prior to Lemma 1, the exact norm $\|\mathbf{x}\| = 1$ is generally unknown, and one often scales the unit-norm directional vector found in (19) by the estimate $\sqrt{\sum_{i=1}^m \psi_i^2/m}$. Next, the approximation error between the estimated norm $\|\mathbf{z}^0\| = \sqrt{\sum_{i=1}^m \psi_i^2/m}$ and the true norm $\|\mathbf{x}\| = 1$ is accounted for. Recall from (19) that the direction of \mathbf{x} is estimated to

be $\tilde{\mathbf{z}}^0$ (of unit norm). Using results similar to those in [10, Lemma 7.8 and Section 7.8], the following holds with high probability, as long as the ratio m/n exceeds some numerical constant

$$\|\mathbf{z}^0 - \tilde{\mathbf{z}}^0\| = \|\|\mathbf{z}^0\| - 1\| \leq (1/20)\|\mathbf{x}\|. \quad (25)$$

Taking the inequalities in (24) and (25) together, it is safe to deduce that

$$\text{dist}(\mathbf{z}^0, \mathbf{x}) \leq \|\mathbf{z}^0 - \tilde{\mathbf{z}}^0\| + \text{dist}(\tilde{\mathbf{z}}^0, \mathbf{x}) \leq (1/10)\|\mathbf{x}\| \quad (26)$$

which confirms that the initial estimate obeys the relative error $\text{dist}(\mathbf{z}^0, \mathbf{x})/\|\mathbf{x}\| \leq 1/10$ for any $\mathbf{x} \in \mathbb{R}^n$ with probability $1 - c_3 e^{-c_2 m}$, provided that $m \geq c_0 |\mathcal{S}| \geq c_1 n$ for some numerical constants $c_0, c_1, c_2, c_3 > 0$.

B. Exact Phase Retrieval From Noiseless Data

It has been demonstrated that the initial estimate \mathbf{z}^0 obtained by means of the weighted maximal correlation initialization strategy has at most a constant relative error to the globally optimal solution \mathbf{x} , i.e., $\text{dist}(\mathbf{z}^0, \mathbf{x}) \leq (1/10)\|\mathbf{x}\|$. We demonstrate in the following that starting from such an initial estimate, the RAF iterates (in Step 4 of Algorithm 1) converge at a linear rate to the global optimum \mathbf{x} ; that is, $\text{dist}(\mathbf{z}^t, \mathbf{x}) \leq (1/10)c^t \|\mathbf{x}\|$ for some constant $0 < c < 1$ depending on the step size $\mu > 0$, the weighting parameter β , and the data $\{(\mathbf{a}_i; \psi_i)\}_{1 \leq i \leq m}$. This constitutes the second part of the proof of Theorem 1. Toward this end, it suffices to show that the iterative update of RAF is locally contractive within a relatively small neighboring region of the true \mathbf{x} . Instead of directly coping with the moments in the weights, we establish a conservative result based directly on [9] and [27]. Recall first that our gradient flow uses the reweighted gradient

$$\nabla \ell_{\text{rw}}(\mathbf{z}) := \frac{1}{m} \sum_{i=1}^m w_i \left(\mathbf{a}_i^* \mathbf{z} - |\mathbf{a}_i^* \mathbf{x}| \frac{\mathbf{a}_i^* \mathbf{z}}{|\mathbf{a}_i^* \mathbf{z}|} \right) \mathbf{a}_i \quad (27)$$

with weights

$$w_i = \frac{1}{1 + \beta/(|\mathbf{a}_i^* \mathbf{z}|/|\mathbf{a}_i^* \mathbf{x}|)}, \quad 1 \leq i \leq m \quad (28)$$

in which the dependence on the iterate index t is ignored for notational brevity.

Proposition 2 (Local error contraction): For an arbitrary $\mathbf{x} \in \mathbb{R}^n$, consider m noise-free measurements $\psi_i = |\mathbf{a}_i^* \mathbf{x}|$, $1 \leq i \leq m$. There exist some numerical constants $c_1, c_2, c_3 > 0$, and $0 < \nu < 1$ such that the following holds with probability exceeding $1 - c_3 e^{-c_2 m}$

$$\text{dist}^2(\mathbf{z} - \mu \nabla \ell_{\text{rw}}(\mathbf{z}), \mathbf{x}) \leq (1 - \nu) \text{dist}^2(\mathbf{z}, \mathbf{x}) \quad (29)$$

for all $\mathbf{x}, \mathbf{z} \in \mathbb{R}^n$ obeying $\text{dist}(\mathbf{z}, \mathbf{x}) \leq (1/10)\|\mathbf{x}\|$, provided that $m \geq c_1 n$ and the constant step size $\mu \leq \mu_0$, where the numerical constant μ_0 depends on the parameter $\beta > 0$ and data $\{(\mathbf{a}_i; \psi_i)\}_{1 \leq i \leq m}$.

Proposition 2 suggests that the distance of RAF's successive iterates to the global optimum \mathbf{x} decreases monotonically once the algorithm's iterate \mathbf{z}^t enters a small neighboring region around the true \mathbf{x} . This small-size neighborhood is commonly known as the *basin of attraction*, and has been widely discussed

in recent non-convex optimization contributions; see e.g., [9], [12], [27]. Expressed differently, RAF's iterates will stay within the region and will be attracted towards \mathbf{x} exponentially fast as soon as they land within the basin of attraction. To substantiate Proposition 2, recall the useful analytical tool of the local regularity condition [10], which plays a key role in establishing linear convergence of iterative procedures to the global optimum in [9], [10], [12], [27], [26], [31], [34].

For RAF, the reweighted gradient $\nabla \ell_{\text{rw}}(\mathbf{z})$ in (27) is said to obey the local regularity condition (LRC), denoted as $\text{LRC}(\mu, \lambda, \epsilon)$ for some constant $\lambda > 0$, if the next inequality

$$\langle \nabla \ell_{\text{rw}}(\mathbf{z}), \mathbf{h} \rangle \geq \frac{\mu}{2} \|\nabla \ell_{\text{rw}}(\mathbf{z})\|^2 + \frac{\lambda}{2} \|\mathbf{h}\|^2 \quad (30)$$

holds for all $\mathbf{z} \in \mathbb{R}^n$ such that $\|\mathbf{h}\| = \|\mathbf{z} - \mathbf{x}\| \leq \epsilon \|\mathbf{x}\|$ for some constant $0 < \epsilon < 1$, where the ball given by $\|\mathbf{z} - \mathbf{x}\| \leq \epsilon \|\mathbf{x}\|$ is the so-termed *basin of attraction*.

Letting $\mathbf{h} := \mathbf{z} - \mathbf{x}$, manipulations in conjunction with the regularity property (30) confirms that

$$\begin{aligned} \text{dist}^2(\mathbf{z} - \mu \nabla \ell_{\text{rw}}(\mathbf{z}), \mathbf{x}) &= \|\mathbf{z} - \mu \nabla \ell_{\text{rw}}(\mathbf{z}) - \mathbf{x}\|^2 \\ &= \|\mathbf{h}\|^2 - 2\mu \langle \mathbf{h}, \nabla \ell_{\text{rw}}(\mathbf{z}) \rangle + \|\mu \nabla \ell_{\text{rw}}(\mathbf{z})\|^2 \end{aligned} \quad (31)$$

$$\begin{aligned} &\leq \|\mathbf{h}\|^2 - 2\mu \left(\frac{\mu}{2} \|\nabla \ell_{\text{rw}}(\mathbf{z})\|^2 + \frac{\lambda}{2} \|\mathbf{h}\|^2 \right) + \|\mu \nabla \ell_{\text{rw}}(\mathbf{z})\|^2 \\ &= (1 - \lambda\mu) \|\mathbf{h}\|^2 = (1 - \lambda\mu) \text{dist}^2(\mathbf{z}, \mathbf{x}) \end{aligned} \quad (32)$$

for all points \mathbf{z} adhering to $\|\mathbf{h}\| \leq \epsilon \|\mathbf{x}\|$. It is evident that if $\text{LRC}(\mu, \lambda, \epsilon)$ can be established for RAF, our ultimate goal of proving the local error contraction in (29) follows straightforwardly upon setting $\nu := \lambda\mu$.

1) *Proof of the Local Regularity Condition in (30)*: The first step to proving the local regularity condition in (30) is to control the size of the reweighted gradient $\nabla \ell_{\text{rw}}(\mathbf{z})$; that is, to upper bound the last term in (31). To start, rewrite the reweighted gradient in a compact matrix-vector representation

$$\nabla \ell_{\text{rw}}(\mathbf{z}) = \frac{1}{m} \sum_{i=1}^m w_i \left(\mathbf{a}_i^* \mathbf{z} - |\mathbf{a}_i^* \mathbf{x}| \frac{\mathbf{a}_i^* \mathbf{z}}{|\mathbf{a}_i^* \mathbf{z}|} \right) \mathbf{a}_i \triangleq \frac{1}{m} \text{dg}(\mathbf{w}) \mathbf{A} \mathbf{v} \quad (33)$$

where $\text{dg}(\mathbf{w}) \in \mathbb{R}^{n \times n}$ is a diagonal matrix holding in order the entries of $\mathbf{w} := [w_1 \cdots w_m]^* \in \mathbb{R}^m$ on its main diagonal, and $\mathbf{v} := [v_1 \cdots v_m]^* \in \mathbb{R}^m$ with $v_i := \mathbf{a}_i^* \mathbf{z} - |\mathbf{a}_i^* \mathbf{x}| \frac{\mathbf{a}_i^* \mathbf{z}}{|\mathbf{a}_i^* \mathbf{z}|}$. Based on the definition of the induced matrix 2-norm (or the matrix spectral norm), it is easy to check that

$$\begin{aligned} \|\nabla \ell_{\text{rw}}(\mathbf{z})\| &= \left\| \frac{1}{m} \text{dg}(\mathbf{w}) \mathbf{A} \mathbf{v} \right\| \leq \frac{1}{m} \|\text{dg}(\mathbf{w})\| \cdot \|\mathbf{A}\| \cdot \|\mathbf{v}\| \\ &\leq \frac{1 + \delta'}{\sqrt{m}} \|\mathbf{v}\| \end{aligned} \quad (34)$$

where we have used the inequalities $\|\text{dg}(\mathbf{w})\| \leq 1$ due to $w_i \leq 1$ for all $1 \leq i \leq m$, and $\|\mathbf{A}\| \leq (1 + \delta')\sqrt{m}$ for some constant $\delta' > 0$ according to [48, Theorem 5.32], provided that m/n is sufficiently large.

The task therefore remains to bound $\|\mathbf{v}\|$ in (34), which is addressed next. To this end, notice that

$$\begin{aligned} \|\mathbf{v}\|^2 &\leq \sum_{i=1}^m (|\mathbf{a}_i^* \mathbf{z}| - |\mathbf{a}_i^* \mathbf{x}|)^2 \leq \sum_{i=1}^m (\mathbf{a}_i^* \mathbf{z} - \mathbf{a}_i^* \mathbf{x})^2 \\ &\leq (1 + \delta'')^2 m \|\mathbf{h}\|^2 \end{aligned} \quad (35)$$

for some numerical constant $\delta'' > 0$, where the last can be obtained using [15, Lemma 3.1], and which holds with probability at least $1 - e^{-c_2 m}$ as long as $m > c_1 n$ holds true.

Combing (34) with (35) and taking $\delta > 0$ larger than the constant $(1 + \delta')(1 + \delta'') - 1$, the size of $\nabla \ell_{\text{rw}}(\mathbf{z})$ can be bounded as

$$\|\nabla \ell_{\text{rw}}(\mathbf{z})\| \leq (1 + \delta) \|\mathbf{h}\| \quad (36)$$

which holds with probability $1 - e^{-c_2 m}$, with a proviso that m/n exceeds some numerical constant $c_7 > 0$. This result indeed asserts that the reweighted gradient of $L(\mathbf{z})$ or the search direction employed in our RAF algorithm is well behaved, implying that the function value along the iterates does not change too much.

In order to prove the LRC, it suffices to show that $\nabla \ell_{\text{rw}}(\mathbf{z})$ ensures sufficient descent, that is, there exists a numerical constant $c > 0$ such that along the search direction $\nabla \ell_{\text{rw}}(\mathbf{z})$ the following uniform lower bound holds

$$\langle \nabla \ell_{\text{rw}}(\mathbf{z}), \mathbf{h} \rangle \geq c \|\mathbf{h}\|^2 \quad (37)$$

which will be addressed next. Formally, this can be summarized in the following proposition, whose proof is deferred to Appendix D.

Proposition 3: For the noise-free measurements $\psi_i = |\mathbf{a}_i^* \mathbf{x}|$, $1 \leq i \leq m$, and any fixed sufficiently small constant $\epsilon > 0$. There exist some numerical constants $c_1, c_2, c_3 > 0$ such that the following holds with probability at least $1 - c_3 e^{-c_2 m}$

$$\langle \mathbf{h}, \nabla \ell_{\text{rw}}(\mathbf{z}) \rangle \geq \zeta_3 \|\mathbf{h}\|^2 \quad (38)$$

for all $\mathbf{x}, \mathbf{z} \in \mathbb{R}^n$ obeying $\|\mathbf{h}\| \leq (1/10)\|\mathbf{x}\|$, provided that $m/n > c_1$, and that $\beta \geq 0$ is small enough. Here, $\zeta_3 := \frac{1 - \zeta_1 - \epsilon}{1 + \beta(1 + \eta)} - 2(\zeta_2 + \epsilon) - \frac{2(0.1271 - \zeta_2 + \epsilon)}{1 + \beta/k}$.

Taking the results in (38) and (36) together back to (30), we deduce that the LRC holds for μ and λ obeying the inequality

$$\zeta_3 \geq \frac{\mu}{2} (1 + \delta)^2 + \frac{\lambda}{2}. \quad (39)$$

For instance, taking $\beta = 2$, $k = 5$, $\eta = 0.5$, and $\epsilon = 0.001$, we have $\zeta_1 = 0.8897$ and $\zeta_2 = 0.0213$, which confirms that $\langle \ell_{\text{rw}}(\mathbf{z}), \mathbf{h} \rangle \geq 0.1065 \|\mathbf{h}\|^2$. Setting further $\delta = 0.001$ leads to

$$0.1065 \geq 0.501\mu + 0.5\lambda \quad (40)$$

which concludes the proof of the LRC in (30). The local error contraction in (29) follows directly after substituting the LRC into (32), hence validating Proposition 2.

VI. CONCLUSION

This paper puts forth a novel linear-time algorithm termed reweighted amplitude flow (RAF) for solving high-dimensional random systems of quadratic equations. Our procedure proceeds in two consecutive stages, namely, a weighted maximal

correlation initialization that entails just a few power or Lanczos iterations, and a sequence of simple iteratively reweighted generalized gradient iterations for the non-convex non-smooth least-squares loss function. Our RAF approach is conceptually simple, easy-to-implement, as well as numerically scalable and effective. It was also proved to achieve the optimal sample and computational complexity orders. Substantial numerical tests using both synthetic data and real-world images corroborated the superior performance of RAF over state-of-the-art iterative solvers. Empirically, RAF solves a set of random quadratic equations in the high-dimensional regime with large probability so long as a unique solution exists, where the number m of equations in the real Gaussian case can be as small as $2n - 1$ with n being the number of unknowns to be recovered.

Future research includes studying robust and/or sparse phase retrieval as well as (semi-definite) matrix recovery by means of (stochastic) reweighted amplitude flow counterparts [17], [37]. Exploiting the possibility of leveraging suitable (re)weighting regularization to improve empirical performance of other non-convex iterative procedures such as [37], [28] is worth investigating as well.

APPENDIX PROOF DETAILS

By homogeneity of (1), we assume without loss of generality that $\|\mathbf{x}\| = 1$ in all proofs.

A. Proof of Lemma 2

Let $\{\mathbf{b}_i^*\}_{1 \leq i \leq |\mathcal{S}|}$ denote rows of $\mathbf{B} \in \mathbb{R}^{|\mathcal{S}| \times n}$, which are obtained by scaling rows of $\mathbf{A}^S := \{\mathbf{a}_i^*\}_{i \in \mathcal{S}} \in \mathbb{R}^{|\mathcal{S}| \times n}$ by weights $\{\psi_i = \psi_i^{\gamma/2}\}_{i \in \mathcal{S}}$ [cf. (17)]. Since $\mathbf{x} = \mathbf{e}_1$, we have $\psi = |\mathbf{A}\mathbf{e}_1| = |\mathbf{A}_1|$, while the index set \mathcal{S} depends solely on the first column of \mathbf{A} , and is independent of the other columns of \mathbf{A} . Using this, partition accordingly $\mathbf{A}^S := [\mathbf{A}_1^S \ \mathbf{A}_r^S]$, where $\mathbf{A}_1^S \in \mathbb{R}^{|\mathcal{S}| \times 1}$ denotes the first column of \mathbf{A}^S , and $\mathbf{A}_r^S \in \mathbb{R}^{|\mathcal{S}| \times (n-1)}$ collects the remaining ones. Likewise, partition $\mathbf{B} = [\mathbf{B}_1 \ \mathbf{B}_r]$ with $\mathbf{B}_1 \in \mathbb{R}^{|\mathcal{S}| \times 1}$ and $\mathbf{B}_r \in \mathbb{R}^{|\mathcal{S}| \times (n-1)}$. By this argument, rows of \mathbf{A}^S are mutually independent, and Gaussian distributed with mean $\mathbf{0}$ and covariance matrix \mathbf{I}_{n-1} . Furthermore, the weights $\psi_i^{\gamma/2} = |\mathbf{a}_i^* \mathbf{e}_1|^{\gamma/2} = |a_{i,1}|^{\gamma/2}$, $\forall i \in \mathcal{S}$ are also independent of the entries in \mathbf{A}^S . As a consequence, rows of \mathbf{B}_r are mutually independent, and one can explicitly write its i -th row as $\mathbf{b}_{r,i} = |\mathbf{a}_{[i]}^* \mathbf{e}_1|^{\gamma/2} \mathbf{a}_{[i],\setminus 1} = |a_{[i],1}|^{\gamma/2} \mathbf{a}_{[i],\setminus 1}$, where $\mathbf{a}_{[i],\setminus 1} \in \mathbb{R}^{n-1}$ is obtained after removing the first entry of $\mathbf{a}_{[i]}$. It is easy to verify that $\mathbb{E}[\mathbf{b}_{r,i}] = \mathbf{0}$, and $\mathbb{E}[\mathbf{b}_{r,i} \mathbf{b}_{r,i}^*] = C_\gamma \mathbf{I}_{n-1}$, where the constant $C_\gamma := \sqrt{2^\gamma/\pi} \Gamma(\gamma + 1/2) \|\mathbf{x}\|^\gamma = \sqrt{2^\gamma/\pi} \Gamma(\gamma + 1/2)$, and $\Gamma(\cdot)$ is the Gamma function.

Given $\mathbf{x}^* \mathbf{x}^\perp = \mathbf{e}_1^* \mathbf{x}^\perp = 0$, one can write $\mathbf{x}^\perp = [0 \ \mathbf{r}^*]^*$ with any unit vector $\mathbf{r} \in \mathbb{R}^{n-1}$; hence,

$$\|\mathbf{B} \mathbf{x}^\perp\|^2 = \|\mathbf{B} [0 \ \mathbf{r}^*]^*\|^2 = \|\mathbf{B}_r \mathbf{r}\|^2 \quad (41)$$

with independent sub-Gaussian rows $\mathbf{b}_{r,i} = |a_{[i],1}|^{\gamma/2} \mathbf{a}_{[i],\setminus 1}$ if $0 \leq \gamma \leq 1$. Standard concentration results on the sum of random positive semi-definite matrices composed of independent

non-isotropic sub-Gaussian rows [48, Remark 5.40.1] assert that

$$\left\| \frac{1}{|\mathcal{S}|} \mathbf{B}_r^* \mathbf{B}_r - C_\gamma \mathbf{I}_{n-1} \right\| \leq \delta \quad (42)$$

holds with probability at least $1 - 2e^{-c_5 n}$ provided that $|\mathcal{S}|/n$ is larger than some positive constant. Here, $\delta > 0$ is a numerical constant that can take arbitrarily small values, and $c_5 > 0$ is a constant depending on δ . With no loss of generality, take $\delta := 0.01 C_\gamma$ in (42). For any unit vector $\mathbf{r} \in \mathbb{R}^{n-1}$, the following holds with probability at least $1 - 2e^{-c_5 n}$

$$\left\| \frac{1}{|\mathcal{S}|} \mathbf{r}^* \mathbf{B}_r^* \mathbf{B}_r \mathbf{r} - C_\gamma \mathbf{r}^* \mathbf{r} \right\| \leq \delta \mathbf{r}^* \mathbf{r} = \delta \quad (43)$$

or

$$\|\mathbf{B}_r \mathbf{r}\|^2 = \mathbf{r}^* \mathbf{B}_r^* \mathbf{B}_r \mathbf{r} \leq 1.01 C_\gamma |\mathcal{S}|. \quad (44)$$

Taking (44) back to (41) confirms that

$$\|\mathbf{B} \mathbf{x}^\perp\|^2 \leq 1.01 C_\gamma |\mathcal{S}| \quad (45)$$

holds with probability at least $1 - 2e^{-c_5 n}$ if $|\mathcal{S}|/n$ exceeds some constant. Note that c_5 depends on the maximum sub-Gaussian norm of the rows \mathbf{b}_i in \mathbf{B}_r , and we assume without loss of generality $c_5 \geq 1/2$. Therefore, one confirms that the numerator $\|\mathbf{B} \mathbf{u}\|^2$ in (18) is upper bounded after replacing \mathbf{x}^\perp with \mathbf{u} in (45).

B. Proof of Lemma 3

This section is devoted to obtaining a meaningful lower bound for the denominator $\|\mathbf{B} \mathbf{x}\|^2$ in (21). Note first that

$$\|\mathbf{B} \mathbf{x}\|^2 = \sum_{i=1}^{|\mathcal{S}|} |\mathbf{b}_i^* \mathbf{x}|^2 = \sum_{i=1}^{|\mathcal{S}|} \psi_i^\gamma |\mathbf{a}_{[i]}^* \mathbf{x}|^2 = \sum_{i=1}^{|\mathcal{S}|} |\mathbf{a}_{[i]}^* \mathbf{x}|^{2+\gamma}.$$

Taking without loss of generality $\mathbf{x} = \mathbf{e}_1$, the term on the right side of the last equality reduces to

$$\|\mathbf{B} \mathbf{x}\|^2 = \sum_{i=1}^{|\mathcal{S}|} |a_{[i],1}|^{2+\gamma}. \quad (46)$$

Since $a_{[i],1}$ follows the standardized normal distribution, the probability density function (pdf) of random variables $|a_{[i],1}|^{2+\gamma}$ can be given in closed form as

$$p(t) = \sqrt{\frac{2}{\pi}} \cdot \frac{1}{2+\gamma} t^{-\frac{1+\gamma}{2+\gamma}} e^{-\frac{1}{2} t^{\frac{2}{2+\gamma}}}, \quad t > 0 \quad (47)$$

which is rather complicated and whose cumulative density function (cdf) does not come in closed form in general. Therefore, instead of dealing with the pdf in (47) directly, we shall take a different route by deriving a lower bound that is a bit looser yet suffices for our purpose.

Since $|a_{[|\mathcal{S}|],1}| \leq \dots \leq |a_{[2],1}| \leq |a_{[1],1}|$, then it holds for all $1 \leq i \leq |\mathcal{S}|$ that $|a_{[i],1}|^{2+\gamma} \geq |a_{[|\mathcal{S}|],1}|^\gamma a_{[i],1}^2$, which yields

$$\|\mathbf{B} \mathbf{x}\|^2 = \sum_{i=1}^{|\mathcal{S}|} |a_{[i],1}|^{2+\gamma} \geq |a_{[|\mathcal{S}|],1}|^\gamma \sum_{i=1}^{|\mathcal{S}|} a_{[i],1}^2. \quad (48)$$

We will next demonstrate that deriving a lower bound for $\|\mathbf{B} \mathbf{x}\|^2$ suffices to derive a lower bound for the summation on the right

hand side (48). The latter can be achieved by appealing to a result in [9, Lemma 3], which for completeness is included in the following.

Lemma 4: For an arbitrary unit-norm vector $\mathbf{x} \in \mathbb{R}^n$, let $\psi_i = |\mathbf{a}_i^* \mathbf{x}|$, $1 \leq i \leq m$ be m noiseless measurements. Then with probability at least $1 - e^{-c_2 m}$, the following holds

$$\sum_{i=1}^{|\mathcal{S}|} a_{[i],1}^2 \geq 0.99|\mathcal{S}|[1 + \log(m/|\mathcal{S}|)] \quad (49)$$

provided that $m \geq c_0|\mathcal{S}| \geq c_1 n$ for some numerical constants $c_0, c_1, c_2 > 0$.

Combining the results in Lemma 4 and (48), one further deduces that

$$\begin{aligned} \|\mathbf{B}\mathbf{x}\|^2 &\geq |a_{[|\mathcal{S}|],1}|^\gamma \sum_{i=1}^{|\mathcal{S}|} a_{[i],1}^2 \\ &\geq |a_{[|\mathcal{S}|],1}|^\gamma \cdot 0.99|\mathcal{S}|[1 + \log(m/|\mathcal{S}|)]. \end{aligned} \quad (50)$$

The task remains to estimate the size of $|a_{[|\mathcal{S}|],1}|$, which we recall is the $|\mathcal{S}|$ -th largest among the m independent realizations $\{\psi_i = |a_{i,1}|\}_{1 \leq i \leq m}$. Taking $\gamma = 1$ in (47) gives the pdf of the half-normal distribution

$$p(t) = \sqrt{\frac{2}{\pi}} e^{-\frac{1}{2}t^2}, \quad t > 0 \quad (51)$$

whose corresponding cdf is

$$F(\tau) = \text{erf}(\tau/\sqrt{2}). \quad (52)$$

Setting $F(\tau_{|\mathcal{S}|}) := 1 - |\mathcal{S}|/m$ or using the complementary cdf $|\mathcal{S}|/m := \text{erfc}(\tau/\sqrt{2})$ based on the complementary error function gives rise to an estimate of the size of the $|\mathcal{S}|$ -th largest (or equivalently, the $(m - |\mathcal{S}|)$ -th smallest) entry in the m realizations, namely

$$\tau_{|\mathcal{S}|} = \sqrt{2} \text{erfc}^{-1}(|\mathcal{S}|/m) \quad (53)$$

where $\text{erfc}^{-1}(\cdot)$ represents the inverse complementary error function. In the sequel, we show that the deviation of the $|\mathcal{S}|$ -th largest realization $\psi_{[|\mathcal{S}|]}$ from its expected value $\tau_{|\mathcal{S}|}$ in (53) is bounded with high probability.

For random variable $\psi = |a|$ with a obeying the standard Gaussian distribution, consider the event $\psi \leq \tau_{|\mathcal{S}|} - \delta$ for a fixed constant $\delta > 0$. Define the indicator random variable $\chi := \mathbb{1}_{\{\psi \leq \tau_{|\mathcal{S}|} - \delta\}}$, whose expectation can be obtained by substituting $\tau = \tau_{|\mathcal{S}|} - \delta$ into the pdf in (52) as

$$\mathbb{E}[\chi_i] = \text{erf}(\tau_{|\mathcal{S}|} - \delta/\sqrt{2}). \quad (54)$$

Considering now the m independent copies $\{\chi_i = \mathbb{1}_{\{\psi_i \leq \tau_{|\mathcal{S}|} - \delta\}}\}_{1 \leq i \leq m}$ of χ , the following holds

$$\begin{aligned} \mathbb{P}(\psi_{[|\mathcal{S}|]} \leq \tau_{|\mathcal{S}|} - \delta) &= \mathbb{P}\left(\sum_{i=1}^m \chi_i \leq m - |\mathcal{S}|\right) \\ &= \mathbb{P}\left(\frac{1}{m} \sum_{i=1}^m (\chi_i - \mathbb{E}[\chi_i]) \leq 1 - \frac{|\mathcal{S}|}{m} - \mathbb{E}[\chi_i]\right) \end{aligned}$$

Clearly, since random variables χ_i are bounded, they are sub-Gaussian [49]. For notational brevity, let $t := 1 - |\mathcal{S}|/m -$

$\mathbb{E}[\chi_i] = 1 - |\mathcal{S}|/m - \text{erf}(\tau_{|\mathcal{S}|} - \delta/\sqrt{2})$. Appealing to a large deviation inequality for sums of independent sub-Gaussian random variables, one establishes that

$$\begin{aligned} &\mathbb{P}(\psi_{[|\mathcal{S}|]} \leq \tau_{|\mathcal{S}|} - \delta) \\ &= \mathbb{P}\left(\frac{1}{m} \sum_{i=1}^m (\chi_i - \mathbb{E}[\chi_i]) \leq 1 - \frac{|\mathcal{S}|}{m} - \mathbb{E}[\chi_i]\right) \leq e^{-c_5 m t^2} \end{aligned} \quad (55)$$

where $c_5 > 0$ is some absolute constant. On the other hand, using the definition of the error function and properties of integration gives rise to

$$\begin{aligned} t &= 1 - |\mathcal{S}|/m - \text{erf}(\tau_{|\mathcal{S}|} - \delta/\sqrt{2}) = \frac{2}{\sqrt{\pi}} \int_{(\tau_{|\mathcal{S}|} - \delta)/\sqrt{2}}^{\tau_{|\mathcal{S}|}/\sqrt{2}} e^{-s^2} ds \\ &\geq \sqrt{\frac{2}{\pi}} \delta e^{-\frac{\tau_{|\mathcal{S}|}^2}{2}} \geq \sqrt{\frac{2}{\pi}} \delta. \end{aligned} \quad (56)$$

Taking the results in (55) and (56) together, one concludes that fixing any constant $\delta > 0$, the following holds with probability at least $1 - e^{-c_2 m}$:

$$\psi_{[|\mathcal{S}|]} \geq \tau_{|\mathcal{S}|} - \delta \geq \sqrt{2} \text{erfc}^{-1}(|\mathcal{S}|/m) - \delta$$

where $c_2 := (2/\pi)c_5\delta^2$. Furthermore, choosing without loss of generality $\delta := 0.01\tau_{|\mathcal{S}|}$ above leads to $\psi_{[|\mathcal{S}|]} \geq 1.4 \text{erfc}^{-1}(|\mathcal{S}|/m)$.

Substituting the last inequality into (50), and under our working assumption $|\mathcal{S}|/m \leq 0.25$, one readily obtains that

$$\begin{aligned} \|\mathbf{B}\mathbf{x}\|^2 &\geq [1.4 \text{erfc}^{-1}(|\mathcal{S}|/m)]^\gamma \cdot 0.99|\mathcal{S}|[1 + \log(m/|\mathcal{S}|)] \\ &\geq 0.99 \cdot 1.14^\gamma |\mathcal{S}|[1 + \log(m/|\mathcal{S}|)] \end{aligned}$$

which holds with probability exceeding $1 - e^{-c_2 m}$ for some constant $c_2 > 0$, thus concluding the proof of Lemma 3.

C. Proof of Proposition 3

To proceed, let us introduce the following events for all $1 \leq i \leq m$:

$$\mathcal{D}_i := \{(\mathbf{a}_i^* \mathbf{x})(\mathbf{a}_i^* \mathbf{z}) < 0\} \quad (57)$$

$$\mathcal{E}_i := \left\{ \frac{|\mathbf{a}_i^* \mathbf{z}|}{|\mathbf{a}_i^* \mathbf{x}|} \geq \frac{1}{1 + \eta} \right\} \quad (58)$$

for some fixed constant $\eta > 0$, in which the former corresponds to the gradients involving wrongly estimated signs, namely $\frac{\mathbf{a}_i^* \mathbf{z}}{|\mathbf{a}_i^* \mathbf{z}|} \neq \frac{\mathbf{a}_i^* \mathbf{x}}{|\mathbf{a}_i^* \mathbf{x}|}$, and the second will be useful for deriving error bounds. Based on the definition of \mathcal{D}_i and with $\mathbb{1}_{\mathcal{D}_i}$ denoting

the indicator function of the event \mathcal{D}_i , we have

$$\begin{aligned}
\langle \ell_{\text{rw}}(\mathbf{z}), \mathbf{h} \rangle &= \frac{1}{m} \sum_{i=1}^m w_i \left(\mathbf{a}_i^* \mathbf{z} - |\mathbf{a}_i^* \mathbf{x}| \frac{\mathbf{a}_i^* \mathbf{z}}{|\mathbf{a}_i^* \mathbf{z}|} \right) (\mathbf{a}_i^* \mathbf{h}) \\
&= \frac{1}{m} \sum_{i=1}^m w_i \left(\mathbf{a}_i^* \mathbf{h} + \mathbf{a}_i^* \mathbf{x} - |\mathbf{a}_i^* \mathbf{x}| \frac{\mathbf{a}_i^* \mathbf{z}}{|\mathbf{a}_i^* \mathbf{z}|} \right) (\mathbf{a}_i^* \mathbf{h}) \\
&= \frac{1}{m} \sum_{i=1}^m w_i (\mathbf{a}_i^* \mathbf{h})^2 + \frac{1}{m} \sum_{i=1}^m 2w_i (\mathbf{a}_i^* \mathbf{x}) (\mathbf{a}_i^* \mathbf{h}) \mathbb{1}_{\mathcal{D}_i} \\
&\geq \frac{1}{m} \sum_{i=1}^m w_i (\mathbf{a}_i^* \mathbf{h})^2 - \frac{1}{m} \sum_{i=1}^m 2w_i |\mathbf{a}_i^* \mathbf{x}| |\mathbf{a}_i^* \mathbf{h}| \mathbb{1}_{\mathcal{D}_i}.
\end{aligned} \tag{59}$$

In the following, we will derive a lower bound for the term on the right hand side of (59). Specifically, a lower bound for the first term $(1/m) \sum_{i=1}^m w_i (\mathbf{a}_i^* \mathbf{h})^2$ and an upper bound for the second term $(1/m) \sum_{i=1}^m 2w_i |\mathbf{a}_i^* \mathbf{x}| |\mathbf{a}_i^* \mathbf{h}| \mathbb{1}_{\mathcal{D}_i}$ will be obtained, based on Lemmas 5 and 6, with their proofs postponed to Appendix E and Appendix F, respectively.

Lemma 5: Fix fixed $\eta, \beta > 0$, and any sufficiently small constant $\epsilon > 0$, the following holds with probability at least $1 - 2e^{-c_5 \epsilon^2 m}$

$$\frac{1}{m} \sum_{i=1}^m w_i (\mathbf{a}_i^* \mathbf{h})^2 \geq \frac{1 - \zeta_1 - \epsilon}{1 + \beta(1 + \eta)} \|\mathbf{h}\|^2 \tag{60}$$

with $w_i = 1/[1 + \beta/(|\mathbf{a}_i^* \mathbf{z}|/|\mathbf{a}_i^* \mathbf{x}|)]$ for all $1 \leq i \leq m$, provided that $m/n > (c_6 \cdot \epsilon^{-2} \log \epsilon^{-1})$ for certain numerical constants $c_5, c_6 > 0$.

Now we turn to the second term in (59). For ease of exposition, let us first introduce the following events

$$\mathcal{B}_i := \{|\mathbf{a}_i^* \mathbf{x}| < |\mathbf{a}_i^* \mathbf{h}| \leq (k+1)|\mathbf{a}_i^* \mathbf{x}|\} \tag{61}$$

$$\mathcal{O}_i := \{(k+1)|\mathbf{a}_i^* \mathbf{x}| < |\mathbf{a}_i^* \mathbf{h}|\} \tag{62}$$

for all $1 \leq i \leq m$ and some fixed constant $k > 0$. The second term can be bounded as follows

$$\begin{aligned}
&\frac{1}{m} \sum_{i=1}^m 2w_i |\mathbf{a}_i^* \mathbf{x}| |\mathbf{a}_i^* \mathbf{h}| \mathbb{1}_{\mathcal{D}_i} \\
&\leq \frac{1}{m} \sum_{i=1}^m w_i [(\mathbf{a}_i^* \mathbf{x})^2 + (\mathbf{a}_i^* \mathbf{h})^2] \mathbb{1}_{\{(\mathbf{a}_i^* \mathbf{z})(\mathbf{a}_i^* \mathbf{x}) < 0\}} \\
&= \frac{1}{m} \sum_{i=1}^m w_i [(\mathbf{a}_i^* \mathbf{x})^2 + (\mathbf{a}_i^* \mathbf{h})^2] \mathbb{1}_{\{(\mathbf{a}_i^* \mathbf{h})(\mathbf{a}_i^* \mathbf{x}) + (\mathbf{a}_i^* \mathbf{x})^2 < 0\}} \\
&\leq \frac{1}{m} \sum_{i=1}^m w_i [(\mathbf{a}_i^* \mathbf{x})^2 + (\mathbf{a}_i^* \mathbf{h})^2] \mathbb{1}_{\{|\mathbf{a}_i^* \mathbf{x}| < |\mathbf{a}_i^* \mathbf{h}|\}} \\
&\leq \frac{2}{m} \sum_{i=1}^m w_i (\mathbf{a}_i^* \mathbf{h})^2 \mathbb{1}_{\{|\mathbf{a}_i^* \mathbf{x}| < |\mathbf{a}_i^* \mathbf{h}|\}} \\
&= \frac{2}{m} \sum_{i=1}^m w_i (\mathbf{a}_i^* \mathbf{h})^2 \mathbb{1}_{\{|\mathbf{a}_i^* \mathbf{x}| < |\mathbf{a}_i^* \mathbf{h}| \leq (k+1)|\mathbf{a}_i^* \mathbf{x}|\}}
\end{aligned}$$

$$\begin{aligned}
&+ \frac{2}{m} \sum_{i=1}^m w_i (\mathbf{a}_i^* \mathbf{h})^2 \mathbb{1}_{\{(k+1)|\mathbf{a}_i^* \mathbf{x}| < |\mathbf{a}_i^* \mathbf{h}|\}} \\
&= \frac{2}{m} \sum_{i=1}^m w_i (\mathbf{a}_i^* \mathbf{h})^2 \mathbb{1}_{\mathcal{B}_i} + \frac{2}{m} \sum_{i=1}^m w_i (\mathbf{a}_i^* \mathbf{h})^2 \mathbb{1}_{\mathcal{O}_i}
\end{aligned} \tag{63}$$

where the first equality is derived by substituting $\mathbf{z} = \mathbf{h} + \mathbf{x}$ according to the definition of \mathbf{h} , the second event suffices for $(\mathbf{a}_i^* \mathbf{h})(\mathbf{a}_i^* \mathbf{x}) + (\mathbf{a}_i^* \mathbf{x})^2 < 0$, and the second equality follows from writing the indicator function $\mathbb{1}_{\{|\mathbf{a}_i^* \mathbf{x}| < |\mathbf{a}_i^* \mathbf{h}|\}}$ as the summation of two indicator functions of two events $\mathbb{1}_{\{|\mathbf{a}_i^* \mathbf{x}| < |\mathbf{a}_i^* \mathbf{h}| \leq (k+1)|\mathbf{a}_i^* \mathbf{x}|\}}$ and $\mathbb{1}_{\{|\mathbf{a}_i^* \mathbf{h}| > (k+1)|\mathbf{a}_i^* \mathbf{x}|\}}$.

The task so far remains to derive upper bounds for the two terms on the right hand side of (63), which leads to Lemma 6.

Lemma 6: Fixing a fixed $k > 0$, define ζ_2 to be the maximum of $\mathbb{E}[w_i]$ in (72) for $\rho = 0.01$ and $\nu = 0.1$, which depends only on k . For any $\epsilon > 0$, if $m/n > c_6 \epsilon^{-2} \log \epsilon^{-1}$, the following hold simultaneously with probability at least $1 - c_3 e^{-c_2 \epsilon^2 m}$

$$\frac{1}{m} \sum_{i=1}^m w_i (\mathbf{a}_i^* \mathbf{h})^2 \mathbb{1}_{\mathcal{O}_i} \leq (\zeta_2 + \epsilon) \|\mathbf{h}\|^2 \tag{64}$$

and

$$\frac{1}{m} \sum_{i=1}^m w_i (\mathbf{a}_i^* \mathbf{h})^2 \mathbb{1}_{\mathcal{B}_i} \leq \frac{0.1271 - \zeta_2 + \epsilon}{1 + \beta/k} \|\mathbf{h}\|^2 \tag{65}$$

for all $\mathbf{h} \in \mathbb{R}^n$ obeying $\|\mathbf{h}\|/\|\mathbf{x}\| \leq 1/10$, where $c_1, c_2, c_3 > 0$ are some universal constants.

Substituting (60), (63), and (64)–(65) established in Lemmas 5 and 6 back into (59), we conclude that

$$\begin{aligned}
\langle \ell_{\text{rw}}(\mathbf{z}), \mathbf{h} \rangle &\geq \frac{1}{m} \sum_{i=1}^m w_i (\mathbf{a}_i^* \mathbf{h})^2 \mathbb{1}_{\mathcal{E}_i} - \frac{1}{m} \sum_{i=1}^m 2w_i |\mathbf{a}_i^* \mathbf{x}| |\mathbf{a}_i^* \mathbf{h}| \mathbb{1}_{\mathcal{D}_i} \\
&= \zeta_e \|\mathbf{h}\|^2
\end{aligned} \tag{66}$$

which will be rendered positive, provided that $\beta > 0$ is small enough, and that parameters $\eta, k > 0$ are suitably chosen.

D. Proof of Lemma 5

Plugging in the weighting parameters $w_i = 1/[1 + \beta/(|\mathbf{a}_i^* \mathbf{z}|/|\mathbf{a}_i^* \mathbf{x}|)]$ and based on the definition of \mathcal{E}_i , the first term in (59) can be lower bounded as

$$\frac{1}{m} \sum_{i=1}^m w_i (\mathbf{a}_i^* \mathbf{h})^2 \geq \frac{1}{m} \sum_{i=1}^m \frac{(\mathbf{a}_i^* \mathbf{h})^2 \mathbb{1}_{\mathcal{E}_i}}{1 + \beta/(|\mathbf{a}_i^* \mathbf{z}|/|\mathbf{a}_i^* \mathbf{x}|)} \tag{67}$$

$$\begin{aligned}
&\geq \frac{1}{m} \sum_{i=1}^m \frac{1}{1 + \beta(1 + \eta)} (\mathbf{a}_i^* \mathbf{h})^2 \mathbb{1}_{\left\{\frac{|\mathbf{a}_i^* \mathbf{z}|}{|\mathbf{a}_i^* \mathbf{x}|} \geq \frac{1}{1 + \eta}\right\}} \\
&= \frac{1}{1 + \beta(1 + \eta)} \cdot \frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^* \mathbf{h})^2 \mathbb{1}_{\mathcal{E}_i}
\end{aligned} \tag{68}$$

where the first inequality arises from dropping some nonnegative terms from the left hand side, and the second one after replacing the ratio $|\mathbf{a}_i^* \mathbf{z}|/|\mathbf{a}_i^* \mathbf{x}|$ in the weights by its lower bound $1/(1 + \eta)$ because the weights are monotonically increasing functions of $|\mathbf{a}_i^* \mathbf{z}|/|\mathbf{a}_i^* \mathbf{x}|$. Using [9, Lemma 5], the last term in

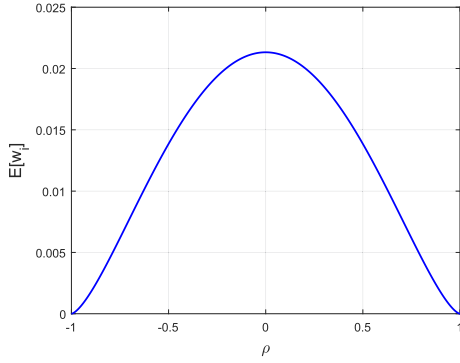


Fig. 8. The expectation $\mathbb{E}[w_i]$ as a function of ρ over $[-1, 1]$.

(68) can be further bounded by

$$\begin{aligned} \frac{1}{m} \sum_{i=1}^m w_i (\mathbf{a}_i^* \mathbf{h})^2 &\geq \frac{1}{1 + \beta(1 + \eta)} \cdot \frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^* \mathbf{h})^2 \mathbb{1}_{\mathcal{E}_i} \\ &\geq \frac{1 - \zeta_1 - \epsilon}{1 + \beta(1 + \eta)} \|\mathbf{h}\|^2 \end{aligned} \quad (69)$$

for any fixed sufficiently small constant $\epsilon > 0$, which holds with probability at least $1 - 2e^{-c_5 \epsilon^2 m}$, if $m > (c_6 \cdot \epsilon^{-2} \log \epsilon^{-1})n$.

E. Proof of Lemma 6

The proof is adapted from [27, Lemma 9]. We first prove the bound (64) for any fixed \mathbf{h} obeying $\|\mathbf{h}\| \leq \|\mathbf{x}\|/10$, and subsequently develop a uniform bound at the end of this section. The bound (65) can be derived directly after subtracting the bound in (64) with k from that bound with $k = 0$, followed by an application of the Bernstein-type sub-exponential tail bound [48]. We only discuss the first bound (64). Because of the discontinuity hence non-Lipschitz of the indicator functions, let us approximate them by a sequence of auxiliary Lipschitz functions. Specifically, with some constant $\varrho > 0$, define for all $1 \leq i \leq m$ the continuous functions $\chi_i(s)$ in (70) as shown at the bottom of this page. Clearly, all $\chi_i(s)$'s are random Lipschitz functions with constant $1/\varrho$. Furthermore, it is easy to verify that

$$\begin{aligned} |\mathbf{a}_i^* \mathbf{h}|^2 \mathbb{1}_{\{(k+1)|\mathbf{a}_i^* \mathbf{x}| < |\mathbf{a}_i^* \mathbf{h}|\}} &\leq \chi_i(|\mathbf{a}_i^* \mathbf{h}|^2) \\ &\leq |\mathbf{a}_i^* \mathbf{h}|^2 \mathbb{1}_{\{\sqrt{1-\varrho}(k+1)|\mathbf{a}_i^* \mathbf{x}| < |\mathbf{a}_i^* \mathbf{h}|\}}. \end{aligned} \quad (71)$$

Since the second term involves the addition event \mathcal{E}_i in (58), define $w_i := \frac{|\mathbf{a}_i^* \mathbf{h}|^2}{\|\mathbf{h}\|^2} \mathbb{1}_{\{\sqrt{1-\varrho}(k+1)|\mathbf{a}_i^* \mathbf{x}| < |\mathbf{a}_i^* \mathbf{h}|\}}$ for $1 \leq i \leq m$, and $\nu := \frac{\|\mathbf{h}\|}{\|\mathbf{x}\|}$ for convenience. If $f(\tau_1, \tau_2)$ denotes the density of two joint Gaussian random variables with correlation coefficient $\rho = \frac{\mathbf{h}^* \mathbf{x}}{\|\mathbf{h}\| \|\mathbf{x}\|} \in (-1, 1)$, then the expectation of w_i can be obtained

using the conditional expectation

$$\begin{aligned} \mathbb{E}[w_i] &= \int_{-\infty}^{\infty} \mathbb{E}[w_i | \mathbf{a}_i^* \mathbf{x} = \tau_1 \|\mathbf{x}\|, \mathbf{a}_i^* \mathbf{h} = \tau_1 \|\mathbf{h}\|] f(\tau_1, \tau_2) d\tau_1 d\tau_2 \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \tau_2^2 \mathbb{1}_{\{\sqrt{1-\varrho}(k+1)|\tau_1| < |\tau_2| \nu\}} f(\tau_1, \tau_2) d\tau_1 d\tau_2 \\ &= \frac{1}{\sqrt{2\pi}} \int_0^{\infty} \tau_2^2 \exp(-\tau_2^2/2) \left[\operatorname{erf}\left(\frac{(\nu/[\sqrt{1-\varrho}(k+1)] - \rho)\tau_2}{\sqrt{2(1-\rho^2)}}\right) \right. \\ &\quad \left. + \operatorname{erf}\left(\frac{(\nu/[\sqrt{1-\varrho}(k+1)] + \rho)\tau_2}{\sqrt{2(1-\rho^2)}}\right) \right] d\tau_2 \end{aligned} \quad (72)$$

$$:= \zeta_2. \quad (73)$$

It is not difficult to see that $\mathbb{E}[w_i] = 0$ for $\rho = \pm 1$, and $\mathbb{E}[w_i]$ is continuous over $\rho \in (-1, 1)$ due to the integration property of continuous functions over a continuous interval. Although the last term in (72) can not be expressed in closed form, it can be evaluated numerically. Note first that for fixed parameters $\varrho > 0$ and $\nu \leq 0.1$, the integration in (72) is monotonically decreasing in $k \geq 0$, and achieves the maximum at $k = 0$. For parameter values $k = 5$, $\nu = 0.1$ and $\varrho = 0.01$, Fig. 8 plots $\mathbb{E}[w_i]$ as a function of ρ , whose maximum $\zeta_2 = 0.0213$ is achieved at $\rho = 0$. Further, from the integration in (72) for fixed $k \geq 0$, $\mathbb{E}[w_i]$ is a monotonically increasing function of both ν and ϱ , and it is therefore safe to conclude that for all $0 < \nu \leq 0.1$, and $\varrho = 0.01$, we have

$$\mathbb{E}[w_i] \leq \zeta_2 = 0.0213. \quad (74)$$

Hence, we can infer that $\mathbb{E}[\chi_i(|\mathbf{a}_i^* \mathbf{h}|^2)] \leq 0.0213 \|\mathbf{h}\|^2$ for $\nu < 0.1$, $\varrho = 0.01$, and $k = 5$. Since the $\chi_i(|\mathbf{a}_i^* \mathbf{h}|^2)$'s are sub-exponential with sub-exponential norm of the order $\mathcal{O}(\|\mathbf{h}\|^2)$, Bernstein-type sub-exponential tail bound [48] confirms that

$$\mathbb{P}\left(\frac{1}{m} \sum_{i=1}^m \frac{\chi_i(|\mathbf{a}_i^* \mathbf{h}|^2)}{\|\mathbf{h}\|^2} > \zeta_2 + \epsilon\right) < e^{-c_7 m \epsilon^2} \quad (75)$$

for some numerical constant $\epsilon > 0$, provided that $\|\mathbf{h}\| \leq \|\mathbf{x}\|/10$. Finally, due to the fact that $w_i \leq 1$ for all $1 \leq i \leq m$, the following holds

$$\frac{1}{m} \sum_{i=1}^m w_i \chi_i(|\mathbf{a}_i^* \mathbf{h}|^2) < (\zeta_2 + \epsilon) \|\mathbf{h}\|^2 \quad (76)$$

with probability at least $1 - e^{-c_7 m \epsilon^2}$.

We have proved the bound in (64) for a fixed vector \mathbf{h} , and the uniform bound for all vectors \mathbf{h} obeying $\|\mathbf{h}\| \leq \|\mathbf{x}\|/10$ can be obtained by similar arguments in the proof [27, Lemma 9] with only minor changes in the constants.

$$\chi_i(s) := \begin{cases} s, & s > (1+k)^2 (\mathbf{a}_i^* \mathbf{x})^2 \\ \frac{1}{\varrho} [s - (k+1)^2 (\mathbf{a}_i^* \mathbf{x})^2] + (k+1)^2 (\mathbf{a}_i^* \mathbf{x})^2, & (1-\varrho)(k+1)^2 (\mathbf{a}_i^* \mathbf{x})^2 \leq s \leq (k+1)^2 (\mathbf{a}_i^* \mathbf{x})^2 \\ 0, & \text{otherwise.} \end{cases} \quad (70)$$

Regarding the second bound (65), it is easy to see that

$$\begin{aligned} & \frac{1}{m} \sum_{i=1}^m |\mathbf{a}_i^* \mathbf{h}|^2 \mathbb{1}_{\{|\mathbf{a}_i^* \mathbf{x}| < |\mathbf{a}_i^* \mathbf{h}|(k+1)|\mathbf{a}_i^* \mathbf{x}|\}} \\ &= \frac{1}{m} \sum_{i=1}^m \left[|\mathbf{a}_i^* \mathbf{h}|^2 \mathbb{1}_{\{|\mathbf{a}_i^* \mathbf{x}| < |\mathbf{a}_i^* \mathbf{h}|\}} - |\mathbf{a}_i^* \mathbf{h}|^2 \mathbb{1}_{\{(k+1)|\mathbf{a}_i^* \mathbf{x}| < |\mathbf{a}_i^* \mathbf{h}|\}} \right] \\ &\leq (0.1271 - \zeta_2 + \epsilon) \|\mathbf{h}\|^2 \end{aligned} \quad (77)$$

where the last inequality follows from subtracting the bound in (64) of k from that corresponding to $k = 0$. To account for the weights $w_i = 1/[1 + \beta/(|\mathbf{a}_i^* \mathbf{z}|/|\mathbf{a}_i^* \mathbf{x}|)]$, first notice that $\mathbf{a}_i^* \mathbf{h} = \mathbf{a}_i^* \mathbf{z} - \mathbf{a}_i^* \mathbf{x}$, and that our second bound works with $(\mathbf{a}_i^* \mathbf{z})(\mathbf{a}_i^* \mathbf{x}) < 0$ in (59), hence $\frac{|\mathbf{a}_i^* \mathbf{z}|}{|\mathbf{a}_i^* \mathbf{x}|} \leq \frac{|\mathbf{a}_i^* \mathbf{h}|}{|\mathbf{a}_i^* \mathbf{x}|} - 1$. Recall that the second bound (65) assumes the event $\{|\mathbf{a}_i^* \mathbf{x}| < |\mathbf{a}_i^* \mathbf{h}| \leq (k+1)|\mathbf{a}_i^* \mathbf{x}|\}$, implying $\frac{|\mathbf{a}_i^* \mathbf{z}|}{|\mathbf{a}_i^* \mathbf{x}|} \leq \frac{|\mathbf{a}_i^* \mathbf{h}|}{|\mathbf{a}_i^* \mathbf{x}|} - 1 \leq k$. Further, because w_i is monotonically increasing in $\frac{|\mathbf{a}_i^* \mathbf{z}|}{|\mathbf{a}_i^* \mathbf{x}|}$, then $w_i \leq \frac{1}{1+\beta/k}$. Taking this result back to (77) yields

$$\begin{aligned} & \frac{1}{m} \sum_{i=1}^m w_i |\mathbf{a}_i^* \mathbf{h}|^2 \mathbb{1}_{\{|\mathbf{a}_i^* \mathbf{x}| < |\mathbf{a}_i^* \mathbf{h}| \leq (k+1)|\mathbf{a}_i^* \mathbf{x}|\}} \\ &\leq \frac{0.1271 - \zeta_2 + \epsilon}{1 + \beta/k} \|\mathbf{h}\|^2 \end{aligned} \quad (78)$$

which proves the second bound in (65).

ACKNOWLEDGMENT

The authors would like to thank the anonymous reviewers for their thorough review and all constructive comments and suggestions, which helped to improve the quality of the manuscript. The authors would also like to thank Prof. J. Duchi for his helpful feedback on our initialization.

REFERENCES

- [1] G. Wang, G. B. Giannakis, Y. Saad, and J. Chen, "Solving most systems of random quadratic equations," in *Proc. Adv. Neural Inf. Process. Syst.*, Long Beach, CA, USA, 2017, pp. 1865–1875.
- [2] R. Balan, P. Casazza, and D. Edidin, "On signal reconstruction without phase," *Appl. Comput. Harmon. Anal.*, vol. 20, no. 3, pp. 345–356, May 2006.
- [3] A. S. Bandeira, J. Cahill, D. G. Mixon, and A. A. Nelson, "Saving phase: Injectivity and stability for phase retrieval," *Appl. Comput. Harmon. Anal.*, vol. 37, no. 1, pp. 106–125, 2014.
- [4] J. R. Fienup, "Phase retrieval algorithms: A comparison," *Appl. Opt.*, vol. 21, no. 15, pp. 2758–2769, Aug. 1982.
- [5] R. W. Gerchberg and W. O. Saxton, "A practical algorithm for the determination of phase from image and diffraction," *Optik*, vol. 35, pp. 237–246, Nov. 1972.
- [6] Y. Shechtman, Y. C. Eldar, O. Cohen, H. N. Chapman, J. Miao, and M. Segev, "Phase retrieval with application to optical imaging: A contemporary overview," *IEEE Signal Process. Mag.*, vol. 32, no. 3, pp. 87–109, May 2015.
- [7] X. Yi, C. Caramanis, and S. Sanghavi, "Alternating minimization for mixed linear regression," in *Proc. Int. Conf. Mach. Learn.*, Beijing, China, 2014, pp. 613–621.
- [8] J. R. Rice, *Numerical Methods in Software and Analysis*. Cambridge, MA, USA: Academic, 1992.
- [9] G. Wang, G. B. Giannakis, and Y. C. Eldar, "Solving systems of random quadratic equations via truncated amplitude flow," *IEEE Trans. Inf. Theory*, vol. 64, no. 2, pp. 773–794, Feb. 2018.
- [10] E. J. Candès, X. Li, and M. Soltanolkotabi, "Phase retrieval via Wirtinger flow: Theory and algorithms," *IEEE Trans. Inf. Theory*, vol. 61, no. 4, pp. 1985–2007, Apr. 2015.
- [11] Y. C. Eldar and S. Mendelson, "Phase retrieval: Stability and recovery guarantees," *Appl. Comput. Harmon. Anal.*, vol. 36, no. 3, pp. 473–494, May 2014.
- [12] Y. Chen and E. J. Candès, "Solving random quadratic systems of equations is nearly as easy as solving linear systems," *Commun. Pure Appl. Math.*, vol. 70, no. 5, pp. 822–883, Dec. 2017.
- [13] A. Ben-Tal and A. Nemirovski, *Lectures on Modern Convex Optimization: Analysis, Algorithms, and Engineering Applications*, vol. 2. Philadelphia, PA, USA: SIAM, 2001.
- [14] A. Chai, M. Moscoso, and G. Papanicolaou, "Array imaging using intensity-only measurements," *Inverse Probl.*, vol. 27, no. 1, Dec. 2011, Art. no. 015005.
- [15] E. J. Candès, T. Strohmer, and V. Voroninski, "PhaseLift: Exact and stable signal recovery from magnitude measurements via convex programming," *Appl. Comput. Harmon. Anal.*, vol. 66, no. 8, pp. 1241–1274, Nov. 2013.
- [16] I. Waldspurger, A. d'Aspremont, and S. Mallat, "Phase recovery, maxcut, and complex semidefinite programming," *Math. Program.*, vol. 149, no. 1, pp. 47–81, 2015.
- [17] Y. Chen, Y. Chi, and A. J. Goldsmith, "Exact and stable covariance estimation from quadratic sampling via convex programming," *IEEE Trans. Inf. Theory*, vol. 61, no. 7, pp. 4034–4059, Jul. 2015.
- [18] E. J. Candès and X. Li, "Solving quadratic equations via PhaseLift when there are about as many equations as unknowns," *Found. Comput. Math.*, vol. 14, no. 5, pp. 1017–1026, 2014.
- [19] Y. M. Lu and M. Vetterli, "Sparse spectral factorization: Unicity and reconstruction algorithms," in *Proc. Int. Conf. Acoust., Speech, Signal Process.*, Prague, Czech Republic, 2011, pp. 5976–5979.
- [20] F. Krahmer and Y.-K. Liu, "Phase retrieval without small-ball probability assumptions," *IEEE Trans. Inf. Theory*, vol. 64, no. 1, pp. 485–500, Jan. 2018.
- [21] T. Goldstein and S. Studer, "PhaseMax: Convex phase retrieval via basis pursuit," *IEEE Trans. Inf. Theory*, vol. 64, no. 4, pp. 2675–2689, Apr. 2018.
- [22] S. Bahmani and J. Romberg, "A flexible convex relaxation for phase retrieval," *Electron. J. Statist.*, vol. 11, no. 2, pp. 5254–5281, Dec. 2017.
- [23] P. Hand and V. Voroninski, "An elementary proof of convex phase retrieval in the natural parameter space via the linear program PhaseMax," arXiv:1611.03935, 2016.
- [24] P. Netrapalli, P. Jain, and S. Sanghavi, "Phase retrieval using alternating minimization," *IEEE Trans. Signal Process.*, vol. 63, no. 18, pp. 4814–4826, Sep. 2015.
- [25] Y. Chen, Y. Chi, J. Fan, and C. Ma, "Gradient descent with random initialization: Fast global convergence for nonconvex phase retrieval," arXiv:1803.07726, 2018.
- [26] Y. Li, C. Ma, Y. Chen, and Y. Chi, "Nonconvex matrix factorization from rank-one measurements," arXiv:1802.06286, 2018.
- [27] H. Zhang, Y. Zhou, Y. Liang, and Y. Chi, "Reshaped Wirtinger flow and incremental algorithm for solving quadratic system of equations," arXiv:1605.07719, 2016.
- [28] H. Zhang, Y. Chi, and Y. Liang, "Provable non-convex phase retrieval with outliers: Median truncated Wirtinger flow," arXiv:1603.03805, 2016.
- [29] T. Bendory, Y. C. Eldar, and N. Boumal, "Non-convex phase retrieval from STFT measurements," *IEEE Trans. Inf. Theory*, vol. 64, no. 1, pp. 467–484, Jan. 2018.
- [30] M. Soltanolkotabi, "Structured signal recovery from quadratic measurements: Breaking sample complexity barriers via nonconvex optimization," arXiv:1702.06175, 2017.
- [31] J. Chen, L. Wang, X. Zhang, and Q. Gu, "Robust Wirtinger flow for phase retrieval with arbitrary corruption," arXiv:1704.06256, 2017.
- [32] Z. Yuan and H. Wang, "Phase retrieval via reweighted Wirtinger flow," *Appl. Opt.*, vol. 56, no. 9, pp. 2418–2427, Mar. 2017.
- [33] G. Wang and G. B. Giannakis, "Solving random systems of quadratic equations via truncated generalized gradient flow," in *Proc. Adv. Neural Inf. Process. Syst.*, Barcelona, Spain, 2016, pp. 568–576.
- [34] G. Wang, G. B. Giannakis, and J. Chen, "Scalable solvers of random quadratic equations via stochastic truncated amplitude flow," *IEEE Trans. Signal Process.*, vol. 65, no. 8, pp. 1961–1974, Apr. 2017.
- [35] G. Wang, L. Zhang, G. B. Giannakis, M. Akçakaya, and J. Chen, "Sparse phase retrieval via truncated amplitude flow," *IEEE Trans. Signal Process.*, vol. 66, no. 2, pp. 479–491, Jan. 2018.
- [36] G. Wang, L. Zhang, G. B. Giannakis, and J. Chen, "Sparse phase retrieval via iteratively reweighted amplitude flow," in *Proc. 26th Eur. Signal Process. Conf.*, Rome, Italy, submitted for publication, 2018.

- [37] J. C. Duchi and F. Ruan, "Solving (most) of a set of quadratic equalities: Composite optimization for robust phase retrieval," arXiv:1705.02356, 2017.
- [38] J. Duchi and F. Ruan, "Stochastic methods for composite optimization problems," arXiv:1703.08570, 2017.
- [39] D. Davis, D. Drusvyatskiy, and C. Paquette, "The nonsmooth landscape of phase retrieval," arXiv:1711.03247, 2017.
- [40] F. H. Clarke, "Generalized gradients and applications," *Trans. Amer. Math. Soc.*, vol. 205, pp. 247–262, 1975.
- [41] J. Sun, Q. Qu, and J. Wright, "A geometric analysis of phase retrieval," *Found. Comput. Math.*, to be published, 2018.
- [42] E. J. Candès, X. Li, and M. Soltanolkotabi, "Phase retrieval from coded diffraction patterns," *Appl. Comput. Harmon. Anal.*, vol. 39, no. 2, pp. 277–299, Sep. 2015.
- [43] R. H. Keshavan, A. Montanari, and S. Oh, "Matrix completion from a few entries," *IEEE Trans. Inf. Theory*, vol. 56, no. 6, pp. 2980–2998, Jun. 2010.
- [44] Y. M. Lu and G. Li, "Phase transitions of spectral initialization for high-dimensional nonconvex estimation," arXiv:1702.06435, 2017.
- [45] M. Mondelli and A. Montanari, "Fundamental limits of weak recovery with applications to phase retrieval," arXiv:1708.05932, 2017.
- [46] Y. Saad, *Numerical Methods for Large Eigenvalue Problems: Revised Ed.* Philadelphia, PA, USA: SIAM, 2011.
- [47] R. Chartrand and W. Yin, "Iteratively reweighted algorithms for compressive sensing," in *Proc. Int. Conf. Acoust., Speech, Signal Process.*, Las Vegas, NV, USA, 2008, pp. 3869–3872.
- [48] P. Chen, A. Fannjiang, and G.-R. Liu, "Phase retrieval with one or two diffraction patterns by alternating projection with null initialization," arXiv:1510.07379, 2015.
- [49] R. Vershynin, "Introduction to the non-asymptotic analysis of random matrices," arXiv:1011.3027, 2010.



Gang Wang (S'12) received the B.Eng. degree in electrical engineering and automation from the Beijing Institute of Technology, Beijing, China, in 2011, and the Ph.D. degree in electrical engineering from the University of Minnesota, Minneapolis, MN, USA, in 2018.

His research interests focus on the areas of statistical signal processing, stochastic and non-convex optimization with applications to smart grids, data analytics, and deep learning. He received a National Scholarship (2014), a Guo Rui Scholarship (2017),

all from China, a Best Student Paper Award at the 2017 European Signal Processing Conference, and the Student Travel Awards from the NSF (2016) and NIPS (2017).



Georgios B. Giannakis (F'97) received the Diploma degree in electrical engineering from the National Technical University of Athens, Athens, Greece, in 1981. From 1982 to 1986, he was with the University of the Southern California, Los Angeles, CA, USA, where he received the M.Sc. degree in electrical engineering in 1983, the M.Sc. degree in mathematics, and the Ph.D. degree in electrical engineering both in 1986.

He was with the University of Virginia from 1987 to 1998 and since 1999, he has been a Professor

with the University of Minnesota, Minneapolis, MN, USA, where he holds an Endowed Chair in Wireless Telecommunications, University of Minnesota McKnight Presidential Chair in ECE, and serves as the Director of the Digital Technology Center. His general interests include the areas of communications, networking, and statistical signal processing—subjects on which he has published more than 400 journal papers, 700 conference papers, 25 book chapters, two edited books, and two research monographs (h-index 129). Current research focuses on learning from big data, wireless cognitive radios, and network science with applications to social, brain, and power networks with renewables. He is the (co-)inventor of 30 patents issued, and the (co-)recipient of nine best paper awards from the IEEE SIGNAL PROCESSING (SP) and Communications Societies, including the G. Marconi Prize Paper Award in Wireless Communications. He is also a recipient of the Technical Achievement Awards from the SP Society (2000), from the EURASIP (2005), a Young Faculty Teaching Award, the G. W. Taylor Award for Distinguished Research from the University of Minnesota, and the IEEE Fourier Technical Field Award (inaugural recipient in 2015). He is a Fellow of the EURASIP, and has served the IEEE in a number of posts, including that of a Distinguished Lecturer for the IEEE-SP Society.



Yousef Saad is a College of Science and Technology Distinguished Professor with the Department of Computer Science and Engineering (CSE) at the University of Minnesota. He received the "Doctorat d'Etat" from the University of Grenoble (France) in 1983. He joined the University of Minnesota in 1990 as a Professor of Computer Science and a Fellow of the Minnesota Supercomputing Institute. He was head of the Department of CSE from January 1997 to June 2000, and became a CSE Distinguished Professor in 2005. From 1981 to 1990, he held positions at the University of California at Berkeley, Yale, the University of Illinois, and the Research Institute for Advanced Computer Science (RIACS).

His current research interests include: numerical linear algebra, sparse matrix computations, iterative methods, parallel computing, numerical methods for electronic structure, and linear algebra methods in data mining. He is the author of two monographs and over 180 journal articles. He is also the developer or co-developer of several software packages for solving sparse linear systems of equations including SPARKIT, pARMS, and ITSOL. Yousef Saad is a SIAM fellow (class of 2010) and a fellow of the AAAS (2011).



Jie Chen (SM'12) received the B.S., M.S., and Ph.D. degrees in control theory and control engineering from the Beijing Institute of Technology, Beijing, China, in 1986, 1996, and 2001, respectively. From 1989 to 1990, he was a Visiting Scholar with the California State University, Long Beach, Long Beach, CA, USA. From 1996 to 1997, he was a Research Fellow with the School of Engineering, the University of Birmingham, Birmingham, U.K.

He is currently a Professor of control science and engineering with the Beijing Institute of Technology.

He is a Member of the Chinese Academy of Engineering, and is also the Head of the State Key Laboratory of Intelligent Control and Decision of Complex Systems, Beijing Institute of Technology. He is currently a Managing Editor for the *Journal of Systems Science & Complexity* (2014–2017), an Associate Editor for the IEEE TRANSACTIONS ON CYBERNETICS (2016–2018) and several other international journals. He has (co-)authored three books and more than 200 research papers. His main research interests include intelligent control and decision in complex systems, multiagent systems, and optimization.