

Median-Truncated Nonconvex Approach for Phase Retrieval with Outliers^{*†}

Huishuai Zhang^s, Yuejie Chi^o and Yingbin Liang^s

^sDepartment of EECS, Syracuse University, Syracuse, NY 13244

^oDepartment of ECE, Ohio State University, Columbus, OH 43210

February 1, 2017

Abstract

This paper investigates the phase retrieval problem, which aims to recover a signal from the magnitudes of its linear measurements. We develop statistically and computationally efficient algorithms for the situation when the measurements are corrupted by sparse outliers that can take arbitrary values. We propose a novel approach to robustify the gradient descent algorithm by using the sample median as a guide for pruning spurious samples in initialization and local search. Adopting the Poisson loss and the reshaped quadratic loss respectively, we obtain two algorithms termed *median-TWF* and *median-RWF*, both of which provably recover the signal from a near-optimal number of measurements when the measurement vectors are composed of i.i.d. Gaussian entries, up to a logarithmic factor, even when a constant fraction of the measurements are adversarially corrupted. We further show that both algorithms are stable in the presence of additional dense bounded noise. Our analysis is accomplished by developing non-trivial concentration results of median-related quantities, which may be of independent interest. We provide numerical experiments to demonstrate the effectiveness of our approach.

1 Introduction

Phase retrieval is a classical problem in signal processing, optics and machine learning that has a wide range of applications such as X-ray crystallography [21], astronomical imaging, and diffraction imaging. Mathematically, it is formulated as recovering a signal $\mathbf{x} \in \mathbb{R}^n/\mathbb{C}^n$ from the magnitudes of its linear measurements:

$$y_i = |\langle \mathbf{a}_i, \mathbf{x} \rangle|^2, \quad i = 1, \dots, m, \quad (1)$$

where m is the total number of measurements, and $\mathbf{a}_i \in \mathbb{R}^n/\mathbb{C}^n$ is the i th known measurement vector, $i = 1, \dots, m$. Phase retrieval is known to be notoriously difficult due to the quadratic form of the measurements. Classical methods based on alternating minimization between the signal of interest and the phase information [22], though computationally simple, are often trapped at local minima and lack rigorous performance guarantees.

There has been, however, a recent line of work that successfully develops provably accurate algorithms for phase retrieval, in particular for the case when the measurement vectors \mathbf{a}_i 's are composed of *independent and identically distributed* (i.i.d.) Gaussian entries. Broadly speaking, two classes of approaches have been proposed based on convex and nonconvex optimization techniques, respectively. Using the lifting trick, the phase retrieval problem can be reformulated as estimating a rank-one positive semidefinite (PSD) matrix $\mathbf{X} = \mathbf{x}\mathbf{x}^T$ from linear measurements [4], for which convex relaxations into semidefinite programs have been studied [8, 10, 16, 19, 35, 50]. In particular, Phaselift [10] perfectly recovers the signal with high probability

^{*}The work of H. Zhang and Y. Liang is supported in part by the grants AFOSR FA9550-16-1-0077 and NSF ECCS 16-09916. The work of Y. Chi is supported in part by the grants NSF ECCS-1650449, AFOSR FA9550-15-1-0205 and ONR N00014-15-1-2387.

[†]The material in this paper was presented in part at the International Conference of Machine Learning (ICML), New York, USA, 2016.

as long as the number of measurements m is on the order of n . However, the computational complexity of Phaselift is at least cubic in n , which becomes expensive when n is large. Very recently, another convex relaxation has been proposed in the natural parameter space without lifting [3, 24, 26], resulting in a linear program that can handle large problem dimensions as long as m is on the order of n .

Another class of approaches aims to find the signal that minimizes a loss function based on certain postulated noise model, which often results in a nonconvex optimization problem due to the quadratic measurements. Despite nonconvexity, it is demonstrated in [9, 42] that the so-called Wirtinger flow (WF) algorithm, based on gradient descent, works remarkably well: it converges to the global optima when properly initialized using the spectral method. Several variants of WF have been proposed thereafter to further improve its performance, including the truncated Wirtinger flow (TWF) algorithm [13], the reshaped Wirtinger flow (RWF) algorithm [55], and the truncated amplitude flow (TAF) algorithm [51]. Notably, TWF, RWF and TAF are shown to converge globally at a linear rate as long as m is on the order of n , and attain ϵ -accuracy within $\mathcal{O}(mn \log(1/\epsilon))$ flops using a constant step size.¹

1.1 Outlier-Robust Phase Retrieval

The aforementioned algorithms are evaluated based on their *statistical* and *computational* performances: statistically, we wish the sample complexity m to be as small as possible; computationally, we wish the run time to be as fast as possible. As can be seen, existing WF-type algorithms are already near-optimal both statistically and computationally. This paper introduces a third consideration, which is the *robustness to outliers*, where we wish the algorithm continues to work well even in the presence of outliers that may take arbitrary magnitudes. This bears great importance in practice, because outliers arise frequently from the phase imaging applications [53] due to various reasons such as detector failures, recording errors, and missing data. Specifically, suppose the set of m measurements are given as

$$y_i = |\langle \mathbf{a}_i, \mathbf{x} \rangle|^2 + \eta_i, \quad i = 1, \dots, m, \quad (2)$$

where $\eta_i \in \mathbb{R}/\mathbb{C}$ for $i = 1, \dots, m$ are outliers that can take arbitrary values. We assume that outliers are sparse with no more than sm nonzero values, i.e., $\|\boldsymbol{\eta}\|_0 \leq sm$, where $\boldsymbol{\eta} = \{\eta_i\}_{i=1}^m \in \mathbb{R}^m/\mathbb{C}^m$. Here, s is a nonzero constant, representing the fraction of measurements that are corrupted by outliers.

The goal of this paper is to develop phase retrieval algorithms with both statistical and computational efficiency, and provable robustness to even a constant proportion of outliers. None of the existing algorithms meet all of the three considerations simultaneously. The performance of WF-type algorithms is very sensitive to outliers which introduce anomalous search directions when their values are excessively deviated. While a form of Phaselift [25] is robust to a constant portion of outliers, it is computationally too expensive.

1.2 Median-Truncated Gradient Descent

A natural idea is to recover the signal as a solution to the following loss minimization problem:

$$\min_{\mathbf{z}} \frac{1}{2m} \sum_{i=1}^m \ell(\mathbf{z}; y_i) \quad (3)$$

where $\ell(\mathbf{z}; y_i)$ is postulated using the negative likelihood of Gaussian or Poisson noise model. Since the measurements are quadratic in \mathbf{x} , the objective function is nonconvex. We consider two choices of $\ell(\mathbf{z}; y_i)$ in this paper. The first one is the Poisson loss function of $|\mathbf{a}_i^T \mathbf{z}|^2$ employed in TWF [13], which is given by

$$\ell(\mathbf{z}; y_i) = |\mathbf{a}_i^T \mathbf{z}|^2 - y_i \log |\mathbf{a}_i^T \mathbf{z}|^2. \quad (4)$$

The second one is the *reshaped*² quadratic loss of $|\mathbf{a}_i^T \mathbf{z}|$ employed in RWF [55], which is given by

$$\ell(\mathbf{z}; y_i) = (|\mathbf{a}_i^T \mathbf{z}| - \sqrt{y_i})^2. \quad (5)$$

¹Notation $f(n) = \mathcal{O}(g(n))$ or $f(n) \lesssim g(n)$ means that there exists a constant $c > 0$ such that $|f(n)| \leq c|g(n)|$.

²It is called “reshaped” in order to distinguish it from the quadratic loss of $|\mathbf{a}_i^T \mathbf{z}|^2$ used in [9].

Though (5) is not smooth everywhere, it resembles more closely the least-squares loss as if the phase information is available than the quadratic loss of $|\mathbf{a}_i^T \mathbf{z}|^2$, resulting in a more amenable curvature.

In the presence of outliers, the signal of interest may no longer be the global optima of (3). Therefore, we wish to only include the clean samples that are not corrupted in the optimization (3), which is, however, impossible as we do not assume any *a priori* knowledge of the outliers. Our key strategy is to prune the bad samples adaptively and iteratively, using a gradient descent procedure that proceeds as follows:

$$\mathbf{z}^{(t+1)} = \mathbf{z}^{(t)} - \frac{\mu}{m} \sum_{i \in T_{t+1}} \nabla \ell(\mathbf{z}^{(t)}; y_i). \quad (6)$$

where $\mathbf{z}^{(t)}$ denotes the t th iterate of the algorithm, $\nabla \ell(\mathbf{z}^{(t)}; y_i)$ is the gradient of $\ell(\mathbf{z}^{(t)}; y_i)$, and μ is the step size, for $t = 0, 1, \dots$. In each iteration, only a subset T_{t+1} of data-dependent and iteration-varying samples contributes to the search direction. But how to select the set T_{t+1} ? Note that the gradient of the loss function typically contains the term $|y_i - |\mathbf{a}_i^T \mathbf{z}^{(t)}||$ (for TWF) or $|\sqrt{y_i} - |\mathbf{a}_i^T \mathbf{z}^{(t)}||$ (for RWF), which measures the residual using the current iterate. With y_i being corrupted by arbitrarily large outliers, the gradient can deviate the search direction from the signal arbitrarily.

Inspired by the utility of *median* to combat outliers in robust statistics [28], we prune samples whose gradient components $\nabla \ell(\mathbf{z}^{(t)}; y_i)$ are much larger than the *sample median* to control the search direction of each update. Hiding some technical details, this gives the main ingredient of our *median-truncated gradient descent* update rule³, i.e., for each iterate $t \geq 0$:

$$T_{t+1} := \{i : |y_i - |\mathbf{a}_i^T \mathbf{z}^{(t)}|| \lesssim \text{med}(\{|y_i - |\mathbf{a}_i^T \mathbf{z}^{(t)}||\}_{i=1}^m)\}, \quad \text{for TWF}, \quad (7)$$

$$T_{t+1} := \{i : |\sqrt{y_i} - |\mathbf{a}_i^T \mathbf{z}^{(t)}|| \lesssim \text{med}(\{|\sqrt{y_i} - |\mathbf{a}_i^T \mathbf{z}^{(t)}||\}_{i=1}^m)\}, \quad \text{for RWF}, \quad (8)$$

where $\text{med}(\cdot)$ denotes the sample median. The robust property of median lies in the fact that the median cannot be arbitrarily perturbed unless the outliers dominate the inliers [28]. This is in sharp contrast to the sample mean, which can be made arbitrarily large even by a single outlier. Thus, using the sample median in the truncation rule can effectively remove the impact of outliers. Finally, there still left the question of initialization, which is critical to the success of the algorithm. We use the spectral method, i.e., initialize $\mathbf{z}^{(0)}$ by a proper rescaling of the top eigenvector of a surrogate matrix

$$\mathbf{Y} = \frac{1}{m} \sum_{i \in T_0} y_i \mathbf{a}_i \mathbf{a}_i^T, \quad (9)$$

where again T_0 includes only a subset of samples whose values are not excessively large compared with the sample median of the measurements, given as

$$T_0 = \{i : y_i \lesssim \text{med}(\{y_i\}_{i=1}^m)\}. \quad (10)$$

Putting things together (the update rule (6) and the initialization (9)), we obtain two new median-truncated gradient descent algorithms, median-TWF and median-RWF, based on applying the median truncation strategy for the loss functions used in TWF and RWF, respectively. The median-TWF and median-RWF algorithms do not assume a priori knowledge of the outliers, such as their existence or the number of outliers, and therefore can be used in an oblivious fashion. Importantly, we establish the following performance guarantees.

Main Result (informal): For the Gaussian measurement model, with high probability, median-TWF and median-RWF recover all signal \mathbf{x} up to the global sign at a linear rate of convergence, even with a constant fraction of outliers, as long as the number of measurements m is on the order of $n \log n$. Furthermore, the reconstruction is stable in the presence of additional bounded dense noise.

Statistically, the sample complexity of both algorithms is near-optimal up to a logarithmic factor, and to reassure, they continue to work even when outliers are absent. Computationally, both algorithms converge linearly, requiring a mere computational cost of $\mathcal{O}(mn \log 1/\epsilon)$ to reach ϵ -accuracy. More importantly, our algorithms now tolerate a constant fraction of arbitrary outliers, without sacrificing performance otherwise. To

³Please see the exact form of the algorithms in Section 2.

the best of our knowledge, this is the first application of the median to robustify high-dimensional statistical estimation in the presence of arbitrary outliers with rigorous non-asymptotic performance guarantees.

To establish the performance guarantees, we first show that the initialization is close enough to the ground truth, and then that within the neighborhood of the ground truth, the gradients satisfy certain *Regularity Condition* [9, 13] that guarantees linear convergence of the descent rule, as long as the fraction of outliers is small enough and the sample complexity is large enough. As a nonlinear operator, the sample median is much more difficult to analyze than the sample mean, which is a linear operator and many existing concentration inequalities are readily applicable. Therefore, considerable technical efforts are devoted to develop novel non-asymptotic concentrations of the sample median, and various statistical properties of the sample median related quantities, which may be of independent interest.

Finally, we note that while median-TWF and median-RWF share similar theoretical performance guarantees, their empirical performances vary under different scenarios, due to the use of different loss functions. Their theoretical analyses also have significant difference that worth separate treatments. While we only consider the loss functions used in TWF and RWF in this paper, we believe the median-truncation technique can be applied to gradient descent algorithms for solving other problems as well.

1.3 Related Work

Our work is closely related to the TWF algorithm [13], which is also a truncated gradient descent algorithm for phase retrieval. However, the truncation rule in TWF is based on the sample mean, which is very sensitive to outliers. In [25, 37, 53], the problem of phase retrieval under outliers is investigated, but the proposed algorithms either lack performance guarantees or are computationally too expensive.

The adoption of median in machine learning is not unfamiliar, for example, K -median clustering [12] and resilient data aggregation for sensor networks [49]. Our work here further extends the applications of median to robustifying high-dimensional estimation problems with theoretical guarantees. Another popular approach in robust estimation is to use the trimmed mean [28], which has found success in robustifying sparse regression [15], subspace clustering [41], etc. However, using the trimmed mean requires knowledge of an upper bound on the number of outliers, whereas median does not require such information.

Developing non-convex algorithms with provable global convergence guarantees has attracted intensive research interest recently. A partial list of these studies include phase retrieval [9, 13, 38, 44, 51], matrix completion [18, 23, 27, 29–31, 45, 57], low-rank matrix recovery [6, 17, 33, 36, 40, 47, 52, 56], robust PCA [39, 54], robust tensor decomposition [1], dictionary learning [2, 43], community detection [5], phase synchronization [7], blind deconvolution [32, 34], joint alignment [14], etc. Our algorithm provides a new instance in this list that emphasizes robust high-dimensional signal estimation under minimal assumptions of outliers.

1.4 Paper Organization and Notations

The rest of this paper is organized as follows. Section 2 describes the proposed two algorithms, median-TWF and median-RWF, in details and their performance guarantees. Section 3 presents numerical experiments. Section 4 provides the preliminaries and the proof road map. Section 5 provides the proofs for median-TWF and Section 6 provides the proofs of median-RWF, respectively. Finally, we conclude in Section 7. Supporting proofs are given in the Appendix.

We adopt the following notations in this paper. Given a set of numbers $\{y_i\}_{i=1}^m$, the sample median is denoted as $\text{med}(\{y_i\}_{i=1}^m)$. The indicator function $\mathbf{1}_A = 1$ if the event A holds, and $\mathbf{1}_A = 0$ otherwise. For a vector \mathbf{y} , $\|\mathbf{y}\|$ denotes the l_2 norm. For two matrices, $\mathbf{A} \preceq \mathbf{B}$ if $\mathbf{B} - \mathbf{A}$ is a positive semidefinite matrix.

2 Algorithms and Performance Guarantees

We consider the following model for phase retrieval, where the measurements are corrupted by not only sparse arbitrary outliers but also dense bounded noise. Under such a model, the measurements are given as

$$y_i = |\langle \mathbf{a}_i, \mathbf{x} \rangle|^2 + w_i + \eta_i, \quad i = 1, \dots, m, \quad (11)$$

Algorithm 1 Median Truncated Wirtinger Flow (Median-TWF)

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

Parameters: thresholds $\alpha_y, \alpha_h, \alpha_l$, and α_u , stepsize μ ;

Initialization: Let $\mathbf{z}^{(0)} = \lambda_0 \tilde{\mathbf{z}}$, where $\lambda_0 = \sqrt{\text{med}(\mathbf{y})/0.455}$ and $\tilde{\mathbf{z}}$ is the leading eigenvector of

$$\mathbf{Y} := \frac{1}{m} \sum_{i=1}^m y_i \mathbf{a}_i \mathbf{a}_i^T \mathbf{1}_{\{|y_i| \leq \alpha_y^2 \lambda_0^2\}}. \quad (14)$$

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

$$\mathbf{z}^{(t+1)} = \mathbf{z}^{(t)} - \frac{\mu}{m} \sum_{i=1}^m \frac{|\mathbf{a}_i^T \mathbf{z}^{(t)}|^2 - y_i}{\mathbf{a}_i^T \mathbf{z}^{(t)}} \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i}, \quad (15)$$

where

$$\begin{aligned} \mathcal{E}_1^i &:= \left\{ \alpha_l \|\mathbf{z}^{(t)}\| \leq |\mathbf{a}_i^T \mathbf{z}^{(t)}| \leq \alpha_u \|\mathbf{z}^{(t)}\| \right\}, \\ \mathcal{E}_2^i &:= \left\{ |y_i - |\mathbf{a}_i^T \mathbf{z}^{(t)}|^2| \leq \alpha_h K_t \frac{|\mathbf{a}_i^T \mathbf{z}^{(t)}|}{\|\mathbf{z}^{(t)}\|} \right\}, \quad \text{and} \quad K_t := \text{med} \left(\{|y_i - |\mathbf{a}_i^T \mathbf{z}^{(t)}|^2|\}_{i=1}^m \right). \end{aligned}$$

Output \mathbf{z}_T .

where $\mathbf{x} \in \mathbb{R}^n$ is the unknown signal⁴, $\mathbf{a}_i \in \mathbb{R}^n$ is the i th measurement vector composed of *i.i.d.* Gaussian entries distributed as $\mathcal{N}(0, 1)$, and $\eta_i \in \mathbb{R}$ for $i = 1, \dots, m$ are outliers with arbitrary values satisfying $\|\boldsymbol{\eta}\|_0 \leq sm$, where s is the fraction of outliers, and $\mathbf{w} = \{w_i\}_{i=1}^m$ is the bounded noise satisfying $\|\mathbf{w}\|_\infty \leq c \|\mathbf{x}\|^2$ for some universal constant c .

It is straightforward that changing the sign of the signal does not affect the measurements. The goal is to recover the signal \mathbf{x} , up to a global sign difference, from the measurements $\mathbf{y} = \{y_i\}_{i=1}^m$ and the measurement vectors $\{\mathbf{a}_i\}_{i=1}^m$. To this end, we define the Euclidean distance between two vectors up to a global sign difference as the performance metric,

$$\text{dist}(\mathbf{z}, \mathbf{x}) := \min\{\|\mathbf{z} + \mathbf{x}\|, \|\mathbf{z} - \mathbf{x}\|\}. \quad (12)$$

We propose two median-truncated gradient descent algorithms, median-TWF in Section 2.1 and median-RWF in Section 2.2, based on different choices of the loss functions. This leads to applying the truncation based on the sample median of $\{|y_i - |\mathbf{a}_i^T \mathbf{z}|^2|\}_{i=1}^m$ in median-TWF, and the sample median of $\{|y_i - |\mathbf{a}_i^T \mathbf{z}|^2|\}_{i=1}^m$ in median-RWF. Section 2.3 provides the theoretical performance guarantees of median-TWF and median-RWF, which turn out to be almost the same at the order level except the choice of constants. The empirical comparisons of median-TWF and median-RWF are demonstrated in Section 3.

2.1 Median-TWF Algorithm

In median-TWF, we adopt the Poisson loss function of $|\mathbf{a}_i^T \mathbf{z}|^2$ employed in TWF [13], given as

$$\ell(\mathbf{z}) := \frac{1}{2m} \sum_{i=1}^m (|\mathbf{a}_i^T \mathbf{z}|^2 - y_i \log |\mathbf{a}_i^T \mathbf{z}|^2). \quad (13)$$

The median-TWF algorithm, as described in Algorithm 1, gradually eliminates the influence of outliers on the way of minimizing (13). Specifically, it comprises an initialization step and a truncated gradient descent step.

1. **Initialization:** As in (14), we initialize $\mathbf{z}^{(0)}$ by the spectral method using a truncated set of samples, where the threshold is determined by $\text{med}(\{y_i\}_{i=1}^m)$. As will be shown in Section 4.2, as long as the fraction

⁴We focus on real signals here, but our analysis can be extended to complex signals.

Algorithm 2 Median Reshaped Wirtinger Flow (median-RWF)

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

Parameters: threshold α'_h , and step size μ ;

Initialization: Same as median-TWF (see Algorithm 1).

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

$$\mathbf{z}^{(t+1)} = \mathbf{z}^{(t)} - \frac{\mu}{m} \sum_{i=1}^m \left(\mathbf{a}_i^T \mathbf{z}^{(t)} - \sqrt{y_i} \cdot \frac{\mathbf{a}_i^T \mathbf{z}^{(t)}}{|\mathbf{a}_i^T \mathbf{z}^{(t)}|} \right) \mathbf{a}_i \mathbf{1}_{\mathcal{T}^i}, \quad (18)$$

where

$$\mathcal{T}^i := \left\{ \left| \sqrt{y_i} - |\mathbf{a}_i^T \mathbf{z}^{(t)}| \right| \leq \alpha'_h M_t \right\}, \quad \text{and} \quad M_t := \text{med} \left(\left\{ \left| \sqrt{y_i} - |\mathbf{a}_i^T \mathbf{z}^{(t)}| \right| \right\}_{i=1}^m \right).$$

Output \mathbf{z}_T .

of outliers is not too large and the sample complexity is large enough, our initialization is guaranteed to be within a small neighborhood of the true signal.

2. **Gradient loop:** for each iteration $0 \leq t \leq T - 1$, median-TWF uses an iteration-varying truncated gradient given as

$$\nabla \ell_{tr}(\mathbf{z}^{(t)}) = \frac{1}{m} \sum_{i=1}^m \frac{|\mathbf{a}_i^T \mathbf{z}^{(t)}|^2 - y_i}{\mathbf{a}_i^T \mathbf{z}^{(t)}} \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i}. \quad (16)$$

In order to remove the contribution of corrupted samples, from the definition of the set \mathcal{E}_2^i (see Algorithm 1), it is clear that samples are truncated if their measurement residuals evaluated using the current iterate are much larger than the sample median. Samples are also truncated according to the set \mathcal{E}_1^i , which removes the contribution of samples outside some confidence interval to better control the search direction, since $\mathbb{E}[|\mathbf{a}_i^T \mathbf{z}|] = \sqrt{2/\pi} \|\mathbf{z}\|$. The median-TWF algorithm closely resembles the TWF algorithm, except that the truncation is guided by the sample median, rather than the sample mean.

We set the step size in median-TWF to be a fixed small constant, i.e., $\mu = 0.4$. The rest of the parameters $\{\alpha_y, \alpha_h, \alpha_l, \alpha_u\}$ are set to satisfy

$$\begin{aligned} \zeta_1 &:= \max \left\{ \mathbb{E} \left[\xi^2 \mathbf{1}_{\{|\xi| < \sqrt{1.01}\alpha_l \text{ or } |\xi| > \sqrt{0.99}\alpha_u\}} \right], \mathbb{E} \left[\mathbf{1}_{\{|\xi| < \sqrt{1.01}\alpha_l \text{ or } |\xi| > \sqrt{0.99}\alpha_u\}} \right] \right\}, \\ \zeta_2 &:= \mathbb{E} \left[\xi^2 \mathbf{1}_{\{|\xi| > 0.248\alpha_h\}} \right], \\ 2(\zeta_1 + \zeta_2) + \sqrt{8/\pi} \alpha_h^{-1} &< 1.99 \\ \alpha_y &\geq 3, \end{aligned} \quad (17)$$

where $\xi \sim \mathcal{N}(0, 1)$. For example, we can set $\alpha_l = 0.3, \alpha_u = 5, \alpha_y = 3$ and $\alpha_h = 12$, and consequently $\zeta_1 \approx 0.24$ and $\zeta_2 \approx 0.032$.

2.2 Median-RWF Algorithm

In median-RWF, we adopt the reshaped quadratic loss function of $|\mathbf{a}_i^T \mathbf{z}|$ employed in RWF [55], given as

$$\mathcal{R}(\mathbf{z}) = \frac{1}{2m} \sum_{i=1}^m \left(\sqrt{y_i} - |\mathbf{a}_i^T \mathbf{z}| \right)^2, \quad (19)$$

which has been shown to be advantageous over other loss functions for phase retrieval [55].

Similarly to median-TWF, the median-RWF algorithm as described in Algorithm 2, gradually eliminates the influence of outliers on the way of minimizing (19). Specifically, it also comprises an initialization step and a truncated gradient descent step.

1. **Initialization:** we initialize in the same manner as in median-TWF (Algorithm 1).

2. **Gradient loop:** for each iteration $0 \leq t \leq T - 1$, median-RWF uses the following iteration-varying truncated gradient:

$$\nabla \mathcal{R}_{tr}(\mathbf{z}^{(t)}) = \frac{1}{m} \sum_{i=1}^m \left(\mathbf{a}_i^T \mathbf{z}^{(t)} - \sqrt{y_i} \cdot \frac{\mathbf{a}_i^T \mathbf{z}^{(t)}}{|\mathbf{a}_i^T \mathbf{z}^{(t)}|} \right) \mathbf{a}_i \mathbf{1}_{\mathcal{T}^i}, \quad (20)$$

From the definition of the set \mathcal{T}^i (see Algorithm 2), samples are truncated by the sample median of gradient components evaluated at the current iteration. We set the step size in median-RWF to be a fixed small constant, i.e., $\mu = 0.8$. Compared with median-TWF, the truncation rule is much simpler with fewer parameters. We simply set the truncation threshold $\alpha'_h = 5$. It is possible that including \mathcal{E}_1^i may further improve the performance, but we wish to highlight that, in this paper, the simple truncation rule is already sufficient to guarantee both robustness and efficiency of median-RWF.

2.3 Performance Guarantees

In this section, we characterize the performance guarantees of median-TWF and median-RWF, which turn out to be very similar though the proofs in fact involve quite different techniques. To avoid repetition, we present the guarantees together for both algorithms. We note that the values of constants in the results can vary for median-TWF and median-RWF.

We first show that median-TWF/median-RWF performs well for the noise-free model in the following proposition, which lends support to the model with outliers. This also justifies that we can run median-TWF/median-RWF without having to know whether the underlying measurements are corrupted.

Proposition 1 (Exact recovery for the noise-free model). *Suppose that the measurements are noise-free, i.e., $\eta_i = 0$ and $w_i = 0$ for $i = 1, \dots, m$ in the model (11). There exist constants $\mu_0 > 0$, $0 < \rho, \nu < 1$ and $c_0, c_1, c_2 > 0$ such that if $m \geq c_0 n \log n$ and $\mu \leq \mu_0$, then with probability at least $1 - c_1 \exp(-c_2 m)$, median-TWF/median-RWF yields*

$$\text{dist}(\mathbf{z}^{(t)}, \mathbf{x}) \leq \nu(1 - \rho)^t \|\mathbf{x}\|, \quad \forall t \in \mathbb{N} \quad (21)$$

simultaneously for all $\mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$.

Proposition 1 suggests that median-TWF/median-RWF allows exact recovery at a linear rate of convergence as long as the sample complexity is on the order of $n \log n$, which is in fact slightly worse, by a logarithmic factor, than existing WF-type algorithms (TWF, RWF and TAF) for the noise-free model. This is a price due to working with the nonlinear operator of median in the proof, and it is not clear whether it is possible to further improve the result. Nonetheless, as the median is quite stable as long as the number of outliers is not so large, the following main theorem indeed establishes that median-TWF/median-RWF still performs well even in the presence of a constant fraction of sparse outliers with the same sample complexity.

Theorem 1 (Exact recovery with sparse arbitrary outliers). *Suppose that the measurements are corrupted by sparse outliers, i.e., $w_i = 0$ for $i = 1, \dots, m$ in the model (11). There exist constants $\mu_0, s_0 > 0$, $0 < \rho, \nu < 1$ and $c_0, c_1, c_2 > 0$ such that if $m \geq c_0 n \log n$, $s < s_0$, $\mu \leq \mu_0$, then with probability at least $1 - c_1 \exp(-c_2 m)$, median-TWF/median-RWF yields*

$$\text{dist}(\mathbf{z}^{(t)}, \mathbf{x}) \leq \nu(1 - \rho)^t \|\mathbf{x}\|, \quad \forall t \in \mathbb{N} \quad (22)$$

simultaneously for all $\mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$.

Theorem 1 indicates that median-TWF/median-RWF admits exact recovery for *all* signals in the presence of sparse outliers with arbitrary magnitudes even when the number of outliers scales linearly with the number of measurements, as long as the sample complexity satisfies $m \gtrsim n \log n$. Moreover, median-TWF/median-RWF converges at a linear rate using a constant step size, with per-iteration cost $\mathcal{O}(mn)$ (note that the median can be computed in linear time [46]). To reach ϵ -accuracy, i.e., $\text{dist}(\mathbf{z}^{(t)}, \mathbf{x}) \leq \epsilon$, only $\mathcal{O}(\log 1/\epsilon)$ iterations are needed, yielding the total computational cost as $\mathcal{O}(mn \log 1/\epsilon)$, which is highly efficient. Empirically in the numerical experiments in Section 3, median-RWF converges faster and tolerates a larger fraction of outliers than median-TWF, which can be due to the use of the reshaped quadratic loss function.

We next consider the model when the measurements are corrupted by both sparse arbitrary outliers and dense bounded noise. Our following theorem characterizes that median-TWF/median-RWF is stable to coexistence of the two types of noises.

Theorem 2 (Stability to sparse arbitrary outliers and dense bounded noises). *Consider the phase retrieval problem given in (11) in which measurements are corrupted by both sparse arbitrary and dense bounded noises. There exist constants $\mu_0, s_0 > 0$, $0 < \rho < 1$ and $c_0, c_1, c_2 > 0$ such that if $m \geq c_0 n \log n$, $s < s_0$, $\mu \leq \mu_0$, then with probability at least $1 - c_1 \exp(-c_2 m)$, median-TWF and median-RWF respectively yield*

$$\text{dist}(\mathbf{z}^{(t)}, \mathbf{x}) \lesssim \frac{\|\mathbf{w}\|_\infty}{\|\mathbf{x}\|} + (1 - \rho)^t \|\mathbf{x}\|, \quad \forall t \in \mathbb{N} \quad \text{for median-TWF}, \quad (23)$$

$$\text{dist}(\mathbf{z}^{(t)}, \mathbf{x}) \lesssim \sqrt{\|\mathbf{w}\|_\infty} + (1 - \rho)^t \|\mathbf{x}\|, \quad \forall t \in \mathbb{N} \quad \text{for median-RWF}, \quad (24)$$

simultaneously for all $\mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$.

Theorem 2 immediately implies the stability of median-TWF/median-RWF when the measurements are only corrupted by dense bounded noise.

Corollary 1. *Consider the phase retrieval problem in which measurements are corrupted only by dense bounded noises, i.e., $\eta_i = 0$ for $i = 1, \dots, m$ in the model (11). There exist constants $\mu_0 > 0$, $0 < \rho < 1$ and $c_0, c_1, c_2 > 0$ such that if $m \geq c_0 n \log n$, $\mu \leq \mu_0$, then with probability at least $1 - c_1 \exp(-c_2 m)$, median-TWF and median-RWF respectively yield*

$$\text{dist}(\mathbf{z}^{(t)}, \mathbf{x}) \lesssim \frac{\|\mathbf{w}\|_\infty}{\|\mathbf{x}\|} + (1 - \rho)^t \|\mathbf{x}\|, \quad \forall t \in \mathbb{N} \quad \text{for median-TWF}, \quad (25)$$

$$\text{dist}(\mathbf{z}^{(t)}, \mathbf{x}) \lesssim \sqrt{\|\mathbf{w}\|_\infty} + (1 - \rho)^t \|\mathbf{x}\|, \quad \forall t \in \mathbb{N} \quad \text{for median-RWF}, \quad (26)$$

simultaneously for all $\mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$.

With both sparse arbitrary outliers and dense bounded noises, Theorem 2 and Corollary 1 imply that median-TWF/median-RWF achieves the same convergence rate and the same level of estimation error as the model with only bounded noise. In fact, together with Theorem 1 and Proposition 1, it can be seen that applying median-TWF/median-RWF does not require the knowledge of the existence of outliers. When there do exist outliers, median-TWF/median-RWF achieves almost the same performance *as if outliers do not exist*.

3 Numerical Experiments

In this section, we provide numerical experiments to demonstrate the effectiveness of median-TWF and median-RWF, which corroborate our theoretical findings.

3.1 Exact Recovery for Noise-free Data

We first show that, in the noise-free case, median-TWF and median-RWF provide similar performance as TWF [13] and RWF [55] for exact recovery. We set the parameters of median-TWF and median-RWF as specified in Section 2.1 and Section 2.2, and those of TWF and RWF as suggested in [13] and [55], respectively. Let the signal length n take values from 1000 to 10000 by a step size of 1000, and the ratio of the number of measurements to the signal dimension, m/n , take values from 2 to 6 by a step size of 0.1. For each pair of $(n, m/n)$, we generate a signal $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{n \times n})$, and the measurement vectors $\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{n \times n})$ i.i.d. for $i = 1, \dots, m$. For all algorithms, a fixed number of iterations $T = 500$ are run, and the trial is declared successful if $\mathbf{z}^{(T)}$, the output of the algorithm, satisfies $\text{dist}(\mathbf{z}^{(T)}, \mathbf{x})/\|\mathbf{x}\| \leq 10^{-8}$. Figure 1 shows the number of successful trials out of 20 trials for all algorithms, with respect to m/n and n . It can be seen that, as soon as m is above $4n$, exact recovery is achieved for all four algorithms. Around the phase transition boundary, the empirical sample complexity of median-TWF is slightly worse than that of TWF, which is

possibly due to the inefficiency of median compared to mean in the noise-free case [28]. Interestingly, the empirical sample complexity of median-RWF is slightly better than RWF because the truncation rule used in median-RWF allows sample pruning that improves the performance.⁵

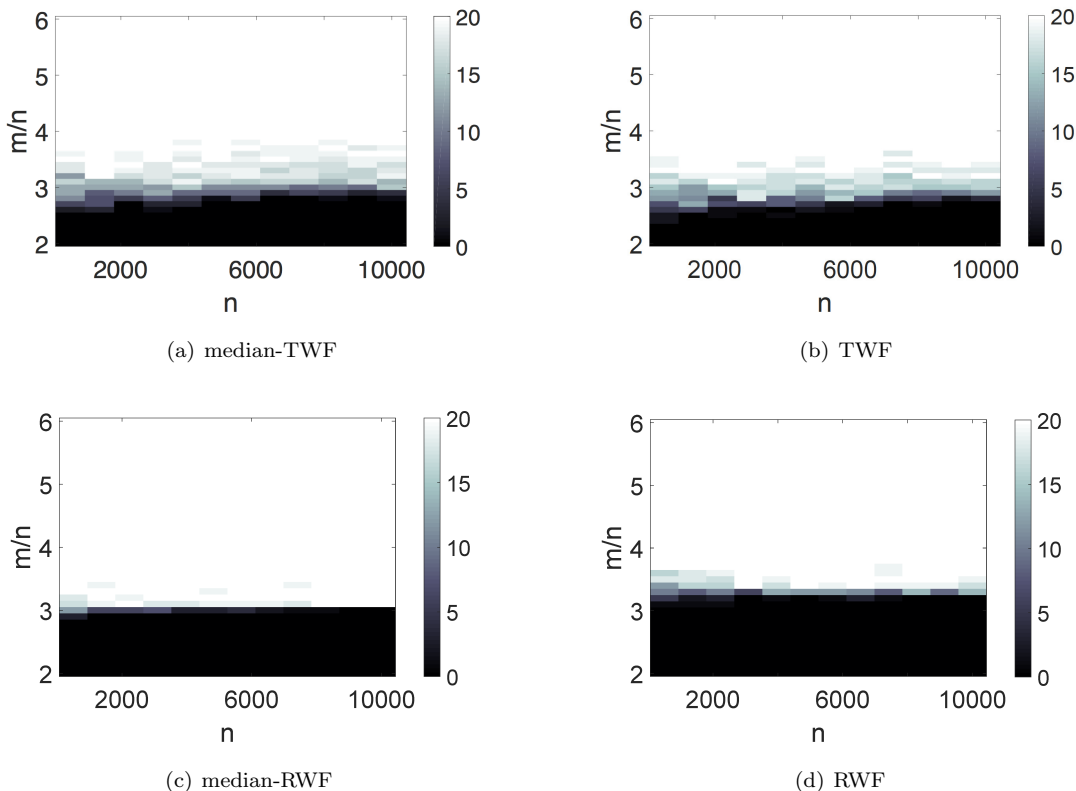


Figure 1: Sample complexity of median-TWF, TWF, median-RWF, and RWF for noise-free data: the gray scale of each cell $(n, m/n)$ indicates the number of successful recovery out of 20 trials.

3.2 Exact Recovery with Sparse Outliers

We next examine the performance of median-TWF and median-RWF in the presence of sparse outliers. We compare the performance of median-TWF and median-RWF with TWF, and also an alternative which we call *trimean-TWF*, based on replacing the sample mean in TWF by the *trimmed mean* for truncation purpose. Specifically, trimean-TWF requires knowing the fraction s of outliers so that samples corresponding to sm largest values in the measurements or gradients are first removed, and truncation is then applied based on the mean of the remaining samples similar to TWF.

We fix the signal length $n = 1000$ and the number of measurements $m = 8000$. Let each measurement y_i be corrupted with probability $s \in [0, 0.4]$ independently, where the corruption value $\eta_i \sim \mathcal{U}(0, \eta_{\max})$ is randomly generated from a uniform distribution. Figure 2 shows the success rate of exact recovery over 100 trials as a function of s at different levels of outlier magnitudes $\eta_{\max}/\|\mathbf{x}\|^2 = 0.1, 1, 10, 100$, for the four algorithms median-TWF, median-RWF, trimean-TWF and TWF.

From Figure 2, it can be seen that median-TWF and median-RWF allow exact recovery as long as s is not too large for all levels of outlier magnitudes, without assuming any knowledge of the outliers, which validates our theoretical analysis. Empirically, median-RWF can tolerate a larger fraction of outliers than median-TWF. This could be due to the fact that the lower-order objective adopted in median-RWF reduces the variance and allows more stable search direction. Unsurprisingly, TWF fails quickly even with a very

⁵The original RWF in [55] does not have sample truncation.

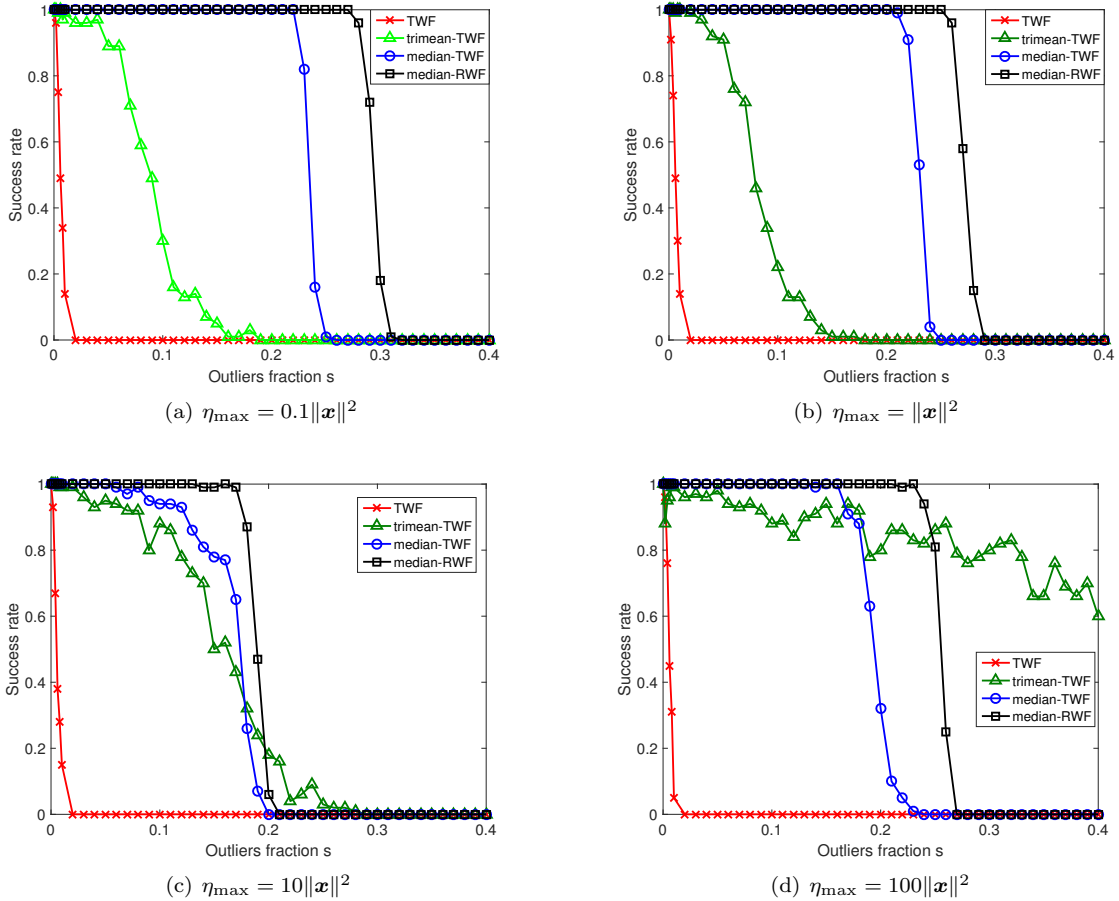


Figure 2: Success rate of exact recovery with respect to the fraction of sparse outliers for median-TWF, median-RWF, trimean-TWF, and TWF at different levels of outlier magnitudes.

small fraction of outliers. No successful instance is observed for TWF when $s \geq 0.02$ irrespective of the value of η_{\max} . Trimean-TWF, even knowing the number of outliers, still does not exhibit a sharp phase transition, and in general underperforms the proposed median-TWF and median-RWF, except when both η_{\max} and s get very large, see Figure 2(d). This is because in this range of large outliers, knowing the fraction s of outliers provides substantial advantage for trimean-TWF to eliminate them, while the sample median can deviate significantly from the true median for large s . Moreover, it is worth mentioning that exact recovery is more challenging for median-TWF and median-RWF when the magnitudes of most outliers are comparable to the measurements, as in Figure 2(c). In such a case, the outliers are more difficult to exclude as opposed to the case with very large outlier magnitudes as in Figure 2(d); and meanwhile, the outlier magnitudes in Figure 2(c) are large enough to affect the accuracy heavily in contrast to the cases in Figures 2(a) and 2(b) where outliers are less prominent.

3.3 Stable Recovery with Sparse Outliers and Dense Noise

We now examine the performance of median-TWF and median-RWF in the presence of both sparse outliers and dense bounded noise. The entries of the dense bounded noise term \mathbf{w} are generated independently from $\mathcal{U}(0, w_{\max})$. The entries of the sparse outlier are then generated as $\eta_i \sim \|\mathbf{w}\| \cdot \text{Bernoulli}(0.1)$ independently. Figure 3(a) and Figure 3(b) depict the relative error $\text{dist}(\mathbf{z}^{(t)}, \mathbf{x})/\|\mathbf{x}\|$ with respect to the iteration count t , when $w_{\max}/\|\mathbf{x}\|^2 = 0.001$ and 0.01 respectively. In the presence of sparse outliers, it can be seen that both median-TWF and median-RWF clearly outperforms TWF under the same situation, and acts as if

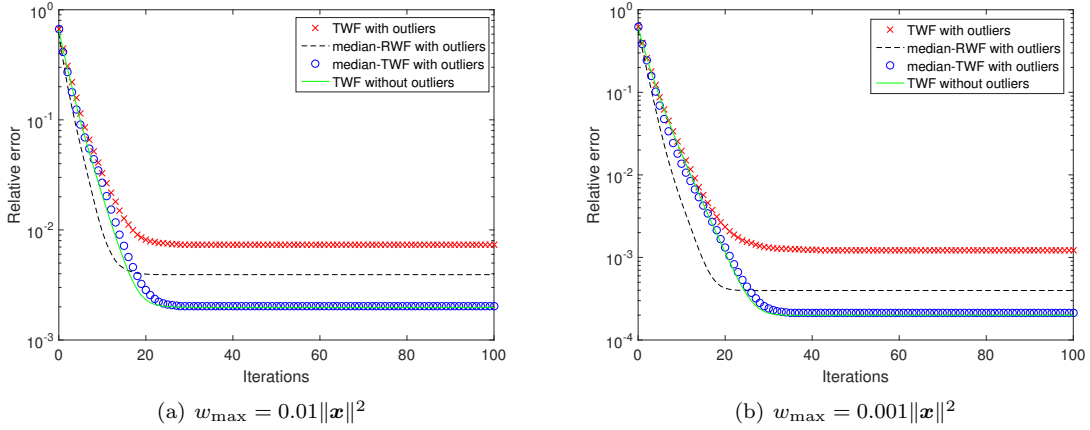


Figure 3: The relative error with respect to the iteration count for median-TWF, median-RWF and TWF with both dense noise and sparse outliers, and TWF with only dense noise. In (a) and (b), the dense noise is generated uniformly at different levels.

the outliers do not exist by achieving almost the same accuracy as TWF without outliers. Moreover, the relative error of the reconstruction using median-TWF or median-RWF has 10 times gain from Figure 3(a) to Figure 3(b) as w_{\max} shrinks by a factor of 10, which corroborates Theorem 2 nicely. Furthermore, it can be seen that median-RWF converges faster than the other algorithms, due to the improved curvature of using low-order objectives, corroborating the result in [55]. On the other hand, median-TWF returns more accurate estimates, due to employing more delicate truncation rules that may help reduce the noise.

Finally, we consider the case when the measurements are corrupted by both Poisson noise and outliers, modeling photon detection in optical imaging applications. We generate each measurement as $y_i \sim \text{Poisson}(|\langle \mathbf{a}_i, \mathbf{x} \rangle|^2)$, for $i = 1, \dots, m$, which is then corrupted with probability $s = 0.1$ by outliers. The entries of the outlier are obtained by first generating $\eta_i \sim \|\mathbf{x}\|^2 \cdot \mathcal{U}(0, 1)$ independently, and then rounding it to the nearest integer. Figure 4 depicts the relative error $\text{dist}(\mathbf{z}^{(t)}, \mathbf{x})/\|\mathbf{x}\|$ with respect to the iteration count t , where median-TWF and median-RWF under both outliers and Poisson noise have almost the same accuracy as, if not better than, TWF under only the Poisson noise.

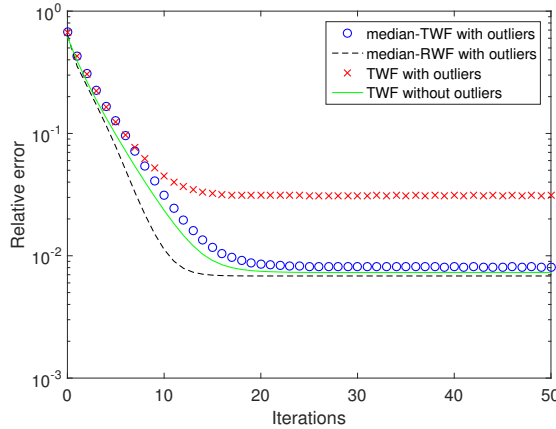


Figure 4: The relative error with respect to the iteration count for median-TWF, median-RWF and TWF with both Poisson noise and sparse outliers, and TWF with only Poisson noise.

4 Preliminaries and Proof Roadmap

Broadly speaking, the proofs for median-TWF and median-RWF follow the same roadmap. The crux is to use the statistical properties of the median to show that the median-truncated gradients satisfy the so-called *Regularity Condition* [9], which guarantees the linear convergence of the update rule, provided the initialization provably lands in a small neighborhood of the true signal.

We first develop a few statistical properties of median that will be useful throughout our analysis in Section 4.1. Section 4.2 analyzes the initialization that is used in both algorithms. We then state the definition of Regularity Condition in Section 4.3 and explain how it leads to the linear convergence rate. We provide separate detailed proofs for two algorithms in Section 5 and Section 6, respectively, because they involve different bounding techniques that may be of independent interest due to different loss functions.

4.1 Properties of Median

We start by the definitions of the quantile of a population distribution and its sample version.

Definition 1 (Generalized quantile function). *Let $0 < p < 1$. For a cumulative distribution function (CDF) F , the generalized quantile function is defined as*

$$F^{-1}(p) = \inf\{x \in \mathbb{R} : F(x) \geq p\}. \quad (27)$$

For simplicity, denote $\theta_p(F) = F^{-1}(p)$ as the p -quantile of F . Moreover for a sample sequence $\{X_i\}_{i=1}^m$, the sample p -quantile $\theta_p(\{X_i\})$ means $\theta_p(\hat{F})$, where \hat{F} is the empirical distribution of the samples $\{X_i\}_{i=1}^m$.

Remark 1. *We note that the median $\text{med}(\{X_i\}) = \theta_{1/2}(\{X_i\})$, and we use both notations interchangeably.*

Next, we show that as long as the sample size is large enough, the sample quantile concentrates around the population quantile (motivated from [11]), as in Lemma 1.

Lemma 1. *Suppose $F(\cdot)$ is cumulative distribution function (i.e., non-decreasing and right-continuous) with continuous density function $F'(\cdot)$. Assume the samples $\{X_i\}_{i=1}^m$ are i.i.d. drawn from F . Let $0 < p < 1$. If $l < F'(\theta) < L$ for all θ in $\{\theta : |\theta - \theta_p| \leq \epsilon\}$, then*

$$|\theta_p(\{X_i\}_{i=1}^m) - \theta_p(F)| < \epsilon \quad (28)$$

holds with probability at least $1 - 2 \exp(-2m\epsilon^2 l^2)$.

Proof. See Appendix A.1. □

Lemma 2 bounds the distance between the median of two sequences.

Lemma 2. *Given a vector $\mathbf{X} = (X_1, X_2, \dots, X_n)$, reorder the entries in a non-decreasing manner*

$$X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n-1)} \leq X_{(n)}.$$

Given another vector $\mathbf{Y} = (Y_1, Y_2, \dots, Y_n)$, then

$$|X_{(k)} - Y_{(k)}| \leq \|\mathbf{X} - \mathbf{Y}\|_\infty, \quad (29)$$

holds for all $k = 1, \dots, n$.

Proof. See Appendix A.2. □

Lemma 3, as a key robustness property of median, suggests that in the presence of outliers, one can bound the sample median from both sides by neighboring quantiles of the corresponding clean samples.

Lemma 3. *Consider clean samples $\{\tilde{X}_i\}_{i=1}^m$. If a fraction s ($s < \frac{1}{2}$) of them are corrupted by outliers, one obtains contaminated samples $\{X_i\}_{i=1}^m$ which contain sm corrupted samples and $(1-s)m$ clean samples. Then for a quantile p such that $s < p < 1-s$, we have*

$$\theta_{p-s}(\{\tilde{X}_i\}) \leq \theta_p(\{X_i\}) \leq \theta_{p+s}(\{\tilde{X}_i\}).$$

Proof. See Appendix A.3. □

Finally, Lemma 4 is related to bound the value of the median, as well as the density at the median for the product of two possibly correlated standard Gaussian random variables.

Lemma 4. *Let $u, v \sim \mathcal{N}(0, 1)$ which can be correlated with the correlation coefficient $|\rho| \leq 1$. Let $r = |uv|$, and $\psi_\rho(x)$ represent the density of r . Denote $\theta_{\frac{1}{2}}(\psi_\rho)$ as the median of r , and the value of $\psi_\rho(x)$ at the median as $\psi_\rho(\theta_{1/2})$. Then for all ρ ,*

$$\begin{aligned} 0.348 &< \theta_{1/2}(\psi_\rho) < 0.455, \\ 0.47 &< \psi_\rho(\theta_{1/2}) < 0.76. \end{aligned}$$

Proof. See Appendix A.4. □

4.2 Robust Initialization with Outliers

Considering the model that the measurements are corrupted by both bounded noise and sparse outliers given by (11), we show that the initialization provided by the median-truncated spectral method in (14) is close enough to the ground truth, i.e., $\text{dist}(\mathbf{z}^{(0)}, \mathbf{x}) \leq \delta \|\mathbf{x}\|$.

Proposition 2. *Fix $\delta > 0$ and $\mathbf{x} \in \mathbb{R}^n$, and consider the model given by (11). Suppose that $\|\mathbf{w}\|_\infty \leq c\|\mathbf{x}\|^2$ for some sufficiently small constant $c > 0$ and that $\|\eta\|_0 \leq sm$ for some sufficiently small constant s . With probability at least $1 - \exp(-\Omega(m))$, the initialization given by the median-truncated spectral method obeys⁶*

$$\text{dist}(\mathbf{z}^{(0)}, \mathbf{x}) \leq \delta \|\mathbf{x}\|, \quad (30)$$

provided that $m > c_0 n$ for some constant $c_0 > 0$.

Proof. See Appendix B.

4.3 Regularity Condition

Once the initialization is guaranteed to be within a small neighborhood of the ground truth, we only need to show that the truncated gradient (16) and (20) satisfy the *Regularity Condition* (RC) [9, 13], which guarantees the geometric convergence of median-TWF/median-RWF once the initialization lands into this neighborhood.

Definition 2. *The gradient $\nabla\ell(\mathbf{z})$ is said to satisfy the Regularity Condition $\text{RC}(\mu, \lambda, c)$ if*

$$\langle \nabla\ell(\mathbf{z}), \mathbf{z} - \mathbf{x} \rangle \geq \frac{\mu}{2} \|\nabla\ell(\mathbf{z})\|^2 + \frac{\lambda}{2} \|\mathbf{z} - \mathbf{x}\|^2 \quad (31)$$

for all \mathbf{z} obeying $\|\mathbf{z} - \mathbf{x}\| \leq c\|\mathbf{x}\|$.

The above RC guarantees that the gradient descent update $\mathbf{z}^{(t+1)} = \mathbf{z}^{(t)} - \mu\nabla\ell(\mathbf{z})$ converges to the true signal \mathbf{x} geometrically [13] if $\mu\lambda < 1$. We repeat this argument below for completeness.

$$\begin{aligned} \text{dist}^2(\mathbf{z} - \mu\nabla\ell(\mathbf{z}), \mathbf{x}) &\leq \|\mathbf{z} - \mu\nabla\ell(\mathbf{z}) - \mathbf{x}\|^2 \\ &= \|\mathbf{z} - \mathbf{x}\|^2 + \|\mu\nabla\ell(\mathbf{z})\|^2 - 2\mu \langle \mathbf{z} - \mathbf{x}, \nabla\ell(\mathbf{z}) \rangle \\ &\leq \|\mathbf{z} - \mathbf{x}\|^2 + \|\mu\nabla\ell(\mathbf{z})\|^2 - \mu^2 \|\nabla\ell(\mathbf{z})\|^2 - \mu\lambda \|\mathbf{z} - \mathbf{x}\|^2 \\ &= (1 - \mu\lambda) \text{dist}^2(\mathbf{z}, \mathbf{x}). \end{aligned}$$

5 Proofs for Median-TWF

We first show that $\nabla\ell_{tr}(\mathbf{z})$ in (16) satisfies the RC for the noise-free case in Section 5.1, and then extend it to the model with only sparse outliers in Section 5.2, thus together with Proposition 2 establishing the global convergence of median-TWF in both cases. Section 5.3 proves Theorem 2 in the presence of both sparse outliers and dense bounded noise.

⁶Notation $f(n) = \Omega(g(n))$ or $f(n) \gtrsim g(n)$ means that there exists a constant $c > 0$ such that $|f(n)| \geq c|g(n)|$.

5.1 Proof of Proposition 1

We consider the noise-free model. The central step to establish the RC is to show that the sample median used in the truncation rule of median-TWF concentrates at the level $\|\mathbf{z} - \mathbf{x}\| \|\mathbf{z}\|$ as stated in the following proposition.

Proposition 3. *If $m > c_0 n \log n$, then with probability at least $1 - c_1 \exp(-c_2 m)$,*

$$0.6 \|\mathbf{z}\| \|\mathbf{z} - \mathbf{x}\| \leq \theta_{0.49}, \theta_{0.5}, \theta_{0.51} (\{ \|\mathbf{a}_i^T \mathbf{x}\|^2 - |\mathbf{a}_i^T \mathbf{z}|^2 \}_{i=1}^m) \leq \|\mathbf{z}\| \|\mathbf{z} - \mathbf{x}\|, \quad (32)$$

holds for all \mathbf{z}, \mathbf{x} satisfying $\|\mathbf{z} - \mathbf{x}\| < 1/11 \|\mathbf{z}\|$.

Proof. Detailed proof is provided in Appendix C.1. \square

We note that a similar property for the sample mean has been shown in [13] as long as the number m of measurements is on the order of n . In fact, the sample median is much more challenging to bound due to its non-linearity, which also causes slightly more measurements compared to the sample mean.

Then we can establish that $\langle \nabla \ell_{tr}(\mathbf{z}), \mathbf{z} - \mathbf{x} \rangle$ is lower bounded on the order of $\|\mathbf{z} - \mathbf{x}\|^2$, as in Proposition 4, and that $\|\nabla \ell_{tr}(\mathbf{z})\|$ is upper bounded on the order of $\|\mathbf{z} - \mathbf{x}\|$, as in Proposition 5.

Proposition 4 (Adapted version of Proposition 2 of [13]). *Consider the noise-free case $y_i = |\mathbf{a}_i^T \mathbf{x}|^2$ for $i = 1, \dots, m$, and any fixed constant $\epsilon > 0$. Under the condition (17), if $m > c_0 n \log n$, then with probability at least $1 - c_1 \exp(-c_2 \epsilon^{-2} m)$,*

$$\langle \nabla \ell_{tr}(\mathbf{z}), \mathbf{z} - \mathbf{x} \rangle \geq \left\{ 1.99 - 2(\zeta_1 + \zeta_2) - \sqrt{8/\pi} \alpha_h^{-1} - \epsilon \right\} \|\mathbf{z} - \mathbf{x}\|^2 \quad (33)$$

holds uniformly over all $\mathbf{x}, \mathbf{z} \in \mathbb{R}^n$ satisfying

$$\frac{\|\mathbf{z} - \mathbf{x}\|}{\|\mathbf{z}\|} \leq \min \left\{ \frac{1}{11}, \frac{\alpha_l}{\alpha_h}, \frac{\alpha_l}{6}, \frac{\sqrt{98/3}(\alpha_l)^2}{2\alpha_u + \alpha_l} \right\}, \quad (34)$$

where $c_0, c_1, c_2 > 0$ are some universal constants, and $\zeta_1, \zeta_2, \alpha_l, \alpha_u$ and α_h are defined in (17).

The proof of Proposition 4 adapts the proof of Proposition 2 of [13], by properly setting parameters based on the properties of sample median. For completeness, we include a short outline of the proof in Appendix C.2.

Proposition 5 (Lemma 7 of [13]). *Under the same condition as in Proposition 4, if $m > c_0 n$, then there exist some constants $c_1, c_2 > 0$ such that with probability at least $1 - c_1 \exp(-c_2 m)$,*

$$\|\nabla \ell_{tr}(\mathbf{z})\| \leq (1 + \delta) \cdot 2\sqrt{1.02 + 2/\alpha_h} \|\mathbf{z} - \mathbf{x}\| \quad (35)$$

holds uniformly over all $\mathbf{x}, \mathbf{z} \in \mathbb{R}^n$ satisfying

$$\frac{\|\mathbf{z} - \mathbf{x}\|}{\|\mathbf{z}\|} \leq \min \left\{ \frac{1}{11}, \frac{\alpha_l}{\alpha_h}, \frac{\alpha_l}{6}, \frac{\sqrt{98/3}(\alpha_l)^2}{2\alpha_u + \alpha_l} \right\}, \quad (36)$$

where δ can be arbitrarily small as long as m/n sufficiently large, and α_l, α_u and α_h are given in (17).

Proof. See the proof of Lemma 7 in [13]. \square

With these two propositions and (17), RC is guaranteed by setting

$$\begin{aligned} \mu < \mu_0 &:= \frac{(1.99 - 2(\zeta_1 + \zeta_2) - \sqrt{8/\pi} \alpha_h^{-1})}{2(1 + \delta)^2 \cdot (1.02 + 2/\alpha_h)}, \\ \lambda + \mu \cdot 4(1 + \delta)^2 \cdot (1.02 + 2/\alpha_h) &< 2 \left\{ 1.99 - 2(\zeta_1 + \zeta_2) - \sqrt{8/\pi} \alpha_h^{-1} - \epsilon \right\}. \end{aligned}$$

5.2 Proof of Theorem 1

We next consider the model (11) with only sparse outliers. It suffices to show that $\nabla \ell_{tr}(\mathbf{z})$ continues to satisfy the RC. The critical step is to bound the sample median of the corrupted measurements. Lemma 3 yields

$$\theta_{\frac{1}{2}-s}(\{|(\mathbf{a}_i^T \mathbf{x})^2 - (\mathbf{a}_i^T \mathbf{z})^2|\}) \leq \theta_{\frac{1}{2}}(\{|y_i - (\mathbf{a}_i^T \mathbf{z})^2|\}) \leq \theta_{\frac{1}{2}+s}(\{|(\mathbf{a}_i^T \mathbf{x})^2 - (\mathbf{a}_i^T \mathbf{z})^2|\}). \quad (37)$$

For simplicity of notation, we let $\mathbf{h} := \mathbf{z} - \mathbf{x}$. Then for the instance of $s = 0.01$, by Proposition 3, we have with probability at least $1 - 2 \exp(-\Omega(m))$,

$$0.6 \|\mathbf{z}\| \|\mathbf{h}\| \leq \theta_{\frac{1}{2}}(\{|y_i - (\mathbf{a}_i^T \mathbf{z})^2|\}) \leq \|\mathbf{z}\| \|\mathbf{h}\|. \quad (38)$$

To differentiate from \mathcal{E}_2^i , we define $\tilde{\mathcal{E}}_2^i := \left\{ |(\mathbf{a}_i^T \mathbf{x})^2 - (\mathbf{a}_i^T \mathbf{z})^2| \leq \alpha_h \text{med} \{|y_i - (\mathbf{a}_i^T \mathbf{z})^2|\} \frac{\|\mathbf{a}_i^T \mathbf{z}\|}{\|\mathbf{z}\|} \right\}$. We then have

$$\begin{aligned} \nabla \ell_{tr}(\mathbf{z}) &= \frac{1}{m} \sum_{i=1}^m \frac{(\mathbf{a}_i^T \mathbf{z})^2 - y_i}{\mathbf{a}_i^T \mathbf{z}} \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i} \\ &= \underbrace{\frac{1}{m} \sum_{i=1}^m \frac{(\mathbf{a}_i^T \mathbf{z})^2 - (\mathbf{a}_i^T \mathbf{x})^2}{\mathbf{a}_i^T \mathbf{z}} \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i \cap \tilde{\mathcal{E}}_2^i}}_{\nabla^{clean} \ell_{tr}(\mathbf{z})} + \underbrace{\frac{1}{m} \sum_{i \in S} \left(\frac{(\mathbf{a}_i^T \mathbf{z})^2 - y_i}{\mathbf{a}_i^T \mathbf{z}} \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i} - \frac{(\mathbf{a}_i^T \mathbf{z})^2 - (\mathbf{a}_i^T \mathbf{x})^2}{\mathbf{a}_i^T \mathbf{z}} \mathbf{1}_{\mathcal{E}_1^i \cap \tilde{\mathcal{E}}_2^i} \right) \mathbf{a}_i}_{\nabla^{extra} \ell_{tr}(\mathbf{z})}. \end{aligned}$$

Choosing ϵ small enough, it is easy to verify that Propositions 4 and 5 are still valid on $\nabla^{clean} \ell_{tr}(\mathbf{z})$. Thus, one has

$$\begin{aligned} \langle \nabla^{clean} \ell_{tr}(\mathbf{z}), \mathbf{h} \rangle &\geq \left\{ 1.99 - 2(\zeta_1 + \zeta_2) - \sqrt{8/\pi} \alpha_h^{-1} - \epsilon \right\} \|\mathbf{h}\|^2, \\ \|\nabla^{clean} \ell_{tr}(\mathbf{z})\| &\leq (1 + \delta) \cdot 2\sqrt{1.02 + 2/\alpha_h} \|\mathbf{h}\|. \end{aligned}$$

We next bound the contribution of $\nabla^{extra} \ell_{tr}(\mathbf{z})$. Introduce $\mathbf{q} = [q_1, \dots, q_m]^T$, where

$$q_i := \left(\frac{(\mathbf{a}_i^T \mathbf{z})^2 - y_i}{\mathbf{a}_i^T \mathbf{z}} \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i} - \frac{(\mathbf{a}_i^T \mathbf{z})^2 - (\mathbf{a}_i^T \mathbf{x})^2}{\mathbf{a}_i^T \mathbf{z}} \mathbf{1}_{\mathcal{E}_1^i \cap \tilde{\mathcal{E}}_2^i} \right) \mathbf{1}_{\{i \in S\}}.$$

It can be seen that $|q_i| \leq 2\alpha_h \|\mathbf{h}\|$. Thus $\|\mathbf{q}\| \leq \sqrt{sm} \cdot 2\alpha_h \|\mathbf{h}\|$, and

$$\begin{aligned} \|\nabla^{extra} \ell_{tr}(\mathbf{z})\| &= \frac{1}{m} \|\mathbf{A}^T \mathbf{q}\| \leq 2(1 + \delta) \sqrt{s} \alpha_h \|\mathbf{h}\|, \\ |\langle \nabla^{extra} \ell_{tr}(\mathbf{z}), \mathbf{h} \rangle| &\leq \|\mathbf{h}\| \cdot \left\| \frac{1}{m} \nabla^{extra} \ell_{tr}(\mathbf{z}) \right\| \leq 2(1 + \delta) \sqrt{s} \alpha_h \|\mathbf{h}\|^2, \end{aligned}$$

where $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_m]^T$. Then, we have

$$\begin{aligned} -\langle \nabla \ell_{tr}(\mathbf{z}), \mathbf{h} \rangle &\geq \langle \nabla^{clean} \ell_{tr}(\mathbf{z}), \mathbf{h} \rangle - |\langle \nabla^{extra} \ell_{tr}(\mathbf{z}), \mathbf{h} \rangle| \\ &\geq \left(1.99 - 2(\zeta_1 + \zeta_2) - \sqrt{8/\pi} \alpha_h^{-1} - \epsilon - 2(1 + \delta) \sqrt{s} \alpha_h \right) \|\mathbf{h}\|^2, \end{aligned}$$

and

$$\begin{aligned} \|\nabla \ell_{tr}(\mathbf{z})\| &\leq \|\nabla^{clean} \ell_{tr}(\mathbf{z})\| + \|\nabla^{extra} \ell_{tr}(\mathbf{z})\| \\ &\leq 2(1 + \delta) \left(\sqrt{1.02 + 2/\alpha_h} + \sqrt{s} \alpha_h \right) \|\mathbf{h}\|. \end{aligned} \quad (39)$$

Therefore, the RC is guaranteed if μ, λ, ϵ are chosen properly and s is sufficiently small.

5.3 Proof of Theorem 2

We consider the model (11), and split our analysis of the gradient loop into two regimes.

- **Regime 1:** $c_4 \|\mathbf{z}\| \geq \|\mathbf{h}\| \geq c_3 \frac{\|\mathbf{w}\|_\infty}{\|\mathbf{z}\|}$. In this regime, error contraction by each gradient step is given by

$$\text{dist}(\mathbf{z} - \mu \nabla \ell_{tr}(\mathbf{z}), \mathbf{x}) \leq (1 - \rho) \text{dist}(\mathbf{z}, \mathbf{x}).$$

It suffices to justify that $\nabla \ell_{tr}(\mathbf{z})$ satisfies the RC. Denote $\tilde{y}_i := (\mathbf{a}_i^T \mathbf{x})^2 + w_i$. Then by Lemma 3, we have

$$\theta_{\frac{1}{2}-s} \{|\tilde{y}_i - (\mathbf{a}_i^T \mathbf{z})^2|\} \leq \text{med} \{|y_i - (\mathbf{a}_i^T \mathbf{z})^2|\} \leq \theta_{\frac{1}{2}+s} \{|\tilde{y}_i - (\mathbf{a}_i^T \mathbf{z})^2|\}.$$

Moreover, by Lemma 2 we have

$$\begin{aligned} \left| \theta_{\frac{1}{2}+s} \{|\tilde{y}_i - (\mathbf{a}_i^T \mathbf{z})^2|\} - \theta_{\frac{1}{2}+s} \{|(\mathbf{a}_i^T \mathbf{x})^2 - (\mathbf{a}_i^T \mathbf{z})^2|\} \right| &\leq \|\mathbf{w}\|_\infty, \\ \left| \theta_{\frac{1}{2}-s} \{|\tilde{y}_i - (\mathbf{a}_i^T \mathbf{z})^2|\} - \theta_{\frac{1}{2}-s} \{|(\mathbf{a}_i^T \mathbf{x})^2 - (\mathbf{a}_i^T \mathbf{z})^2|\} \right| &\leq \|\mathbf{w}\|_\infty. \end{aligned}$$

Assume that $s = 0.01$. By Proposition 3, if c_3 is sufficiently large (i.e., $c_3 > 100$), we still have

$$0.6 \|\mathbf{x} - \mathbf{z}\| \|\mathbf{z}\| \leq \text{med} \{|y_i - (\mathbf{a}_i^T \mathbf{z})^2|\} \leq \|\mathbf{x} - \mathbf{z}\| \|\mathbf{z}\|. \quad (40)$$

Furthermore, recall $\tilde{\mathcal{E}}_2^i := \{ |(\mathbf{a}_i^T \mathbf{x})^2 - (\mathbf{a}_i^T \mathbf{z})^2| \leq \alpha_h \text{med} \{ |(\mathbf{a}_i^T \mathbf{z})^2 - y_i | \} \frac{|\mathbf{a}_i^T \mathbf{z}|}{\|\mathbf{z}\|} \}$. Then,

$$\begin{aligned} \nabla \ell_{tr}(\mathbf{z}) &= \frac{1}{m} \sum_{i=1}^m \frac{(\mathbf{a}_i^T \mathbf{z})^2 - y_i}{\mathbf{a}_i^T \mathbf{z}} \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i} \\ &= \frac{1}{m} \left(\underbrace{\sum_{i \notin S} \frac{(\mathbf{a}_i^T \mathbf{z})^2 - (\mathbf{a}_i^T \mathbf{x})^2}{\mathbf{a}_i^T \mathbf{z}} \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i}}_{\nabla^{clean} \ell_{tr}(\mathbf{z})} + \sum_{i \in S} \frac{(\mathbf{a}_i^T \mathbf{z})^2 - (\mathbf{a}_i^T \mathbf{x})^2}{\mathbf{a}_i^T \mathbf{z}} \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i \cap \tilde{\mathcal{E}}_2^i} \right) \\ &\quad - \underbrace{\frac{1}{m} \sum_{i \notin S} \frac{w_i}{\mathbf{a}_i^T \mathbf{z}} \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i}}_{\nabla^{noise} \ell_{tr}(\mathbf{z})} + \underbrace{\frac{1}{m} \sum_{i \in S} \left(\frac{(\mathbf{a}_i^T \mathbf{z})^2 - y_i}{\mathbf{a}_i^T \mathbf{z}} \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i} - \frac{(\mathbf{a}_i^T \mathbf{z})^2 - (\mathbf{a}_i^T \mathbf{x})^2}{\mathbf{a}_i^T \mathbf{z}} \mathbf{1}_{\mathcal{E}_1^i \cap \tilde{\mathcal{E}}_2^i} \right) \mathbf{a}_i}_{\nabla^{extra} \ell_{tr}(\mathbf{z})}. \end{aligned}$$

For $i \notin S$, the inclusion property (i.e. $\mathcal{E}_3^i \subseteq \mathcal{E}_2^i \subseteq \mathcal{E}_4^i$) holds because

$$|y_i - (\mathbf{a}_i^T \mathbf{z})^2| \in |(\mathbf{a}_i^T \mathbf{x})^2 - (\mathbf{a}_i^T \mathbf{z})^2| \pm |w_i|$$

and $|w_i| \leq \frac{1}{c_3} \|\mathbf{h}\| \|\mathbf{z}\|$ for some sufficient large c_3 . For $i \in S$, the inclusion $\mathcal{E}_3^i \subseteq \tilde{\mathcal{E}}_2^i \subseteq \mathcal{E}_4^i$ holds because of (40). All the proof arguments for Propositions 4 and 5 are also valid for $\nabla^{clean} \ell_{tr}(\mathbf{z})$, and thus we have

$$\begin{aligned} \langle \nabla^{clean} \ell_{tr}(\mathbf{z}), \mathbf{h} \rangle &\geq \left\{ 1.99 - 2(\zeta_1 + \zeta_2) - \sqrt{8/\pi} \alpha_h^{-1} - \epsilon \right\} \|\mathbf{h}\|^2, \\ \|\nabla^{clean} \ell_{tr}(\mathbf{z})\| &\leq (1 + \delta) \cdot 2\sqrt{1.02 + 2/\alpha_h} \|\mathbf{h}\|. \end{aligned}$$

Next, we turn to control the contribution of the noise. Let $\tilde{w}_i = \frac{w_i}{\mathbf{a}_i^T \mathbf{z}} \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i}$, and then we have

$$\|\nabla^{noise} \ell_{tr}(\mathbf{z})\| = \left\| \frac{1}{m} \mathbf{A}^T \tilde{\mathbf{w}} \right\| \leq \left\| \frac{1}{\sqrt{m}} \mathbf{A}^T \right\| \left\| \frac{\tilde{\mathbf{w}}}{\sqrt{m}} \right\| \leq (1 + \delta) \|\tilde{\mathbf{w}}\|_\infty \leq (1 + \delta) \frac{\|\mathbf{w}\|_\infty}{\alpha_l \|\mathbf{z}\|},$$

when m/n is sufficiently large. Given the regime condition $\|\mathbf{h}\| \geq c_3 \frac{\|\mathbf{w}\|_\infty}{\|\mathbf{z}\|}$, we further have

$$\begin{aligned} \|\nabla^{noise} \ell_{tr}(\mathbf{z})\| &\leq \frac{(1 + \delta)}{c_3 \alpha_l} \|\mathbf{h}\|, \\ |\langle \nabla^{noise} \ell_{tr}(\mathbf{z}), \mathbf{h} \rangle| &\leq \|\nabla^{noise} \ell_{tr}(\mathbf{z})\| \cdot \|\mathbf{h}\| \leq \frac{(1 + \delta)}{c_3 \alpha_l} \|\mathbf{h}\|^2. \end{aligned}$$

We next bound the contribution of $\nabla^{extra} \ell_{tr}(\mathbf{z})$. Introduce $\mathbf{q} = [q_1, \dots, q_m]^T$, where

$$q_i := \left(\frac{(\mathbf{a}_i^T \mathbf{z})^2 - y_i}{\mathbf{a}_i^T \mathbf{z}} \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i} - \frac{(\mathbf{a}_i^T \mathbf{z})^2 - (\mathbf{a}_i^T \mathbf{x})^2}{\mathbf{a}_i^T \mathbf{z}} \mathbf{1}_{\mathcal{E}_1^i \cap \bar{\mathcal{E}}_2^i} \right) \mathbf{1}_{\{i \in \mathcal{S}\}}.$$

Then $|q_i| \leq 2\alpha_h \|\mathbf{h}\|$, and $\|\mathbf{q}\| \leq \sqrt{sm} \cdot 2\alpha_h \|\mathbf{h}\|$. We thus have

$$\begin{aligned} \|\nabla^{extra} \ell_{tr}(\mathbf{z})\| &= \frac{1}{m} \|\mathbf{A}^T \mathbf{q}\| \leq 2(1 + \delta) \sqrt{s} \alpha_h \|\mathbf{h}\|, \\ |\langle \nabla^{extra} \ell_{tr}(\mathbf{z}), \mathbf{h} \rangle| &\leq \|\mathbf{h}\| \cdot \|\nabla^{extra} \ell_{tr}(\mathbf{z})\| \leq 2(1 + \delta) \sqrt{s} \alpha_h \|\mathbf{h}\|^2. \end{aligned}$$

Putting these together, one has

$$\begin{aligned} \langle \nabla \ell_{tr}(\mathbf{z}), \mathbf{h} \rangle &\geq \langle \nabla^{clean} \ell_{tr}(\mathbf{z}), \mathbf{h} \rangle - |\langle \nabla^{noise} \ell_{tr}(\mathbf{z}), \mathbf{h} \rangle| - |\langle \nabla^{extra} \ell_{tr}(\mathbf{z}), \mathbf{h} \rangle| \\ &\geq \left(1.99 - 2(\zeta_1 + \zeta_2) - \sqrt{8/\pi} \alpha_h^{-1} - \epsilon - (1 + \delta)(1/(c_3 \alpha_z^l) + 2\sqrt{s} \alpha_h) \right) \|\mathbf{h}\|^2, \end{aligned} \quad (41)$$

and

$$\begin{aligned} \|\nabla \ell_{tr}(\mathbf{z})\| &\leq \|\nabla^{clean} \ell_{tr}(\mathbf{z})\| + \|\nabla^{noise} \ell_{tr}(\mathbf{z})\| + \|\nabla^{extra} \ell_{tr}(\mathbf{z})\| \\ &\leq (1 + \delta) \left(2\sqrt{1.02 + 2/\alpha_h} + 1/(c_3 \alpha_z^l) + 2\sqrt{s} \alpha_h \right) \|\mathbf{h}\|. \end{aligned} \quad (42)$$

The RC is guaranteed if μ, λ, ϵ are chosen properly, c_3 is sufficiently large and s is sufficiently small.

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

$$\text{dist}(\mathbf{z} - \mu \nabla \ell_{tr}(\mathbf{z}), \mathbf{x}) \leq c_5 \frac{\|\mathbf{w}\|_\infty}{\|\mathbf{x}\|}$$

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

6 Proofs for Median-RWF

We first show that $\nabla \mathcal{R}_{tr}(\mathbf{z})$ in (20) satisfies the RC for the noise-free case in Section 6.1, and then extend it to the model with only sparse outliers in Section 6.2, thus together with Proposition 2 establishing the global convergence of median-RWF in both cases. Section 6.3 proves Theorem 2 in the presence of both sparse outliers and dense bounded noise.

6.1 Proof of Proposition 1

The central step to establish the RC is to show that the sample median used in the truncation rule of median-RWF concentrates on the order of $\|\mathbf{z} - \mathbf{x}\|$ as stated in the following proposition.

Proposition 6. *If $m > c_0 n \log n$, then with probability at least $1 - c_1 \exp(-c_2 m)$,*

$$0.5 \|\mathbf{z} - \mathbf{x}\| \leq \theta_{0.49}, \theta_{1/2}, \theta_{0.51} \left(\left\{ \|\mathbf{a}_i^T \mathbf{z}\| - \|\mathbf{a}_i^T \mathbf{x}\| \right\}_{i=1}^m \right) \leq 0.8 \|\mathbf{z} - \mathbf{x}\|, \quad (43)$$

holds for all \mathbf{z}, \mathbf{x} satisfying $\|\mathbf{z} - \mathbf{x}\| < 1/11 \|\mathbf{z}\|$.

Proof. See Appendix D.1. □

Next we give a bound on the left hand side of RC.

Proposition 7 (Adapted version of Proposition 2 of [13]). *Consider the noise-free measurements $y_i = |\mathbf{a}_i^T \mathbf{x}|$ and any fixed constant $\epsilon > 0$. If $m > c_0 n \log n$, then with probability at least $1 - c_1 \exp(-c_2 m)$,*

$$\langle \nabla \mathcal{R}_{tr}(\mathbf{z}), \mathbf{z} - \mathbf{x} \rangle \geq \{0.88 - \zeta'_1 - \zeta'_2 - \epsilon\} \|\mathbf{z} - \mathbf{x}\|^2 \quad (44)$$

holds uniformly over all $\mathbf{x}, \mathbf{z} \in \mathbb{R}^n$ satisfying $\frac{\|\mathbf{z} - \mathbf{x}\|}{\|\mathbf{z}\|} \leq \frac{1}{20}$, where $c_0, c_1, c_2 > 0$ are some universal constants, and ζ'_1, ζ'_2 are given by

$$\begin{aligned} \zeta'_1 &:= 1 - \min \left\{ \mathbb{E} \left[\xi^2 \mathbf{1}_{\{\xi \geq 0.5\sqrt{1.01}\alpha'_h \frac{\|\mathbf{z} - \mathbf{x}\|}{\|\mathbf{z}\|}\}} \right], \mathbb{E} \left[\mathbf{1}_{\{\xi \geq 0.5\sqrt{1.01}\alpha'_h \frac{\|\mathbf{z} - \mathbf{x}\|}{\|\mathbf{z}\|}\}} \right] \right\} \\ \zeta'_2 &:= \mathbb{E} \left[\xi^2 \mathbf{1}_{\{|\xi| > 0.5\sqrt{0.99}\alpha'_h\}} \right] \end{aligned}$$

for some $\xi \sim \mathcal{N}(0, 1)$ and $\alpha'_h = 5$.

Proof. See Appendix D.2. □

Proposition 7 indicates that $\langle \nabla \mathcal{R}_{tr}(\mathbf{z}), \mathbf{z} - \mathbf{x} \rangle$ is lower bounded by $\|\mathbf{z} - \mathbf{x}\|^2$ with some positive constant coefficient. In order to prove the RC, it suffices to show that $\|\nabla \mathcal{R}_{tr}(\mathbf{z})\|$ is upper bounded by the order of $\|\mathbf{z} - \mathbf{x}\|$ when \mathbf{z} is within the neighborhood of true signal \mathbf{x} .

Proposition 8 (Lemma 7 of [13]). *If $m > c_0 n$, then there exist some constants $c_1, c_2 > 0$ such that with probability at least $1 - c_1 \exp(-c_2 m)$,*

$$\|\nabla \mathcal{R}_{tr}(\mathbf{z})\| \leq (1.8 + \delta) \|\mathbf{z} - \mathbf{x}\| \quad (45)$$

holds uniformly over all $\mathbf{x}, \mathbf{z} \in \mathbb{R}^n$ satisfying $\|\mathbf{x} - \mathbf{z}\| \leq \frac{1}{11} \|\mathbf{x}\|$ where δ can be arbitrarily small as long as c_0 sufficiently large.

Proof. See Appendix D.3. □

With the above two propositions, RC is guaranteed by setting $\mu < \mu_0 := \frac{2(0.88 - \zeta'_1 - \zeta'_2 - \epsilon)}{(1.8 + \delta)^2}$ and $\lambda + \mu \cdot (1.8 + \delta)^2 < 2(0.88 - \zeta'_1 - \zeta'_2 - \epsilon)$.

6.2 Proof of Theorem 1

We consider the model (11) with only outliers, i.e., $y_i = |\langle \mathbf{a}_i, \mathbf{x} \rangle|^2 + \eta_i$ for $i = 1, \dots, m$. It suffices to show that $\nabla \mathcal{R}_{tr}(\mathbf{z})$ satisfies the RC. The critical step is to lower and upper bound the sample median of the corrupted measurements. Lemma 3 yields

$$\theta_{\frac{1}{2}-s}(\{|\mathbf{a}_i^T \mathbf{x}| - |\mathbf{a}_i^T \mathbf{z}|\}) \leq \theta_{\frac{1}{2}}(\{|\sqrt{y_i} - |\mathbf{a}_i^T \mathbf{z}|\}) \leq \theta_{\frac{1}{2}+s}(\{|\mathbf{a}_i^T \mathbf{x}| - |\mathbf{a}_i^T \mathbf{z}|\}). \quad (46)$$

For the simplicity of notation, we let $\mathbf{h} := \mathbf{z} - \mathbf{x}$. Then for the instance of $s = 0.01$, Proposition 6 yields that if $m > c_0 n \log n$, then

$$0.5 \|\mathbf{h}\| \leq \theta_{\frac{1}{2}}(\{|\sqrt{y_i} - |\mathbf{a}_i^T \mathbf{z}|\}) \leq 0.8 \|\mathbf{h}\| \quad (47)$$

holds with probability at least $1 - 2 \exp(-\Omega(m))$.

To differentiate from \mathcal{T}^i , we define $\tilde{\mathcal{T}}^i := \{|\mathbf{a}_i^T \mathbf{x}| - |\mathbf{a}_i^T \mathbf{z}| \leq \alpha'_h \text{med} \{|\sqrt{y_i} - |\mathbf{a}_i^T \mathbf{z}|\}\}$. We then have

$$\begin{aligned} \nabla \mathcal{R}_{tr}(\mathbf{z}) &= \frac{1}{m} \sum_{i=1}^m (|\mathbf{a}_i^T \mathbf{z}| - \sqrt{y_i}) \mathbf{a}_i \mathbf{1}_{\mathcal{T}^i} \\ &= \underbrace{\frac{1}{m} \sum_{i=1}^m (|\mathbf{a}_i^T \mathbf{z}| - \sqrt{y_i}) \mathbf{a}_i \mathbf{1}_{\tilde{\mathcal{T}}^i}}_{\nabla^{clean} \mathcal{R}_{tr}(\mathbf{z})} + \underbrace{\frac{1}{m} \sum_{i \in S} ((|\mathbf{a}_i^T \mathbf{z}| - \sqrt{y_i}) \mathbf{1}_{\mathcal{T}^i} - (|\mathbf{a}_i^T \mathbf{z}| - |\mathbf{a}_i^T \mathbf{x}|) \mathbf{1}_{\tilde{\mathcal{T}}^i}) \mathbf{a}_i}_{\nabla^{extra} \mathcal{R}_{tr}(\mathbf{z})}. \end{aligned}$$

Under the condition (47), the inclusion property (i.e., $\mathcal{T}_1^i \subseteq \tilde{\mathcal{T}}^i \subseteq \mathcal{T}_2^i$) holds, and all the proof arguments for Propositions 7 and 8 are also valid to $\nabla^{clean}\mathcal{R}_{tr}(\mathbf{z})$. Thus, one has

$$\begin{aligned} \langle \nabla^{clean}\mathcal{R}_{tr}(\mathbf{z}), \mathbf{h} \rangle &\geq (0.88 - \zeta'_1 - \zeta'_2 - \epsilon) \|\mathbf{h}\|^2 \\ \|\nabla^{clean}\mathcal{R}_{tr}(\mathbf{z})\| &\leq (1.8 + \delta) \|\mathbf{h}\|. \end{aligned}$$

We next bound the contribution of $\nabla^{extra}\mathcal{R}_{tr}(\mathbf{z})$. Introduce $\mathbf{q} = [q_1, \dots, q_m]^T$, where

$$q_i := ((|\mathbf{a}_i^T \mathbf{z}| - \sqrt{y_i}) \mathbf{1}_{\mathcal{T}^i} - (|\mathbf{a}_i^T \mathbf{z}| - |\mathbf{a}_i^T \mathbf{x}|) \mathbf{1}_{\tilde{\mathcal{T}}^i}) \mathbf{1}_{\{i \in S\}},$$

and then $|q_i| \leq 1.6\alpha'_h \|\mathbf{h}\|$. Thus, $\|\mathbf{q}\| \leq \sqrt{sm} \cdot 1.6\alpha'_h \|\mathbf{h}\|$, and

$$\begin{aligned} \|\nabla^{extra}\mathcal{R}_{tr}(\mathbf{z})\| &= \frac{1}{m} \|\mathbf{A}^T \mathbf{q}\| \leq 1.6(1 + \delta) \sqrt{s} \alpha'_h \|\mathbf{h}\|, \\ |\langle \nabla^{extra}\mathcal{R}_{tr}(\mathbf{z}), \mathbf{h} \rangle| &\leq \|\mathbf{h}\| \cdot \|\nabla^{extra}\mathcal{R}_{tr}(\mathbf{z})\| \leq 1.6(1 + \delta) \sqrt{s} \alpha'_h \|\mathbf{h}\|^2, \end{aligned}$$

where $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_m]^T$. Then, we have

$$\begin{aligned} \langle \nabla \mathcal{R}_{tr}(\mathbf{z}), \mathbf{h} \rangle &\geq \langle \nabla^{clean}\mathcal{R}_{tr}(\mathbf{z}), \mathbf{h} \rangle - |\langle \nabla^{extra}\mathcal{R}_{tr}(\mathbf{z}), \mathbf{h} \rangle| \\ &\geq (0.88 - \zeta'_1 - \zeta'_2 - \epsilon - 1.6(1 + \delta) \sqrt{s} \alpha'_h) \|\mathbf{h}\|^2, \end{aligned}$$

and

$$\begin{aligned} \|\nabla \mathcal{R}_{tr}(\mathbf{z})\| &\leq \|\nabla^{clean}\mathcal{R}_{tr}(\mathbf{z})\| + \|\nabla^{extra}\mathcal{R}_{tr}(\mathbf{z})\| \\ &\leq (1.8 + \delta + 1.6(1 + \delta) \sqrt{s} \alpha'_h) \|\mathbf{h}\|. \end{aligned}$$

Therefore the RC is guaranteed if μ, λ are chosen properly, δ is chosen sufficiently small and s is sufficiently small.

6.3 Proof of Theorem 2

We consider the model (11) with outliers and bounded noise. We split our analysis of the gradient loop into two regimes.

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

$$\text{dist}(\mathbf{z} - \mu \nabla \mathcal{R}_{tr}(\mathbf{z}), \mathbf{x}) \leq (1 - \rho) \text{dist}(\mathbf{z}, \mathbf{x}). \quad (48)$$

It suffices to justify that $\nabla \mathcal{R}_{tr}(\mathbf{z})$ satisfies the RC. Denote $\tilde{y}_i := (\mathbf{a}_i^T \mathbf{x})^2 + w_i$. Then by Lemma 3, we have

$$\theta_{\frac{1}{2}-s} \left\{ \left| \sqrt{\tilde{y}_i} - |\mathbf{a}_i^T \mathbf{z}| \right| \right\} \leq \text{med} \left\{ \left| \sqrt{\tilde{y}_i} - |\mathbf{a}_i^T \mathbf{x}| \right| \right\} \leq \theta_{\frac{1}{2}+s} \left\{ \left| \sqrt{\tilde{y}_i} - |\mathbf{a}_i^T \mathbf{z}| \right| \right\}.$$

Moreover, by Lemma 2 we have

$$\begin{aligned} \left| \theta_{\frac{1}{2}+s} \left\{ \left| \sqrt{\tilde{y}_i} - |\mathbf{a}_i^T \mathbf{z}| \right| \right\} - \theta_{\frac{1}{2}+s} \left\{ \left| |\mathbf{a}_i^T \mathbf{x}| - |\mathbf{a}_i^T \mathbf{z}| \right| \right\} \right| &\leq \sqrt{\|\mathbf{w}\|_\infty}, \\ \left| \theta_{\frac{1}{2}-s} \left\{ \left| \sqrt{\tilde{y}_i} - |\mathbf{a}_i^T \mathbf{z}| \right| \right\} - \theta_{\frac{1}{2}-s} \left\{ \left| |\mathbf{a}_i^T \mathbf{x}| - |\mathbf{a}_i^T \mathbf{z}| \right| \right\} \right| &\leq \sqrt{\|\mathbf{w}\|_\infty}. \end{aligned}$$

Assume that $s = 0.01$. By Proposition 6, if c_3 is sufficiently large (i.e., $c_3 > 100$), we still have

$$0.5 \|\mathbf{h}\| \leq \text{med} \left\{ \left| \sqrt{\tilde{y}_i} - |\mathbf{a}_i^T \mathbf{z}| \right| \right\} \leq 0.8 \|\mathbf{h}\|. \quad (49)$$

Furthermore, recall $\tilde{\mathcal{T}}^i := \{ \|\mathbf{a}_i^T \mathbf{x} - \mathbf{a}_i^T \mathbf{z}\| \leq \alpha'_h \text{med} \{ \|\mathbf{a}_i^T \mathbf{z} - \sqrt{y_i}\| \} \}$. Then,

$$\begin{aligned} \nabla \mathcal{R}_{tr}(\mathbf{z}) &= \frac{1}{m} \sum_{i=1}^m (|\mathbf{a}_i^T \mathbf{z} - \sqrt{y_i}|) \mathbf{a}_i \mathbf{1}_{\mathcal{T}^i} \\ &= \frac{1}{m} \left(\underbrace{\sum_{i \notin S} (|\mathbf{a}_i^T \mathbf{z} - \mathbf{a}_i^T \mathbf{x}|) \mathbf{a}_i \mathbf{1}_{\mathcal{T}^i} + \sum_{i \in S} (|\mathbf{a}_i^T \mathbf{z} - \mathbf{a}_i^T \mathbf{x}|) \mathbf{a}_i \mathbf{1}_{\tilde{\mathcal{T}}^i}}_{\nabla^{clean} \mathcal{R}_{tr}(\mathbf{z})} \right) \\ &\quad - \underbrace{\frac{1}{m} \sum_{i \notin S} (\sqrt{y_i} - \mathbf{a}_i^T \mathbf{x}) \mathbf{a}_i \mathbf{1}_{\mathcal{T}^i}}_{\nabla^{noise} \mathcal{R}_{tr}(\mathbf{z})} + \underbrace{\frac{1}{m} \sum_{i \in S} ((|\mathbf{a}_i^T \mathbf{z} - \sqrt{y_i}|) \mathbf{1}_{\mathcal{T}^i} - (|\mathbf{a}_i^T \mathbf{z} - \mathbf{a}_i^T \mathbf{x}|) \mathbf{1}_{\tilde{\mathcal{T}}^i}) \mathbf{a}_i}_{\nabla^{extra} \mathcal{R}_{tr}(\mathbf{z})}. \end{aligned}$$

For $i \notin S$, the inclusion property (i.e. $\mathcal{T}_1^i \subseteq \mathcal{T}^i \subseteq \mathcal{T}_2^i$) holds because

$$|\sqrt{y_i} - \mathbf{a}_i^T \mathbf{z}| \in \|\mathbf{a}_i^T \mathbf{x} - \mathbf{a}_i^T \mathbf{z}\| \pm \sqrt{|w_i|}$$

and $\sqrt{|w_i|} \leq \frac{1}{c_3} \|\mathbf{h}\|$ for some sufficient large c_3 . For $i \in S$, the inclusion $\mathcal{T}_1^i \subseteq \tilde{\mathcal{T}}^i \subseteq \mathcal{T}_2^i$ holds because of (49). All the proof arguments for Propositions 7 and 8 are also valid for $\nabla^{clean} \mathcal{R}_{tr}(\mathbf{z})$, and thus we have

$$\begin{aligned} \langle \nabla^{clean} \mathcal{R}_{tr}(\mathbf{z}), \mathbf{h} \rangle &\geq (0.88 - \zeta'_1 - \zeta'_2 - \epsilon) \|\mathbf{h}\|^2, \\ \|\nabla^{clean} \mathcal{R}_{tr}(\mathbf{z})\| &\leq (1.8 + \delta) \|\mathbf{h}\|. \end{aligned}$$

Next, we turn to control the contribution of the noise. Let $\tilde{\mathbf{w}}_i = (\sqrt{y_i} - \mathbf{a}_i^T \mathbf{x}) \mathbf{1}_{\mathcal{T}^i}$. Then $|\tilde{w}_i| < \sqrt{|w_i|}$ and we have

$$\|\nabla^{noise} \mathcal{R}_{tr}(\mathbf{z})\| = \left\| \frac{1}{m} \mathbf{A}^T \tilde{\mathbf{w}} \right\| \leq \left\| \frac{1}{\sqrt{m}} \mathbf{A}^T \right\| \left\| \frac{\tilde{\mathbf{w}}}{\sqrt{m}} \right\| \leq (1 + \delta) \|\tilde{\mathbf{w}}\|_\infty \leq (1 + \delta) \sqrt{\|\mathbf{w}\|_\infty},$$

when m/n is sufficiently large. Given the regime condition $\|\mathbf{h}\| \geq c_3 \sqrt{\|\mathbf{w}\|_\infty}$, we further have

$$\begin{aligned} \|\nabla^{noise} \mathcal{R}_{tr}(\mathbf{z})\| &\leq \frac{(1 + \delta)}{c_3} \|\mathbf{h}\|, \\ |\langle \nabla^{noise} \mathcal{R}_{tr}(\mathbf{z}), \mathbf{h} \rangle| &\leq \|\nabla^{noise} \mathcal{R}_{tr}(\mathbf{z})\| \cdot \|\mathbf{h}\| \leq \frac{(1 + \delta)}{c_3} \|\mathbf{h}\|^2. \end{aligned}$$

We next bound the contribution of $\nabla^{extra} \mathcal{R}_{tr}(\mathbf{z})$. Introduce $\mathbf{q} = [q_1, \dots, q_m]^T$, where

$$q_i := ((|\mathbf{a}_i^T \mathbf{z} - \sqrt{y_i}|) \mathbf{1}_{\mathcal{T}^i} - (|\mathbf{a}_i^T \mathbf{z} - \mathbf{a}_i^T \mathbf{x}|) \mathbf{1}_{\tilde{\mathcal{T}}^i}) \mathbf{1}_{\{i \in S\}}.$$

Then $|q_i| \leq 1.6\alpha'_h \|\mathbf{h}\|$, and $\|\mathbf{q}\| \leq \sqrt{sm} \cdot 1.6\alpha'_h \|\mathbf{h}\|$. We thus have

$$\begin{aligned} \|\nabla^{extra} \mathcal{R}_{tr}(\mathbf{z})\| &= \frac{1}{m} \|\mathbf{A}^T \mathbf{q}\| \leq 1.6(1 + \delta) \sqrt{s} \alpha'_h \|\mathbf{h}\|, \\ |\langle \nabla^{extra} \mathcal{R}_{tr}(\mathbf{z}), \mathbf{h} \rangle| &\leq \|\mathbf{h}\| \cdot \|\nabla^{extra} \mathcal{R}_{tr}(\mathbf{z})\| \leq 1.6(1 + \delta) \sqrt{s} \alpha'_h \|\mathbf{h}\|^2. \end{aligned}$$

Putting these together, one has

$$\begin{aligned} \langle \nabla \mathcal{R}_{tr}(\mathbf{z}), \mathbf{h} \rangle &\geq \langle \nabla^{clean} \mathcal{R}_{tr}(\mathbf{z}), \mathbf{h} \rangle - |\langle \nabla^{noise} \mathcal{R}_{tr}(\mathbf{z}), \mathbf{h} \rangle| - |\langle \nabla^{extra} \mathcal{R}_{tr}(\mathbf{z}), \mathbf{h} \rangle| \\ &\geq (0.88 - \zeta'_1 - \zeta'_2 - \epsilon - (1 + \delta)(1/c_3 - 1.6\sqrt{s}\alpha'_h)) \|\mathbf{h}\|^2, \end{aligned}$$

and

$$\begin{aligned} \|\nabla \mathcal{R}_{tr}(\mathbf{z})\| &\leq \|\nabla^{clean} \mathcal{R}_{tr}(\mathbf{z})\| + \|\nabla^{noise} \mathcal{R}_{tr}(\mathbf{z})\| + \|\nabla^{extra} \mathcal{R}_{tr}(\mathbf{z})\| \\ &\leq (1.8 + \delta + (1 + \delta) \cdot (1/c_3 + 1.6\sqrt{s}\alpha'_h)) \|\mathbf{h}\|. \end{aligned} \tag{50}$$

Thus, the RC is guaranteed if μ, λ, ϵ are chosen properly, c_0, c_3 are sufficiently large and s is sufficiently small.

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

$$\text{dist}(\mathbf{z} - \mu\nabla\mathcal{R}_{tr}(\mathbf{z}), \mathbf{x}) \leq c_5\sqrt{\|\mathbf{w}\|_\infty} \quad (51)$$

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

7 Conclusions

In this paper, we propose provably effective approaches, median-TWF and median-RWF, for phase retrieval when the measurements are corrupted by sparse outliers that can take arbitrary values. Our strategy is to apply gradient descent with respect to carefully chosen loss functions, where both the initialization and the search directions are pruned guided by the sample median. We show that both algorithms allow exact recovery even with a constant proportion of arbitrary outliers for robust phase retrieval using a near-optimal number of measurements up to a logarithmic factor. Our algorithm performs well for phase retrieval problem under sparse corruptions. We anticipate that the technique developed in this paper will be useful for designing provably robust algorithms for other inference problems under sparse corruptions.

Appendix

A Proof of Properties of Median

A.1 Proof of Lemma 1

For simplicity, denote $\theta_p := \theta_p(F)$ and $\hat{\theta}_p := \theta_p(\{X_i\}_{i=1}^m)$. Since F' is continuous and positive, for an ϵ , there exists a constant δ_1 such that $\mathbb{P}(X \leq \theta_p - \epsilon) = p - \delta_1$, where $\delta_1 \in (\epsilon l, \epsilon L)$. Then one has

$$\begin{aligned} \mathbb{P}(\hat{\theta}_p < \theta_p - \epsilon) &\stackrel{(a)}{=} \mathbb{P}\left(\sum_{i=1}^m \mathbf{1}_{\{X_i \leq \theta_p - \epsilon\}} \geq pm\right) = \mathbb{P}\left(\frac{1}{m} \sum_{i=1}^m \mathbf{1}_{\{X_i \leq \theta_p - \epsilon\}} \geq (p - \delta_1) + \delta_1\right) \\ &\stackrel{(b)}{\leq} \exp(-2m\delta_1^2) \leq \exp(-2m\epsilon^2 l^2), \end{aligned}$$

where (a) is due to the definition of the quantile function in (27) and (b) is due to the fact that $\mathbf{1}_{\{X_i \leq \theta_p - \epsilon\}} \sim \text{Bernoulli}(p - \delta_1)$ i.i.d., followed by the Hoeffding inequality. Similarly, one can show for some $\delta_2 \in (\epsilon l, \epsilon L)$,

$$\mathbb{P}(\hat{\theta}_p > \theta_p + \epsilon) \leq \exp(-2m\delta_2^2) \leq \exp(-2m\epsilon^2 l^2).$$

Combining these two inequalities, one has the conclusion.

A.2 Proof of Lemma 2

It suffices to show that

$$|X_{(k)} - Y_{(k)}| \leq \max_l |X_l - Y_l|, \quad \forall k = 1, \dots, n. \quad (52)$$

Case 1: $k = n$, suppose $X_{(n)} = X_i$ and $Y_{(n)} = Y_j$, i.e., X_i is the largest among $\{X_l\}_{l=1}^n$ and Y_j is the largest among $\{Y_l\}_{l=1}^n$. Then we have either $X_j \leq X_i \leq Y_j$ or $Y_i \leq Y_j \leq X_i$. Hence,

$$|X_{(n)} - Y_{(n)}| = |X_i - Y_j| \leq \max\{|X_i - Y_i|, |X_j - Y_j|\}.$$

Case 2: $k = 1$, suppose that $X_{(1)} = X_i$ and $Y_{(1)} = Y_j$. Similarly

$$|X_{(1)} - Y_{(1)}| = |X_i - Y_j| \leq \max\{|X_i - Y_i|, |X_j - Y_j|\}.$$

Case 3: $1 < k < n$, suppose that $X_{(k)} = X_i$, $Y_{(k)} = Y_j$, and without loss of generality assume that $X_i < Y_j$ (if $X_i = Y_j$, $0 = |X_{(k)} - Y_{(k)}| \leq \max_l |X_l - Y_l|$ holds trivially). We show the conclusion by contradiction.

Assume $|X_{(k)} - Y_{(k)}| > \max_l |X_l - Y_l|$. Then one must have $Y_i < Y_j$ and $X_j > X_i$ and $i \neq j$. Moreover for any $p < k$ and $q > k$, the index of $X_{(p)}$ cannot be equal to the index of $Y_{(q)}$; otherwise the assumption is violated.

Thus, all $Y_{(q)}$ for $q > k$ must share the same index set with $X_{(p)}$ for $p > k$. However, X_j , which is larger than X_i (thus if $X_j = X_{(k')}$, then $k' > k$), shares the same index with Y_j , where $Y_j = Y_{(k)}$. This yields contradiction.

A.3 Proof of Lemma 3

Assume that sm is an integer. Since there are sm corrupted samples in total, one can select at least $\lceil (p-s)m \rceil$ clean samples from the left p portion of ordered contaminated samples $\{\theta_{1/m}(\{X_i\}), \theta_{2/m}(\{X_i\}), \dots, \theta_p(\{X_i\})\}$. Thus one has the left inequality. Furthermore, one can also select out at least $\lceil (1-p-s)m \rceil$ clean samples from the right $1-p$ portion of ordered contaminated samples $\{\theta_p(\{X_i\}), \dots, \theta_1(\{X_i\})\}$. One has the right inequality.

A.4 Proof of Lemma 4

First we introduce some general facts for the distribution of the product of two correlated standard Gaussian random variables [20]. Let $u \sim \mathcal{N}(0, 1)$, $v \sim \mathcal{N}(0, 1)$, and their correlation coefficient be $\rho \in [-1, 1]$. Then the density of uv is given by

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

where $K_0(\cdot)$ is the modified Bessel function of the second kind. Thus the density of $r = |uv|$ is

$$\psi_\rho(x) = \frac{1}{\pi\sqrt{1-\rho^2}} \left[\exp\left(\frac{\rho x}{1-\rho^2}\right) + \exp\left(-\frac{\rho x}{1-\rho^2}\right) \right] K_0\left(\frac{|x|}{1-\rho^2}\right), \quad x > 0, \quad (53)$$

for $|\rho| < 1$. If $|\rho| = 1$, r becomes a χ_1^2 random variable, with the density

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

It can be seen from (53) that the density of r only relates to the correlation coefficient $\rho \in [-1, 1]$.

Let $\theta_{1/2}(\psi_\rho)$ be the 1/2 quantile (median) of the distribution $\psi_\rho(x)$, and $\psi_\rho(\theta_{1/2})$ be the value of the function ψ_ρ at the point $\theta_{1/2}(\psi_\rho)$. Although it is difficult to derive the analytical expressions of $\theta_{1/2}(\psi_\rho)$ and $\psi_\rho(\theta_{1/2})$ due to the complicated form of ψ_ρ in (53), due to the continuity of $\psi_\rho(x)$ and $\theta_{1/2}(\psi_\rho)$, we can calculate them numerically, as illustrated in Figure 5. From the numerical calculation, one can see that both $\psi_\rho(\theta_{1/2})$ and $\theta_{1/2}(\psi_\rho)$ are bounded from below and above for all $\rho \in [0, 1]$ ($\psi_\rho(\cdot)$ is symmetric over ρ , hence it is sufficient to consider $\rho \in [0, 1]$), satisfying

$$0.348 < \theta_{1/2}(\psi_\rho) < 0.455, \quad 0.47 < \psi_\rho(\theta_{1/2}) < 0.76. \quad (54)$$

B Proof of Proposition 2

Denote $\tilde{y}_i := |\mathbf{a}_i^T \mathbf{x}|^2 + w_i$ for convenience. We first bound the concentration of $\text{med}(\{y_i\})$, also denoted by $\theta_{\frac{1}{2}}(\{y_i\})$. Lemma 3 yields

$$\theta_{\frac{1}{2}-s}(\{\tilde{y}_i\}) < \theta_{\frac{1}{2}}(\{y_i\}) < \theta_{\frac{1}{2}+s}(\{\tilde{y}_i\}). \quad (55)$$

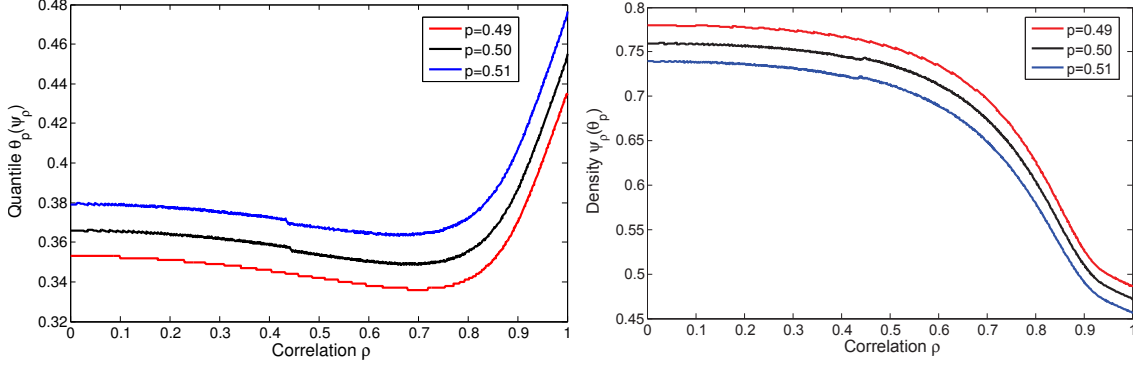


Figure 5: Quantiles and density at quantiles of $\psi_\rho(x)$ across ρ .

Moreover, Lemma 2 indicates that

$$\theta_{\frac{1}{2}-s}(\{\tilde{y}_i\}) \geq \theta_{\frac{1}{2}-s}(\{\mathbf{a}_i^T \mathbf{x}^2\}) - \|\mathbf{w}\|_\infty, \quad (56)$$

$$\theta_{\frac{1}{2}+s}(\{\tilde{y}_i\}) \leq \theta_{\frac{1}{2}+s}(\{\mathbf{a}_i^T \mathbf{x}^2\}) + \|\mathbf{w}\|_\infty. \quad (57)$$

Observe that $\mathbf{a}_i^T \mathbf{x} = \tilde{a}_{i1} \|\mathbf{x}\|$, where $\tilde{a}_{i1} = \mathbf{a}_i^T \mathbf{x} / \|\mathbf{x}\|$ is a standard Gaussian random variable. Thus $|\tilde{a}_{i1}|^2$ is a χ_1^2 random variable, whose cumulative distribution function is denoted as $K(x)$. Moreover by Lemma 1, for a small ϵ , one has $|\theta_{\frac{1}{2}-s}(\{|\tilde{a}_{i1}|^2\}) - \theta_{\frac{1}{2}-s}(K)| < \epsilon$ and $|\theta_{\frac{1}{2}+s}(\{|\tilde{a}_{i1}|^2\}) - \theta_{\frac{1}{2}+s}(K)| < \epsilon$ with probability $1 - 2\exp(-cm\epsilon^2)$ and c is a constant around 2×0.47^2 (see Figure 5). We note that $\theta_{\frac{1}{2}}(K) = 0.455$ and both $\theta_{\frac{1}{2}-s}(K)$ and $\theta_{\frac{1}{2}+s}(K)$ can be arbitrarily close to $\theta_{\frac{1}{2}}(K)$ simultaneously as long as s is small enough (independent of n). Thus, one has

$$\left(\theta_{\frac{1}{2}-s}(K) - \epsilon - c\right) \|\mathbf{x}\|^2 < \theta_{\frac{1}{2}}(\{y_i\}) < \left(\theta_{\frac{1}{2}+s}(K) + \epsilon + c\right) \|\mathbf{x}\|^2, \quad (58)$$

with probability at least $1 - \exp(-cm\epsilon^2)$. For the sake of simplicity, we introduce two new notations $\zeta_s := \theta_{\frac{1}{2}-s}(K)$ and $\zeta^s := \theta_{\frac{1}{2}+s}(K)$. Specifically for the instance of $s = 0.01$, one has $\zeta_s = 0.434$ and $\zeta^s = 0.477$. It is easy to see that $\zeta^s - \zeta_s$ can be arbitrarily small if s is small enough.

We next estimate the direction of \mathbf{x} , assuming $\|\mathbf{x}\| = 1$. On the event that (58) holds, the truncation function has the following bounds,

$$\begin{aligned} \mathbf{1}_{\{y_i \leq \alpha_y^2 \theta_{1/2}(\{y_i\})/0.455\}} &\leq \mathbf{1}_{\{y_i \leq \alpha_y^2 (\zeta^s + \epsilon)/0.455\}} \leq \mathbf{1}_{\{\mathbf{a}_i^T \mathbf{x}^2 \leq \alpha_y^2 (\zeta^s + \epsilon + c)/0.455\}} \\ \mathbf{1}_{\{y_i \leq \alpha_y^2 \theta_{1/2}(\{y_i\})/0.455\}} &\geq \mathbf{1}_{\{y_i \leq \alpha_y^2 (\zeta_s - \epsilon)/0.455\}} \geq \mathbf{1}_{\{\mathbf{a}_i^T \mathbf{x}^2 \leq \alpha_y^2 (\zeta_s - \epsilon - c)/0.455\}}. \end{aligned}$$

On the other hand, denote the support of the outliers as S , and we have

$$\mathbf{Y} = \frac{1}{m} \sum_{i \notin S} \mathbf{a}_i \mathbf{a}_i^T \tilde{y}_i \mathbf{1}_{\{\mathbf{a}_i^T \mathbf{x}^2 \leq \alpha_y^2 \theta_{1/2}(\{y_i\})/0.455\}} + \frac{1}{m} \sum_{i \in S} \mathbf{a}_i \mathbf{a}_i^T y_i \mathbf{1}_{\{y_i \leq \alpha_y^2 \theta_{1/2}(\{y_i\})/0.455\}}.$$

Consequently, one can bound \mathbf{Y} as

$$\begin{aligned} \mathbf{Y}_1 &:= \frac{1}{m} \sum_{i \notin S} \mathbf{a}_i \mathbf{a}_i^T (\mathbf{a}_i^T \mathbf{x})^2 \mathbf{1}_{\{\mathbf{a}_i^T \mathbf{x}^2 \leq \alpha_y^2 (\zeta_s - \epsilon - c)/0.455\}} - c \cdot \frac{1}{m} \sum_{i \notin S} \mathbf{a}_i \mathbf{a}_i^T \preceq \mathbf{Y} \\ &\preceq \frac{1}{m} \sum_{i \notin S} \mathbf{a}_i \mathbf{a}_i^T (\mathbf{a}_i^T \mathbf{x})^2 \mathbf{1}_{\{\mathbf{a}_i^T \mathbf{x}^2 \leq \alpha_y^2 (\zeta^s + \epsilon + c)/0.455\}} + c \cdot \frac{1}{m} \sum_{i \notin S} \mathbf{a}_i \mathbf{a}_i^T + \frac{1}{m} \sum_{i \in S} \mathbf{a}_i \mathbf{a}_i^T \alpha_y^2 (\zeta^s + \epsilon + c)/0.455 =: \mathbf{Y}_2, \end{aligned}$$

where we have

$$\mathbb{E}[\mathbf{Y}_1] = (1-s)(\beta_1 \mathbf{x} \mathbf{x}^T + \beta_2 \mathbf{I} - c \mathbf{I}), \quad \mathbb{E}[\mathbf{Y}_2] = (1-s)(\beta_3 \mathbf{x} \mathbf{x}^T + \beta_4 \mathbf{I} + c \mathbf{I}) + s \alpha_y^2 \frac{(\zeta^s + \epsilon)}{0.455} \mathbf{I}, \quad (59)$$

with

$$\begin{aligned}
\beta_1 &:= \mathbb{E} \left[\xi^4 \mathbf{1}_{\{|\xi| \leq \alpha_y \sqrt{(\zeta_s - \epsilon - c)/0.455}\}} \right] - \mathbb{E} \left[\xi^2 \mathbf{1}_{\{|\xi| \leq \alpha_y \sqrt{(\zeta_s - \epsilon - c)/0.455}\}} \right] \\
\beta_2 &:= \mathbb{E} \left[\xi^2 \mathbf{1}_{\{|\xi| \leq \alpha_y \sqrt{(\zeta_s - \epsilon - c)/0.455}\}} \right] \\
\beta_3 &:= \mathbb{E} \left[\xi^4 \mathbf{1}_{\{|\xi| \leq \alpha_y \sqrt{(\zeta^s + \epsilon + c)/0.455}\}} \right] - \mathbb{E} \left[\xi^2 \mathbf{1}_{\{|\xi| \leq \alpha_y \sqrt{(\zeta^s + \epsilon + c)/0.455}\}} \right] \\
\beta_4 &:= \mathbb{E} \left[\xi^2 \mathbf{1}_{\{|\xi| \leq \alpha_y \sqrt{(\zeta^s + \epsilon + c)/0.455}\}} \right]
\end{aligned}$$

where $\xi \sim \mathcal{N}(0, 1)$.

Applying standard results on random matrices with non-isotropic sub-Gaussian rows [48, equation (5.26)] and noticing that $\mathbf{a}_i \mathbf{a}_i^T (\mathbf{a}_i^T \mathbf{x})^2 \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{x}| \leq c\}}$ can be rewritten as $\mathbf{b}_i \mathbf{b}_i^T$ where $\mathbf{b}_i := \mathbf{a}_i (\mathbf{a}_i^T \mathbf{x}) \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{x}| \leq c\}}$ is sub-Gaussian, one can obtain

$$\|\mathbf{Y}_1 - \mathbb{E}[\mathbf{Y}_1]\| \leq \delta, \quad \|\mathbf{Y}_2 - \mathbb{E}[\mathbf{Y}_2]\| \leq \delta \quad (60)$$

with probability $1 - \exp(-\Omega(m))$, provided that m/n exceeds some large constant. Furthermore, when ϵ, c and s are sufficiently small, one further has $\|\mathbb{E}[\mathbf{Y}_1] - \mathbb{E}[\mathbf{Y}_2]\| \leq \delta$. Putting these together, one has

$$\|\mathbf{Y} - (1-s)(\beta_1 \mathbf{x} \mathbf{x}^T + \beta_2 \mathbf{I} - c \mathbf{I})\| \leq 3\delta. \quad (61)$$

Let $\tilde{\mathbf{z}}^{(0)}$ be the normalized leading eigenvector of \mathbf{Y} . Repeating the same argument as in [9, Section 7.8] and taking δ, ϵ to be sufficiently small, one has

$$\text{dist}(\tilde{\mathbf{z}}^{(0)}, \mathbf{x}) \leq \tilde{\delta}, \quad (62)$$

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

Furthermore let $\mathbf{z}^{(0)} = \sqrt{\text{med}\{y_i\}/0.455} \tilde{\mathbf{z}}^{(0)}$ to handle cases $\|\mathbf{x}\| \neq 1$. By the bound (58), one has

$$\left| \frac{\text{med}\{y_i\}}{0.455} - \|\mathbf{x}\|^2 \right| \leq \max \left\{ \left| \frac{\zeta_s - \epsilon - c}{0.455} - 1 \right|, \left| \frac{\zeta^s + \epsilon + c}{0.455} - 1 \right| \right\} \|\mathbf{x}\|^2 \leq \frac{\zeta^s - \zeta_s + 2\epsilon + 2c}{0.455} \|\mathbf{x}\|^2. \quad (63)$$

Thus

$$\text{dist}(\mathbf{z}^{(0)}, \mathbf{x}) \leq \frac{\zeta^s - \zeta_s + 2\epsilon + 2c}{0.455} \|\mathbf{x}\| + \tilde{\delta} \|\mathbf{x}\| \leq \frac{1}{11} \|\mathbf{x}\|$$

as long as s and c are small enough constants. \square

C Supporting Proofs for median-TWF

C.1 Proof of Proposition 3

We show that the sample median used in the truncation rule concentrates at the level $\|\mathbf{z} - \mathbf{x}\| \|\mathbf{z}\|$. Along the way, we also establish that the sample quantiles around the median are also concentrated at the level $\|\mathbf{z} - \mathbf{x}\| \|\mathbf{z}\|$.

We first show that for a fixed pair \mathbf{z} and \mathbf{x} , (32) holds with high probability. For simplicity of notation, we let $\mathbf{h} := \mathbf{z} - \mathbf{x}$. Let $r_i = |(\mathbf{a}_i^T \mathbf{x})^2 - (\mathbf{a}_i^T \mathbf{z})^2|$. Then r_i 's are i.i.d. copies of a random variable r , where $r = |(\mathbf{a}^T \mathbf{x})^2 - (\mathbf{a}^T \mathbf{z})^2|$ with the entries of \mathbf{a} composed of i.i.d. standard Gaussian random variables. Note that the distribution of r is fixed once given \mathbf{h} and \mathbf{z} . Let $\mathbf{x}(1)$ denote the first element of a generic vector \mathbf{x} , and \mathbf{x}_{-1} denote the remaining vector of \mathbf{x} after eliminating the first element. Let \mathbf{U}_h be an orthonormal

matrix with first row being $\mathbf{h}^T/\|\mathbf{h}\|$, $\tilde{\mathbf{a}} = \mathbf{U}_h \mathbf{a}$, and $\tilde{\mathbf{z}} = \mathbf{U}_h \mathbf{z}$. Similarly, define $\mathbf{U}_{\tilde{\mathbf{z}}_{-1}}$ and let $\tilde{\mathbf{b}} = \mathbf{U}_{\tilde{\mathbf{z}}_{-1}} \tilde{\mathbf{a}}_{-1}$. Then $\tilde{\mathbf{a}}(1)$ and $\tilde{\mathbf{b}}(1)$ are independent standard normal random variables. We further express r as follows.

$$\begin{aligned}
r &= |(\mathbf{a}^T \mathbf{z})^2 - (\mathbf{a}^T \mathbf{x})^2| \\
&= |(2\mathbf{a}^T \mathbf{z} - \mathbf{a}^T \mathbf{h})(\mathbf{a}^T \mathbf{h})| \\
&= |(2\tilde{\mathbf{a}}^T \tilde{\mathbf{z}} - \tilde{\mathbf{a}}(1)\|\mathbf{h}\|)(\tilde{\mathbf{a}}(1)\|\mathbf{h}\|)| \\
&= |(2\mathbf{h}^T \mathbf{z} - \|\mathbf{h}\|^2)\tilde{\mathbf{a}}(1)^2 + 2(\tilde{\mathbf{a}}_{-1}^T \tilde{\mathbf{z}}_{-1})(\tilde{\mathbf{a}}(1)\|\mathbf{h}\|)| \\
&= |(2\mathbf{h}^T \mathbf{z} - \|\mathbf{h}\|^2)\tilde{\mathbf{a}}(1)^2 + 2\tilde{\mathbf{b}}(1)\|\tilde{\mathbf{z}}_{-1}\|\tilde{\mathbf{a}}(1)\|\mathbf{h}\|| \\
&= |(2\mathbf{h}^T \mathbf{z} - \|\mathbf{h}\|^2)\tilde{\mathbf{a}}(1)^2 + 2\sqrt{\|\mathbf{z}\|^2 - \tilde{\mathbf{z}}(1)^2}\tilde{\mathbf{a}}(1)\tilde{\mathbf{b}}(1)\|\mathbf{h}\|| \\
&= \left| \left(2\frac{\mathbf{h}^T \mathbf{z}}{\|\mathbf{h}\|\|\mathbf{z}\|} - \frac{\|\mathbf{h}\|}{\|\mathbf{z}\|} \right) \tilde{\mathbf{a}}(1)^2 + 2\sqrt{1 - \left(\frac{\mathbf{h}^T \mathbf{z}}{\|\mathbf{h}\|\|\mathbf{z}\|} \right)^2} \tilde{\mathbf{a}}(1)\tilde{\mathbf{b}}(1) \right| \cdot \|\mathbf{h}\|\|\mathbf{z}\| \\
&=: |(2\cos(\omega) - t)\tilde{\mathbf{a}}(1)^2 + 2\sqrt{1 - \cos^2(\omega)}\tilde{\mathbf{a}}(1)\tilde{\mathbf{b}}(1)| \cdot \|\mathbf{h}\|\|\mathbf{z}\| \\
&=: |u\tilde{v}| \cdot \|\mathbf{h}\|\|\mathbf{z}\|
\end{aligned}$$

where ω is the angle between \mathbf{h} and \mathbf{z} , and $t = \|\mathbf{h}\|/\|\mathbf{z}\| < 1/11$. Consequently, $u = \tilde{\mathbf{a}}(1) \sim \mathcal{N}(0, 1)$ and $\tilde{v} = (2\cos(\omega) - t)\tilde{\mathbf{a}}(1) + 2|\sin(\omega)|\tilde{\mathbf{b}}(1)$ is also a Gaussian random variable with variance $3.6 < \text{Var}(\tilde{v}) < 4$ under the assumption $t < 1/11$.

Let $v = \tilde{v}/\sqrt{\text{Var}(\tilde{v})}$, and then $v \sim \mathcal{N}(0, 1)$. Furthermore, let $r' = |uv|$. Denote the density function of r' as $\psi_\rho(\cdot)$ and the 1/2-quantile point of r' as $\theta_{1/2}(\psi_\rho)$. By Lemma 4, we have

$$0.47 < \psi_\rho(\theta_{1/2}) < 0.76, \quad 0.348 < \theta_{1/2}(\psi_\rho) < 0.455. \quad (64)$$

By Lemma 1, we have with probability at least $1 - 2\exp(-c\epsilon^2)$ (here c is around 2×0.47^2),

$$0.348 - \epsilon < \text{med}(\{r'_i\}_{i=1}^m) < 0.455 + \epsilon. \quad (65)$$

The same arguments carry over to other quantiles $\theta_{0.49}(\{r'_i\})$ and $\theta_{0.51}(\{r'_i\})$. From Figure. 5, we observe that for $\rho \in [0, 1]$

$$0.45 < \psi_\rho(\theta_{0.49}), \psi_\rho(\theta_{0.51}) < 0.78, \quad 0.336 < \theta_{0.49}(\psi_\rho), \theta_{0.51}(\psi_\rho) < 0.477 \quad (66)$$

and then we have with probability at least $1 - 2\exp(-c\epsilon^2)$ (here c is around 2×0.45^2),

$$0.336 - \epsilon < \theta_{0.49}(\{r'_m\}), \theta_{0.51}(\{r'_m\}) < 0.477 + \epsilon. \quad (67)$$

Hence, by multiplying by $\sqrt{\text{Var}(\tilde{v})}$, we have with probability $1 - 2\exp(-c\epsilon^2)$,

$$(0.65 - \epsilon)\|\mathbf{z} - \mathbf{x}\|\|\mathbf{z}\| \leq \text{med}(\{ |(\mathbf{a}_i^T \mathbf{z})^2 - (\mathbf{a}_i^T \mathbf{x})^2| \}) \leq (0.91 + \epsilon)\|\mathbf{z} - \mathbf{x}\|\|\mathbf{z}\|, \quad (68)$$

$$(0.63 - \epsilon)\|\mathbf{z} - \mathbf{x}\|\|\mathbf{z}\| \leq \theta_{0.49}, \theta_{0.51}(\{ |(\mathbf{a}_i^T \mathbf{z})^2 - (\mathbf{a}_i^T \mathbf{x})^2| \}) \leq (0.95 + \epsilon)\|\mathbf{z} - \mathbf{x}\|\|\mathbf{z}\|. \quad (69)$$

We note that, to keep notation simple, c and ϵ may vary line by line within constant factors.

Up to now, we prove that for any fixed \mathbf{z} and \mathbf{x} , the median or neighboring quantiles of $\{ |(\mathbf{a}_i^T \mathbf{z})^2 - (\mathbf{a}_i^T \mathbf{x})^2| \}$ are upper and lower bounded by $\|\mathbf{z} - \mathbf{x}\|\|\mathbf{z}\|$ times constant factors. To prove (32) for all \mathbf{z} and \mathbf{x} with $\|\mathbf{z} - \mathbf{x}\| \leq \frac{1}{11}\|\mathbf{z}\|$, we use the net covering argument. Still we argue for median first and the same arguments carry over to other quantiles.

To proceed, we restate (68) as

$$(0.65 - \epsilon) \leq \text{med} \left(\left\{ \left| \left(\frac{2(\mathbf{a}_i^T \mathbf{z})}{\|\mathbf{z}\|} - \frac{\mathbf{a}_i^T \mathbf{h} \|\mathbf{h}\|}{\|\mathbf{h}\| \|\mathbf{z}\|} \right) \frac{\mathbf{a}_i^T \mathbf{h}}{\|\mathbf{h}\|} \right| \right\} \right) \leq (0.91 + \epsilon) \quad (70)$$

holds with probability at least $1 - 2\exp(-c\epsilon^2)$ for a given pair \mathbf{h}, \mathbf{z} satisfying $\|\mathbf{h}\|/\|\mathbf{z}\| \leq 1/11$.

Let $\tau = \epsilon/(6n + 6m)$, let \mathcal{S}_τ be a τ -net covering the unit sphere, \mathcal{L}_τ be a τ -net covering a line with length $1/11$, and set

$$\mathcal{N}_\tau = \{(z_0, \mathbf{h}_0, t_0) : (z_0, \mathbf{h}_0, t_0) \in \mathcal{S}_\tau \times \mathcal{S}_\tau \times \mathcal{L}_\tau\}. \quad (71)$$

One has cardinality bound (i.e., the upper bound on the covering number) $|\mathcal{N}_\tau| \leq (1 + 2/\tau)^{2n}/(11\tau) < (1 + 2/\tau)^{2n+1}$. Taking the union bound, we have

$$(0.65 - \epsilon) \leq \text{med}(\{|2(\mathbf{a}_i^T z_0) - (\mathbf{a}_i^T \mathbf{h}_0)t_0| |\mathbf{a}_i^T \mathbf{h}_0|\}) \leq (0.91 + \epsilon), \quad \forall (z_0, \mathbf{h}_0, t_0) \in \mathcal{N}_\epsilon \quad (72)$$

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

We next argue that (72) holds with probability $1 - c_1 \exp(-c_2 m \epsilon^2)$ for some constants c_1, c_2 as long as $m \geq c_0(\epsilon^{-2} \log \epsilon^{-1})n \log n$ for sufficiently large constant c_0 . To prove this claim, we first observe

$$(1 + 2/\tau)^{2n+1} \asymp \exp(2n(\log(n + m) + \log 12 + \log(1/\epsilon))) \asymp \exp(2n(\log m)).$$

We note that once ϵ is chosen, it is fixed in the whole proof and does not scale with m or n . For simplicity, assume that $\epsilon < 1/e$. Fix some positive constant $c' < c - c_2$. It then suffices to show that there exists a large constant c_0 such that if $m \geq c_0(\epsilon^{-2} \log \epsilon^{-1})n \log n$, then

$$2n \log m < c' m \epsilon^2. \quad (73)$$

For any fixed n , if (73) holds for some m and $m > (2/c')\epsilon^{-2}n$, then (73) always holds for larger m , because

$$\begin{aligned} 2n \log(m + 1) &= 2n \log m + 2n(\log(m + 1) - \log m) = 2n \log m + \frac{2n}{m} \log(1 + \frac{1}{m})^m \\ &\leq 2n \log m + \frac{2n}{m} \leq c' m \epsilon^2 + c' \epsilon^2 = c'(m + 1)\epsilon^2. \end{aligned}$$

Next, for any n , we can always find a constant c_0 such that (73) holds for $m = c_0(\epsilon^{-2} \log \epsilon^{-1})n \log n$. Such c_0 can be easily found for large n . For example, $c_0 = 4/c'$ is a valid option if

$$(4/c')(\epsilon^{-2} \log \epsilon^{-1})n \log n < n^2. \quad (74)$$

Moreover, since the number of n that violates (74) is finite, the maximum over all such c_0 serves the purpose.

Next, one needs to bound

$$|\text{med}(\{|2(\mathbf{a}_i^T z_0) - (\mathbf{a}_i^T \mathbf{h}_0)t_0| |\mathbf{a}_i^T \mathbf{h}_0|\}) - \text{med}(\{|2(\mathbf{a}_i^T z) - (\mathbf{a}_i^T \mathbf{h})t| |\mathbf{a}_i^T \mathbf{h}|\})|$$

for any $\|z - z_0\| < \tau$, $\|\mathbf{z} - \mathbf{z}_0\| < \tau$ and $\|t - t_0\| < \tau$.

By Lemma 2 and the inequality $||x| - |y|| \leq |x - y|$, we have

$$\begin{aligned} &|\text{med}(\{|2(\mathbf{a}_i^T z_0) - (\mathbf{a}_i^T \mathbf{h}_0)t_0| |\mathbf{a}_i^T \mathbf{h}_0|\}) - \text{med}(\{|2(\mathbf{a}_i^T z) - (\mathbf{a}_i^T \mathbf{h})t| |\mathbf{a}_i^T \mathbf{h}|\})| \\ &\leq \max_{i \in [m]} |(2(\mathbf{a}_i^T z_0) - (\mathbf{a}_i^T \mathbf{h}_0)t_0) (\mathbf{a}_i^T \mathbf{h}_0) - (2(\mathbf{a}_i^T z) - (\mathbf{a}_i^T \mathbf{h})t) (\mathbf{a}_i^T \mathbf{h})| \\ &\leq \max_{i \in [m]} |(2(\mathbf{a}_i^T z_0) - (\mathbf{a}_i^T \mathbf{h}_0)t_0) (\mathbf{a}_i^T \mathbf{h}_0) - (2(\mathbf{a}_i^T z) - (\mathbf{a}_i^T \mathbf{h})t) (\mathbf{a}_i^T \mathbf{h}_0)| \\ &\quad + \max_{i \in [m]} |(2(\mathbf{a}_i^T z) - (\mathbf{a}_i^T \mathbf{h})t) (\mathbf{a}_i^T \mathbf{h}_0) - (2(\mathbf{a}_i^T z) - (\mathbf{a}_i^T \mathbf{h})t) (\mathbf{a}_i^T \mathbf{h})| \\ &\leq \max_{i \in [m]} (|2\mathbf{a}_i^T(z_0 - z)| + |(\mathbf{a}_i^T \mathbf{h}_0)t_0 - (\mathbf{a}_i^T \mathbf{h})t|) |\mathbf{a}_i^T \mathbf{h}_0| + \max_{i \in [m]} |2(\mathbf{a}_i^T z) - (\mathbf{a}_i^T \mathbf{h})t| |\mathbf{a}_i^T(\mathbf{h}_0 - \mathbf{h})| \\ &\leq \max_{i \in [m]} \|\mathbf{a}_i\|^2(3 + t)\tau + \max_{i \in [m]} \|\mathbf{a}_i\|^2(2 + t)\tau \\ &\leq \max_{i \in [m]} \|\mathbf{a}_i\|^2(5 + 2t)\tau \end{aligned}$$

On the event $E_1 := \{\max_{i \in [m]} \|\mathbf{a}_i\|^2 \leq m + n\}$, one can show that

$$|\text{med}(\{|2(\mathbf{a}_i^T z_0) - (\mathbf{a}_i^T \mathbf{h}_0)t_0| |\mathbf{a}_i^T \mathbf{h}_0|\}) - \text{med}(\{|2(\mathbf{a}_i^T z) - (\mathbf{a}_i^T \mathbf{h})t| |\mathbf{a}_i^T \mathbf{h}|\})| < 6(m + n)\tau < \epsilon. \quad (75)$$

We claim that E_1 holds with probability at least $1 - m \exp(-m/8)$ if $m > n$. This can be argued as follows. Note that $\|\mathbf{a}_i\|^2 = \sum_{j=1}^n \mathbf{a}_i(j)^2$, where $\mathbf{a}_i(j)$ is the j -th element of \mathbf{a}_i . Hence, $\|\mathbf{a}_i\|^2$ is a sum of n i.i.d. χ_1^2 random variables. Applying the Bernstein-type inequality [48, Corollary 5.17] and observing that the sub-exponential norm of χ_1^2 is smaller than 2, we have

$$\mathbb{P} \{ \|\mathbf{a}_i\|^2 \geq m + n \} \leq \exp(-m/8). \quad (76)$$

Then a union bound concludes the claim.

Further note that (72) holds on an event E_2 , which has probability $1 - c_1 \exp(-c_2 m \epsilon^2)$ as long as $m \geq c_0 (\epsilon^{-2} \log \frac{1}{\epsilon}) n \log n$. On the intersection of E_1 and E_2 , inequality for $\theta_{\frac{1}{2}}$ (i.e., median) in (32) holds. Such net covering arguments can also carry over to show that inequalities of $\theta_{0.49}$ and $\theta_{0.51}$ in (32) also hold for all \mathbf{x} and \mathbf{z} obeying $\|\mathbf{x} - \mathbf{z}\| \leq \frac{1}{11} \|\mathbf{z}\|$.

C.2 Proof of Proposition 4

The proof adapts that of [13, Proposition 2]. We outline the main steps for completeness. Observe that for the noise-free case, $y_i = (\mathbf{a}_i^T \mathbf{x})^2$. We obtain

$$\begin{aligned} \nabla \ell_{tr}(\mathbf{z}) &= \frac{1}{m} \sum_{i=1}^m \frac{(\mathbf{a}_i^T \mathbf{z})^2 - (\mathbf{a}_i^T \mathbf{x})^2}{\mathbf{a}_i^T \mathbf{z}} \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i} \\ &= \frac{1}{m} \sum_{i=1}^m 2(\mathbf{a}_i^T \mathbf{h}) \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i} - \frac{1}{m} \sum_{i=1}^m \frac{(\mathbf{a}_i^T \mathbf{h})^2}{\mathbf{a}_i^T \mathbf{z}} \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i}. \end{aligned} \quad (77)$$

One expects the contribution of the second term in (77) to be small as $\|\mathbf{h}\|/\|\mathbf{z}\|$ decreases.

For each i , we introduce two new events

$$\begin{aligned} \mathcal{E}_3^i &:= \{ |(\mathbf{a}_i^T \mathbf{x})^2 - (\mathbf{a}_i^T \mathbf{z})^2| \leq 0.6 \alpha_h \|\mathbf{h}\| \cdot |\mathbf{a}_i^T \mathbf{z}| \}, \\ \mathcal{E}_4^i &:= \{ |(\mathbf{a}_i^T \mathbf{x})^2 - (\mathbf{a}_i^T \mathbf{z})^2| \leq 1.0 \alpha_h \|\mathbf{h}\| \cdot |\mathbf{a}_i^T \mathbf{z}| \}. \end{aligned}$$

One the event that Proposition 3 holds, the following inclusion property

$$\mathcal{E}_3^i \subseteq \mathcal{E}_2^i \subseteq \mathcal{E}_4^i \quad (78)$$

is true for all i , where \mathcal{E}_2^i is defined in Algorithm 1. It is easier to work with these new events because \mathcal{E}_3^i 's (resp. \mathcal{E}_4^i 's) are statistically independent across i for any fixed \mathbf{x} and \mathbf{z} . To further decouple the quadratic inequalities in \mathcal{E}_3^i and \mathcal{E}_4^i into linear inequalities, we introduce two more events and state their properties in the following lemma.

Lemma 5 (Lemma 3 in [13]). *For any $\gamma > 0$, define*

$$\mathcal{D}_\gamma^i := \{ |(\mathbf{a}_i^T \mathbf{x})^2 - (\mathbf{a}_i^T \mathbf{z})^2| \leq \gamma \|\mathbf{h}\| |\mathbf{a}_i^T \mathbf{z}| \}, \quad (79)$$

$$\mathcal{D}_\gamma^{i,1} := \left\{ \frac{|\mathbf{a}_i^T \mathbf{h}|}{\|\mathbf{h}\|} \leq \gamma \right\}, \quad (80)$$

$$\mathcal{D}_\gamma^{i,2} := \left\{ \left| \frac{\mathbf{a}_i^T \mathbf{h}}{\|\mathbf{h}\|} - \frac{2\mathbf{a}_i^T \mathbf{z}}{\|\mathbf{h}\|} \right| \leq \gamma \right\}. \quad (81)$$

On the event \mathcal{E}_1^i defined in Algorithm 1, the quadratic inequality specifying \mathcal{D}_γ^i implicates that $\mathbf{a}_i^T \mathbf{h}$ belongs to two intervals centered around 0 and $2\mathbf{a}_i^T \mathbf{z}$, respectively, i.e., $\mathcal{D}_\gamma^{i,1}$ and $\mathcal{D}_\gamma^{i,2}$. The following inclusion property holds

$$\left(\mathcal{D}_\gamma^{i,1} \cap \mathcal{E}_1^i \right) \cup \left(\mathcal{D}_\gamma^{i,2} \cap \mathcal{E}_1^i \right) \subseteq \mathcal{D}_\gamma^i \cap \mathcal{E}_1^i \subseteq \left(\mathcal{D}_\gamma^{i,1} \cap \mathcal{E}_1^i \right) \cup \left(\mathcal{D}_\gamma^{i,2} \cap \mathcal{E}_1^i \right). \quad (82)$$

Specifically, following the two inclusion properties (78) and (82), we have

$$\mathcal{D}_{\gamma_3}^{i,1} \cap \mathcal{E}_{1,\gamma_3}^i \subseteq \mathcal{E}_3^i \cap \mathcal{E}_1^i \subseteq \mathcal{E}_2^i \cap \mathcal{E}_1^i \subseteq \mathcal{E}_4^i \cap \mathcal{E}_1^i \subseteq (\mathcal{D}_{\gamma_4}^{i,1} \cup \mathcal{D}_{\gamma_4}^{i,2}) \cap \mathcal{E}_1^i \quad (83)$$

where the parameters γ_3, γ_4 are given by

$$\gamma_3 := 0.248\alpha_h, \quad \text{and} \quad \gamma_4 := \alpha_h.$$

Further using the identity (77), we have the following lower bound

$$\langle \nabla \ell_{tr}(\mathbf{z}), \mathbf{h} \rangle \geq \frac{2}{m} \sum_{i=1}^m (\mathbf{a}_i^T \mathbf{h})^2 \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{D}_{\gamma_3}^{i,1}} - \frac{1}{m} \sum_{i=1}^m \frac{|\mathbf{a}_i^T \mathbf{h}|^3}{|\mathbf{a}_i^T \mathbf{z}|} \mathbf{1}_{\mathcal{D}_{\gamma_4}^{i,1} \cap \mathcal{E}_1^i} - \frac{1}{m} \sum_{i=1}^m \frac{|\mathbf{a}_i^T \mathbf{h}|^3}{|\mathbf{a}_i^T \mathbf{z}|} \mathbf{1}_{\mathcal{D}_{\gamma_4}^{i,2} \cap \mathcal{E}_1^i}. \quad (84)$$

The three terms in (84) can be bounded following Lemmas 4, 5, and 6 in [13], which concludes the proof.

D Supporting Proofs for Median-RWF

D.1 Proof of Proposition 6

Observe that

$$\| |\mathbf{a}_i^T \mathbf{x}| - |\mathbf{a}_i^T \mathbf{z}| \| = \begin{cases} |\mathbf{a}_i^T \mathbf{h}|, & \text{if } (\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{z}) \geq 0; \\ |2\mathbf{a}_i^T \mathbf{x} + \mathbf{a}_i^T \mathbf{h}|, & \text{if } (\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{z}) < 0. \end{cases}$$

The following lemma states that $\{(\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{z}) < 0\}$ are rare events when $\|\mathbf{x} - \mathbf{z}\|$ is small. Hence, $\text{med}(\{ \| |\mathbf{a}_i^T \mathbf{x}| - |\mathbf{a}_i^T \mathbf{z}| \| \}_{i=1}^m)$ can be viewed as $\text{med}(\{ |\mathbf{a}_i^T \mathbf{h}| \}_{i=1}^m)$ with a small perturbation.

Lemma 6. *If $m > c_0 n$, then with probability at least $1 - c_1 \exp(-c_2 m)$,*

$$\frac{1}{m} \sum_{i=1}^m \mathbf{1}_{\{(\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{z}) < 0\}} < 0.05 \quad (85)$$

holds for all \mathbf{z}, \mathbf{x} satisfying $\|\mathbf{z} - \mathbf{x}\| < \frac{1}{11} \|\mathbf{x}\|$.

Proof. See Appendix D.4. □

By Lemma 3 and Lemma 6, we have

$$\theta_{p-0.05}(\{ |\mathbf{a}_i^T \mathbf{h}| \}) \leq \theta_p(\{ \| |\mathbf{a}_i^T \mathbf{x}| - |\mathbf{a}_i^T \mathbf{z}| \| \}) \leq \theta_{p+0.05}(\{ |\mathbf{a}_i^T \mathbf{h}| \}) \quad (86)$$

for all \mathbf{x} and \mathbf{z} satisfying $\|\mathbf{x} - \mathbf{z}\| \leq \frac{1}{11} \|\mathbf{x}\|$ with high probability.

For the model (2) with a fraction s of outliers, due to Lemma 3, we have that

$$\theta_{\frac{1}{2}-s}(\{ \| |\mathbf{a}_i^T \mathbf{x}| - |\mathbf{a}_i^T \mathbf{z}| \| \}) \leq \theta_{\frac{1}{2}}(\{ |\sqrt{y_i} - |\mathbf{a}_i^T \mathbf{z}| \| \}) \leq \theta_{\frac{1}{2}+s}(\{ \| |\mathbf{a}_i^T \mathbf{x}| - |\mathbf{a}_i^T \mathbf{z}| \| \}). \quad (87)$$

Combining with (86), we obtain that

$$\theta_{0.45-s}(\{ |\mathbf{a}_i^T \mathbf{h}| \}) \leq \theta_{\frac{1}{2}}(\{ |\sqrt{y_i} - |\mathbf{a}_i^T \mathbf{z}| \| \}) \leq \theta_{0.55+s}(\{ |\mathbf{a}_i^T \mathbf{h}| \}). \quad (88)$$

Next it suffices to show that $\theta_{0.45-s}, \theta_{0.55+s}(\{ |\mathbf{a}_i^T \mathbf{h}| \})$ are on the order of $\|\mathbf{h}\|$ for small s .

Let $\tilde{a}_i = |\mathbf{a}_i^T \mathbf{h}| / \|\mathbf{h}\|$. Then \tilde{a}_i 's are i.i.d. copies of a *folded standard Gaussian* random variable (i.e., $|\xi|$ where $\xi \sim \mathcal{N}(0, 1)$). We use $\phi(\cdot)$ to denote the density of folded standard Gaussian distribution.

For $s = 0.01$, we calculate that

$$\phi(\theta_{0.44}) = 0.67, \quad \phi(\theta_{0.45}) = 0.67, \quad \phi(\theta_{0.55}) = 0.60, \quad \phi(\theta_{0.56}) = 0.59 \quad (89)$$

$$\theta_{0.44}(\phi) = 0.58, \quad \theta_{0.45}(\phi) = 0.6, \quad \theta_{0.55}(\phi) = 0.76, \quad \theta_{0.56}(\phi) = 0.78. \quad (90)$$

By Lemma 1, the sample quantiles concentrate on population quantiles. Thus, for any fixed pair (\mathbf{x}, \mathbf{z}) ,

$$(0.6 - \epsilon) \|\mathbf{h}\| \leq \theta_{1/2}(\{ \| |\mathbf{a}_i^T \mathbf{x}| - |\mathbf{a}_i^T \mathbf{z}| \| \}_{i=1}^m) \leq (0.76 + \epsilon) \|\mathbf{h}\|, \quad (91)$$

holds with probability at least $1 - 2 \exp(-cm\epsilon^{-2})$.

Following the argument of net covering similarly to that in Appendix C.1, the proposition is proved.

D.2 Proof of Proposition 7

The proof adapts the proof of Proposition 2 in [13]. We outline the main steps for completeness. Observe that for the noise-free case, $y_i = |\mathbf{a}_i^T \mathbf{x}|$. We obtain

$$\nabla \mathcal{R}_{tr}(\mathbf{z}) = \frac{1}{m} \sum_{i=1}^m \left((\mathbf{a}_i^T \mathbf{z}) - |\mathbf{a}_i^T \mathbf{x}| \cdot \frac{\mathbf{a}_i^T \mathbf{z}}{|\mathbf{a}_i^T \mathbf{z}|} \right) \mathbf{a}_i \mathbf{1}_{\mathcal{T}^i} = \frac{1}{m} \sum_{i \notin \mathcal{B}} (\mathbf{a}_i^T \mathbf{h}) \mathbf{a}_i \mathbf{1}_{\mathcal{T}^i} + \frac{1}{m} \sum_{i \in \mathcal{B}} (\mathbf{a}_i^T \mathbf{z} + \mathbf{a}_i^T \mathbf{x}) \mathbf{a}_i \mathbf{1}_{\mathcal{T}^i}, \quad (92)$$

where $\mathcal{B} := \{i : (\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{z}) < 0\}$. If $\|\mathbf{h}\|/\|\mathbf{x}\|$ is small enough, the cardinality of \mathcal{B} is small and thus one expects the contribution of the second term in (92) to be negligible.

We note that events \mathcal{T}^i are not statistically independent. To remove such dependency, we introduce two new series of events

$$\mathcal{T}_1^i := \{|\mathbf{a}_i^T \mathbf{x}| - |\mathbf{a}_i^T \mathbf{z}| \leq 0.5\alpha'_h \|\mathbf{h}\|\}, \quad (93)$$

$$\mathcal{T}_2^i := \{|\mathbf{a}_i^T \mathbf{x}| - |\mathbf{a}_i^T \mathbf{z}| \leq 0.8\alpha'_h \|\mathbf{h}\|\}. \quad (94)$$

Due to Proposition 6, the following inclusion property

$$\mathcal{T}_1^i \subseteq \mathcal{T}^i \subseteq \mathcal{T}_2^i \quad (95)$$

holds for all i , where \mathcal{T}^i is defined in Algorithm 2. It is easier to work with these new events because $\mathcal{T}_1^{i'}$'s (resp. $\mathcal{T}_2^{i'}$'s) are statistically independent for any fixed \mathbf{x} and \mathbf{z} . Because of the inclusion property (95), we have

$$\langle \nabla \mathcal{R}_{tr}(\mathbf{z}), \mathbf{h} \rangle \geq \frac{1}{m} \sum_{i \notin \mathcal{B}} (\mathbf{a}_i^T \mathbf{h})^2 \mathbf{1}_{\mathcal{T}_1^i} - \frac{1}{m} \sum_{i \in \mathcal{B}} |\mathbf{a}_i^T \mathbf{z} + \mathbf{a}_i^T \mathbf{x}| \cdot |\mathbf{a}_i^T \mathbf{h}| \mathbf{1}_{\mathcal{T}_2^i}. \quad (96)$$

Under the condition $i \notin \mathcal{B}$, we have $\mathcal{T}_1^i = \{|\mathbf{a}_i^T \mathbf{h}| \leq 0.5\alpha'_h \|\mathbf{h}\|\}$. Under the condition $i \in \mathcal{B}$, we have $\mathcal{T}_2^i = \{|\mathbf{a}_i^T \mathbf{x} + \mathbf{a}_i^T \mathbf{z}| \leq 0.8\alpha'_h \|\mathbf{h}\|\}$. For convenience, we introduce two parameters $\gamma_1 = 0.5\alpha'_h$ and $\gamma_2 = 0.8\alpha'_h$.

We next bound the two terms in (96) respectively. For the first term, because of the inclusion $\mathcal{B} \subseteq \{i : |\mathbf{a}_i^T \mathbf{x}| < |\mathbf{a}_i^T \mathbf{h}|\}$, we have

$$\begin{aligned} \frac{1}{m} \sum_{i \notin \mathcal{B}} (\mathbf{a}_i^T \mathbf{h})^2 \mathbf{1}_{\mathcal{T}_1^i} &= \frac{1}{m} \sum_{i \notin \mathcal{B}} (\mathbf{a}_i^T \mathbf{h})^2 \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{h}| \leq \gamma_1 \|\mathbf{h}\|\}} \\ &\geq \frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^T \mathbf{h})^2 \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{h}| \leq \gamma_1 \|\mathbf{h}\|\}} \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{x}| \geq |\mathbf{a}_i^T \mathbf{h}|\}} \\ &\geq \frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^T \mathbf{h})^2 \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{h}| \leq \gamma_1 \|\mathbf{h}\|\}} \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{x}| \geq \gamma_1 \|\mathbf{h}\|\}}. \end{aligned}$$

A simpler version of Lemma 4 in [13] gives that if $m > c_0 n$, with probability at least $1 - c_1 \exp(-c_2 m \epsilon^2)$

$$\frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^T \mathbf{h})^2 \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{h}| \leq \gamma_1 \|\mathbf{h}\|\}} \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{x}| \geq \gamma_1 \|\mathbf{h}\|\}} \geq (1 - \zeta'_1 - \zeta'_2 - \epsilon) \|\mathbf{h}\|^2 \quad (97)$$

holds for all $\mathbf{h} \in \mathbb{R}^n$, where $\zeta'_1 := 1 - \min \left\{ \mathbb{E} \left[\xi^2 \mathbf{1}_{\{\xi \geq \sqrt{1.01} \gamma_1 \frac{\|\mathbf{h}\|}{\|\mathbf{x}\|}\}} \right], \mathbb{E} \left[\mathbf{1}_{\{\xi \geq \sqrt{1.01} \gamma_1 \frac{\|\mathbf{h}\|}{\|\mathbf{x}\|}\}} \right] \right\}$ and $\zeta'_2 := \mathbb{E} \left[\xi^2 \mathbf{1}_{\{|\xi| > \sqrt{0.99} \gamma_1\}} \right]$ for $\xi \sim \mathcal{N}(0, 1)$.

For the second term, we have

$$\frac{1}{m} \sum_{i \in \mathcal{B}} |\mathbf{a}_i^T \mathbf{z} + \mathbf{a}_i^T \mathbf{x}| \cdot |\mathbf{a}_i^T \mathbf{h}| \mathbf{1}_{\mathcal{T}_2^i} \leq \gamma_2 \|\mathbf{h}\| \frac{1}{m} \sum_{i \in \mathcal{B}} |\mathbf{a}_i^T \mathbf{h}| \leq \gamma_2 \|\mathbf{h}\| \frac{1}{m} \sum_{i=1}^m |\mathbf{a}_i^T \mathbf{h}| \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{x}| < |\mathbf{a}_i^T \mathbf{h}|\}}, \quad (98)$$

where the second inequality is due to the inclusion property $\mathcal{B} \subseteq \{i : |\mathbf{a}_i^T \mathbf{x}| < |\mathbf{a}_i^T \mathbf{h}|\}$.

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

$$\frac{1}{m} \sum_{i=1}^m |\mathbf{a}_i^T \mathbf{h}| \cdot \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{x}| < |\mathbf{a}_i^T \mathbf{h}|\}} \leq (0.12 + \epsilon) \|\mathbf{h}\| \quad (99)$$

holds for all non-zero vectors $\mathbf{x}, \mathbf{h} \in \mathbb{R}^n$ satisfying $\|\mathbf{h}\| \leq \frac{1}{20} \|\mathbf{x}\|$. Here, $c_0, c_1, C > 0$ are some universal constants.

Proof. See Appendix D.5. □

Thus, putting together (97), (98) and Lemma 7 concludes the proof.

D.3 Proof of Proposition 8

This proof adapts the proof of Lemma 7 in [13]. Denote $v_i := (\mathbf{a}_i^T \mathbf{z} - |\mathbf{a}_i^T \mathbf{x}| \operatorname{sgn}(\mathbf{a}_i^T \mathbf{z})) \mathbf{1}_{\mathcal{T}^i}$. Then

$$\nabla \mathcal{R}_{tr}(\mathbf{z}) = \frac{1}{m} \mathbf{A}^T \mathbf{v},$$

where \mathbf{A} is a matrix with each row being \mathbf{a}_i^T and \mathbf{v} is a m -dimensional vector with each entry being v_i . Thus, for sufficiently large m/n , we have

$$\|\nabla \mathcal{R}_{tr}(\mathbf{z})\| = \left\| \frac{1}{m} \mathbf{A}^T \mathbf{v} \right\| \leq \frac{1}{m} \|\mathbf{A}\| \cdot \|\mathbf{v}\| \leq (1 + \delta) \frac{\|\mathbf{v}\|}{\sqrt{m}}$$

where the last inequality is due to the spectral norm bound $\|\mathbf{A}\| \leq \sqrt{m}(1 + \delta)$ following from [48, Theorem 5.32].

We next bound $\|\mathbf{v}\|$. Let $\mathbf{v} = \mathbf{v}^{(1)} + \mathbf{v}^{(2)}$, where $v_i^{(1)} = \mathbf{a}_i^T \mathbf{h} \mathbf{1}_{\mathcal{T}^i \setminus \mathcal{B}^i}$ and $v_i^{(2)} = (\mathbf{a}_i^T \mathbf{x} + \mathbf{a}_i^T \mathbf{z}) \mathbf{1}_{\mathcal{T}^i \cap \mathcal{B}^i}$, where $\mathcal{B}^i := \{(\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{z}) < 0\}$. By triangle inequality, we have $\|\mathbf{v}\| \leq \|\mathbf{v}^{(1)}\| + \|\mathbf{v}^{(2)}\|$. Furthermore, given $m > c_0 n$, by [10, Lemma 3.1] with probability $1 - \exp(-cm)$, we have

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

By Lemma 6, we have with probability $1 - C \exp(-c_1 m)$

$$\frac{1}{m} \|\mathbf{v}^{(2)}\|^2 \leq (0.8 \alpha'_h \|\mathbf{h}\|)^2 \cdot \left(\frac{1}{m} \sum_{i=1}^m \mathbf{1}_{\{(\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{z}) < 0\}} \right) \leq 0.8 \|\mathbf{h}\|^2$$

holds, where the last inequality is due to Lemma 6. Hence,

$$\frac{\|\mathbf{v}\|}{\sqrt{m}} \leq \left(\sqrt{1 + \delta} + \sqrt{0.8} \right) \|\mathbf{h}\|.$$

This concludes the proof.

D.4 Proof of Lemma 6

Denote correlation $\rho := \frac{\mathbf{z}^T \mathbf{x}}{\|\mathbf{z}\| \|\mathbf{x}\|}$. Under the condition $\|\mathbf{z} - \mathbf{x}\| \leq \frac{1}{11} \|\mathbf{x}\|$, simple calculation yields $0.995 < \rho \leq 1$. It suffices to show that the result holds with high probability for all \mathbf{x} and \mathbf{z} satisfying $\rho > 0.995$. Since now the claim is invariant with the norms of \mathbf{x} and \mathbf{z} , we assume that both \mathbf{x} and \mathbf{z} have unit length without loss of generality.

We first establish the result for any fixed \mathbf{x} and \mathbf{z} and then develop a uniform bound by covering net argument in the end. We introduce a Lipschitz function to approximate the indicator function. Define

$$\chi(t) := \begin{cases} 1, & \text{if } t < 0; \\ -\frac{1}{\delta} \cdot t + 1, & \text{if } 0 \leq t \leq \delta; \\ 0, & \text{else;} \end{cases}$$

and then $\chi(t)$ is a Lipschitz function with Lipschitz constant $\frac{1}{\delta}$. In the following proof, we set $\delta = 0.001$. We further have

$$\mathbf{1}_{\{(\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{z}) < 0\}} \leq \chi((\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{z})) \leq \mathbf{1}_{\{(\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{z}) < \delta\}}. \quad (100)$$

For convenience, we denote $b_i := \mathbf{a}_i^T \mathbf{x}$ and $\tilde{b}_i := \mathbf{a}_i^T \mathbf{z}$. Then (b_i, \tilde{b}_i) takes the jointly Gaussian distribution with mean $\mu = (0, 0)^T$ and correlation ρ (b_i and \tilde{b}_i have unit variance). We next estimate the expectation of $\mathbf{1}_{\{(\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{z}) < \delta\}}$ as follows.

$$\mathbb{E}[\mathbf{1}_{\{(\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{z}) < \delta\}}] = \mathbb{P}\{(\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{z}) < \delta\} = \iint_{\tau_1, \tau_2 < \delta} f(\tau_1, \tau_2) d\tau_1 d\tau_2, \quad (101)$$

where $f(\tau_1, \tau_2)$ is the density of the jointly Gaussian random variables (b_i, \tilde{b}_i) . Note that $\mathbb{E}[\mathbf{1}_{\{(\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{z}) < \delta\}}]$ is decreasing on ρ and for the case $\rho = 0.995$ we calculate $\mathbb{E}[\mathbf{1}_{\{(\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{z}) < \delta\}}] = 0.045$ numerically. This implies that

$$\mathbb{E}[\chi((\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{z}))] \leq 0.045$$

for $\delta = 0.001$. Furthermore, $\chi((\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{z}))$ for all i are bounded and hence sub-Gaussian. By Hoeffding type inequality for sub-Gaussian tail [48], we have

$$\mathbb{P}\left[\frac{1}{m} \sum_{i=1}^m \chi((\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{z})) > (0.045 + \epsilon)\right] < \exp(-cm\epsilon^2), \quad (102)$$

for some universal constant c , as long as $\rho \geq 0.995$.

We have proved so far that the claim holds for fixed \mathbf{x} and \mathbf{z} . We next obtain a uniform bound over all \mathbf{x} and \mathbf{z} with unit length. Let \mathcal{N}'_ϵ be an ϵ -net covering the unit sphere in \mathbb{R}^n and set

$$\mathcal{N}_\epsilon = \{(\mathbf{x}_0, \mathbf{z}_0) : (\mathbf{x}_0, \mathbf{z}_0) \in \mathcal{N}'_\epsilon \times \mathcal{N}'_\epsilon\}. \quad (103)$$

One has cardinality bound (i.e., the upper bound on the covering number) $|\mathcal{N}_\epsilon| \leq (1 + 2/\epsilon)^{2n}$. Then for any pair (\mathbf{x}, \mathbf{z}) with $\|\mathbf{x}\| = \|\mathbf{z}\| = 1$, there exists a pair $(\mathbf{x}_0, \mathbf{z}_0) \in \mathcal{N}_\epsilon$ such that $\|\mathbf{x} - \mathbf{x}_0\| \leq \epsilon$ and $\|\mathbf{z} - \mathbf{z}_0\| \leq \epsilon$. Taking the union bound for all the points on the net, we claim that

$$\frac{1}{m} \sum_{i=1}^m \chi((\mathbf{a}_i^T \mathbf{x}_0)(\mathbf{a}_i^T \mathbf{z}_0)) \leq 0.045 + \epsilon, \quad \forall (\mathbf{x}_0, \mathbf{z}_0) \in \mathcal{N}_\epsilon \quad (104)$$

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

Since $\chi(t)$ is Lipschitz with constant $1/\delta$, we have

$$|\chi((\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{z})) - \chi((\mathbf{a}_i^T \mathbf{x}_0)(\mathbf{a}_i^T \mathbf{z}_0))| \leq \frac{1}{\delta} |(\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{z}) - (\mathbf{a}_i^T \mathbf{x}_0)(\mathbf{a}_i^T \mathbf{z}_0)|. \quad (105)$$

Moreover, by [13, Lemma 1],

$$\frac{1}{m} \|\mathcal{A}(\mathbf{M})\|_1 \leq c_2 \|\mathbf{M}\|_F, \quad \text{for all symmetric rank-2 matrices } \mathbf{M} \in \mathbb{R}^{n \times n}, \quad (106)$$

holds with probability at least $1 - C \exp(-c_1 m)$ as long as $m > c_0 n$ for some constants $C, c_0, c_1, c_2 > 0$. Consequently, on the event that (106) holds, we have

$$\begin{aligned} & \left| \frac{1}{m} \sum_{i=1}^m \chi((\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{z})) - \frac{1}{m} \sum_{i=1}^m \chi((\mathbf{a}_i^T \mathbf{x}_0)(\mathbf{a}_i^T \mathbf{z}_0)) \right| \\ & \leq \frac{1}{m} \sum_{i=1}^m |\chi((\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{z})) - \chi((\mathbf{a}_i^T \mathbf{x}_0)(\mathbf{a}_i^T \mathbf{z}_0))| \\ & \leq \frac{1}{\delta} \cdot \frac{1}{m} \|\mathcal{A}(\mathbf{x}\mathbf{z}^T - \mathbf{x}_0\mathbf{z}_0^T)\|_1 \quad \text{due to (105)} \\ & \leq \frac{1}{\delta} \cdot c_2 \|\mathbf{x}\mathbf{z}^T - \mathbf{x}_0\mathbf{z}_0^T\|_F \quad \text{due to (106)} \\ & \leq \frac{1}{\delta} \cdot c_2 (\|\mathbf{x} - \mathbf{x}_0\| \cdot \|\mathbf{z}\| + \|\mathbf{z} - \mathbf{z}_0\| \cdot \|\mathbf{x}_0\|) \leq 2c_3 \epsilon / \delta. \end{aligned}$$

On the intersection of events that (104) and (106) hold, we have

$$\frac{1}{m} \sum_{i=1}^m \chi((\mathbf{a}_i^T \mathbf{x})(\mathbf{a}_i^T \mathbf{z})) \leq (0.045 + \epsilon + 2c_3\epsilon/\delta), \quad (107)$$

for all \mathbf{x} and \mathbf{z} with unit length and $\rho \geq 0.995$. Since ϵ can be arbitrarily small, the proof is completed.

D.5 Proof of Lemma 7

We first observe that for any γ ,

$$\mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{x}| < |\mathbf{a}_i^T \mathbf{h}|\}} \leq \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{x}| < \gamma \|\mathbf{x}\|\}} + \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{h}| \geq \gamma \|\mathbf{x}\|\}} \leq \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{x}| < \gamma \|\mathbf{x}\|\}} + \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{h}| \geq 20\gamma \|\mathbf{h}\|\}} \quad (108)$$

where the last inequality is due to the assumption $\frac{\|\mathbf{h}\|}{\|\mathbf{x}\|} \leq \frac{1}{20}$.

To establish the lemma, we set $\gamma = 0.15$ and denote $\gamma' := 20\gamma = 3$. We next respectively show that

$$\frac{1}{m} \sum_{i=1}^m |\mathbf{a}_i^T \mathbf{h}| \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{x}| < \gamma \|\mathbf{x}\|\}} \leq (0.11 + \epsilon) \|\mathbf{h}\| \quad (109)$$

for all $\mathbf{x}, \mathbf{h} \in \mathbb{R}^n$, and

$$\frac{1}{m} \sum_{i=1}^m |\mathbf{a}_i^T \mathbf{h}| \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{h}| > \gamma' \|\mathbf{h}\|\}} \leq (0.01 + \epsilon) \|\mathbf{h}\| \quad (110)$$

for all $\mathbf{h} \in \mathbb{R}^n$.

We first prove (109). Without loss of generality, we assume that \mathbf{h} and \mathbf{x} have unit length. We introduce a Lipschitz function to approximate the indicator functions, which is defined as

$$\chi_x(t) := \begin{cases} 1, & \text{if } |t| < \gamma; \\ \frac{1}{\delta}(\gamma - |t|) + 1, & \text{if } \gamma \leq |t| \leq \gamma + \delta; \\ 0, & \text{else.} \end{cases}$$

Then $\chi_x(t)$ is a Lipschitz function with constant $\frac{1}{\delta}$. We further have

$$\mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{x}| < \gamma\}} \leq \chi_x(\mathbf{a}_i^T \mathbf{x}) \leq \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{x}| < \gamma + \delta\}}. \quad (111)$$

We first prove bounds for any fixed pair \mathbf{h}, \mathbf{x} , and then develop a uniform bound later on.

We next estimate the expectation of $|\mathbf{a}_i^T \mathbf{h}| \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{x}| < \gamma + \delta\}}$,

$$\mathbb{E}[|\mathbf{a}_i^T \mathbf{h}| \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{x}| < \gamma + \delta\}}] = \iint_{-\infty}^{\infty} |\tau_1| \mathbf{1}_{\{|\tau_2| < \gamma + \delta\}} \cdot f(\tau_1, \tau_2) d\tau_1 d\tau_2, \quad (112)$$

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

$$\mathbb{E}[|\mathbf{a}_i^T \mathbf{h}| \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{x}| < \gamma + \delta\}}] = \frac{1}{2\pi\sqrt{1-\rho^2}} \int_{-\infty}^{\infty} |\tau_1| \exp\left(-\frac{\tau_1^2}{2}\right) \cdot \int_{-(\gamma+\delta)}^{\gamma+\delta} \exp\left(-\frac{(\tau_2 - \rho\tau_1)^2}{2(1-\rho^2)}\right) d\tau_2 d\tau_1 \quad (113)$$

$$= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} |\tau_1| \exp\left(-\frac{\tau_1^2}{2}\right) \cdot \int_{\frac{-\gamma-\delta-\rho\tau_1}{\sqrt{2(1-\rho^2)}}}^{\frac{\gamma+\delta-\rho\tau_1}{\sqrt{2(1-\rho^2)}}} \exp(-\tau^2) d\tau d\tau_1 \quad \text{by changing variables}$$

$$= \frac{1}{\sqrt{8\pi}} \int_{-\infty}^{\infty} |\tau_1| \exp\left(-\frac{\tau_1^2}{2}\right) \cdot \left(\operatorname{erf}\left(\frac{\gamma+\delta-\rho\tau_1}{\sqrt{2(1-\rho^2)}}\right) - \operatorname{erf}\left(\frac{-\gamma-\delta-\rho\tau_1}{\sqrt{2(1-\rho^2)}}\right) \right) d\tau_1 \quad (114)$$

For $|\rho| < 1$, $\mathbb{E}[|\mathbf{a}_i^T \mathbf{h}| \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{x}| < \gamma + \delta\}}]$ is a continuous function of ρ . The last integral (114) can be calculated numerically. Figure 6 plots $\mathbb{E}[|\mathbf{a}_i^T \mathbf{h}| \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{x}| < \gamma + \delta\}}]$ for $\gamma = 0.15$ and $\delta = 0.01$ over $\rho \in (-1, 1)$. Furthermore, (113) indicates that $\mathbb{E}[|\mathbf{a}_i^T \mathbf{h}| \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{x}| < \gamma + \delta\}}]$ is monotonically increasing with both θ and δ . Thus, we obtain a universal bound

$$\mathbb{E}[|\mathbf{a}_i^T \mathbf{h}| \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{x}| < \gamma + \delta\}}] \leq 0.11 \|\mathbf{h}\| \quad \text{for } \gamma < 0.15 \text{ and } \delta = 0.01, \quad (115)$$

which further implies $\mathbb{E}[|\mathbf{a}_i^T \mathbf{h}| \chi_x(\mathbf{a}_i^T \mathbf{x})] \leq 0.11 \|\mathbf{h}\|$ for $\gamma < 0.15$ and $\delta = 0.01$.

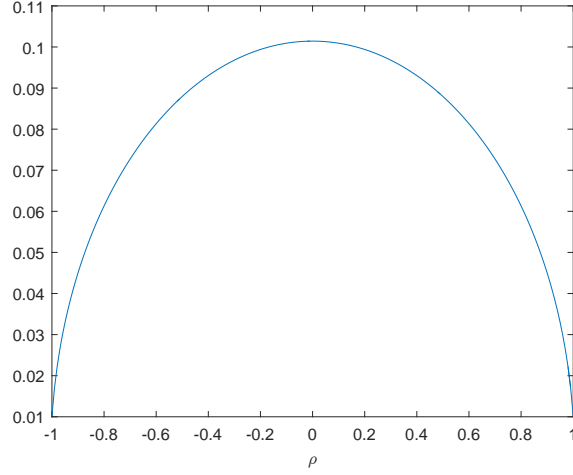


Figure 6: $\mathbb{E}[|\mathbf{a}_i^T \mathbf{h}| \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{x}| < \gamma + \delta\}}]$ with respect to ρ

Furthermore, $|\mathbf{a}_i^T \mathbf{h}| \chi_x(\mathbf{a}_i^T \mathbf{x})$'s are sub-Gaussian with sub-Gaussian norm $\mathcal{O}(\|\mathbf{h}\|)$. By the Hoeffding type of sub-Gaussian tail bound [48], we have

$$\mathcal{P} \left[\frac{1}{m} \sum_{i=1}^m |\mathbf{a}_i^T \mathbf{h}| \chi_x(\mathbf{a}_i^T \mathbf{x}) > (0.11 + \epsilon) \|\mathbf{h}\| \right] < \exp(-cm\epsilon^2), \quad (116)$$

for some universal constant c .

We have proved so far that the claim holds for a fixed pair \mathbf{h}, \mathbf{x} . We next obtain a uniform bound over all \mathbf{x} and \mathbf{h} with unit length. Let \mathcal{N}'_ϵ be a ϵ -net covering the unit sphere in \mathbb{R}^n and set

$$\mathcal{N}_\epsilon = \{(\mathbf{x}_0, \mathbf{h}_0) : (\mathbf{x}_0, \mathbf{h}_0) \in \mathcal{N}'_\epsilon \times \mathcal{N}'_\epsilon\}.$$

One has cardinality bound (i.e., the upper bound on the covering number) $|\mathcal{N}_\epsilon| \leq (1 + 2/\epsilon)^{2n}$. Then for any pair (\mathbf{x}, \mathbf{h}) with $\|\mathbf{x}\| = \|\mathbf{h}\| = 1$, there exists a pair $(\mathbf{x}_0, \mathbf{h}_0) \in \mathcal{N}_\epsilon$ such that $\|\mathbf{x} - \mathbf{x}_0\| \leq \epsilon$ and $\|\mathbf{h} - \mathbf{h}_0\| \leq \epsilon$. Taking the union bound for all the points on the net, one can show

$$\frac{1}{m} \sum_{i=1}^m |\mathbf{a}_i^T \mathbf{h}_0| \chi_x(\mathbf{a}_i^T \mathbf{x}_0) \leq 0.11 + \epsilon, \quad \forall (\mathbf{x}_0, \mathbf{h}_0) \in \mathcal{N}_\epsilon \quad (117)$$

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

Since $\chi_x(t)$ is Lipschitz with constant $1/\delta$, we have the following bound

$$|\chi_x(\mathbf{a}_i^T \mathbf{x}) - \chi_x(\mathbf{a}_i^T \mathbf{x}_0)| \leq \frac{1}{\delta} |\mathbf{a}_i^T (\mathbf{x} - \mathbf{x}_0)|. \quad (118)$$

Consequently, on the event that (106) holds, we have

$$\begin{aligned}
& \left| \frac{1}{m} \sum_{i=1}^m |\mathbf{a}_i^T \mathbf{h}|_{\chi_x}(\mathbf{a}_i^T \mathbf{x}) - \frac{1}{m} \sum_{i=1}^m |\mathbf{a}_i^T \mathbf{h}_0|_{\chi_x}(\mathbf{a}_i^T \mathbf{x}_0) \right| \\
& \leq \frac{1}{m} \sum_{i=1}^m \left| |\mathbf{a}_i^T \mathbf{h}|_{\chi_x}(\mathbf{a}_i^T \mathbf{x}) - |\mathbf{a}_i^T \mathbf{h}_0|_{\chi_x}(\mathbf{a}_i^T \mathbf{x}_0) \right| \\
& \leq \frac{1}{m} \sum_{i=1}^m |\mathbf{a}_i^T (\mathbf{h} - \mathbf{h}_0)| + \frac{1}{m} \sum_{i=1}^m \frac{1}{\delta} |\mathbf{a}_i^T \mathbf{h}_0| \cdot |\mathbf{a}_i^T \mathbf{x} - \mathbf{a}_i^T \mathbf{x}_0| \quad \text{due to (118)} \\
& \leq c'_2 \|\mathbf{h} - \mathbf{h}_0\| + \frac{1}{\delta} \cdot c_2 \|\mathbf{h}_0(\mathbf{x} - \mathbf{x}_0)^T\|_F \quad \text{due to (106)} \\
& \leq c_3 \epsilon / \delta.
\end{aligned}$$

On the intersection of events that (117) and (106) hold, we have

$$\frac{1}{m} \sum_{i=1}^m |\mathbf{a}_i^T \mathbf{h}|_{\chi_x}(\mathbf{a}_i^T \mathbf{x}_0) \leq (0.11 + \epsilon + 2c_3 \epsilon / \delta), \quad (119)$$

for all \mathbf{x} and \mathbf{h} with unit length.

We next prove (110). Without loss of generality, we assume that \mathbf{h} has unit length. We introduce a Lipschitz function to approximate the indicator functions, which is defined as

$$\chi_h(t) := \begin{cases} |t|, & \text{if } |t| > \gamma'; \\ \frac{1}{\delta}(|t| - \gamma') + \gamma', & \text{if } \gamma'(1 - \delta) \leq |t| \leq \gamma'; \\ 0, & \text{else.} \end{cases}$$

Then, $\chi_h(t)$ is a Lipschitz function with constant $\frac{1}{\delta}$. We further have

$$|\mathbf{a}_i^T \mathbf{h}| \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{h}| > \gamma' \|\mathbf{h}\|\}} \leq \chi_h(\mathbf{a}_i^T \mathbf{h}) \leq |\mathbf{a}_i^T \mathbf{h}| \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{h}| > \gamma'(1 - \delta) \|\mathbf{h}\|\}}. \quad (120)$$

We first prove bounds for any fixed \mathbf{h} , and then develop a uniform bound later on.

We next estimate the expectation of $|\mathbf{a}_i^T \mathbf{h}| \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{h}| > \gamma'(1 - \delta) \|\mathbf{h}\|\}}$ as follows:

$$\begin{aligned}
\mathbb{E}[|\mathbf{a}_i^T \mathbf{h}| \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{h}| > \gamma'(1 - \delta) \|\mathbf{h}\|\}}] &= \int_{-\infty}^{\infty} |\tau| \mathbf{1}_{\{|\tau| > \gamma'(1 - \delta)\}} \cdot f(\tau) d\tau, \\
&= 2 \cdot \frac{1}{\sqrt{2\pi}} \int_{\gamma'(1 - \delta)}^{\infty} \tau \exp\left(-\frac{\tau^2}{2}\right) d\tau \\
&= \sqrt{\frac{2}{\pi}} \exp(-\gamma'^2(1 - \delta)^2/2) < 0.01 \quad \text{for } \gamma' = 3, \delta = 0.01, \quad (121)
\end{aligned}$$

where $f(\tau)$ is the density of the standard Gaussian distribution. We note that $\mathbb{E}[|\mathbf{a}_i^T \mathbf{h}| \mathbf{1}_{\{|\mathbf{a}_i^T \mathbf{h}| > \gamma'(1 - \delta) \|\mathbf{h}\|\}}]$ is monotonically increasing with δ and decreasing with γ' . Furthermore, $\mathbb{E}[\chi_h(\mathbf{a}_i^T \mathbf{h})] \leq 0.01 \|\mathbf{h}\|$ for $\gamma' \geq 3$ and $\delta \leq 0.01$.

Moreover, $\chi_h(\mathbf{a}_i^T \mathbf{h})$ for all i are sub-Gaussian with sub-Gaussian norm $\mathcal{O}(\|\mathbf{h}\|)$. By the Hoeffding type sub-Gaussian tail bound [48], we have

$$\mathcal{P} \left[\frac{1}{m} \sum_{i=1}^m \chi_h(\mathbf{a}_i^T \mathbf{h}) > (0.01 + \epsilon) \|\mathbf{h}\| \right] < \exp(-cm\epsilon^2), \quad (122)$$

for some universal constant c .

We have proved so far that the claim holds for a fixed \mathbf{h} . We next obtain a uniform bound over all \mathbf{h} with unit length. Let \mathcal{N}'_ϵ be an ϵ -net covering the unit sphere in \mathbb{R}^n . One has cardinality bound (i.e., the

upper bound on the covering number) $|\mathcal{N}_\epsilon| \leq (1 + 2/\epsilon)^n$. Then for any \mathbf{h} with unit length, there exists a $\mathbf{h}_0 \in \mathcal{N}_\epsilon$ such that $\|\mathbf{h} - \mathbf{h}_0\| \leq \epsilon$. Taking the union bound for all the points on the net, one can show

$$\frac{1}{m} \sum_{i=1}^m \chi_h(\mathbf{a}_i^T \mathbf{h}_0) \leq 0.01 + \epsilon, \quad \forall \mathbf{h}_0 \in \mathcal{N}_\epsilon \quad (123)$$

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

Consequently, we have

$$\begin{aligned} & \left| \frac{1}{m} \sum_{i=1}^m \chi_h(\mathbf{a}_i^T \mathbf{h}) - \frac{1}{m} \sum_{i=1}^m \chi_h(\mathbf{a}_i^T \mathbf{h}_0) \right| \\ & \leq \frac{1}{m} \sum_{i=1}^m |\chi_h(\mathbf{a}_i^T \mathbf{h}) - \chi_h(\mathbf{a}_i^T \mathbf{h}_0)| \\ & \leq \frac{1}{\delta} \cdot \frac{1}{m} \sum_{i=1}^m |\mathbf{a}_i^T (\mathbf{h} - \mathbf{h}_0)| \\ & \leq \frac{1}{\delta} c'_2 \|\mathbf{h} - \mathbf{h}_0\| \leq c_3 \epsilon / \delta, \end{aligned}$$

where the second inequality is because $\chi_h(t)$ is Lipschitz continuous with constant $1/\delta$.

On the intersection of events that (123) and (106) hold, we have

$$\frac{1}{m} \sum_{i=1}^m \chi_h(\mathbf{a}_i^T \mathbf{h}) \leq (0.01 + \epsilon + c_3 \epsilon / \delta), \quad (124)$$

for all \mathbf{h} with unit length.

Putting together (119) and (124), and since ϵ can be arbitrarily small, the proof is completed.

References

- [1] A. Anandkumar, P. Jain, Y. Shi, and U. Niranjan. Tensor vs matrix methods: Robust tensor decomposition under block sparse perturbations. *arXiv preprint arXiv:1510.04747*, 2015.
- [2] S. Arora, R. Ge, T. Ma, and A. Moitra. Simple, efficient, and neural algorithms for sparse coding. *arXiv preprint arXiv:1503.00778*, 2015.
- [3] S. Bahmani and J. Romberg. Phase retrieval meets statistical learning theory: A flexible convex relaxation. *arXiv preprint arXiv:1610.04210*, 2016.
- [4] R. Balan, P. Casazza, and D. Edidin. On signal reconstruction without phase. *Applied and Computational Harmonic Analysis*, 20(3):345–356, 2006.
- [5] A. S. Bandeira, N. Boumal, and V. Voroninski. On the low-rank approach for semidefinite programs arising in synchronization and community detection. In *29th Annual Conference on Learning Theory*, 2016.
- [6] S. Bhojanapalli, B. Neyshabur, and N. Srebro. Global optimality of local search for low rank matrix recovery. *arXiv preprint arXiv:1605.07221*, 2016.
- [7] N. Boumal. Nonconvex phase synchronization. *arXiv preprint arXiv:1601.06114*, 2016.
- [8] E. J. Candès and X. Li. Solving quadratic equations via phaselift when there are about as many equations as unknowns. *Foundations of Computational Mathematics*, 14(5):1017–1026, 2014.
- [9] E. J. Candès, X. Li, and M. Soltanolkotabi. Phase retrieval via wirtinger flow: Theory and algorithms. *IEEE Transactions on Information Theory*, 61(4):1985–2007, 2015.

- [10] E. J. Candès, T. Strohmer, and V. Voroninski. Phaselift: Exact and stable signal recovery from magnitude measurements via convex programming. *Communications on Pure and Applied Mathematics*, 66(8):1241–1274, 2013.
- [11] M. Charikar, K. Chen, and M. Farach-Colton. Finding frequent items in data streams. In *Automata, languages and programming*, pages 693–703. Springer, 2002.
- [12] K. Chen. On k-median clustering in high dimensions. In *Proceedings of the seventeenth annual ACM-SIAM symposium on Discrete algorithm*, 2006.
- [13] Y. Chen and E. Candes. Solving random quadratic systems of equations is nearly as easy as solving linear systems. In *Advances in Neural Information Processing Systems (NIPS)*. 2015.
- [14] Y. Chen and E. Candes. The projected power method: An efficient algorithm for joint alignment from pairwise differences. *arXiv preprint arXiv:1609.05820*, 2016.
- [15] Y. Chen, C. Caramanis, and S. Mannor. Robust sparse regression under adversarial corruption. In *Proceedings of the 30th International Conference on Machine Learning (ICML)*, 2013.
- [16] Y. Chen, Y. Chi, and A. Goldsmith. Exact and stable covariance estimation from quadratic sampling via convex programming. *IEEE Transactions on Information Theory*, 61(7):4034–4059, July 2015.
- [17] Y. Chen and M. J. Wainwright. Fast low-rank estimation by projected gradient descent: General statistical and algorithmic guarantees. *arXiv preprint arXiv:1509.03025*, 2015.
- [18] C. De Sa, K. Olukotun, and C. Ré. Global convergence of stochastic gradient descent for some non-convex matrix problems. *arXiv preprint arXiv:1411.1134v3*, 2015.
- [19] L. Demanet and P. Hand. Stable optimizationless recovery from phaseless linear measurements. *Journal of Fourier Analysis and Applications*, 20(1):199–221, 2014.
- [20] J. D. Donahue. Products and quotients of random variables and their applications. Technical report, DTIC Document, 1964.
- [21] J. Drenth. *X-Ray Crystallography*. Wiley Online Library, 2007.
- [22] J. R. Fienup. Phase retrieval algorithms: a comparison. *Applied Optics*, 21(15):2758–2769, 1982.
- [23] R. Ge, J. D. Lee, and T. Ma. Matrix completion has no spurious local minimum. *arXiv preprint arXiv:1605.07272*, 2016.
- [24] T. Goldstein and C. Studer. Phasemax: Convex phase retrieval via basis pursuit. *arXiv preprint arXiv:1610.07531*, 2016.
- [25] P. Hand. Phaselift is robust to a constant fraction of arbitrary errors. *arXiv preprint arXiv:1502.04241*, 2015.
- [26] P. Hand and V. Voroninski. An elementary proof of convex phase retrieval in the natural parameter space via the linear program phasemax. *arXiv preprint arXiv:1611.03935*, 2016.
- [27] M. Hardt. Understanding alternating minimization for matrix completion. In *IEEE 55th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 651–660. IEEE, 2014.
- [28] P. J. Huber. *Robust statistics*. Springer, 2011.
- [29] P. Jain, P. Netrapalli, and S. Sanghavi. Low-rank matrix completion using alternating minimization. In *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*, 2013.
- [30] C. Jin, S. M. Kakade, and P. Netrapalli. Provable efficient online matrix completion via non-convex stochastic gradient descent. *arXiv preprint arXiv:1605.08370*, 2016.

- [31] R. H. Keshavan, A. Montanari, and S. Oh. Matrix completion from a few entries. *IEEE Transactions on Information Theory*, 56(6):2980–2998, June 2010.
- [32] K. Lee, Y. Li, M. Junge, and Y. Bresler. Blind recovery of sparse signals from subsampled convolution. *arXiv preprint arXiv:1511.06149*, 2015.
- [33] Q. Li and G. Tang. The nonconvex geometry of low-rank matrix optimizations with general objective functions. *arXiv preprint arXiv:1611.03060*, 2016.
- [34] X. Li, S. Ling, T. Strohmer, and K. Wei. Rapid, robust, and reliable blind deconvolution via nonconvex optimization. *arXiv preprint arXiv:1606.04933*, 2016.
- [35] X. Li and V. Voroninski. Sparse signal recovery from quadratic measurements via convex programming. *SIAM Journal on Mathematical Analysis*, 2013.
- [36] X. Li, Z. Wang, J. Lu, R. Arora, J. Haupt, H. Liu, and T. Zhao. Symmetry, saddle points, and global geometry of nonconvex matrix factorization. *arXiv preprint arXiv:1612.09296*, 2016.
- [37] Y. Li, Y. Sun, and Y. Chi. Low-rank positive semidefinite matrix recovery from corrupted rank-one measurements. *IEEE Transactions on Signal Processing*, 65(2):397–408, Jan 2017.
- [38] P. Netrapalli, P. Jain, and S. Sanghavi. Phase retrieval using alternating minimization. *Advances in Neural Information Processing Systems (NIPS)*, 2013.
- [39] P. Netrapalli, U. Niranjan, S. Sanghavi, A. Anandkumar, and P. Jain. Non-convex robust pca. In *Advances in Neural Information Processing Systems (NIPS)*, 2014.
- [40] D. Park, A. Kyrillidis, S. Bhojanapalli, C. Caramanis, and S. Sanghavi. Provable non-convex projected gradient descent for a class of constrained matrix optimization problems. *arXiv preprint arXiv:1606.01316*, 2016.
- [41] C. Qu and H. Xu. Subspace clustering with irrelevant features via robust dantzig selector. In *Advances in Neural Information Processing Systems (NIPS)*. 2015.
- [42] M. Soltanolkotabi. *Algorithms and theory for clustering and nonconvex quadratic programming*. PhD thesis, Stanford University, 2014.
- [43] J. Sun, Q. Qu, and J. Wright. Complete dictionary recovery using nonconvex optimization. In *Proceedings of the 32nd International Conference on Machine Learning (ICML)*, 2015.
- [44] J. Sun, Q. Qu, and J. Wright. A geometric analysis of phase retrieval. *arXiv preprint arXiv:1602.06664*, 2016.
- [45] R. Sun and Z.-Q. Luo. Guaranteed matrix completion via nonconvex factorization. In *IEEE 56th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 270–289, 2015.
- [46] R. J. Tibshirani. Fast computation of the median by successive binning. *arXiv preprint arXiv:0806.3301*, 2008.
- [47] S. Tu, R. Boczar, M. Soltanolkotabi, and B. Recht. Low-rank solutions of linear matrix equations via procrustes flow. *arXiv preprint arXiv:1507.03566*, 2015.
- [48] R. Vershynin. Introduction to the non-asymptotic analysis of random matrices. *Compressed Sensing, Theory and Applications*, pages 210 – 268, 2012.
- [49] D. Wagner. Resilient aggregation in sensor networks. In *Proceedings of the 2nd ACM workshop on Security of ad hoc and sensor networks*, pages 78–87. ACM, 2004.
- [50] I. Waldspurger, A. d’Aspremont, and S. Mallat. Phase recovery, maxcut and complex semidefinite programming. *Mathematical Programming*, 149(1-2):47–81, 2015.

- [51] G. Wang, G. B. Giannakis, and Y. C. Eldar. Solving systems of random quadratic equations via truncated amplitude flow. *arXiv preprint arXiv:1605.08285*, 2016.
- [52] K. Wei, J.-F. Cai, T. F. Chan, and S. Leung. Guarantees of riemannian optimization for low rank matrix recovery. *arXiv preprint arXiv:1511.01562*, 2015.
- [53] D. Weller, A. Pnueli, G. Divon, O. Radzyner, Y. Eldar, and J. Fessler. Undersampled phase retrieval with outliers. *IEEE Transactions on Computational Imaging*, 1(4):247–258, Dec 2015.
- [54] X. Yi, D. Park, Y. Chen, and C. Caramanis. Fast algorithms for robust pca via gradient descent. *arXiv preprint arXiv:1605.07784*, 2016.
- [55] H. Zhang, Y. Zhou, Y. Liang, and Y. Chi. Reshaped Wirtinger Flow and Incremental Algorithm for Solving Quadratic System of Equations. *ArXiv 1605.07719*, May 2016.
- [56] Q. Zheng and J. Lafferty. A convergent gradient descent algorithm for rank minimization and semidefinite programming from random linear measurements. In *Advances in Neural Information Processing Systems (NIPS)*, 2015.
- [57] Q. Zheng and J. Lafferty. Convergence analysis for rectangular matrix completion using burer-monteiro factorization and gradient descent. *arXiv preprint arXiv:1605.07051*, 2016.