

# Solving Random Quadratic Systems of Equations Is Nearly as Easy as Solving Linear Systems

Yuxin Chen <sup>\*</sup>      Emmanuel J. Candès <sup>\*†</sup>

May 2015; Revised January 2016

## Abstract

We consider the fundamental problem of solving quadratic systems of equations in  $n$  variables, where  $y_i = |\langle \mathbf{a}_i, \mathbf{x} \rangle|^2$ ,  $i = 1, \dots, m$  and  $\mathbf{x} \in \mathbb{R}^n$  is unknown. We propose a novel method, which starting with an initial guess computed by means of a spectral method, proceeds by minimizing a nonconvex functional as in the Wirtinger flow approach [1]. There are several key distinguishing features, most notably, a distinct objective functional and novel update rules, which operate in an adaptive fashion and drop terms bearing too much influence on the search direction. These careful selection rules provide a tighter initial guess, better descent directions, and thus enhanced practical performance. On the theoretical side, we prove that for certain unstructured models of quadratic systems, our algorithms return the correct solution in linear time, i.e. in time proportional to reading the data  $\{\mathbf{a}_i\}$  and  $\{y_i\}$  as soon as the ratio  $m/n$  between the number of equations and unknowns exceeds a fixed numerical constant. We extend the theory to deal with noisy systems in which we only have  $y_i \approx |\langle \mathbf{a}_i, \mathbf{x} \rangle|^2$  and prove that our algorithms achieve a statistical accuracy, which is nearly un-improvable. We complement our theoretical study with numerical examples showing that solving random quadratic systems is both computationally and statistically not much harder than solving linear systems of the same size—hence the title of this paper. For instance, we demonstrate empirically that the computational cost of our algorithm is about four times that of solving a least-squares problem of the same size.

## 1 Introduction

### 1.1 Problem formulation

Imagine we are given a set of  $m$  quadratic equations taking the form

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

where the data  $\mathbf{y} = [y_i]_{1 \leq i \leq m}$  and design vectors  $\mathbf{a}_i \in \mathbb{R}^n / \mathbb{C}^n$  are known whereas  $\mathbf{x} \in \mathbb{R}^n / \mathbb{C}^n$  is unknown. Having information about  $|\langle \mathbf{a}_i, \mathbf{x} \rangle|^2$ —or, equivalently,  $|\langle \mathbf{a}_i, \mathbf{x} \rangle|$ —means that we a priori know nothing about the phases or signs of the linear products  $\langle \mathbf{a}_i, \mathbf{x} \rangle$ . The problem is this: can we hope to identify a solution, if any, compatible with this nonlinear system of equations?

This problem is combinatorial in nature as one can alternatively pose it as recovering the missing signs of  $\langle \mathbf{a}_i, \mathbf{x} \rangle$  from magnitude-only observations. As is well known, many classical combinatorial problems with Boolean variables may be cast as special instances of (1). As an example, consider the NP-hard *stone problem* [2] in which we have  $n$  stones each of weight  $w_i > 0$  ( $1 \leq i \leq n$ ), which we would like to divide into two groups of equal sum weight. Letting  $x_i \in \{-1, 1\}$  indicate which of the two groups the  $i$ th stone belongs to, one can formulate this problem as solving the following quadratic system

$$\begin{cases} x_i^2 = 1, & i = 1, \dots, n, \\ (w_1 x_1 + \dots + w_n x_n)^2 = 0. \end{cases} \quad (2)$$

---

<sup>\*</sup>Department of Statistics, Stanford University, Stanford, CA 94305, U.S.A.

<sup>†</sup>Department of Mathematics, Stanford University, Stanford, CA 94305, U.S.A.

However simple this formulation may seem, even checking whether a solution to (2) exists or not is known to be NP-hard.

Moving from combinatorial optimization to the physical sciences, one application of paramount importance is the *phase retrieval* [3, 4] problem, which permeates through a wide spectrum of techniques including X-ray crystallography, diffraction imaging, microscopy, and even quantum mechanics. In a nutshell, the problem of phase retrieval arises due to the physical limitation of optical sensors, which are often only able to record the intensities of the diffracted waves scattered by an object under study. Notably, upon illuminating an object  $\mathbf{x}$ , the diffraction pattern is of the form of  $\mathbf{Ax}$ ; however, it is only possible to obtain intensity measurements  $\mathbf{y} = |\mathbf{Ax}|^2$  leading to the quadratic system (1).<sup>1</sup> In the Fraunhofer regime where data is collected in the far-field zone,  $\mathbf{A}$  is given by the spatial Fourier transform. We refer to [5] for in-depth reviews of this subject.

Continuing this motivating line of thought, in any real-world application recorded intensities are always corrupted by at least a small amount of noise so that observed data are only about  $|\langle \mathbf{a}_i, \mathbf{x} \rangle|^2$ ; i.e.

$$y_i \approx |\langle \mathbf{a}_i, \mathbf{x} \rangle|^2, \quad i = 1, \dots, m. \quad (3)$$

Although we present results for arbitrary noise distributions—even for non-stochastic noise—we shall pay particular attention to the Poisson data model, which assumes

$$y_i \stackrel{\text{ind.}}{\sim} \text{Poisson}(|\langle \mathbf{a}_i, \mathbf{x} \rangle|^2), \quad i = 1, \dots, m. \quad (4)$$

The reason why this statistical model is of special interest is that it naturally describes the variation in the number of photons detected by an optical sensor in various imaging applications.

## 1.2 Nonconvex optimization

Under a stochastic noise model with independent samples, a first impulse for solving (3) is to seek the maximum likelihood estimate (MLE), namely,

$$\text{minimize}_{\mathbf{z}} \quad - \sum_{i=1}^m \ell(\mathbf{z}; y_i), \quad (5)$$

where  $\ell(\mathbf{z}; y_i)$  denotes the log-likelihood of a candidate solution  $\mathbf{z}$  given the outcome  $y_i$ . For instance, under the Poisson data model (4) one can write

$$\ell(\mathbf{z}; y_i) = y_i \log(|\mathbf{a}_i^* \mathbf{z}|^2) - |\mathbf{a}_i^* \mathbf{z}|^2 \quad (6)$$

modulo some constant offset. Unfortunately, the log-likelihood is usually not concave, thus making the problem of computing the MLE NP-hard in general.

To alleviate this computational intractability, several convex surrogates have been proposed that work particularly well when the design vectors  $\{\mathbf{a}_i\}$  are chosen at random [6–20]. The basic idea is to introduce a rank-one matrix  $\mathbf{X} = \mathbf{x}\mathbf{x}^*$  to linearize the quadratic constraints, and then relax the rank-one constraint. Suppose we have Poisson data, then this strategy converts the problem into a convenient convex program:

$$\begin{aligned} & \text{minimize}_{\mathbf{X}} && \sum_{i=1}^m (\mu_i - y_i \log \mu_i) + \lambda \text{Tr}(\mathbf{X}) \\ & \text{subject to} && \mu_i = \mathbf{a}_i^T \mathbf{X} \mathbf{a}_i, \quad 1 \leq i \leq m, \\ & && \mathbf{X} \succeq \mathbf{0}. \end{aligned}$$

Note that the log-likelihood function is augmented by the trace functional  $\text{Tr}(\cdot)$  whose role is to promote low-rank solutions. While such convex relaxation schemes enjoy intriguing performance guarantees in many aspects (e.g. they achieve minimal sample complexity and near-optimal error bounds for certain noise models), the computational cost typically far exceeds the order of  $n^3$ . This limits applicability to large-dimensional data.

This paper follows another route: rather than lifting the problem into higher dimensions by introducing matrix variables, this paradigm maintains its iterates within the vector domain and optimize the nonconvex

<sup>1</sup>Here and below, for  $\mathbf{z} \in \mathbb{C}^n$ ,  $|\mathbf{z}|$  (resp.  $|\mathbf{z}|^2$ ) represents the vector of magnitudes  $(|z_1|, \dots, |z_n|)^\top$  (resp. squared magnitudes  $(|z_1|^2, \dots, |z_n|^2)^\top$ ).

objective directly (e.g. [1, 3, 21–28]). One promising approach along this line is the recently proposed two-stage algorithm called *Wirtinger Flow* (WF) [1]. Simply put, WF starts by computing a suitable initial guess  $\mathbf{z}^{(0)}$  using a spectral method, and then successively refines the estimate via an update rule that bears a strong resemblance to a gradient descent scheme, namely,

$$\mathbf{z}^{(t+1)} = \mathbf{z}^{(t)} + \frac{\mu_t}{m} \sum_{i=1}^m \nabla \ell(\mathbf{z}^{(t)}; y_i),$$

where  $\mathbf{z}^{(t)}$  denotes the  $t$ th iterate of the algorithm, and  $\mu_t$  is the step size (or learning rate). Here,  $\nabla \ell(\mathbf{z}; y_i)$  stands for the Wirtinger derivative w.r.t.  $\mathbf{z}$ , which in the real-valued case reduces to the ordinary gradient. The main results of [1] demonstrate that WF is surprisingly accurate for independent Gaussian design. Specifically, when  $\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  or  $\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) + j\mathcal{N}(\mathbf{0}, \mathbf{I})$ :

1. WF achieves exact recovery from  $m = O(n \log n)$  quadratic equations when there is no noise;<sup>2</sup>
2. WF attains  $\epsilon$ -accuracy—in a relative sense—within  $O(mn^2 \log(1/\epsilon))$  time (or flops);
3. In the presence of Gaussian noise, WF is stable and converges to the MLE as shown in [29].

While these results formalize the advantages of WF, the computational complexity of WF is still much larger than the best one can hope for. Moreover, the statistical guarantee in terms of the sample complexity is weaker than that achievable by convex relaxations.<sup>3</sup>

### 1.3 This paper: Truncated Wirtinger Flow

This paper develops an efficient linear-time algorithm, which also enjoys near-optimal statistical guarantees. Following the spirit of WF, we propose a novel procedure called *Truncated Wirtinger Flow* (TWF) adopting a more adaptive gradient flow. Informally, TWF proceeds in two stages:

1. **Initialization:** compute an initial guess  $\mathbf{z}^{(0)}$  by means of a spectral method applied to a subset  $\mathcal{T}_0$  of the observations  $\{y_i\}$ ;
2. **Loop:** for  $0 \leq t < T$ ,

$$\mathbf{z}^{(t+1)} = \mathbf{z}^{(t)} + \frac{\mu_t}{m} \sum_{i \in \mathcal{T}_{t+1}} \nabla \ell(\mathbf{z}^{(t)}; y_i) \tag{7}$$

for some index subset  $\mathcal{T}_{t+1} \subseteq \{1, \dots, m\}$  determined by  $\mathbf{z}^{(t)}$ .

Three remarks are in order.

- Firstly, we regularize both the initialization and the gradient flow in a data-dependent fashion by operating only upon some iteration-varying index subsets  $\mathcal{T}_t$ . This is a distinguishing feature of TWF in comparison to WF and other gradient descent variants. In words,  $\mathcal{T}_t$  corresponds to those data  $\{y_i\}$  whose resulting spectral or gradient components are in some sense not excessively large; see Sections 2 and 3 for details. As we shall see later, the main point is that this careful data trimming procedure gives us a tighter initial guess and more stable search directions.
- Secondly, we recommend that the step size  $\mu_t$  is either taken as some appropriate constant or determined by a backtracking line search. For instance, under appropriate conditions, we can take  $\mu_t = 0.2$  for all  $t$ .

---

<sup>2</sup> The standard notation  $f(n) = O(g(n))$  or  $f(n) \lesssim g(n)$  (resp.  $f(n) = \Omega(g(n))$  or  $f(n) \gtrsim g(n)$ ) means that there exists a constant  $c > 0$  such that  $|f(n)| \leq c|g(n)|$  (resp.  $|f(n)| \geq c|g(n)|$ ).  $f(n) \asymp g(n)$  means that there exist constants  $c_1, c_2 > 0$  such that  $c_1|g(n)| \leq |f(n)| \leq c_2|g(n)|$ .

<sup>3</sup> M. Soltanolkotabi recently informed us that the sample complexity of WF may be improved if one employs a better initialization procedure.

- Finally, the most expensive part of the gradient stage consists in computing  $\nabla\ell(\mathbf{z}; y_i)$ ,  $1 \leq i \leq m$ , which can often be performed in an efficient manner. More concretely, under the *real-valued* Poisson data model (4) one has

$$\nabla\ell(\mathbf{z}; y_i) = 2 \left\{ \frac{y_i}{|\mathbf{a}_i^\top \mathbf{z}|^2} \mathbf{a}_i \mathbf{a}_i^\top \mathbf{z} - \mathbf{a}_i \mathbf{a}_i^\top \mathbf{z} \right\} = 2 \left( \frac{y_i - |\mathbf{a}_i^\top \mathbf{z}|^2}{\mathbf{a}_i^\top \mathbf{z}} \right) \mathbf{a}_i,$$

Thus, calculating  $\{\nabla\ell(\mathbf{z}; y_i)\}$  essentially amounts to two matrix-vector products. Letting  $\mathbf{A} := [\mathbf{a}_1, \dots, \mathbf{a}_m]^\top$  as before, we have

$$\sum_{i \in \mathcal{T}_{t+1}} \nabla\ell(\mathbf{z}^{(t)}; y_i) = \mathbf{A}^\top \mathbf{v}, \quad v_i = \begin{cases} 2 \frac{y_i - |\mathbf{a}_i^\top \mathbf{z}|^2}{\mathbf{a}_i^\top \mathbf{z}}, & i \in \mathcal{T}_{t+1}, \\ 0, & \text{otherwise.} \end{cases}$$

Hence,  $\mathbf{A}\mathbf{z}$  gives  $\mathbf{v}$  and  $\mathbf{A}^\top \mathbf{v}$  the desired regularized gradient.

A detailed specification of the algorithm is deferred to Section 2.

## 1.4 Numerical surprises

To give the readers a sense of the practical power of TWF, we present here three illustrative numerical examples. Since it is impossible to recover the global sign—i.e. we cannot distinguish  $\mathbf{x}$  from  $-\mathbf{x}$ —we will evaluate our solutions to the quadratic equations through the distance measure put forth in [1] representing the Euclidean distance modulo a global sign: for complex-valued signals,

$$\text{dist}(\mathbf{z}, \mathbf{x}) := \min_{\varphi \in [0, 2\pi)} \|e^{-j\varphi} \mathbf{z} - \mathbf{x}\|, \quad (8)$$

while it is simply  $\min \|\mathbf{z} \pm \mathbf{x}\|$  in the real-valued case. We shall use  $\text{dist}(\hat{\mathbf{x}}, \mathbf{x})/\|\mathbf{x}\|$  throughout to denote the relative error of an estimate  $\hat{\mathbf{x}}$ . In the sequel, TWF proceeds by attempting to maximize the Poisson log-likelihood (6). Standalone Matlab implementations of TWF are available at <http://statweb.stanford.edu/~candes/publications.html> (see [30] for straight WF implementations).

The first numerical example concerns the following two problems under noiseless real-valued data:

$$\begin{aligned} \text{(a)} \quad & \text{find } \mathbf{x} \in \mathbb{R}^n && \text{s.t. } b_i = \mathbf{a}_i^\top \mathbf{x}, && 1 \leq i \leq m; \\ \text{(b)} \quad & \text{find } \mathbf{x} \in \mathbb{R}^n && \text{s.t. } b_i = |\mathbf{a}_i^\top \mathbf{x}|, && 1 \leq i \leq m. \end{aligned}$$

Apparently, (a) involves solving a linear system of equations (or a linear least squares problem), while (b) is tantamount to solving a quadratic system. Arguably the most popular method for solving large-scale least squares problems is the conjugate gradient (CG) method [31] applied to the normal equations. We are going to compare the computational efficiency between CG (for solving least squares) and TWF with a step size  $\mu_t \equiv 0.2$  (for solving a quadratic system). Set  $m = 8n$  and generate  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and  $\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ,  $1 \leq i \leq m$ , independently. This gives a matrix  $\mathbf{A}^\top \mathbf{A}$  with a low condition number equal to about  $(1 + \sqrt{1/8})^2 / (1 - \sqrt{1/8})^2 \approx 4.38$  by the Marchenko-Pastur law. Therefore, this is an ideal setting for CG as it converges extremely rapidly [32, Theorem 38.5]. Fig. 1 shows the relative estimation error of each method as a function of the iteration count, where TWF is seeded through 10 power iterations. For ease of comparison, we illustrate the iteration counts in different scales so that 4 TWF iterations are equivalent to 1 CG iteration.

Recognizing that each iteration of CG and TWF involves two matrix vector products  $\mathbf{A}\mathbf{z}$  and  $\mathbf{A}^\top \mathbf{v}$ , for such a design we reach a suprising observation:

*Even when all phase information is missing, TWF is capable of solving a quadratic system of equations only about 4 times slower than solving a least squares problem of the same size!*

To illustrate the applicability of TWF on real images, we turn to testing our algorithm on a digital photograph of Stanford main quad containing  $320 \times 1280$  pixels. We consider a type of measurements that falls under the category of coded diffraction patterns (CDP) [33] and set

$$\mathbf{y}^{(l)} = |\mathbf{F}\mathbf{D}^{(l)}\mathbf{x}|^2, \quad 1 \leq l \leq L. \quad (9)$$

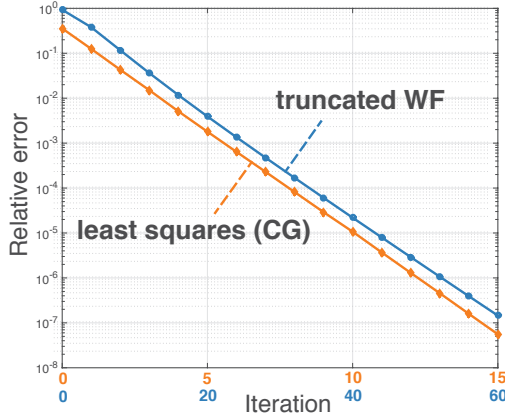


Figure 1: Relative errors of CG and TWF vs. iteration count. Here,  $n = 1000$ ,  $m = 8n$ , and TWF is seeded using just 10 power iterations.

Here,  $\mathbf{F}$  stands for a discrete Fourier transform (DFT) matrix, and  $\mathbf{D}^{(l)}$  is a diagonal matrix whose diagonal entries are independently and uniformly drawn from  $\{1, -1, j, -j\}$  (phase delays). In phase retrieval, each  $\mathbf{D}^{(l)}$  represents a random mask placed after the object so as to modulate the illumination patterns. When  $L$  masks are employed, the total number of quadratic measurements is  $m = nL$ . In this example,  $L = 12$  random coded patterns are generated to measure each color band (i.e. red, green, or blue) separately. The experiment is carried out on a MacBook Pro equipped with a 3 GHz Intel Core i7 and 16GB of memory. We run 50 iterations of the truncated power method for initialization, and 50 regularized gradient iterations, which in total costs 43.9 seconds or 2400 FFTs for each color band. The relative error after regularized spectral initialization and after 50 TWF iterations are 0.4773 and  $2.16 \times 10^{-5}$ , respectively, with the recovered images displayed in Fig. 2. In comparison, the spectral initialization using 50 untruncated power iterations returns an image of relative error 1.409, which is almost like a random guess and extremely far from the truth.

While the above experiments concern noiseless data, the numerical surprise extends to the noisy realm. Suppose the data are drawn according to the Poisson noise model (4), with  $\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  independently generated. Fig. 3 displays the empirical relative mean-square error (MSE) of TWF as a function of the signal-to-noise ratio (SNR), where the relative MSE for an estimate  $\hat{\mathbf{x}}$  and the SNR are defined as<sup>4</sup>

$$\text{MSE} := \frac{\text{dist}^2(\hat{\mathbf{x}}, \mathbf{x})}{\|\mathbf{x}\|^2}, \quad \text{and} \quad \text{SNR} := 3\|\mathbf{x}\|^2. \quad (10)$$

Both SNR and MSE are displayed on a dB scale (i.e. the values of  $10 \log_{10}(\text{SNR})$  and  $10 \log_{10}(\text{rel. MSE})$  are plotted). To evaluate the accuracy of the TWF solutions, we consider the performance achieved by MLE applied to an *ideal* problem in which the true phases are revealed. In this ideal scenario, in addition to the data  $\{y_i\}$  we are further given exact phase information  $\{\varphi_i = \text{sign}(\mathbf{a}_i^\top \mathbf{x})\}$ . Such precious information gives away the phase retrieval problem and makes the MLE efficiently solvable since the MLE problem with side information

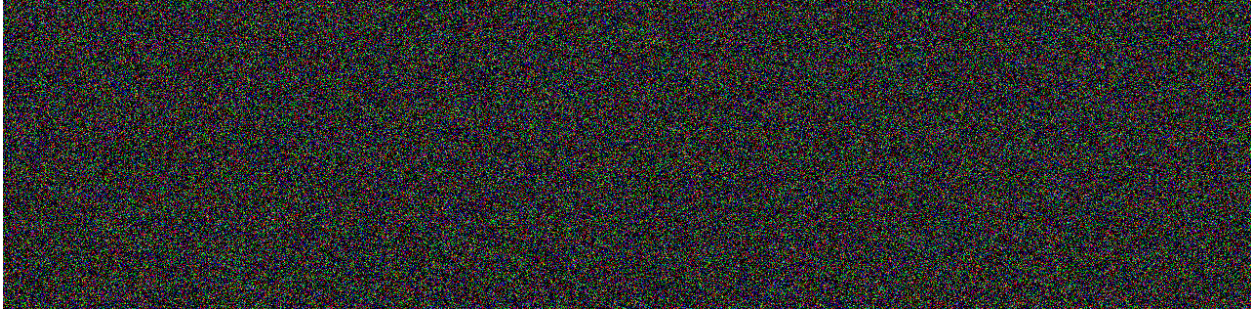
$$\begin{aligned} & \underset{\mathbf{z} \in \mathbb{R}^n}{\text{minimize}} && - \sum_{i=1}^m y_i \log(|\mathbf{a}_i^\top \mathbf{z}|^2) + (\mathbf{a}_i^\top \mathbf{z})^2 \\ & \text{subject to} && \varphi_i = \text{sign}(\mathbf{a}_i^\top \mathbf{z}) \end{aligned}$$

can be cast as a convex program

$$\underset{\mathbf{z} \in \mathbb{R}^n}{\text{minimize}} \quad - \sum_{i=1}^m 2y_i \log(\varphi_i \mathbf{a}_i^\top \mathbf{z}) + (\mathbf{a}_i^\top \mathbf{z})^2.$$

Fig. 3 illustrates the empirical performance for this ideal problem. The plots demonstrate that even when all phases are erased, TWF yields a solution of nearly the best possible quality, since it only incurs an extra 1.5

<sup>4</sup>To justify the definition of SNR, note that the signals and noise are captured by  $\mu_i = (\mathbf{a}_i^\top \mathbf{x})^2$  and  $y_i - \mu_i$ ,  $1 \leq i \leq m$ , respectively. The ratio of the signal power to the noise power is therefore  $\frac{\sum_{i=1}^m \mu_i^2}{\sum_{i=1}^m \text{Var}[y_i]} = \frac{\sum_{i=1}^m |\mathbf{a}_i^\top \mathbf{x}|^4}{\sum_{i=1}^m |\mathbf{a}_i^\top \mathbf{x}|^2} \approx \frac{3m\|\mathbf{x}\|^4}{m\|\mathbf{x}\|^2} = 3\|\mathbf{x}\|^2$ .



(a)



(b)



(c)

Figure 2: The recovered images after (a) spectral initialization; (b) regularized spectral initialization; and (c) 50 TWF gradient iterations following the regularized initialization.

dB loss compared to ideal MLE computed with all true phases revealed. This phenomenon arises regardless of the SNR!

## 1.5 Main results

The preceding numerical discoveries unveil promising features of TWF in three aspects: (1) exponentially fast convergence; (2) exact recovery from noiseless data with sample complexity  $O(n)$ ; (3) nearly minimal mean-square loss in the presence of noise. This paper offers a formal substantiation of all these findings. To this end, we assume a tractable model in which the design vectors  $\mathbf{a}_i$ 's are independent Gaussian:

$$\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n). \quad (11)$$

For concreteness, our results are concerned with TWF designed based on the Poisson log-likelihood function

$$\ell_i(\mathbf{z}) := \ell(\mathbf{z}; y_i) = y_i \log(|\mathbf{a}_i^\top \mathbf{z}|^2) - |\mathbf{a}_i^\top \mathbf{z}|^2, \quad (12)$$

where we shall use  $\ell_i(\mathbf{z})$  as a shorthand for  $\ell(\mathbf{z}; y_i)$  from now on. We begin with the performance guarantees of TWF in the absence of noise.

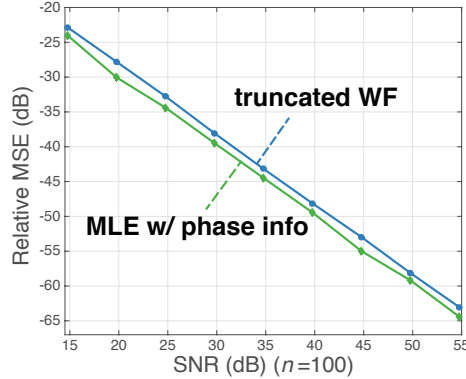


Figure 3: Relative MSE vs. SNR in dB. The curves are shown for two settings: TWF for solving quadratic equations (blue), and MLE had we observed additional phase information (green). The results are shown for  $n = 100$ , and each point is averaged over 50 Monte Carlo trials.

**Theorem 1 (Exact recovery).** *Consider the noiseless case (1) with an arbitrary signal  $\mathbf{x} \in \mathbb{R}^n$ . Suppose that the step size  $\mu_t$  is either taken to be a positive constant  $\mu_t \equiv \mu$  or chosen via a backtracking line search. Then there exist some universal constants  $0 < \rho, \nu < 1$  and  $\mu_0, c_0, c_1, c_2 > 0$  such that with probability exceeding  $1 - c_1 \exp(-c_2 m)$ , the truncated Wirtinger Flow estimates (Algorithm 1 with parameters specified in Table 1) obey*

$$\text{dist}(\mathbf{z}^{(t)}, \mathbf{x}) \leq \nu(1 - \rho)^t \|\mathbf{x}\|, \quad \forall t \in \mathbb{N}, \quad (13)$$

provided that

$$m \geq c_0 n \quad \text{and} \quad 0 < \mu \leq \mu_0.$$

As explained below, we can often take  $\mu_0 \approx 0.3$ .

**Remark 1.** As will be made precise in Section 5 (and in particular Proposition 1), one can take

$$\mu_0 = \frac{0.994 - \zeta_1 - \zeta_2 - \sqrt{2/(9\pi)}\alpha_h^{-1}}{2(1.02 + 0.665/\alpha_h)}$$

for some small quantities  $\zeta_1, \zeta_2$  and some predetermined threshold  $\alpha_h$  that is usually taken to be  $\alpha_h \geq 5$ . Under appropriate conditions, one can treat  $\mu_0$  as  $\mu_0 \approx 0.3$ .

Theorem 1 justifies at least two appealing features of TWF: (i) *minimal sample complexity* and (ii) *linear-time computational cost*. Specifically, TWF allows exact recovery from  $O(n)$  quadratic equations, which is optimal since one needs at least  $n$  measurements to have a well-posed problem. Also, because of the geometric convergence rate, TWF achieves  $\epsilon$ -accuracy (i.e.  $\text{dist}(\mathbf{z}^{(t)}, \mathbf{x}) \leq \epsilon \|\mathbf{x}\|$ ) within at most  $O(\log(1/\epsilon))$  iterations. The total computational cost is therefore  $O(mn \log(1/\epsilon))$ , which is linear in the problem size. These outperform the performance guarantees of WF [1], which runs in  $O(mn^2 \log(1/\epsilon))$  time and requires  $O(n \log n)$  sample complexity.

We emphasize that enhanced performance vis-à-vis WF is not the result of a sharper analysis, but rather, the result of key algorithmic changes. In both the initialization and iterative refinement stages, TWF proceeds in a more prudent manner by means of proper regularization, which effectively trims away those components that are too influential on either the initial guess or search directions, thus reducing the volatility of each movement. With a tighter initialization and better-controlled search directions in place, we take the step size in a far more liberal fashion—which is some constant bounded away from 0—compared to a step size which is  $O(1/n)$  as explained in [1]. In fact, what enables the movement to be more aggressive is exactly the cautious choice of  $\mathcal{T}_t$ , which precludes adverse effects from high-leverage samples.

To be broadly applicable, the proposed algorithm must guarantee reasonably faithful estimates in the presence of noise. Suppose that

$$y_i = |\langle \mathbf{a}_i, \mathbf{x} \rangle|^2 + \eta_i, \quad 1 \leq i \leq m, \quad (14)$$

where  $\eta_i$  represents an error term. We claim that TWF is stable against additive noise, as demonstrated in the theorem below.

**Theorem 2 (Stability).** *Consider the noisy case (14). Suppose that the step size  $\mu_t$  is either taken to be a positive constant  $\mu_t \equiv \mu$  or chosen via a backtracking line search. If*

$$m \geq c_0 n, \quad \mu \leq \mu_0, \quad \text{and} \quad \|\boldsymbol{\eta}\|_\infty \leq c_1 \|\mathbf{x}\|^2, \quad (15)$$

*then with probability at least  $1 - c_2 \exp(-c_3 m)$ , the truncated Wirtinger Flow estimates (Algorithm 1 with parameters specified in Table 1) satisfy*

$$\text{dist}(\mathbf{z}^{(t)}, \mathbf{x}) \lesssim \frac{\|\boldsymbol{\eta}\|}{\sqrt{m}\|\mathbf{x}\|} + (1 - \rho)^t \|\mathbf{x}\|, \quad \forall t \in \mathbb{N} \quad (16)$$

*simultaneously for all  $\mathbf{x} \in \mathbb{R}^n$ . Here,  $0 < \rho < 1$  and  $\mu_0, c_0, c_1, c_2, c_3 > 0$  are some universal constants.*

*Under the Poisson noise model (4), one has*

$$\text{dist}(\mathbf{z}^{(t)}, \mathbf{x}) \lesssim 1 + (1 - \rho)^t \|\mathbf{x}\|, \quad \forall t \in \mathbb{N} \quad (17)$$

*with probability approaching one, provided that  $\|\mathbf{x}\| \geq \log^{1.5} m$ .*

**Remark 2.** In the main text, we will prove Theorem 2 only for the case where  $\mathbf{x}$  is fixed and independent of the design vectors  $\{\mathbf{a}_i\}$ . Interested readers are referred to the supplemental materials [34] for the proof of the universal theory (i.e. the case simultaneously accommodating all  $\mathbf{x} \in \mathbb{R}^n$ ). Note that when there is no noise ( $\boldsymbol{\eta} = \mathbf{0}$ ), this stronger result guarantees the universality of the noiseless recovery.

**Remark 3.** [29] establishes stability estimates using the WF approach under Gaussian noise. There, the sample and computational complexities are still on the order of  $n \log n$  and  $mn^2$  respectively whereas the computational complexity in Theorem 2 is linear, i.e. on the order of  $mn$ .

Theorem 2 essentially reveals that the estimation error of TWF rapidly shrinks to  $O\left(\frac{\|\boldsymbol{\eta}\|/\sqrt{m}}{\|\mathbf{x}\|}\right)$  within logarithmic iterations. Put another way, since the SNR for the model (14) is captured by

$$\text{SNR} := \frac{\sum_{i=1}^m |\langle \mathbf{a}_i, \mathbf{x} \rangle|^4}{\|\boldsymbol{\eta}\|^2} \approx \frac{3m\|\mathbf{x}\|^4}{\|\boldsymbol{\eta}\|^2}, \quad (18)$$

we immediately arrive at an alternative form of the performance guarantee:

$$\text{dist}(\mathbf{z}^{(t)}, \mathbf{x}) \lesssim \frac{1}{\sqrt{\text{SNR}}} \|\mathbf{x}\| + (1 - \rho)^t \|\mathbf{x}\|, \quad \forall t \in \mathbb{N}, \quad (19)$$

revealing the stability of TWF as a function of SNR. We emphasize that this estimate holds for any error term  $\boldsymbol{\eta}$ —i.e. any noise structure, even deterministic. This being said, specializing this estimate to the Poisson noise model (4) with  $\|\mathbf{x}\| \gtrsim \log^{1.5} m$  gives an estimation error that will eventually approach a numerical constant, independent of  $n$  and  $m$ .

Encouragingly, this is already the best statistical guarantee any algorithm can achieve. We formalize this claim by deriving a fundamental lower bound on the minimax estimation error.

**Theorem 3 (Lower bound on the minimax risk).** *Suppose that  $\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ,  $m = \kappa n$  for some fixed  $\kappa$  independent of  $n$ , and  $n$  is sufficiently large. For any  $K \geq \log^{1.5} m$ , define<sup>5</sup>*

$$\Upsilon(K) := \{\mathbf{x} \in \mathbb{R}^n \mid \|\mathbf{x}\| \in (1 \pm 0.1)K\}.$$

*With probability approaching one, the minimax risk under the Poisson model (4) obeys*

$$\inf_{\hat{\mathbf{x}}} \sup_{\mathbf{x} \in \Upsilon(K)} \mathbb{E}[\text{dist}(\hat{\mathbf{x}}, \mathbf{x}) \mid \{\mathbf{a}_i\}_{1 \leq i \leq m}] \geq \frac{\varepsilon_1}{\sqrt{\kappa}}, \quad (20)$$

*where the infimum is over all estimator  $\hat{\mathbf{x}}$ . Here,  $\varepsilon_1 > 0$  is a numerical constant independent of  $n$  and  $m$ .*

<sup>5</sup>Here, 0.1 can be replaced by any positive constant within  $(0, 1/2)$ .



When the number  $m$  of measurements is proportional to  $n$  and the energy of the planted solution exceeds  $\log^3 m$ , Theorem 3 asserts that there exists absolutely no estimator that can achieve an estimation error that vanishes as  $n$  increases. This lower limit matches the estimation error of TWF, which corroborates the optimality of TWF under noisy data.

Recall that in many optical imaging applications, the output data we collect are the intensities of the diffractive waves scattered by the sample or specimen under study. The Poisson noise model employs the input  $\mathbf{x}$  and output  $\mathbf{y}$  to describe the numbers of photons diffracted by the specimen and detected by the optical sensor, respectively. Each specimen needs to be sufficiently illuminated in order for the receiver to sense the diffracted light. In such settings, the low-intensity regime  $\|\mathbf{x}\| \leq \log^{1.5} m$  is of little practical interest as it corresponds to an illumination with just very few photons. We forego the details.

It is worth noting that apart from WF, various other nonconvex procedures have been proposed as well for phase retrieval, including the error reduction schemes dating back to Gerchberg-Saxton and Fienup [3, 4], iterated projections [22], alternating minimization [21], generalized approximate message passing [23], Kaczmarz method [35], and greedy methods that exploit additional sparsity constraint [27], to name just a few. While these paradigms enjoy favorable empirical behavior, most of them fall short of theoretical support, except for a version of alternating minimization (called AltMinPhase) [21] that requires fresh samples for each iteration. In comparison, AltMinPhase attains  $\epsilon$ -accuracy when the sample complexity exceeds the order of  $n \log^3 n + n \log^2 n \log(1/\epsilon)$ , which is at least a factor of  $\log^3 n$  from optimal and is empirically largely outperformed by the variant that reuses all samples. In contrast, our algorithm uses the same set of samples all the time and is therefore practically appealing. Furthermore, none of these algorithms come with provable stability guarantees, which are particularly important in most realistic scenarios. Numerically, each iteration of Fienup’s algorithm (or alternating minimization) involves solving a least squares problem, and the algorithm converges in tens or hundreds of iterations. This is computationally more expensive than TWF, whose computational complexity is merely about 4 times that of solving a least squares problem. Interesting readers are referred to [1] for a comparison of several non-convex schemes, and [33] for a discussion of other alternative approaches (e.g. [36, 37]) and performance lower bounds (e.g. [38, 39]).

## 2 Algorithm: Truncated Wirtinger Flow

This section describes the two stages of TWF in details, presented in a reverse order. For each stage, we start with some algorithmic issues encountered by WF, which is then used to motivate and explain the basic principles of TWF. Here and throughout, we let  $\mathcal{A} : \mathbb{R}^{n \times n} \mapsto \mathbb{R}^m$  be the linear map

$$\mathbf{M} \in \mathbb{R}^{n \times n} \quad \mapsto \quad \mathcal{A}(\mathbf{M}) := \{\mathbf{a}_i^\top \mathbf{M} \mathbf{a}_i\}_{1 \leq i \leq m}$$

and  $\mathbf{A}$  the design matrix

$$\mathbf{A} := [\mathbf{a}_1, \dots, \mathbf{a}_m]^\top.$$

### 2.1 Regularized gradient stage

For independent samples, the gradient of the real-valued Poisson log-likelihood obeys

$$\sum_{i=1}^m \nabla \ell_i(\mathbf{z}) = \sum_{i=1}^m 2 \underbrace{\frac{y_i - |\mathbf{a}_i^\top \mathbf{z}|^2}{\mathbf{a}_i^\top \mathbf{z}}}_{:= \nu_i} \mathbf{a}_i, \tag{21}$$

where  $\nu_i$  represents the weight assigned to each  $\mathbf{a}_i$ . This forms the descent direction of WF updates.

Unfortunately, WF moving along the preceding direction might not come close to the truth unless  $\mathbf{z}$  is already very close to  $\mathbf{x}$ . To see this, it is helpful to consider any fixed vector  $\mathbf{z} \in \mathbb{R}^n$  independent of the design vectors. The typical size of  $\min_{1 \leq i \leq m} |\mathbf{a}_i^\top \mathbf{z}|$  is about on the order of  $\frac{1}{m} \|\mathbf{z}\|$ , introducing some unreasonably large weights  $\nu_i$ , which can be as large as  $m \|\mathbf{x}\|^2 / \|\mathbf{z}\|$ . Consequently, the iterative updates based on (1.2) often overshoot, and this arises starting from the very initial stage<sup>6</sup>.

<sup>6</sup>For complex-valued data where  $\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) + j\mathcal{N}(\mathbf{0}, \mathbf{I})$ , WF converges empirically, as  $\min_i |\mathbf{a}_i^* \mathbf{z}|$  is much larger than the real-valued case.

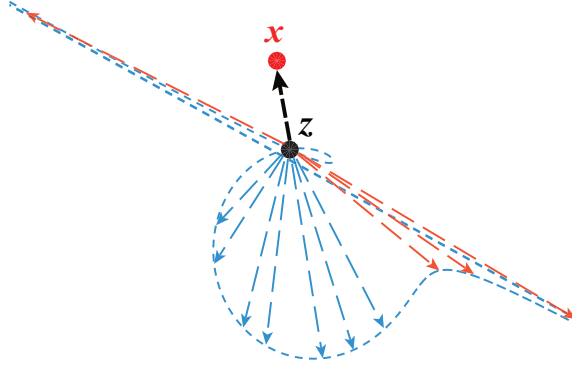


Figure 4: The locus of  $-\frac{1}{2}\nabla\ell_i(\mathbf{z}) = \frac{|\mathbf{a}_i^\top \mathbf{z}|^2 - |\mathbf{a}_i^\top \mathbf{x}|^2}{\mathbf{a}_i^\top \mathbf{z}} \mathbf{a}_i$  when  $\mathbf{a}_i$  ranges over all unit vectors, where  $\mathbf{x} = (2.7, 8)$  and  $\mathbf{z} = (3, 6)$ . For each direction  $\mathbf{a}_i$ ,  $-\frac{1}{2}\nabla\ell_i(\mathbf{z})$  is aligned with  $\mathbf{a}_i$ , and its length represents the weight assigned to this component. In particular, the red arrows depict a few directions that behave like outliers, whereas the blue arrows depict several directions whose resulting gradients take typical sizes.

Fig. 4 illustrates this phenomenon by showing the locus of  $-\nabla\ell_i(\mathbf{z})$  when  $\mathbf{a}_i$  has unit norm and ranges over all possible directions. Examination of the figure seems to suggest that most of the gradient components  $\nabla\ell_i(\mathbf{z})$  are more or less pointing towards the truth  $\mathbf{x}$  and forming reasonable search directions. But there exist a few outlier components that are excessively large, which lead to unstable search directions. Notably, an underlying premise for a nonconvex procedure to succeed is to ensure all iterates reside within a *basin of attraction*, that is, a neighborhood surrounding  $\mathbf{x}$  within which  $\mathbf{x}$  is the unique stationary point of the objective. When a gradient is not well-controlled, the iterative procedure might overshoot and end up leaving this basin of attraction. This intuition is corroborated by numerical experiments under *real-valued* data. As illustrated in Fig. 5, the solutions returned by the WF (designed for a real-valued Poisson log-likelihood and  $m = 8n$ ) are very far from the ground truth.

Hence, to remedy the aforementioned stability issue, it would be natural to separate the small fraction of abnormal gradient components by regularizing the weights  $\nu_i$ , possibly via data-dependent trimming rules. This gives rise to the update rule of TWF:

$$\mathbf{z}^{(t+1)} = \mathbf{z}^{(t)} + \frac{\mu t}{m} \nabla\ell_{\text{tr}}(\mathbf{z}^{(t)}), \quad \forall t \in \mathbb{N}, \quad (22)$$

where  $\nabla\ell_{\text{tr}}(\cdot)$  denotes the regularized gradient given by<sup>7</sup>

$$\nabla\ell_{\text{tr}}(\mathbf{z}) := \sum_{i=1}^m 2 \frac{y_i - |\mathbf{a}_i^\top \mathbf{z}|^2}{\mathbf{a}_i^\top \mathbf{z}} \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i(\mathbf{z}) \cap \mathcal{E}_2^i(\mathbf{z})}. \quad (23)$$

for some trimming criteria specified by  $\mathcal{E}_1^i(\cdot)$  and  $\mathcal{E}_2^i(\cdot)$ . In our algorithm, we take  $\mathcal{E}_1^i(\mathbf{z})$  and  $\mathcal{E}_2^i(\mathbf{z})$  to be two collections of events given by

$$\mathcal{E}_1^i(\mathbf{z}) := \left\{ \alpha_z^{\text{lb}} \leq \frac{|\mathbf{a}_i^\top \mathbf{z}|}{\|\mathbf{z}\|} \leq \alpha_z^{\text{ub}} \right\}, \quad (24)$$

$$\mathcal{E}_2^i(\mathbf{z}) := \left\{ |y_i - |\mathbf{a}_i^\top \mathbf{z}|^2| \leq \frac{\alpha_h}{m} \|\mathbf{y} - \mathcal{A}(\mathbf{z}\mathbf{z}^\top)\|_1 \frac{|\mathbf{a}_i^\top \mathbf{z}|}{\|\mathbf{z}\|} \right\}, \quad (25)$$

where  $\alpha_z^{\text{lb}}$ ,  $\alpha_z^{\text{ub}}$ ,  $\alpha_z$  are predetermined thresholds. To keep notation light, we shall use  $\mathcal{E}_1^i$  and  $\mathcal{E}_2^i$  rather than  $\mathcal{E}_1^i(\mathbf{z})$  and  $\mathcal{E}_2^i(\mathbf{z})$  whenever it is clear from context.

<sup>7</sup>In the complex-valued case, the trimming rule is enforced upon the Wirtinger derivative, which reads  $\nabla\ell_{\text{tr}}(\mathbf{z}) := \sum_{i=1}^m 2 \frac{y_i - |\mathbf{z}^* \mathbf{a}_i|^2}{\mathbf{z}^* \mathbf{a}_i} \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i(\mathbf{z}) \cap \mathcal{E}_2^i(\mathbf{z})}$ .

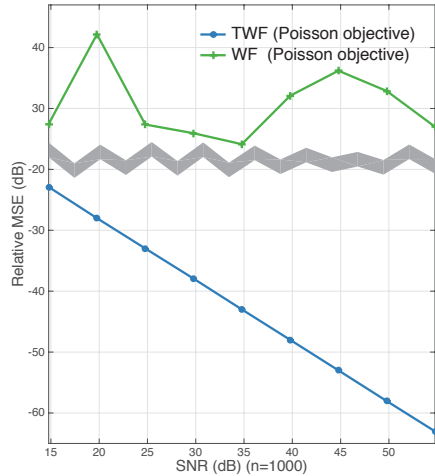


Figure 5: Relative MSE vs. SNR in dB. The curves are shown for WF and TWF, both employing the Poisson log-likelihood. Here,  $\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ,  $n = 1000$ ,  $m = 8n$ , and each point is averaged over 100 Monte Carlo trials.

We emphasize that the above trimming procedure simply throws away those components whose weights  $\nu_i$ 's fall outside some confidence range, so as to remove the influence of outlier components. To achieve this, we regularize both the numerator and denominator of  $\nu_i$  by enforcing separate trimming rules. Recognize that for any fixed  $\mathbf{z}$ , the denominator obeys

$$\mathbb{E} [|\mathbf{a}_i^\top \mathbf{z}|] = \sqrt{2/\pi} \|\mathbf{z}\|,$$

leading up to the rule (24). Regarding the numerator, by the law of large numbers one would expect

$$\mathbb{E} [ |y_i - |\mathbf{a}_i^\top \mathbf{z}|^2 | ] \approx \frac{1}{m} \|\mathbf{y} - \mathcal{A}(\mathbf{z}\mathbf{z}^\top)\|_1,$$

and hence it is natural to regularize the numerator by ensuring

$$|y_i - |\mathbf{a}_i^\top \mathbf{z}|^2| \lesssim \frac{1}{m} \|\mathbf{y} - \mathcal{A}(\mathbf{z}\mathbf{z}^\top)\|_1.$$

As a remark, we include an extra term  $|\mathbf{a}_i^\top \mathbf{z}|/\|\mathbf{z}\|$  in (25) to sharpen the theory, but all our results continue to hold (up to some modification of constants) if we drop this term in (25). Detailed procedures are summarized in Algorithm 1<sup>8</sup>.

The proposed paradigm could be counter-intuitive at first glance, since one might expect the larger terms to be better aligned with the desired search direction. The issue, however, is that the large terms are extremely volatile and could have too high of a leverage on the descent directions. In contrast, TWF discards these high-leverage data, which slightly increases the bias but remarkably reduces the variance of the descent direction. We expect such gradient regularization and variance reduction schemes to be beneficial for solving a broad family of nonconvex problems.

## 2.2 Truncated spectral initialization

In order for the gradient stage to converge rapidly, we need to seed it with a suitable initialization. One natural alternative is the spectral method adopted in [1, 21], which amounts to computing the leading

<sup>8</sup>Careful readers might note that we include some extra factor  $\frac{\sqrt{n}}{\|\mathbf{a}_i\|}$  (which is approximately 1 in the Gaussian model) in Algorithm 1. This occurs since we present Algorithm 1 in a more general fashion that applies beyond the model  $\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , but all results / proofs continue to hold in the presence of this extra term.

---

**Algorithm 1** Truncated Wirtinger Flow.

**Input:** Measurements  $\{y_i \mid 1 \leq i \leq m\}$  and sampling vectors  $\{\mathbf{a}_i \mid 1 \leq i \leq m\}$ ; trimming thresholds  $\alpha_z^{\text{lb}}$ ,  $\alpha_z^{\text{ub}}$ ,  $\alpha_h$ , and  $\alpha_y$  (see default values in Table 1).

**Initialize**  $\mathbf{z}^{(0)}$  to be  $\sqrt{\frac{mn}{\sum_{i=1}^m \|\mathbf{a}_i\|^2}} \lambda_0 \tilde{\mathbf{z}}$ , where  $\lambda_0 = \sqrt{\frac{1}{m} \sum_{i=1}^m y_i}$  and  $\tilde{\mathbf{z}}$  is the leading eigenvector of

$$\mathbf{Y} = \frac{1}{m} \sum_{i=1}^m y_i \mathbf{a}_i \mathbf{a}_i^* \mathbf{1}_{\{|y_i| \leq \alpha_y^2 \lambda_0^2\}}. \quad (26)$$

**Loop:** for  $t = 0 : T$  do

$$\mathbf{z}^{(t+1)} = \mathbf{z}^{(t)} + \frac{2\mu_t}{m} \sum_{i=1}^m \frac{y_i - |\mathbf{a}_i^* \mathbf{z}^{(t)}|^2}{\mathbf{z}^{(t)*} \mathbf{a}_i} \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i}, \quad (27)$$

where

$$\mathcal{E}_1^i := \left\{ \alpha_z^{\text{lb}} \leq \frac{\sqrt{n} |\mathbf{a}_i^* \mathbf{z}^{(t)}|}{\|\mathbf{a}_i\| \|\mathbf{z}^{(t)}\|} \leq \alpha_z^{\text{ub}} \right\}, \quad \mathcal{E}_2^i := \left\{ |y_i - |\mathbf{a}_i^* \mathbf{z}^{(t)}|^2| \leq \alpha_h K_t \frac{\sqrt{n} |\mathbf{a}_i^* \mathbf{z}^{(t)}|}{\|\mathbf{a}_i\| \|\mathbf{z}^{(t)}\|} \right\}, \quad (28)$$

$$\text{and } K_t := \frac{1}{m} \sum_{l=1}^m |y_l - |\mathbf{a}_l^* \mathbf{z}^{(t)}|^2|.$$

**Output**  $\mathbf{z}_T$ .

---

eigenvector of  $\tilde{\mathbf{Y}} := \frac{1}{m} \sum_{i=1}^m y_i \mathbf{a}_i \mathbf{a}_i^\top$ . This arises from the observation that when  $\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and  $\|\mathbf{x}\| = 1$ ,

$$\mathbb{E}[\tilde{\mathbf{Y}}] = \mathbf{I} + 2\mathbf{x}\mathbf{x}^\top,$$

whose leading eigenvector is exactly  $\mathbf{x}$  with an eigenvalue of 3.

Unfortunately, this spectral technique converges to a good initial point only when  $m \gtrsim n \log n$ , due to the fact that  $(\mathbf{a}_i^\top \mathbf{x})^2 \mathbf{a}_i \mathbf{a}_i^\top$  is heavy-tailed, a random quantity which does not have a moment generating function. To be more precise, consider the noiseless case  $y_i = |\mathbf{a}_i^\top \mathbf{x}|^2$  and recall that  $\max_i y_i \approx 2 \log m$ . Letting  $k = \arg \max_i y_i$ , one can calculate

$$\left( \frac{\mathbf{a}_k}{\|\mathbf{a}_k\|} \right)^\top \tilde{\mathbf{Y}} \frac{\mathbf{a}_k}{\|\mathbf{a}_k\|} \geq \left( \frac{\mathbf{a}_k}{\|\mathbf{a}_k\|} \right)^\top \left( \frac{1}{m} \mathbf{a}_k \mathbf{a}_k^\top \right) (\mathbf{a}_k^\top \mathbf{x})^2 \left( \frac{\mathbf{a}_k}{\|\mathbf{a}_k\|} \right) \approx \frac{2n \log m}{m},$$

which is much larger than  $\mathbf{x}^\top \tilde{\mathbf{Y}} \mathbf{x} = 3$  unless  $m/n$  is very large. This tells us that in the regime where  $m \asymp n$ , there exists some unit vector  $\mathbf{a}_k / \|\mathbf{a}_k\|$  that is closer to the leading eigenvector of  $\tilde{\mathbf{Y}}$  than  $\mathbf{x}$ . This phenomenon happens because the summands of  $\tilde{\mathbf{Y}}$  have huge tails so that even one large term could end up dominating the empirical sum, thus preventing the spectral method from returning a meaningful initial guess.

To address this issue, we propose a more robust version of spectral method, which discards those observations  $y_i$  that are several times larger than the mean during spectral initialization. Specifically, the initial estimate is obtained by computing the leading eigenvector  $\tilde{\mathbf{z}}$  of the truncated sum

$$\mathbf{Y} := \frac{1}{m} \sum_{i=1}^m y_i \mathbf{a}_i \mathbf{a}_i^\top \mathbf{1}_{\{|y_i| \leq \alpha_y^2 (\frac{1}{m} \sum_{i=1}^m y_i)\}} \quad (29)$$

for some predetermined threshold  $\alpha_y$ , and then rescaling  $\tilde{\mathbf{z}}$  so as to have roughly the same norm as  $\mathbf{x}$  (which is estimated to be  $\frac{1}{m} \sum_{l=1}^m y_l$ ); see Algorithm 1 for the detailed procedure.

Notably, the aforementioned drawback of the spectral method is not merely a theoretical concern but rather a substantial practical issue. We have seen this in Fig. 2 (main quad example) showing the enormous advantage of truncated spectral initialization. This is also further illustrated in Fig. 6, which compares the empirical efficiency of both methods with  $\alpha_y = 3$  set to be the truncation threshold. For both Gaussian designs and CDP models, the empirical loss incurred by the original spectral method increases as  $n$  grows,

Table 1: Range of algorithmic parameters

(a) **When a fixed step size  $\mu_t \equiv \mu$  is employed:**  $(\alpha_z^{\text{lb}}, \alpha_z^{\text{ub}}, \alpha_h, \alpha_y)$  obeys

$$\begin{cases} \zeta_1 := \max \left\{ \mathbb{E} \left[ \xi^2 \mathbf{1}_{\{|\xi| \leq \sqrt{1.01} \alpha_z^{\text{lb}} \text{ or } |\xi| \geq \sqrt{0.99} \alpha_z^{\text{ub}}\}} \right], \mathbb{P} \left( |\xi| \leq \sqrt{1.01} \alpha_z^{\text{lb}} \text{ or } |\xi| \geq \sqrt{0.99} \alpha_z^{\text{ub}} \right) \right\} \\ \zeta_2 := \mathbb{E} \left[ \xi^2 \mathbf{1}_{\{|\xi| > 0.473 \alpha_h\}} \right], \\ 2(\zeta_1 + \zeta_2) + \sqrt{8/(9\pi)} \alpha_h^{-1} < 1.99, \\ \alpha_y \geq 3, \end{cases} \quad (30)$$

where  $\xi \sim \mathcal{N}(0, 1)$ . By default,  $\alpha_z^{\text{lb}} = 0.3$ ,  $\alpha_z^{\text{ub}} = \alpha_h = 5$ , and  $\alpha_y = 3$ .

(b) **When  $\mu_t$  is chosen by a backtracking line search:**  $(\alpha_z^{\text{lb}}, \alpha_z^{\text{ub}}, \alpha_h, \alpha_y, \alpha_p)$  obeys

$$0 < \alpha_z^{\text{lb}} \leq 0.1, \quad \alpha_z^{\text{ub}} \geq 5, \quad \alpha_h \geq 6, \quad \alpha_y \geq 3, \quad \text{and} \quad \alpha_p \geq 5. \quad (31)$$

By default,  $\alpha_z^{\text{lb}} = 0.1$ ,  $\alpha_z^{\text{ub}} = 5$ ,  $\alpha_h = 6$ ,  $\alpha_y = 3$ , and  $\alpha_p = 5$ .

which is in stark contrast to the truncated spectral method that achieves almost identical accuracy over the same range of  $n$ .

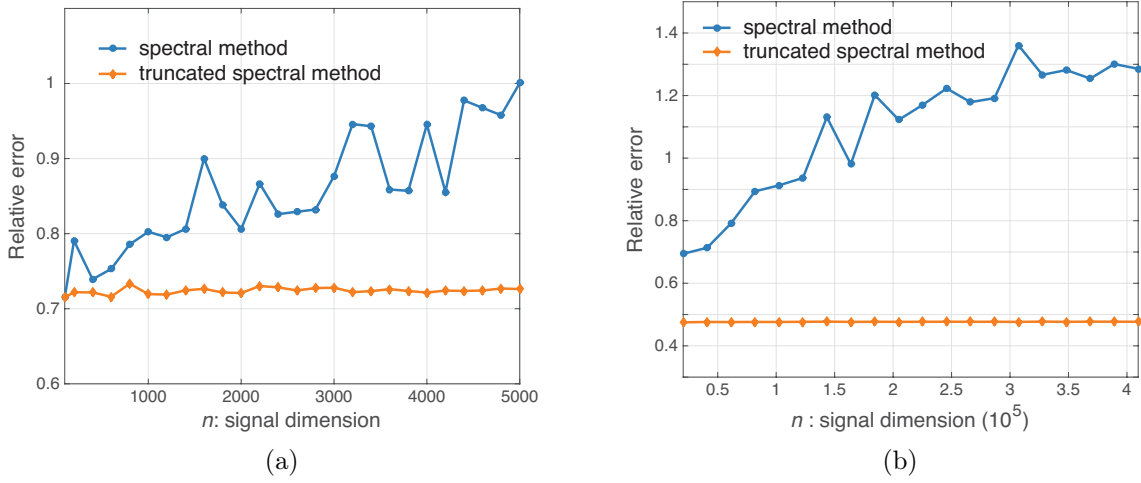


Figure 6: The empirical relative error for both the spectral and the truncated spectral methods. The results are averaged over 50 Monte Carlo runs, and are shown for: (a) 1-D Gaussian measurement where  $\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and  $m = 6n$ ; (b) 2-D CDP model (9) where the diagonal entries of  $\mathbf{D}^{(l)}$  are uniformly drawn from  $\{1, -1, j, -j\}$ ,  $n = n_1 \times n_2$  with  $n_1 = 300$  and  $n_2$  ranging from 64 to 1280, and  $m = 12n$ .

### 2.3 Choice of algorithmic parameters

One implementation detail to specify is the step size  $\mu_t$  at each iteration  $t$ . There are two alternatives that work well in both theory and practice:

1. **Fixed step size.** Take  $\mu_t \equiv \mu$  ( $\forall t \in \mathbb{N}$ ) for some constant  $\mu > 0$ . As long as  $\mu$  is not too large, our main results state that this strategy always works—although the convergence rate depends on  $\mu$ . Under appropriate conditions, our theorems hold for any constant  $0 < \mu < 0.28$ .

2. **Backtracking line search with truncated objective.** This strategy performs a line search along the descent direction

$$\mathbf{p}_t := \frac{1}{m} \nabla \ell_{\text{tr}}(\mathbf{z}_t)$$

and determines an appropriate step size that guarantees a sufficient improvement. In contrast to the conventional search strategy that determines the sufficient progress with respect to the true objective function, we propose to evaluate instead a regularized version of the objective function. Specifically, put

$$\widehat{\ell}(\mathbf{z}) := \sum_{i \in \widehat{\mathcal{T}}(\mathbf{z})} \{y_i \log(|\mathbf{a}_i^\top \mathbf{z}|^2) - |\mathbf{a}_i^\top \mathbf{z}|^2\}, \quad (32)$$

where

$$\widehat{\mathcal{T}}(\mathbf{z}) := \{i \mid |\mathbf{a}_i^\top \mathbf{z}| \geq \alpha_z^{\text{lb}} \|\mathbf{z}\| \text{ and } |\mathbf{a}_i^\top \mathbf{p}| \leq \alpha_p \|\mathbf{p}\|\}.$$

Then the backtracking line search proceeds as

- (a) Start with  $\tau = 1$ ;
- (b) Repeat  $\tau \leftarrow \beta\tau$  until

$$\frac{1}{m} \widehat{\ell}(\mathbf{z}^{(t)} + \tau \mathbf{p}^{(t)}) \geq \frac{1}{m} \widehat{\ell}(\mathbf{z}^{(t)}) + \frac{1}{2} \tau \|\mathbf{p}^{(t)}\|^2, \quad (33)$$

where  $\beta \in (0, 1)$  is some pre-determined constant;

- (c) Set  $\mu_t = \tau$ .

By definition (32), evaluating  $\widehat{\ell}(\mathbf{z}^{(t)} + \tau \mathbf{p}^{(t)})$  mainly consists in calculating the matrix-vector product  $\mathbf{A}(\mathbf{z}^{(t)} + \tau \mathbf{p}^{(t)})$ . In total, we are going to evaluate  $\widehat{\ell}(\mathbf{z}^{(t)} + \tau \mathbf{p}^{(t)})$  for  $O(\log(1/\beta))$  different  $\tau$ 's, and hence the total cost amounts to computing  $\mathbf{A}\mathbf{z}^{(t)}$ ,  $\mathbf{A}\mathbf{p}^{(t)}$  as well as  $O(m \log(1/\beta))$  additional flops. Note that the matrix-vector products  $\mathbf{A}\mathbf{z}^{(t)}$  and  $\mathbf{A}\mathbf{p}^{(t)}$  need to be computed even when one adopts a pre-determined step size. Hence, the extra cost incurred by a backtracking line search, which is  $O(m \log(1/\beta))$  flops, is negligible compared to that of computing the gradient even once.

Another set of important algorithmic parameters to determine is the trimming thresholds  $\alpha_h$ ,  $\alpha_z^{\text{lb}}$ ,  $\alpha_z^{\text{ub}}$ ,  $\alpha_y$ , and  $\alpha_p$  (for a backtracking line search only). The present paper isolates the set of  $(\alpha_h, \alpha_z^{\text{lb}}, \alpha_z^{\text{ub}}, \alpha_y)$  obeying (30) as given in Table 1 when a fixed step size is employed. More concretely, this range subsumes as special cases all parameters obeying the following constraints:

$$0 < \alpha_z^{\text{lb}} \leq 0.5, \quad \alpha_z^{\text{ub}} \geq 5, \quad \alpha_h \geq 5, \quad \text{and} \quad \alpha_y \geq 3. \quad (34)$$

When a backtracking line search is adopted, an extra parameter  $\alpha_p$  is needed, which we take to be  $\alpha_p \geq 5$ . In all theory presented herein, we assume that the parameters fall within the range singled out in Table 1.

### 3 Why TWF works?

Before proceeding, it is best to develop an intuitive understanding of the TWF iterations. We start with a notation representing the (unrecoverable) global phase [1] for real-valued data

$$\phi(\mathbf{z}) := \begin{cases} 0, & \text{if } \|\mathbf{z} - \mathbf{x}\| \leq \|\mathbf{z} + \mathbf{x}\|, \\ \pi, & \text{else.} \end{cases} \quad (35)$$

It is self-evident that

$$(-\mathbf{z}) + \frac{\mu}{m} \nabla_{\text{tr}} \ell(-\mathbf{z}) = - \left\{ \mathbf{z} + \frac{\mu}{m} \nabla_{\text{tr}} \ell(\mathbf{z}) \right\},$$

and hence (cf. Definition (8))

$$\text{dist} \left( (-\mathbf{z}) + \frac{\mu}{m} \nabla_{\text{tr}} \ell(-\mathbf{z}), \mathbf{x} \right) = \text{dist} \left( \mathbf{z} + \frac{\mu}{m} \nabla_{\text{tr}} \ell(\mathbf{z}), \mathbf{x} \right)$$

despite the global phase uncertainty. For simplicity of presentation, we shall drop the phase term by letting  $\mathbf{z}$  be  $e^{-j\phi(\mathbf{z})}\mathbf{z}$  and setting  $\mathbf{h} = \mathbf{z} - \mathbf{x}$ , whenever it is clear from context.

The first object to consider is the descent direction. To this end, we find it convenient to work with a fixed  $\mathbf{z}$  independent of the design vectors  $\mathbf{a}_i$ , which is of course heuristic but helpful in developing some intuition. Rewrite

$$\begin{aligned}\nabla\ell_i(\mathbf{z}) &= 2\frac{(\mathbf{a}_i^\top\mathbf{x})^2 - (\mathbf{a}_i^\top\mathbf{z})^2}{\mathbf{a}_i^\top\mathbf{z}}\mathbf{a}_i \stackrel{(i)}{=} -2\frac{(\mathbf{a}_i^\top\mathbf{h})(2\mathbf{a}_i^\top\mathbf{z} - \mathbf{a}_i^\top\mathbf{h})}{\mathbf{a}_i^\top\mathbf{z}}\mathbf{a}_i \\ &= -4(\mathbf{a}_i^\top\mathbf{h})\mathbf{a}_i + 2\underbrace{\frac{(\mathbf{a}_i^\top\mathbf{h})^2}{\mathbf{a}_i^\top\mathbf{z}}\mathbf{a}_i}_{:=\mathbf{r}_i},\end{aligned}\tag{36}$$

where (i) follows from the identity  $a^2 - b^2 = (a + b)(a - b)$ . The first component of (36), which on average gives  $-4\mathbf{h}$ , makes a good search direction when averaged over all the observations  $i = 1, \dots, m$ . The issue is that the other term  $\mathbf{r}_i$ —which is in general non-integrable—could be devastating. The reason is that  $\mathbf{a}_i^\top\mathbf{z}$  could be arbitrarily small, thus resulting in an unbounded  $\mathbf{r}_i$ . As a consequence, a non-negligible portion of the  $\mathbf{r}_i$ 's may exert a very strong influence on the descent direction in an undesired manner.

Such an issue can be prevented if one can detect and separate those gradient components bearing abnormal  $\mathbf{r}_i$ 's. Since we cannot observe the individual components of the decomposition (36), we cannot reject indices with large values of  $\mathbf{r}_i$  directly. Instead, we examine each gradient component as a whole and discard it if its size is not absolutely controlled. Fortunately, such a strategy is sufficient to ensure that most of the contribution from the regularized gradient comes from the first component of (36), namely,  $-4(\mathbf{a}_i^\top\mathbf{h})\mathbf{a}_i$ . As will be made precise in Proposition 2 and Lemma 7, the regularized gradient obeys

$$-\left\langle\frac{1}{m}\nabla\ell_{\text{tr}}(\mathbf{z}), \mathbf{h}\right\rangle \geq (4 - \epsilon)\|\mathbf{h}\|^2 - O\left(\frac{\|\mathbf{h}\|^3}{\|\mathbf{z}\|}\right)\tag{37}$$

$$\text{and } \left\|\frac{1}{m}\nabla\ell_{\text{tr}}(\mathbf{z})\right\| \lesssim \|\mathbf{h}\|.\tag{38}$$

Here, one has  $(4 - \epsilon)\|\mathbf{h}\|^2$  in (37) instead of  $4\|\mathbf{h}\|^2$  to account for the bias introduced by adaptive trimming, where  $\epsilon$  is small as long as we only throw away a small fraction of data. Looking at (37) and (38) we see that the search direction is sufficiently aligned with the deviation  $-\mathbf{h} = \mathbf{x} - \mathbf{z}$  of the current iterate; i.e. they form a reasonably good angle that is bounded away from  $90^\circ$ . Consequently,  $\mathbf{z}$  is expected to be dragged towards  $\mathbf{x}$  provided that the step size is appropriately chosen.

The observations (37) and (38) are reminiscent of a (local) regularity condition given in [1], which is a fundamental criterion that dictates rapid convergence of iterative procedures (including WF and other gradient descent schemes). When specialized to TWF, we say that  $-\frac{1}{m}\nabla\ell_{\text{tr}}(\cdot)$  satisfies the *regularity condition*, denoted by RC  $(\mu, \lambda, \epsilon)$ , if

$$\left\langle\mathbf{h}, -\frac{1}{m}\nabla\ell_{\text{tr}}(\mathbf{z})\right\rangle \geq \frac{\mu}{2}\left\|\frac{1}{m}\nabla\ell_{\text{tr}}(\mathbf{z})\right\|^2 + \frac{\lambda}{2}\|\mathbf{h}\|^2\tag{39}$$

holds for all  $\mathbf{z}$  obeying  $\|\mathbf{z} - \mathbf{x}\| \leq \epsilon\|\mathbf{x}\|$ , where  $0 < \epsilon < 1$  is some constant. Such an  $\epsilon$ -ball around  $\mathbf{x}$  forms a basin of attraction. Formally, under RC  $(\mu, \lambda, \epsilon)$ , a little algebra gives

$$\begin{aligned}\text{dist}^2\left(\mathbf{z} + \frac{\mu}{m}\nabla\ell_{\text{tr}}(\mathbf{z}), \mathbf{x}\right) &\leq \left\|\mathbf{z} + \frac{\mu}{m}\nabla\ell_{\text{tr}}(\mathbf{z}) - \mathbf{x}\right\|^2 \\ &= \|\mathbf{h}\|^2 + \left\|\frac{\mu}{m}\nabla\ell_{\text{tr}}(\mathbf{z})\right\|^2 + 2\mu\left\langle\mathbf{h}, \frac{1}{m}\nabla\ell_{\text{tr}}(\mathbf{z})\right\rangle \\ &\leq \|\mathbf{h}\|^2 + \left\|\frac{\mu}{m}\nabla\ell_{\text{tr}}(\mathbf{z})\right\|^2 - \mu^2\left\|\frac{1}{m}\nabla\ell_{\text{tr}}(\mathbf{z})\right\|^2 - \mu\lambda\|\mathbf{h}\|^2 \\ &= (1 - \mu\lambda)\text{dist}^2(\mathbf{z}, \mathbf{x})\end{aligned}\tag{40}$$

for any  $\mathbf{z}$  with  $\|\mathbf{z} - \mathbf{x}\| \leq \epsilon$ . In words, the TWF update rule is locally contractive around the planted solution, provided that RC  $(\mu, \lambda, \epsilon)$  holds for some nonzero  $\mu$  and  $\lambda$ . Apparently, Conditions (37) and (38)

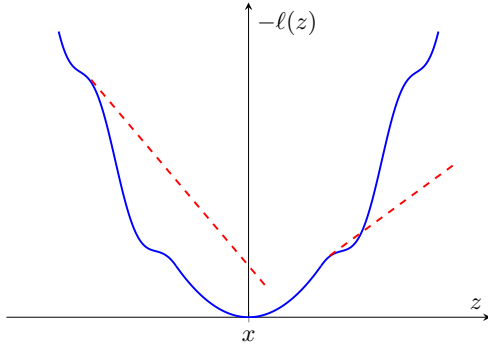


Figure 7: A function  $-\ell(z)$  satisfying RC:  $-\ell(z) = z^2$  for any  $z \in [-6, 6]$ , and  $-\ell(z) = z^2 + 1.5|z|(\cos(|z| - 6) - 1)$  otherwise.

already imply the validity of RC for some constants  $\mu, \lambda \asymp 1$  when  $\|\mathbf{h}\|/\|\mathbf{z}\|$  is reasonably small, which in turn allows us to take a constant step size  $\mu$  and enables a constant contraction rate  $1 - \mu\lambda$ .

Finally, caution must be exercised when connecting RC with strong convexity, since the former does not necessarily guarantee the latter within the basin of attraction. As an illustration, Fig. 7 plots the graph of a non-convex function obeying RC. The distinction stems from the fact that RC is stated only for those pairs  $\mathbf{z}$  and  $\mathbf{h} = \mathbf{z} - \mathbf{x}$  with  $\mathbf{x}$  being a fixed component, rather than simultaneously accommodating all possible  $\mathbf{z}$  and  $\mathbf{h} = \mathbf{z} - \tilde{\mathbf{z}}$  with  $\tilde{\mathbf{z}}$  being an arbitrary vector. In contrast, RC says that the only stationary point of the truncated objective in a neighborhood of  $\mathbf{x}$  is  $\mathbf{x}$ , which often suffices for a gradient-descent type scheme to succeed.

## 4 Numerical experiments

In this section, we report additional numerical results to verify the practical applicability of TWF. In all numerical experiments conducted in the current paper, we set

$$\alpha_z^{\text{lb}} = 0.3, \quad \alpha_z^{\text{ub}} = 5, \quad \alpha_h = 5, \quad \text{and} \quad \alpha_y = 3. \quad (41)$$

This is a concrete combination of parameters satisfying our condition (30). Unless otherwise noted, we employ 50 power iterations for initialization, adopt a fixed step size  $\mu_t \equiv 0.2$  when updating TWF iterates, and set the maximum number of iterations to be  $T = 1000$  for the iterative refinement stage.

The first series of experiments concerns exact recovery from noise-free data. Set  $n = 1000$  and generate a real-valued signal  $\mathbf{x}$  at random. Then for  $m$  varying between  $2n$  and  $6n$ , generate  $m$  design vectors  $\mathbf{a}_i$  independently drawn from  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ . An experiment is claimed to succeed if the returned estimate  $\hat{\mathbf{x}}$  satisfies  $\text{dist}(\hat{\mathbf{x}}, \mathbf{x})/\|\mathbf{x}\| \leq 10^{-5}$ . Fig. 8 illustrates the empirical success rate of TWF (over 100 Monte Carlo trials for each  $m$ ) revealing that exact recovery is practically guaranteed from fewer than 1000 iterations when the number of quadratic constraints is about 5 times the ambient dimension.

To see how special the real-valued Gaussian designs are to our theoretical finding, we perform experiments on two other types of measurement models. In the first, TWF is applied to complex-valued data by generating  $\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \frac{1}{2}\mathbf{I}) + j\mathcal{N}(\mathbf{0}, \frac{1}{2}\mathbf{I})$ . The other is the model of coded diffraction patterns described in (9). Fig. 9 depicts the average success rate for both types of measurements over 100 Monte Carlo trials, indicating that  $m > 4.5n$  and  $m \geq 6n$  are often sufficient under complex-valued Gaussian and CDP models, respectively.

For the sake of comparison, we also report the empirical performance of WF in all the above settings, where the step size is set to be the default choice of [1], that is,  $\mu_t = \min\{1 - e^{-t/330}, 0.2\}$ . As can be seen, the empirical success rates of TWF outperform WF when  $T = 1000$  under Gaussian models, suggesting that TWF either converges faster or exhibits better phase transition behavior.

Another series of experiments has been carried out to demonstrate the stability of TWF when the number  $m$  of quadratic equations varies. We consider the case where  $n = 1000$ , and vary the SNR (cf. (10)) from 15 dB to 55dB. The design vectors are real-valued independent Gaussian  $\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , while the measurements



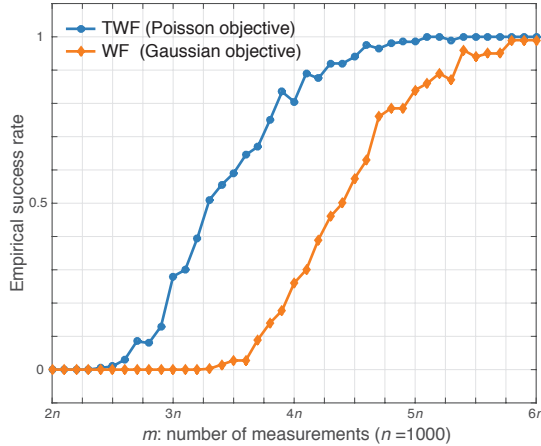


Figure 8: Empirical success rate under real-valued Gaussian sampling  $\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ .

$y_i$  are generated according to the Poisson noise model (4). Fig. 10 shows the relative mean square error—in the dB scale—as a function of SNR, when averaged over 100 independent runs. For all choices of  $m$ , the numerical experiments demonstrate that the relative MSE scales inversely proportional to SNR, which matches our stability guarantees in Theorem 2 (since we observe that on the dB scale, the slope is about -1 as predicted by the theory (19)).

## 5 Exact recovery from noiseless data

This section proves the theoretical guarantees of TWF in the absence of noise (i.e. Theorem 1). We separate the noiseless case mainly out of pedagogical reasons, as most of the steps carry over to the noisy case with slight modification.

The analysis for the truncated spectral method relies on the celebrated Davis-Kahan  $\sin \Theta$  theorem [56], which we defer to Appendix C. In short, for any fixed  $\delta > 0$  and  $\mathbf{x} \in \mathbb{R}^n$ , the initial point  $\mathbf{z}^{(0)}$  returned by the truncated spectral method obeys

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

with high probability, provided that  $m/n$  exceeds some numerical constant. With this in place, it suffices to demonstrate that the TWF update rule is locally contractive, as stated in the following proposition.

**Proposition 1 (Local error contraction).** *Consider the noiseless case (1). Under the condition (30), there exist some universal constants  $0 < \rho_0 < 1$  and  $c_0, c_1, c_2 > 0$  such that with probability exceeding  $1 - c_1 \exp(-c_2 m)$ ,*

$$\text{dist}^2\left(\mathbf{z} + \frac{\mu}{m} \nabla \ell_{\text{tr}}(\mathbf{z}), \mathbf{x}\right) \leq (1 - \rho_0) \text{dist}^2(\mathbf{z}, \mathbf{x}) \quad (42)$$

holds simultaneously for all  $\mathbf{x}, \mathbf{z} \in \mathbb{R}^n$  obeying

$$\frac{\text{dist}(\mathbf{z}, \mathbf{x})}{\|\mathbf{z}\|} \leq \min \left\{ \frac{1}{11}, \frac{\alpha_z^{\text{lb}}}{3\alpha_h}, \frac{\alpha_z^{\text{lb}}}{6}, \frac{5.7(\alpha_z^{\text{lb}})^2}{2\alpha_z^{\text{ub}} + \alpha_z^{\text{lb}}} \right\}, \quad (43)$$

provided that  $m \geq c_0 n$  and that  $\mu$  is some constant obeying

$$0 < \mu \leq \mu_0 := \frac{0.994 - \zeta_1 - \zeta_2 - \sqrt{2/(9\pi)} \alpha_h^{-1}}{2(1.02 + 0.665/\alpha_h)}.$$

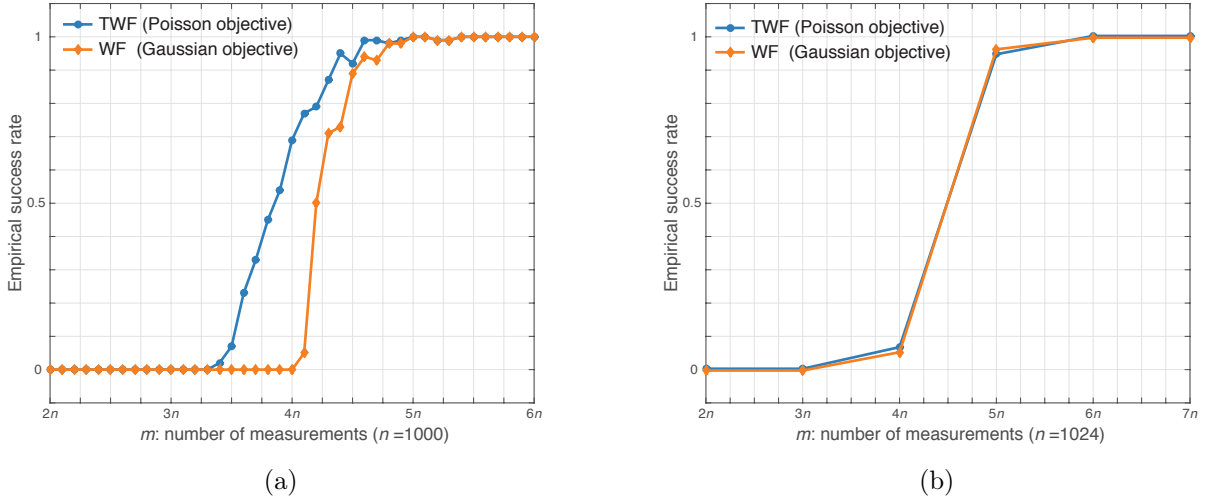


Figure 9: Empirical success rate for exact recovery using TWF. The results are shown for (a) complex-valued Gaussian sampling  $\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \frac{1}{2}\mathbf{I}_n) + j\mathcal{N}(\mathbf{0}, \frac{1}{2}\mathbf{I}_n)$ , and (b) CDP with masks uniformly drawn from  $\{1, -1, j, -j\}$ .

Proposition 1 reveals the monotonicity of the estimation error: once entering a neighborhood around  $\mathbf{x}$  of a reasonably small size, the iterative updates will remain within this neighborhood all the time and be attracted towards  $\mathbf{x}$  at a geometric rate.

As shown in Section 3, under the hypothesis  $\text{RC}(\mu, \lambda, \epsilon)$  one can conclude

$$\text{dist}^2\left(\mathbf{z} + \frac{\mu}{m}\nabla\ell_{\text{tr}}(\mathbf{z}), \mathbf{x}\right) \leq (1 - \mu\lambda)\text{dist}^2(\mathbf{z}, \mathbf{x}), \quad \forall(\mathbf{z}, \mathbf{x}) \text{ with } \text{dist}(\mathbf{z}, \mathbf{x}) \leq \epsilon. \quad (44)$$

Thus, everything now boils down to showing  $\text{RC}(\mu, \lambda, \epsilon)$  for some constants  $\mu, \lambda, \epsilon > 0$ . This occupies the rest of this section.

### 5.1 Preliminary facts about $\{\mathcal{E}_1^i\}$ and $\{\mathcal{E}_2^i\}$

Before proceeding, we gather a few properties of the events  $\mathcal{E}_1^i$  and  $\mathcal{E}_2^i$ , which will prove crucial in establishing  $\text{RC}(\mu, \lambda, \epsilon)$ . To begin with, recall that the truncation level given in  $\mathcal{E}_2^i$  depends on  $\frac{1}{m}\|\mathcal{A}(\mathbf{x}\mathbf{x}^\top - \mathbf{z}\mathbf{z}^\top)\|_1$ . Instead of working with this random variable directly, we use deterministic quantities that are more amenable to analysis. Specifically, we claim that  $\frac{1}{m}\|\mathcal{A}(\mathbf{x}\mathbf{x}^\top - \mathbf{z}\mathbf{z}^\top)\|_1$  offers a uniform and orderwise tight estimate on  $\|\mathbf{h}\|\|\mathbf{z}\|$ , which can be seen from the following two facts.

**Lemma 1.** Fix  $\zeta \in (0, 1)$ . If  $m > c_0 n \zeta^{-2} \log \frac{1}{\zeta}$ , then with probability at least  $1 - C \exp(-c_1 \zeta^2 m)$ ,

$$0.9(1 - \zeta)\|\mathbf{M}\|_{\text{F}} \leq \frac{1}{m}\|\mathcal{A}(\mathbf{M})\|_1 \leq (1 + \zeta)\sqrt{2}\|\mathbf{M}\|_{\text{F}} \quad (45)$$

holds for all symmetric rank-2 matrices  $\mathbf{M} \in \mathbb{R}^{n \times n}$ . Here,  $c_0, c_1, C > 0$  are some universal constants.

*Proof.* Since [6, Lemma 3.1] already establishes the upper bound, it suffices to prove the lower tail bound. Consider all symmetric rank-2 matrices  $\mathbf{M}$  with eigenvalues 1 and  $-t$  for some  $-1 \leq t \leq 1$ . When  $t \in [0, 1]$ , it has been shown in the proof of [6, Lemma 3.2] that with high probability,

$$\frac{1}{m}\|\mathcal{A}(\mathbf{M})\|_1 \geq (1 - \zeta)f(t), \quad (46)$$

for all such rank-2 matrices  $\mathbf{M}$ , where  $f(t) := \frac{2}{\pi}\{2\sqrt{t} + (1 - t)(\pi/2 - 2\arctan(\sqrt{t}))\}$ . The lower bound in this case can then be justified by recognizing that  $f(t)/\sqrt{1+t^2} \geq 0.9$  for all  $t \in [0, 1]$ , as illustrated in Fig. 11. The case where  $t \in [-1, 0]$  is an immediate consequence from [6, Lemma 3.1].  $\square$

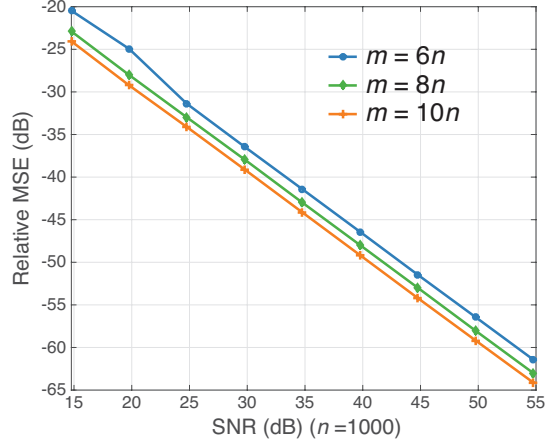


Figure 10: Relative MSE vs. SNR when the  $y_i$ 's follow the Poisson model.

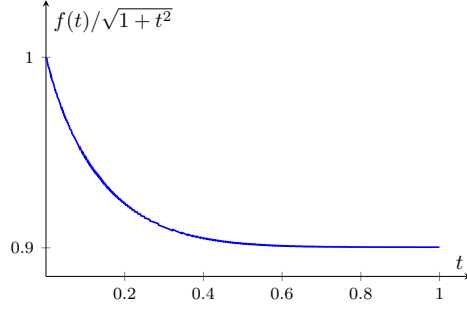


Figure 11:  $\frac{f(t)}{\sqrt{1+t^2}}$  as a function of  $t$ .

**Lemma 2.** Consider any  $\mathbf{x}, \mathbf{z} \in \mathbb{R}^n$  obeying  $\|\mathbf{z} - \mathbf{x}\| \leq \delta \|\mathbf{z}\|$  for some  $\delta < \frac{1}{2}$ . Then one has

$$\sqrt{2-4\delta} \|\mathbf{z} - \mathbf{x}\| \|\mathbf{z}\| \leq \|\mathbf{x}\mathbf{x}^\top - \mathbf{z}\mathbf{z}^\top\|_{\text{F}} \leq (2+\delta) \|\mathbf{z} - \mathbf{x}\| \|\mathbf{z}\|. \quad (47)$$

*Proof.* Take  $\mathbf{h} = \mathbf{z} - \mathbf{x}$  and write

$$\begin{aligned} \|\mathbf{x}\mathbf{x}^\top - \mathbf{z}\mathbf{z}^\top\|_{\text{F}}^2 &= \|\mathbf{h}\mathbf{z}^\top - \mathbf{z}\mathbf{h}^\top + \mathbf{h}\mathbf{h}^\top\|_{\text{F}}^2 \\ &= \|\mathbf{h}\mathbf{z}^\top + \mathbf{z}\mathbf{h}^\top\|_{\text{F}}^2 + \|\mathbf{h}\mathbf{h}^\top\|_{\text{F}}^2 - 2\langle \mathbf{h}\mathbf{z}^\top + \mathbf{z}\mathbf{h}^\top, \mathbf{h}\mathbf{h}^\top \rangle \\ &= 2\|\mathbf{z}\|^2 \|\mathbf{h}\|^2 + 2|\mathbf{h}^\top \mathbf{z}|^2 + \|\mathbf{h}\mathbf{h}^\top\|_{\text{F}}^2 - 2\|\mathbf{h}\|^2 (\mathbf{h}^\top \mathbf{z} + \mathbf{z}^\top \mathbf{h}). \end{aligned}$$

When  $\|\mathbf{h}\| < \frac{1}{2}\|\mathbf{z}\|$ , the Cauchy-Schwarz inequality gives

$$2\|\mathbf{z}\|^2 \|\mathbf{h}\|^2 - 4\|\mathbf{z}\| \|\mathbf{h}\|^3 \leq \|\mathbf{x}\mathbf{x}^\top - \mathbf{z}\mathbf{z}^\top\|_{\text{F}}^2 \leq 4\|\mathbf{z}\|^2 \|\mathbf{h}\|^2 + 4\|\mathbf{h}\|^3 \|\mathbf{z}\| + \|\mathbf{h}\mathbf{h}^\top\|_{\text{F}}^2, \quad (48)$$

$$\Rightarrow \sqrt{(2\|\mathbf{z}\| - 4\|\mathbf{h}\|) \|\mathbf{z}\|} \cdot \|\mathbf{h}\| \leq \|\mathbf{x}\mathbf{x}^\top - \mathbf{z}\mathbf{z}^\top\|_{\text{F}} \leq (2\|\mathbf{z}\| + \|\mathbf{h}\|) \cdot \|\mathbf{h}\| \quad (49)$$

as claimed.  $\square$

Taken together the above two facts demonstrate that with probability  $1 - \exp(-\Omega(m))$ ,

$$1.15 \|\mathbf{z} - \mathbf{x}\| \|\mathbf{z}\| \leq \frac{1}{m} \|\mathcal{A}(\mathbf{x}\mathbf{x}^\top - \mathbf{z}\mathbf{z}^\top)\|_1 \leq 3 \|\mathbf{z} - \mathbf{x}\| \|\mathbf{z}\| \quad (50)$$

holds simultaneously for all  $\mathbf{z}$  and  $\mathbf{x}$  satisfying  $\|\mathbf{h}\| \leq \frac{1}{11} \|\mathbf{z}\|$ . Conditional on (50), the inclusion

$$\mathcal{E}_3^i \subseteq \mathcal{E}_2^i \subseteq \mathcal{E}_4^i \quad (51)$$

holds with respect to the following events

$$\mathcal{E}_3^i : = \left\{ \left| |\mathbf{a}_i^\top \mathbf{x}|^2 - |\mathbf{a}_i^\top \mathbf{z}|^2 \right| \leq 1.15\alpha_h \|\mathbf{h}\| \cdot |\mathbf{a}_i^\top \mathbf{z}| \right\}, \quad (52)$$

$$\mathcal{E}_4^i : = \left\{ \left| |\mathbf{a}_i^\top \mathbf{x}|^2 - |\mathbf{a}_i^\top \mathbf{z}|^2 \right| \leq 3\alpha_h \|\mathbf{h}\| \cdot |\mathbf{a}_i^\top \mathbf{z}| \right\}. \quad (53)$$

The point of introducing these new events is that the  $\mathcal{E}_3^i$ 's (resp.  $\mathcal{E}_4^i$ 's) are statistically independent for any fixed  $\mathbf{x}$  and  $\mathbf{z}$  and are, therefore, easier to work with.

Note that each  $\mathcal{E}_3^i$  (resp.  $\mathcal{E}_4^i$ ) is specified by a quadratic inequality. A closer inspection reveals that in order to satisfy these quadratic inequalities, the quantity  $\mathbf{a}_i^\top \mathbf{h}$  must fall within two intervals centered around 0 and  $2\mathbf{a}_i^\top \mathbf{z}$ , respectively. One can thus facilitate analysis by decoupling each quadratic inequality of interest into two simple linear inequalities, as stated in the following lemma.

**Lemma 3.** *For any  $\gamma > 0$ , define*

$$\mathcal{D}_\gamma^i := \left\{ \left| |\mathbf{a}_i^\top \mathbf{x}|^2 - |\mathbf{a}_i^\top \mathbf{z}|^2 \right| \leq \gamma \|\mathbf{h}\| |\mathbf{a}_i^\top \mathbf{z}| \right\}, \quad (54)$$

$$\mathcal{D}_\gamma^{i,1} := \left\{ \frac{|\mathbf{a}_i^\top \mathbf{h}|}{\|\mathbf{h}\|} \leq \gamma \right\}, \quad (55)$$

$$\text{and } \mathcal{D}_\gamma^{i,2} := \left\{ \left| \frac{\mathbf{a}_i^\top \mathbf{h}}{\|\mathbf{h}\|} - \frac{2\mathbf{a}_i^\top \mathbf{z}}{\|\mathbf{h}\|} \right| \leq \gamma \right\}. \quad (56)$$

Thus,  $\mathcal{D}_\gamma^{i,1}$  and  $\mathcal{D}_\gamma^{i,2}$  represent the two intervals on  $\mathbf{a}_i^\top \mathbf{h}$  centered around 0 and  $2\mathbf{a}_i^\top \mathbf{z}$ . If  $\frac{\|\mathbf{h}\|}{\|\mathbf{z}\|} \leq \frac{\alpha^{\text{lb}}}{\gamma}$ , then the following inclusion holds

$$\left( \mathcal{D}_{\frac{\gamma}{1+\sqrt{2}}}^{i,1} \cap \mathcal{E}_1^i \right) \cup \left( \mathcal{D}_{\frac{\gamma}{1+\sqrt{2}}}^{i,2} \cap \mathcal{E}_1^i \right) \subseteq \mathcal{D}_\gamma^i \cap \mathcal{E}_1^i \subseteq \left( \mathcal{D}_\gamma^{i,1} \cap \mathcal{E}_1^i \right) \cup \left( \mathcal{D}_\gamma^{i,2} \cap \mathcal{E}_1^i \right). \quad (57)$$

## 5.2 Proof of the regularity condition

By definition, one step towards proving the regularity condition (39) is to control the norm of the regularized gradient. In fact, a crude argument already reveals that  $\|\frac{1}{m} \nabla \ell_{\text{tr}}(\mathbf{z})\| \lesssim \|\mathbf{h}\|$ . To see this, introduce  $\mathbf{v} = [v_i]_{1 \leq i \leq m}$  with  $v_i := 2 \frac{|\mathbf{a}_i^\top \mathbf{x}|^2 - |\mathbf{a}_i^\top \mathbf{z}|^2}{\mathbf{a}_i^\top \mathbf{z}} \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i}$ . It comes from the trimming rule  $\mathcal{E}_1^i$  as well as the inclusion property (51) that

$$|\mathbf{a}_i^\top \mathbf{z}| \gtrsim \|\mathbf{z}\| \quad \text{and} \quad \left| y_i - |\mathbf{a}_i^\top \mathbf{z}|^2 \right| \leq \frac{1}{m} \|\mathcal{A}(\mathbf{x}\mathbf{x}^\top - \mathbf{z}\mathbf{z}^\top)\|_1 \asymp \|\mathbf{h}\| \|\mathbf{z}\|,$$

implying  $|v_i| \lesssim \|\mathbf{h}\|$  and hence  $\|\mathbf{v}\| \lesssim \sqrt{m} \|\mathbf{h}\|$ . The Marchenko–Pastur law gives  $\|\mathbf{A}\| \lesssim \sqrt{m}$ , whence

$$\frac{1}{m} \|\nabla \ell_{\text{tr}}(\mathbf{z})\| = \frac{1}{m} \|\mathbf{A}^\top \mathbf{v}\| \leq \frac{1}{m} \|\mathbf{A}\| \cdot \|\mathbf{v}\| \lesssim \|\mathbf{h}\|. \quad (58)$$

A more refined estimate will be provided in Lemma 7.

The above argument essentially tells us that to establish RC, it suffices to verify a uniform lower bound of the form

$$-\left\langle \mathbf{h}, \frac{1}{m} \nabla \ell_{\text{tr}}(\mathbf{z}) \right\rangle \gtrsim \|\mathbf{h}\|^2, \quad (59)$$

as formally derived in the following proposition.

**Proposition 2.** *Consider the noise-free measurements  $y_i = |\mathbf{a}_i^\top \mathbf{x}|^2$  and any fixed constant  $\epsilon > 0$ . Under the condition (30), if  $m > c_1 n$ , then with probability exceeding  $1 - C \exp(-c_0 m)$ ,*

$$-\left\langle \mathbf{h}, \frac{1}{m} \nabla \ell_{\text{tr}}(\mathbf{z}) \right\rangle \geq 2 \left\{ 1.99 - 2(\zeta_1 + \zeta_2) - \sqrt{8/(9\pi)} \alpha_h^{-1} - \epsilon \right\} \|\mathbf{h}\|^2 \quad (60)$$

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

$$\frac{\|\mathbf{h}\|}{\|\mathbf{z}\|} \leq \min \left\{ \frac{1}{11}, \frac{\alpha_z^{\text{lb}}}{3\alpha_h}, \frac{\alpha_z^{\text{lb}}}{6}, \frac{5.7(\alpha_z^{\text{lb}})^2}{2\alpha_z^{\text{ub}} + \alpha_z^{\text{lb}}} \right\}. \quad (61)$$

Here,  $c_0, c_1, C > 0$  are some universal constants, and  $\zeta_1$  and  $\zeta_2$  are defined in (30).

The basic starting point is the observation that  $(\mathbf{a}_i^\top \mathbf{z}) - (\mathbf{a}_i^\top \mathbf{x})^2 = (\mathbf{a}_i^\top \mathbf{h})(2\mathbf{a}_i^\top \mathbf{z} - \mathbf{a}_i^\top \mathbf{h})$  and hence

$$\begin{aligned} -\frac{1}{2m} \nabla \ell_{\text{tr}}(\mathbf{z}) &= \frac{1}{m} \sum_{i=1}^m \frac{(\mathbf{a}_i^\top \mathbf{z})^2 - (\mathbf{a}_i^\top \mathbf{x})^2}{\mathbf{a}_i^\top \mathbf{z}} \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i} \\ &= \frac{1}{m} \sum_{i=1}^m 2(\mathbf{a}_i^\top \mathbf{h}) \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i} - \frac{1}{m} \sum_{i=1}^m \frac{(\mathbf{a}_i^\top \mathbf{h})^2}{\mathbf{a}_i^\top \mathbf{z}} \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i}. \end{aligned} \quad (62)$$

One would expect the contribution of the second term of (62) (which is a second-order quantity) to be small as  $\|\mathbf{h}\| / \|\mathbf{z}\|$  decreases.

To facilitate analysis, we rewrite (62) in terms of the more convenient events  $\mathcal{D}_\gamma^{i,1}$  and  $\mathcal{D}_\gamma^{i,2}$ . Specifically, the inclusion property (51) together with Lemma 3 reveals that

$$\mathcal{D}_{\gamma_3}^{i,1} \cap \mathcal{E}_1^i \subseteq \mathcal{E}_3^i \cap \mathcal{E}_1^i \subseteq \mathcal{E}_2^i \cap \mathcal{E}_1^i \subseteq \mathcal{E}_4^i \cap \mathcal{E}_1^i \subseteq (\mathcal{D}_{\gamma_4}^{i,1} \cup \mathcal{D}_{\gamma_4}^{i,2}) \cap \mathcal{E}_1^i, \quad (63)$$

where the parameters  $\gamma_3, \gamma_4$  are given by

$$\gamma_3 := 0.476\alpha_h, \quad \text{and} \quad \gamma_4 := 3\alpha_h. \quad (64)$$

This taken collectively with the identity (62) leads to a lower estimate

$$\begin{aligned} -\left\langle \frac{1}{2m} \nabla \ell_{\text{tr}}(\mathbf{z}), \mathbf{h} \right\rangle &\geq \\ &\frac{2}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h})^2 \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{D}_{\gamma_3}^{i,1}} - \frac{1}{m} \sum_{i=1}^m \frac{|\mathbf{a}_i^\top \mathbf{h}|^3}{|\mathbf{a}_i^\top \mathbf{z}|} \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{D}_{\gamma_4}^{i,1}} - \frac{1}{m} \sum_{i=1}^m \frac{|\mathbf{a}_i^\top \mathbf{h}|^3}{|\mathbf{a}_i^\top \mathbf{z}|} \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{D}_{\gamma_4}^{i,2}}, \end{aligned} \quad (65)$$

leaving us with three quantities in the right-hand side to deal with. We pause here to explain and compare the influences of these three terms.

To begin with, as long as the trimming step does not discard too many data, the first term should be close to  $\frac{2}{m} \sum_i |\mathbf{a}_i^\top \mathbf{h}|^2$ , which approximately gives  $2\|\mathbf{h}\|^2$  from the law of large numbers. This term turns out to be dominant in the right-hand side of (65) as long as  $\|\mathbf{h}\|/\|\mathbf{z}\|$  is reasonably small. To see this, please recognize that the second term in the right-hand side is  $O(\|\mathbf{h}\|^3/\|\mathbf{z}\|)$ , simply because both  $\mathbf{a}_i^\top \mathbf{h}$  and  $\mathbf{a}_i^\top \mathbf{z}$  are absolutely controlled on  $\mathcal{D}_{\gamma_4}^{i,1} \cap \mathcal{E}_1^i$ . However,  $\mathcal{D}_{\gamma_4}^{i,2}$  does not share such a desired feature. By the very definition of  $\mathcal{D}_{\gamma_4}^{i,2}$ , each nonzero summand of the last term of (65) must obey  $|\mathbf{a}_i^\top \mathbf{h}| \approx 2|\mathbf{a}_i^\top \mathbf{z}|$  and, therefore,  $\frac{|\mathbf{a}_i^\top \mathbf{h}|^3}{|\mathbf{a}_i^\top \mathbf{z}|} \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{D}_{\gamma_4}^{i,2}}$  is roughly of the order of  $\|\mathbf{z}\|^2$ ; this could be much larger than our target level  $\|\mathbf{h}\|^2$ . Fortunately,  $\mathcal{D}_{\gamma_4}^{i,2}$  is a rare event, thus precluding a noticeable influence upon the descent direction. All of this is made rigorous in Lemma 4 (first term), Lemma 5 (second term) and Lemma 6 (third term) together with subsequent analysis.

**Lemma 4.** Fix  $\gamma > 0$ , and let  $\mathcal{E}_1^i$  and  $\mathcal{D}_\gamma^{i,1}$  be defined in (24) and (55), respectively. Set

$$\zeta_1 := 1 - \min \left\{ \mathbb{E} \left[ \xi^2 \mathbf{1}_{\{\sqrt{1.01}\alpha_z^{\text{lb}} \leq |\xi| \leq \sqrt{0.99}\alpha_z^{\text{ub}}\}} \right], \mathbb{E} \left[ \mathbf{1}_{\{\sqrt{1.01}\alpha_z^{\text{lb}} \leq |\xi| \leq \sqrt{0.99}\alpha_z^{\text{ub}}\}} \right] \right\} \quad (66)$$

$$\text{and} \quad \zeta_2 := \mathbb{E} \left[ \xi^2 \mathbf{1}_{\{|\xi| > \sqrt{0.99}\gamma\}} \right], \quad (67)$$

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

$$\frac{1}{m} \sum_{i=1}^m |\mathbf{a}_i^\top \mathbf{h}|^2 \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{D}_\gamma^{i,1}} \geq (1 - \zeta_1 - \zeta_2 - \epsilon) \|\mathbf{h}\|^2 \quad (68)$$

holds for all non-zero vectors  $\mathbf{h}, \mathbf{z} \in \mathbb{R}^n$ . Here,  $c_0, c_1, C > 0$  are some universal constants.

We now move on to the second term in the right-hand side of (65). For any fixed  $\gamma > 0$ , the definition of  $\mathcal{E}_1^i$  gives rise to an upper estimate

$$\frac{1}{m} \sum_{i=1}^m \frac{|\mathbf{a}_i^\top \mathbf{h}|^3}{|\mathbf{a}_i^\top \mathbf{z}|} \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{D}_\gamma^{i,1}} \leq \frac{1}{\alpha_z^{\text{lb}} \|\mathbf{z}\|} \frac{1}{m} \sum_{i=1}^m |\mathbf{a}_i^\top \mathbf{h}|^3 \mathbf{1}_{\mathcal{D}_\gamma^{i,1}} \leq \frac{(1+\epsilon) \sqrt{8/\pi} \|\mathbf{h}\|^3}{\alpha_z^{\text{lb}} \|\mathbf{z}\|}, \quad (69)$$

where  $\sqrt{8/\pi} \|\mathbf{h}\|^3$  is exactly the untruncated moment  $\mathbb{E}[|\mathbf{a}_i^\top \mathbf{h}|^3]$ . The second inequality is a consequence of the lemma below, which arises by observing that the summands  $|\mathbf{a}_i^\top \mathbf{h}|^3 \mathbf{1}_{\mathcal{D}_\gamma^{i,1}}$  are independent sub-Gaussian random variables.

**Lemma 5.** *For any constant  $\gamma > 0$ , if  $m/n \geq c_0 \cdot \epsilon^{-2} \log \epsilon^{-1}$ , then*

$$\frac{1}{m} \sum_{i=1}^m |\mathbf{a}_i^\top \mathbf{h}|^3 \mathbf{1}_{\mathcal{D}_\gamma^{i,1}} \leq (1+\epsilon) \sqrt{8/\pi} \|\mathbf{h}\|^3, \quad \forall \mathbf{h} \in \mathbb{R}^n \quad (70)$$

with probability at least  $1 - C \exp(-c_1 \epsilon^2 m)$  for some universal constants  $c_0, c_1, C > 0$ .

It remains to control the last term of (65). As mentioned above, the influence of this term is small since the set of  $\mathbf{a}_i$ 's satisfying  $\mathcal{D}_\gamma^{i,2}$  accounts for a small fraction of measurements. Put formally, the number of equations satisfying  $|\mathbf{a}_i^\top \mathbf{h}| \geq \gamma \|\mathbf{h}\|$  decays rapidly for large  $\gamma$  (at least at a quadratic rate), as stated below.

**Lemma 6.** *For any  $0 < \epsilon < 1$ , there exist some universal constants  $c_0, c_1, C > 0$  such that*

$$\frac{1}{m} \sum_{i=1}^m \mathbf{1}_{\{|\mathbf{a}_i^\top \mathbf{h}| \geq \gamma \|\mathbf{h}\|\}} \leq \frac{1}{0.49\gamma} \exp(-0.485\gamma^2) + \frac{\epsilon}{\gamma^2}, \quad \forall \mathbf{h} \in \mathbb{R}^n \setminus \{\mathbf{0}\} \text{ and } \gamma \geq 2 \quad (71)$$

with probability at least  $1 - C \exp(-c_0 \epsilon^2 m)$ . This holds with the proviso  $m/n \geq c_1 \cdot \epsilon^{-2} \log \epsilon^{-1}$ .

To connect this lemma with the last term of (65), we recognize that when  $\gamma \leq \frac{\alpha_z^{\text{lb}} \|\mathbf{z}\|}{\|\mathbf{h}\|}$ , one has

$$\mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{D}_\gamma^{i,2}} \leq \mathbf{1}_{\{|\mathbf{a}_i^\top \mathbf{h}| \geq \alpha_z^{\text{lb}} \|\mathbf{z}\|\}}. \quad (72)$$

The constraint  $\left| \frac{\mathbf{a}_i^\top \mathbf{h}}{\|\mathbf{h}\|} - \frac{2\mathbf{a}_i^\top \mathbf{z}}{\|\mathbf{h}\|} \right| \leq \gamma$  of  $\mathcal{D}_\gamma^{i,2}$  necessarily requires

$$\frac{|\mathbf{a}_i^\top \mathbf{h}|}{\|\mathbf{h}\|} \geq \frac{2|\mathbf{a}_i^\top \mathbf{z}|}{\|\mathbf{h}\|} - \gamma \geq \frac{2\alpha_z^{\text{lb}} \|\mathbf{z}\|}{\|\mathbf{h}\|} - \gamma \geq \frac{\alpha_z^{\text{lb}} \|\mathbf{z}\|}{\|\mathbf{h}\|}, \quad (73)$$

where the last inequality comes from our assumption on  $\gamma$ . With Lemma 6 in place, (72) immediately gives

$$\begin{aligned} \sum_{i=1}^m \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{D}_\gamma^{i,2}} &\leq \frac{\|\mathbf{h}\|}{0.49\alpha_z^{\text{lb}} \|\mathbf{z}\|} \exp\left(-0.485 \left(\frac{\alpha_z^{\text{lb}} \|\mathbf{z}\|}{\|\mathbf{h}\|}\right)^2\right) + \frac{\epsilon \|\mathbf{h}\|^2}{(\alpha_z^{\text{lb}})^2 \|\mathbf{z}\|^2} \\ &\leq \frac{1}{9800} \left(\frac{\|\mathbf{h}\|}{\alpha_z^{\text{lb}} \|\mathbf{z}\|}\right)^4 + \frac{\epsilon}{(\alpha_z^{\text{lb}})^2} \left(\frac{\|\mathbf{h}\|}{\|\mathbf{z}\|}\right)^2 \end{aligned} \quad (74)$$

as long as  $\frac{\|\mathbf{h}\|}{\|\mathbf{z}\|} \leq \frac{\alpha_z^{\text{lb}}}{6}$ , where the last inequality uses the majorization  $\frac{1}{20000x^4} \geq \frac{1}{x} \exp(-0.485x^2)$  holding for any  $x \geq 6$ .

In addition, on  $\mathcal{E}_1^i \cap \mathcal{D}_\gamma^{i,2}$ , the amplitude of each summand can be bounded in such a way that

$$\frac{|\mathbf{a}_i^\top \mathbf{h}|^3}{|\mathbf{a}_i^\top \mathbf{z}|} \leq \frac{|2\mathbf{a}_i^\top \mathbf{z}| + \gamma \|\mathbf{h}\|}{|\mathbf{a}_i^\top \mathbf{z}|} (2\alpha_z^{\text{ub}} \|\mathbf{z}\| + \gamma \|\mathbf{h}\|)^2 \quad (75)$$

$$\leq \left(2 + \frac{\gamma \|\mathbf{h}\|}{\alpha_z^{\text{lb}} \|\mathbf{z}\|}\right) \left(2\alpha_z^{\text{ub}} + \gamma \frac{\|\mathbf{h}\|}{\|\mathbf{z}\|}\right)^2 \|\mathbf{z}\|^2, \quad (76)$$

where both inequalities are immediate consequences from the definitions of  $\mathcal{D}_\gamma^{i,2}$  and  $\mathcal{E}_1^i$  (see (56) and (24)). Taking this together with the cardinality bound (74) and picking  $\epsilon$  appropriately, we get

$$\frac{1}{m} \sum_{i=1}^m \frac{|\mathbf{a}_i^\top \mathbf{h}|^3}{|\mathbf{a}_i^\top \mathbf{z}|} \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{D}_\gamma^{i,2}} \leq \left\{ \underbrace{\frac{\left(2 + \frac{\gamma}{\alpha_z^{\text{lb}}}\frac{\|\mathbf{h}\|}{\|\mathbf{z}\|}\right) \left(2\alpha_z^{\text{ub}} + \gamma\frac{\|\mathbf{h}\|}{\|\mathbf{z}\|}\right)^2}{9800(\alpha_z^{\text{lb}})^4}}_{\vartheta_1} \|\mathbf{z}\|^2 + \epsilon \right\} \|\mathbf{h}\|^2. \quad (77)$$

Furthermore, under the condition that

$$\gamma \leq \alpha_z^{\text{lb}} \frac{\|\mathbf{z}\|}{\|\mathbf{h}\|} \quad \text{and} \quad \frac{\|\mathbf{h}\|}{\|\mathbf{z}\|} \leq \frac{\sqrt{98}(\alpha_z^{\text{lb}})^2}{\sqrt{3}(2\alpha_z^{\text{ub}} + \alpha_z^{\text{lb}})},$$

one can simplify (77) by observing that  $\vartheta_1 \leq \frac{1}{100}$ , which results in

$$\frac{1}{m} \sum_{i=1}^m \frac{|\mathbf{a}_i^\top \mathbf{h}|^3}{|\mathbf{a}_i^\top \mathbf{z}|} \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{D}_\gamma^{i,2}} \leq \left( \frac{1}{100} + \epsilon \right) \|\mathbf{h}\|^2. \quad (78)$$

Putting all preceding results in this subsection together reveals that with probability exceeding  $1 - \exp(-\Omega(m))$ ,

$$\begin{aligned} - \left\langle \mathbf{h}, \frac{1}{2m} \nabla \ell_{\text{tr}}(\mathbf{z}) \right\rangle &\geq \left\{ 1.99 - 2(\zeta_1 + \zeta_2) - \sqrt{8/\pi} \frac{\|\mathbf{h}\|}{\alpha_z^{\text{lb}} \|\mathbf{z}\|} - 3\epsilon \right\} \|\mathbf{h}\|^2 \\ &\geq \left\{ 1.99 - 2(\zeta_1 + \zeta_2) - \sqrt{8/\pi} (3\alpha_h)^{-1} - 3\epsilon \right\} \|\mathbf{h}\|^2 \end{aligned} \quad (79)$$

holds simultaneously over all  $\mathbf{x}$  and  $\mathbf{z}$  satisfying

$$\frac{\|\mathbf{h}\|}{\|\mathbf{z}\|} \leq \min \left\{ \frac{\alpha_z^{\text{lb}}}{3\alpha_h}, \frac{\alpha_z^{\text{lb}}}{6}, \frac{\sqrt{98/3}(\alpha_z^{\text{lb}})^2}{2\alpha_z^{\text{ub}} + \alpha_z^{\text{lb}}}, \frac{1}{11} \right\} \quad (80)$$

as claimed in Proposition 2.

To conclude this section, we provide a tighter estimate about the norm of the regularized gradient.

**Lemma 7.** Fix  $\delta > 0$ , and assume that  $y_i = (\mathbf{a}_i^\top \mathbf{x})^2$ . Suppose that  $m \geq c_0 n$  for some large constant  $c_0 > 0$ . There exist some universal constants  $c, C > 0$  such that with probability at least  $1 - C \exp(-cm)$ ,

$$\frac{1}{m} \|\nabla \ell_{\text{tr}}(\mathbf{z})\| \leq (1 + \delta) \cdot 4\sqrt{1.02 + 0.665/\alpha_h} \|\mathbf{h}\| \quad (81)$$

holds simultaneously for all  $\mathbf{x}, \mathbf{z} \in \mathbb{R}^n$  satisfying  $\frac{\|\mathbf{h}\|}{\|\mathbf{z}\|} \leq \min \left\{ \frac{\alpha_z^{\text{lb}}}{3\alpha_h}, \frac{\alpha_z^{\text{lb}}}{6}, \frac{\sqrt{98/3}(\alpha_z^{\text{lb}})^2}{2\alpha_z^{\text{ub}} + \alpha_z^{\text{lb}}}, \frac{1}{11} \right\}$ .

Lemma 7 complements the preceding arguments by allowing us to identify a concrete plausible range for the step size. Specifically, putting Lemma 7 and Proposition 2 together suggests that

$$- \left\langle \mathbf{h}, \frac{1}{m} \nabla \ell_{\text{tr}}(\mathbf{z}) \right\rangle \geq \frac{2 \left\{ 1.99 - 2(\zeta_1 + \zeta_2) - \sqrt{8/(9\pi)} \alpha_h^{-1} - \epsilon \right\}}{(1 + \delta)^2 \cdot 16(1.02 + 0.665/\alpha_h)} \left\| \frac{1}{m} \nabla \ell_{\text{tr}}(\mathbf{z}) \right\|^2. \quad (82)$$

Taking  $\epsilon$  and  $\delta$  to be sufficiently small we arrive at a feasible range (cf. Definition (39))

$$\mu \leq \frac{0.994 - \zeta_1 - \zeta_2 - \sqrt{2/(9\pi)} \alpha_h^{-1}}{2(1.02 + 0.665/\alpha_h)} := \mu_0. \quad (83)$$

This establishes Proposition 1 and in turn Theorem 1 when  $\mu_t$  is taken to be a fixed constant.

To justify the contraction under backtracking line search, it suffices to prove that the resulting step size falls within this range (83), which we defer to Appendix D.

## 6 Stability

This section goes in the direction of establishing stability guarantees of TWF. We concentrate on the iterative gradient stage, and defer the analysis for the initialization stage to Appendix C.

Before continuing, we collect two bounds that we shall use several times. The first is the observation that

$$\begin{aligned} \frac{1}{m} \|\mathbf{y} - \mathcal{A}(\mathbf{z}\mathbf{z}^\top)\|_1 &\leq \frac{1}{m} \|\mathcal{A}(\mathbf{x}\mathbf{x}^\top - \mathbf{z}\mathbf{z}^\top)\|_1 + \frac{1}{m} \|\boldsymbol{\eta}\|_1 \\ &\lesssim \|\mathbf{h}\| \|\mathbf{z}\| + \frac{1}{m} \|\boldsymbol{\eta}\|_1 \lesssim \|\mathbf{h}\| \|\mathbf{z}\| + \frac{1}{\sqrt{m}} \|\boldsymbol{\eta}\|, \end{aligned} \quad (84)$$

where the last inequality follows from Cauchy-Schwarz. Setting

$$v_i := 2 \frac{y_i - |\mathbf{a}_i^\top \mathbf{z}|^2}{\mathbf{a}_i^\top \mathbf{z}} \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i}$$

as usual, this inequality together with the trimming rules  $\mathcal{E}_1^i$  and  $\mathcal{E}_2^i$  gives

$$\begin{aligned} |v_i| &\lesssim \|\mathbf{h}\| + \frac{\|\boldsymbol{\eta}\|}{\sqrt{m} \|\mathbf{z}\|} \\ \implies \left\| \frac{1}{m} \nabla \ell_{\text{tr}}(\mathbf{z}) \right\| &= \frac{1}{m} \|\mathbf{A}^\top \mathbf{v}\| \leq \left\| \frac{1}{\sqrt{m}} \mathbf{A} \right\| \frac{1}{\sqrt{m}} \|\mathbf{v}\| \stackrel{(i)}{\lesssim} \frac{1}{\sqrt{m}} \|\mathbf{v}\| \lesssim \|\mathbf{h}\| + \frac{\|\boldsymbol{\eta}\|}{\sqrt{m} \|\mathbf{z}\|}, \end{aligned} \quad (85)$$

where (i) arises from [40, Corollary 5.35].

As discussed in Section 3, the estimation error is contractive if  $-\frac{1}{m} \nabla \ell_{\text{tr}}(\mathbf{z})$  satisfies the regularity condition. With (85) in place, RC reduces to

$$-\frac{1}{m} \langle \nabla \ell_{\text{tr}}(\mathbf{z}), \mathbf{h} \rangle \gtrsim \|\mathbf{h}\|^2. \quad (86)$$

Unfortunately, (86) does not hold for all  $\mathbf{z}$  within the neighborhood of  $\mathbf{x}$  due to the existence of noise. Instead we establish the following:

- The condition (86) holds for all  $\mathbf{h}$  obeying

$$c_3 \frac{\|\boldsymbol{\eta}\| / \sqrt{m}}{\|\mathbf{z}\|} \leq \|\mathbf{h}\| \leq c_4 \|\mathbf{x}\| \quad (87)$$

for some constants  $c_3, c_4 > 0$  (we shall call it *Regime 1*); this will be proved later. In this regime, the reasoning in Section 3 gives

$$\text{dist}\left(\mathbf{z} + \frac{\mu}{m} \nabla \ell_{\text{tr}}(\mathbf{z}), \mathbf{x}\right) \leq (1 - \rho) \text{dist}(\mathbf{z}, \mathbf{x}) \quad (88)$$

for some appropriate constants  $\mu, \rho > 0$  and, hence, error contraction occurs as in the noiseless setting.

- However, once the iterate enters *Regime 2* where

$$\|\mathbf{h}\| \leq \frac{c_3 \|\boldsymbol{\eta}\|}{\sqrt{m} \|\mathbf{z}\|}, \quad (89)$$

the estimation error might no longer be contractive. Fortunately, in this regime each move by  $\frac{\mu}{m} \nabla \ell_{\text{tr}}(\mathbf{z})$  is of size at most  $O(\frac{\|\boldsymbol{\eta}\|}{\sqrt{m} \|\mathbf{z}\|})$ , compare (85). As a result, at each iteration the estimation error cannot increase by more than a numerical constant times  $\frac{\|\boldsymbol{\eta}\|}{\sqrt{m} \|\mathbf{z}\|}$  before possibly jumping out (of this regime). Therefore,

$$\text{dist}\left(\mathbf{z} + \frac{\mu}{m} \nabla \ell_{\text{tr}}(\mathbf{z}), \mathbf{x}\right) \leq c_5 \frac{\|\boldsymbol{\eta}\|}{\sqrt{m} \|\mathbf{x}\|} \quad (90)$$

for some constant  $c_5 > 0$ . Moreover, as long as  $\|\boldsymbol{\eta}\|_\infty / \|\mathbf{x}\|^2$  is sufficiently small, one can guarantee that  $c_5 \frac{\|\boldsymbol{\eta}\|}{\sqrt{m} \|\mathbf{x}\|} \leq c_5 \frac{\|\boldsymbol{\eta}\|_\infty}{\|\mathbf{x}\|} \leq c_4 \|\mathbf{x}\|$ . In other words, if the iterate jumps out of Regime 2, it will still fall within Regime 1.



To summarize, suppose the initial guess  $\mathbf{z}^{(0)}$  obeys  $\text{dist}(\mathbf{z}^{(0)}, \mathbf{x}) \leq c_4 \|\mathbf{x}\|$ . Then the estimation error will shrink at a geometric rate  $1 - \rho$  before it enters Regime 2. Afterwards,  $\mathbf{z}^{(t)}$  will either stay within Regime 2 or jump back and forth between Regimes 1 and 2. Because of the bounds (90) and (88), the estimation errors will never exceed the order of  $\frac{\|\boldsymbol{\eta}\|}{\sqrt{m}\|\mathbf{x}\|}$  from then on. Putting these together establishes (16), namely, the first part of the theorem.

Below we justify the condition (86) for Regime 1, for which we start by gathering additional properties of the trimming rules. By Cauchy-Schwarz,  $\frac{1}{m} \|\boldsymbol{\eta}\|_1 \leq \frac{1}{\sqrt{m}} \|\boldsymbol{\eta}\| \leq \frac{1}{c_3} \|\mathbf{h}\| \|\mathbf{z}\|$ . When  $c_3$  is sufficiently large, applying Lemmas 1 and 2 gives

$$\begin{aligned} \frac{1}{m} \sum_{l=1}^m \left| y_l - |\mathbf{a}_l^\top \mathbf{z}|^2 \right| &\leq \frac{1}{m} \|\mathcal{A}(\mathbf{x}\mathbf{x}^\top - \mathbf{z}\mathbf{z}^\top)\|_1 + \frac{1}{m} \|\boldsymbol{\eta}\|_1 \leq 2.98 \|\mathbf{h}\| \|\mathbf{z}\|; \\ \frac{1}{m} \sum_{l=1}^m \left| y_l - |\mathbf{a}_l^\top \mathbf{z}|^2 \right| &\geq \frac{1}{m} \|\mathcal{A}(\mathbf{x}\mathbf{x}^\top - \mathbf{z}\mathbf{z}^\top)\|_1 - \frac{1}{m} \|\boldsymbol{\eta}\|_1 \geq 1.151 \|\mathbf{h}\| \|\mathbf{z}\|. \end{aligned} \quad (91)$$

From now on, we shall denote  $\tilde{\mathcal{E}}_2^i := \left\{ \left| |\mathbf{a}_i^\top \mathbf{x}|^2 - |\mathbf{a}_i^\top \mathbf{z}|^2 \right| \leq \frac{\alpha_h}{m} \|\mathbf{y} - \mathcal{A}(\mathbf{z}\mathbf{z}^\top)\|_1 \frac{|\mathbf{a}_i^\top \mathbf{z}|}{\|\mathbf{z}\|} \right\}$  to differentiate from  $\mathcal{E}_2^i$ . For any small constant  $\epsilon > 0$ , we introduce the index set  $\mathcal{G} := \{i : |\eta_i| \leq C_\epsilon \|\boldsymbol{\eta}\| / \sqrt{m}\}$  that satisfies  $|\mathcal{G}| = (1 - \epsilon)m$ . Note that  $C_\epsilon$  must be bounded as  $n$  scales, since

$$\|\boldsymbol{\eta}\|^2 \geq \sum_{i \notin \mathcal{G}} \eta_i^2 \geq (m - |\mathcal{G}|) \cdot C_\epsilon^2 \|\boldsymbol{\eta}\|^2 / m \geq \epsilon C_\epsilon^2 \|\boldsymbol{\eta}\|^2 \Rightarrow C_\epsilon \leq 1/\sqrt{\epsilon}. \quad (92)$$

We are now ready to analyze the regularized gradient, which we separate into several components as follows

$$\begin{aligned} \nabla_{\text{tr}} \ell(\mathbf{z}) &= \underbrace{2 \sum_{i \in \mathcal{G}} \frac{|\mathbf{a}_i^\top \mathbf{x}|^2 - |\mathbf{a}_i^\top \mathbf{z}|^2}{\mathbf{a}_i^\top \mathbf{z}} \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i} + 2 \sum_{i \notin \mathcal{G}} \frac{|\mathbf{a}_i^\top \mathbf{x}|^2 - |\mathbf{a}_i^\top \mathbf{z}|^2}{\mathbf{a}_i^\top \mathbf{z}} \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i \cap \tilde{\mathcal{E}}_2^i}}_{:= \nabla_{\text{tr}}^{\text{clean}} \ell(\mathbf{z})} \\ &+ \underbrace{2 \sum_{i \in \mathcal{G}} \frac{\eta_i}{\mathbf{a}_i^\top \mathbf{z}} \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i}}_{:= \nabla_{\text{tr}}^{\text{noise}} \ell(\mathbf{z})} + \underbrace{2 \sum_{i \notin \mathcal{G}} \left( \frac{y_i - |\mathbf{a}_i^\top \mathbf{z}|^2}{\mathbf{a}_i^\top \mathbf{z}} \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i} - \frac{|\mathbf{a}_i^\top \mathbf{x}|^2 - |\mathbf{a}_i^\top \mathbf{z}|^2}{\mathbf{a}_i^\top \mathbf{z}} \mathbf{1}_{\mathcal{E}_1^i \cap \tilde{\mathcal{E}}_2^i} \right) \mathbf{a}_i}_{:= \nabla_{\text{tr}}^{\text{extra}} \ell(\mathbf{z})}. \end{aligned} \quad (93)$$

- For each index  $i \in \mathcal{G}$ , the inclusion property (51) (i.e.  $\mathcal{E}_3^i \subseteq \mathcal{E}_2^i \subseteq \mathcal{E}_4^i$ ) holds. To see this, observe that

$$|y_i - |\mathbf{a}_i^\top \mathbf{z}|^2| \in \left[ \left| |\mathbf{a}_i^\top \mathbf{x}|^2 - |\mathbf{a}_i^\top \mathbf{z}|^2 \right| \pm |\eta_i| \right].$$

Since  $|\eta_i| \leq C_\epsilon \|\boldsymbol{\eta}\| / \sqrt{m} \ll \|\mathbf{h}\| \|\mathbf{z}\|$  when  $c_3$  is sufficiently large, one can derive the inclusion (51) immediately from (91). As a result, all the proof arguments for Proposition 2 carry over to  $\nabla_{\text{tr}}^{\text{clean}} \ell(\mathbf{z})$ , suggesting that

$$-\left\langle \mathbf{h}, \frac{1}{m} \nabla_{\text{tr}}^{\text{clean}} \ell(\mathbf{z}) \right\rangle \geq 2 \left\{ 1.99 - 2(\zeta_1 + \zeta_2) - \sqrt{8/(9\pi)} \alpha_h^{-1} - \epsilon \right\} \|\mathbf{h}\|^2. \quad (94)$$

- Next, letting  $w_i = \frac{2\eta_i}{\mathbf{a}_i^\top \mathbf{z}} \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i} \mathbf{1}_{\{i \in \mathcal{G}\}}$ , we see that for any constant  $\delta > 0$ , the noise component obeys

$$\begin{aligned} \left\| \frac{1}{m} \nabla_{\text{tr}}^{\text{noise}} \ell(\mathbf{z}) \right\| &= \left\| \frac{1}{m} \mathbf{A}^\top \mathbf{w} \right\| \leq \left\| \frac{1}{\sqrt{m}} \mathbf{A} \right\| \left\| \frac{1}{\sqrt{m}} \mathbf{w} \right\| \\ &\stackrel{\text{(ii)}}{\leq} \frac{1 + \delta}{\sqrt{m}} \|\mathbf{w}\| \leq (1 + \delta) \frac{2\|\boldsymbol{\eta}\|/\sqrt{m}}{\alpha_z^{\text{lb}} \|\mathbf{z}\|}, \end{aligned} \quad (95)$$

provided that  $m/n$  is sufficiently large. Here, (ii) arises from [40, Corollary 5.35], and the last inequality is a consequence of the upper estimate

$$\|\mathbf{w}\|^2 \leq 4 \sum_{i=1}^m \frac{|\eta_i|^2}{(\mathbf{a}_i^\top \mathbf{z})^2} \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i} \leq 4 \sum_{i=1}^m \frac{|\eta_i|^2}{(\alpha_z^{\text{lb}} \|\mathbf{z}\|)^2} = \frac{4\|\boldsymbol{\eta}\|^2}{(\alpha_z^{\text{lb}} \|\mathbf{z}\|)^2}. \quad (96)$$

In turn, this immediately gives

$$\left| \left\langle \mathbf{h}, \frac{1}{m} \nabla_{\text{tr}}^{\text{noise}} \ell(\mathbf{z}) \right\rangle \right| \leq \|\mathbf{h}\| \left\| \frac{1}{m} \nabla_{\text{tr}}^{\text{noise}} \ell(\mathbf{z}) \right\| \leq \frac{2(1+\delta)}{\alpha_z^{\text{lb}}} \frac{\|\boldsymbol{\eta}\|}{\sqrt{m}\|\mathbf{z}\|} \|\mathbf{h}\|. \quad (97)$$

- We now turn to the last term  $\nabla_{\text{tr}}^{\text{extra}} \ell(\mathbf{z})$ . According to the definition of  $\mathcal{E}_2^i$  and  $\tilde{\mathcal{E}}_2^i$  as well as the property (91), the weight  $q_i := 2 \left( \frac{y_i - |\mathbf{a}_i^\top \mathbf{z}|^2}{\alpha_i^\top \mathbf{z}} \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i} - \frac{|\mathbf{a}_i^\top \mathbf{x}|^2 - |\mathbf{a}_i^\top \mathbf{z}|^2}{\alpha_i^\top \mathbf{z}} \mathbf{1}_{\tilde{\mathcal{E}}_1^i \cap \tilde{\mathcal{E}}_2^i} \right) \mathbf{1}_{\{i \notin \mathcal{G}\}}$  is bounded in magnitude by  $6\|\mathbf{h}\|$ . This gives

$$\|\mathbf{q}\| \leq \sqrt{m - |\mathcal{G}|} \cdot 6\|\mathbf{h}\| \leq 6\sqrt{\epsilon m} \|\mathbf{h}\|,$$

and hence

$$\left| \left\langle \frac{1}{m} \nabla_{\text{tr}}^{\text{extra}} \ell(\mathbf{z}), \mathbf{h} \right\rangle \right| \leq \|\mathbf{h}\| \cdot \left\| \frac{1}{m} \nabla_{\text{tr}}^{\text{extra}} \ell(\mathbf{z}) \right\| = \frac{1}{m} \|\mathbf{h}\| \cdot \|\mathbf{A}^\top \mathbf{q}\| \leq 6(1+\delta) \sqrt{\epsilon} \|\mathbf{h}\|^2. \quad (98)$$

Taking the above bounds together yields

$$\begin{aligned} -\frac{1}{m} \langle \nabla \ell_{\text{tr}}(\mathbf{z}), \mathbf{h} \rangle &\geq 2 \left\{ 1.99 - 2(\zeta_1 + \zeta_2) - \sqrt{\frac{8}{9\pi}} \frac{1}{\alpha_h} - 6(1+\delta) \sqrt{\epsilon} - \epsilon \right\} \|\mathbf{h}\|^2 \\ &\quad - \frac{2(1+\delta)}{\alpha_z^{\text{lb}}} \frac{\|\boldsymbol{\eta}\|}{\sqrt{m}\|\mathbf{z}\|} \|\mathbf{h}\|. \end{aligned}$$

Since  $\|\mathbf{h}\| \geq c_3 \frac{\|\boldsymbol{\eta}\|}{\sqrt{m}\|\mathbf{z}\|}$  for some large constant  $c_3 > 0$ , setting  $\epsilon$  to be small one obtains

$$-\frac{1}{m} \langle \nabla \ell_{\text{tr}}(\mathbf{z}), \mathbf{h} \rangle \geq 2 \left\{ 1.95 - 2(\zeta_1 + \zeta_2) - \sqrt{8/(9\pi)} \alpha_h^{-1} \right\} \|\mathbf{h}\|^2 \quad (99)$$

for all  $\mathbf{h}$  obeying

$$\frac{c_3 \|\boldsymbol{\eta}\| / \sqrt{m}}{\|\mathbf{z}\|} \leq \|\mathbf{h}\| \leq \min \left\{ \frac{1}{11}, \frac{\alpha_z^{\text{lb}}}{3\alpha_h}, \frac{\alpha_z^{\text{lb}}}{6}, \frac{\sqrt{98/3} (\alpha_z^{\text{lb}})^2}{2\alpha_z^{\text{ub}} + \alpha_z^{\text{lb}}} \right\} \|\mathbf{z}\|,$$

which finishes the proof of Theorem 2 for general  $\boldsymbol{\eta}$ .

Up until now, we have established the theorem for general  $\boldsymbol{\eta}$ , and it remains to specialize it to the Poisson model. Standard concentration results, which we omit, give

$$\frac{1}{m} \|\boldsymbol{\eta}\|^2 \approx \frac{1}{m} \sum_{i=1}^m \mathbb{E} [\eta_i^2] = \frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{x})^2 \approx \|\mathbf{x}\|^2 \quad (100)$$

with high probability. Substitution into (16) completes the proof.

## 7 Minimax lower bound

The goal of this section is to establish the minimax lower bound given in Theorem 3. For notational simplicity, we denote by  $\mathbb{P}(\mathbf{y} | \mathbf{w})$  the likelihood of  $y_i \stackrel{\text{ind.}}{\sim} \text{Poisson}(|\mathbf{a}_i^\top \mathbf{w}|^2)$ ,  $1 \leq i \leq m$  conditional on  $\{\mathbf{a}_i\}$ . For any two probability measures  $P$  and  $Q$ , we denote by  $\text{KL}(P||Q)$  the Kullback–Leibler (KL) divergence between them:

$$\text{KL}(P||Q) := \int \log \left( \frac{dP}{dQ} \right) dP, \quad (101)$$

The basic idea is to adopt the general reduction scheme discussed in [41, Section 2.2], which amounts to finding a finite collection of hypotheses that are minimally separated. Below we gather one result useful for constructing and analyzing such hypotheses.

**Lemma 8.** *Suppose that  $\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ ,  $n$  is sufficiently large, and  $m = \kappa n$  for some sufficiently large constant  $\kappa > 0$ . Consider any  $\mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$ . On an event  $\mathcal{B}$  of probability approaching one, there exists a collection  $\mathcal{M}$  of  $M = \exp(n/30)$  distinct vectors obeying the following properties:*

(i)  $\mathbf{x} \in \mathcal{M}$ ;

(ii) for all  $\mathbf{w}^{(l)}, \mathbf{w}^{(j)} \in \mathcal{M}$ ,

$$1/\sqrt{8} - (2n)^{-1/2} \leq \|\mathbf{w}^{(l)} - \mathbf{w}^{(j)}\| \leq 3/2 + n^{-1/2}; \quad (102)$$

(iii) for all  $\mathbf{w} \in \mathcal{M}$ ,

$$\frac{|\mathbf{a}_i^\top (\mathbf{w} - \mathbf{x})|^2}{|\mathbf{a}_i^\top \mathbf{x}|^2} \leq \frac{\|\mathbf{w} - \mathbf{x}\|^2}{\|\mathbf{x}\|^2} \{2 + 17 \log^3 m\}, \quad 1 \leq i \leq m. \quad (103)$$

In words, Lemma 8 constructs a set  $\mathcal{M}$  of exponentially many vectors/hypotheses scattered around  $\mathbf{x}$  and yet well separated. From (ii) we see that each pair of hypotheses in  $\mathcal{M}$  is separated by a distance roughly on the order of 1, and all hypotheses reside within a spherical ball centered at  $\mathbf{x}$  of radius  $3/2 + o(1)$ . When  $\|\mathbf{x}\| \geq \log^{1.5} m$ , every hypothesis  $\mathbf{w} \in \mathcal{M}$  satisfies  $\|\mathbf{w}\| \approx \|\mathbf{x}\| \gg 1$ . In addition, (iii) says that the quantities  $|\mathbf{a}_i^\top (\mathbf{w} - \mathbf{x})|/|\mathbf{a}_i^\top \mathbf{x}|$  are all very well controlled (modulo some logarithmic factor). In particular, when  $\|\mathbf{x}\| \geq \log^{1.5} m$ , one must have

$$\frac{|\mathbf{a}_i^\top (\mathbf{w} - \mathbf{x})|^2}{|\mathbf{a}_i^\top \mathbf{x}|^2} \lesssim \frac{\|\mathbf{w} - \mathbf{x}\|^2}{\|\mathbf{x}\|^2} \log^3 m \lesssim \frac{1}{\log^3 m} \log^3 m \lesssim 1. \quad (104)$$

In the Poisson model, such a quantity turns out to be crucial in controlling the information divergence between two hypotheses, as demonstrated in the following lemma.

**Lemma 9.** *Fix a family of design vectors  $\{\mathbf{a}_i\}$ . Then for any  $\mathbf{w}$  and  $\mathbf{r} \in \mathbb{R}^n$ ,*

$$\text{KL}(\mathbb{P}(\mathbf{y} | \mathbf{w} + \mathbf{r}) \parallel \mathbb{P}(\mathbf{y} | \mathbf{w})) \leq \sum_{i=1}^m |\mathbf{a}_i^\top \mathbf{r}|^2 \left(8 + \frac{2|\mathbf{a}_i^\top \mathbf{r}|^2}{|\mathbf{a}_i^\top \mathbf{w}|^2}\right). \quad (105)$$

Lemma 9 and (104) taken collectively suggest that on the event  $\mathcal{B} \cap \mathcal{C}$  ( $\mathcal{B}$  is in Lemma 8 and  $\mathcal{C} := \{\|\mathbf{A}\| \leq \sqrt{2m}\}$ ), the conditional KL divergence (we condition on the  $\mathbf{a}_i$ 's) obeys

$$\text{KL}(\mathbb{P}(\mathbf{y} | \mathbf{w}) \parallel \mathbb{P}(\mathbf{y} | \mathbf{x})) \leq c_3 \sum_{i=1}^m |\mathbf{a}_i^\top (\mathbf{w} - \mathbf{x})|^2 \leq 2c_3 m \|\mathbf{w} - \mathbf{x}\|^2, \quad \forall \mathbf{w} \in \mathcal{M}; \quad (106)$$

here, the inequality holds for some constant  $c_3 > 0$  provided that  $\|\mathbf{x}\| \geq \log^{1.5} m$ , and the last inequality is a result of  $\mathcal{C}$  (which occurs with high probability). We now use hypotheses as in Lemma 8 but rescaled in such a way that

$$\|\mathbf{w} - \mathbf{x}\| \asymp \delta, \quad \text{and} \quad \|\mathbf{w} - \tilde{\mathbf{w}}\| \asymp \delta, \quad \forall \mathbf{w}, \tilde{\mathbf{w}} \in \mathcal{M} \text{ with } \mathbf{w} \neq \tilde{\mathbf{w}}. \quad (107)$$

for some  $0 < \delta < 1$ . This is achieved via the substitution  $\mathbf{w} \leftarrow \mathbf{x} + \delta(\mathbf{w} - \mathbf{x})$ ; with a slight abuse of notation,  $\mathcal{M}$  denotes the new set.

The hardness of a minimax estimation problem is known to be dictated by information divergence inequalities such as (106). Indeed, suppose that

$$\frac{1}{M-1} \sum_{\mathbf{w} \in \mathcal{M} \setminus \{\mathbf{x}\}} \text{KL}(\mathbb{P}(\mathbf{y} | \mathbf{w}) \parallel \mathbb{P}(\mathbf{y} | \mathbf{x})) \leq \frac{1}{10} \log(M-1) \quad (108)$$

holds, then the Fano-type minimax lower bound [41, Theorem 2.7] asserts that

$$\inf_{\hat{\mathbf{x}}} \sup_{\mathbf{x} \in \mathcal{M}} \mathbb{E}[\|\hat{\mathbf{x}} - \mathbf{x}\| \mid \{\mathbf{a}_i\}] \gtrsim \min_{\mathbf{w}, \tilde{\mathbf{w}} \in \mathcal{M}, \mathbf{w} \neq \tilde{\mathbf{w}}} \|\mathbf{w} - \tilde{\mathbf{w}}\|. \quad (109)$$

Since  $M = \exp(n/30)$ , (108) would follow from

$$2c_3 \|\mathbf{w} - \mathbf{x}\|^2 \leq n/(300m). \quad \mathbf{w} \in \mathcal{M}. \quad (110)$$

Hence, we just need to select  $\delta$  to be a small multiple of  $\sqrt{n/m}$ . This in turn gives

$$\inf_{\hat{\mathbf{x}}} \sup_{\mathbf{x} \in \mathcal{M}} \mathbb{E} [\|\hat{\mathbf{x}} - \mathbf{x}\| \mid \{\mathbf{a}_i\}] \gtrsim \min_{\mathbf{w}, \tilde{\mathbf{w}} \in \mathcal{M}, \mathbf{w} \neq \tilde{\mathbf{w}}} \|\mathbf{w} - \tilde{\mathbf{w}}\| \gtrsim \sqrt{n/m}. \quad (111)$$

Finally, it remains to connect  $\|\hat{\mathbf{x}} - \mathbf{x}\|$  with  $\text{dist}(\hat{\mathbf{x}}, \mathbf{x})$ . Since all the  $\mathbf{w} \in \mathcal{M}$  are clustered around  $\mathbf{x}$  and are at a mutual distance about  $\delta$  that is much smaller than  $\|\mathbf{x}\|$ , we can see that for any reasonable estimator,  $\text{dist}(\hat{\mathbf{x}}, \mathbf{x}) = \|\hat{\mathbf{x}} - \mathbf{x}\|$ . This finishes the proof.

## 8 Discussion

To keep our treatment concise, this paper does not strive to explore all possible generalizations of the theory. There are nevertheless a few extensions worth pointing out.

- **More general objective functions.** For concreteness, we restrict our analysis to the Poisson log-likelihood function, but the analysis framework we laid out easily carries over to a broad class of (nonconvex) objective functions. For instance, all results continue to hold if we replace the Poisson log-likelihood by the Gaussian log-likelihood; that is, the polynomial function  $-\sum_{i=1}^m (y_i - |\mathbf{a}_i^\top \mathbf{z}|^2)^2$  studied in [1]. A general guideline is to first check whether the expected regularity condition

$$\mathbb{E} \left[ -\left\langle \frac{1}{m} \nabla \ell_{\text{tr}}(\mathbf{z}), \mathbf{h} \right\rangle \right] \gtrsim \|\mathbf{h}\|^2$$

holds for any fixed  $\mathbf{z}$  within a neighborhood around  $\mathbf{x}$ . If so, then often times RC holds uniformly within this neighborhood due to sharp concentration of measure ensured by the regularization procedure.

- **Sub-Gaussian measurements.** The theory extends to the situation where the  $\mathbf{a}_i$ 's are i.i.d. sub-Gaussian random vectors, although the truncation threshold might need to be tweaked based on the sub-Gaussian norm of  $\mathbf{a}_i$ . A more challenging scenario, however, is the case where the  $\mathbf{a}_i$ 's are generated according to the CDP model, since there is much less randomness to exploit in the mathematical analysis. We leave this to future research.

Having demonstrated the power of TWF in recovering a rank-one matrix  $\mathbf{x}\mathbf{x}^*$  from quadratic equations, we remark on the potential of TWF towards recovering low-rank matrices from rank-one measurements. Imagine that we wish to estimate a rank- $r$  matrix  $\mathbf{X} \succeq \mathbf{0}$  and that all we know about  $\mathbf{X}$  is

$$y_i = \mathbf{a}_i^\top \mathbf{X} \mathbf{a}_i, \quad 1 \leq i \leq m.$$

It is known that this problem can be efficiently solved by using more computational-intensive semidefinite programs [9, 14]. With the hope of developing a linear-time algorithm, one might consider a modified TWF scheme, which would maintain a rank- $r$  matrix variable and operate as follows: perform truncated spectral initialization, and then successively update the current guess via a regularized gradient descent rule applied to a presumed log-likelihood function.

Moving away from i.i.d. sub-Gaussian measurements, there is a proliferation of problems that involve completion of a low-rank matrix  $\mathbf{X}$  from partial entries, where the rank is known *a priori*. It is self-evident that such entry-wise observations can also be cast as rank-one measurements of  $\mathbf{X}$ . Therefore, the preceding modified TWF may add to recent literature in applying non-convex schemes for low-rank matrix completion [42–45], robust PCA [46], or even a broader family of latent-variable models (e.g. dictionary learning [47, 48], sparse coding [49], and mixture problems [50, 51]). A concrete application of this flavor is a simple form of the fundamental alignment/matching problem [52–54]. Imagine a collection of  $n$  instances, each representing an image of the same physical object but with different shift  $r_i \in \{0, \dots, M-1\}$ . The goal is to align all these instances from observations on the relative shift between pairs of them. Denoting by  $\mathbf{X}_i$  the cyclic shift by an amount  $r_i$  of  $\mathbf{I}_M$ , one sees that the collection matrix  $\mathbf{X} := [\mathbf{X}_i^\top \mathbf{X}_j]_{1 \leq i, j \leq k}$  is a rank- $M$  matrix, and the relative shift observations can be treated as rank-one measurements of  $\mathbf{X}$ . Running TWF over this problem instance might result in a statistically and computationally efficient solution. This would be of great practical interest.

## A Proofs for Section 5

### A.1 Proof of Lemma 3

First, we make the observation that  $(\mathbf{a}_i^\top \mathbf{z})^2 - (\mathbf{a}_i^\top \mathbf{x})^2 = (2\mathbf{a}_i^\top \mathbf{z} - \mathbf{a}_i^\top \mathbf{h}) \mathbf{a}_i^\top \mathbf{h}$  is a quadratic function in  $\mathbf{a}_i^\top \mathbf{h}$ . If we assume  $\gamma \leq \frac{\alpha_z^{\text{lb}} \|\mathbf{z}\|}{\|\mathbf{h}\|}$ , then on the event  $\mathcal{E}_1^i$  one has

$$(\mathbf{a}_i^\top \mathbf{z})^2 \geq \alpha_z^{\text{lb}} \|\mathbf{z}\| \cdot |\mathbf{a}_i^\top \mathbf{z}| \geq \gamma \|\mathbf{h}\| |\mathbf{a}_i^\top \mathbf{z}|. \quad (112)$$

Solving the quadratic inequality that specifies  $\mathcal{D}_\gamma^i$  gives

$$\begin{aligned} \mathbf{a}_i^\top \mathbf{h} &\in \left[ \mathbf{a}_i^\top \mathbf{z} - \sqrt{(\mathbf{a}_i^\top \mathbf{z})^2 + \gamma \|\mathbf{h}\| |\mathbf{a}_i^\top \mathbf{z}|}, \mathbf{a}_i^\top \mathbf{z} + \sqrt{(\mathbf{a}_i^\top \mathbf{z})^2 - \gamma \|\mathbf{h}\| |\mathbf{a}_i^\top \mathbf{z}|} \right], \\ \text{or } \mathbf{a}_i^\top \mathbf{h} &\in \left[ \mathbf{a}_i^\top \mathbf{z} + \sqrt{(\mathbf{a}_i^\top \mathbf{z})^2 - \gamma \|\mathbf{h}\| |\mathbf{a}_i^\top \mathbf{z}|}, \mathbf{a}_i^\top \mathbf{z} - \sqrt{(\mathbf{a}_i^\top \mathbf{z})^2 + \gamma \|\mathbf{h}\| |\mathbf{a}_i^\top \mathbf{z}|} \right], \end{aligned}$$

which we will simplify in the sequel.

Suppose for the moment that  $\mathbf{a}_i^\top \mathbf{z} \geq 0$ , then the preceding two intervals are respectively equivalent to

$$\begin{aligned} \mathbf{a}_i^\top \mathbf{h} &\in \underbrace{\left[ \frac{-\gamma \|\mathbf{h}\| |\mathbf{a}_i^\top \mathbf{z}|}{\mathbf{a}_i^\top \mathbf{z} + \sqrt{(\mathbf{a}_i^\top \mathbf{z})^2 + \gamma \|\mathbf{h}\| |\mathbf{a}_i^\top \mathbf{z}|}}, \frac{\gamma \|\mathbf{h}\| |\mathbf{a}_i^\top \mathbf{z}|}{\mathbf{a}_i^\top \mathbf{z} + \sqrt{(\mathbf{a}_i^\top \mathbf{z})^2 - \gamma \|\mathbf{h}\| |\mathbf{a}_i^\top \mathbf{z}|}} \right]}_{:=I_1}; \\ \mathbf{a}_i^\top \mathbf{h} - 2\mathbf{a}_i^\top \mathbf{z} &\in \underbrace{\left[ \frac{-\gamma \|\mathbf{h}\| |\mathbf{a}_i^\top \mathbf{z}|}{\mathbf{a}_i^\top \mathbf{z} + \sqrt{(\mathbf{a}_i^\top \mathbf{z})^2 - \gamma \|\mathbf{h}\| |\mathbf{a}_i^\top \mathbf{z}|}}, \frac{\gamma \|\mathbf{h}\| |\mathbf{a}_i^\top \mathbf{z}|}{\mathbf{a}_i^\top \mathbf{z} + \sqrt{(\mathbf{a}_i^\top \mathbf{z})^2 + \gamma \|\mathbf{h}\| |\mathbf{a}_i^\top \mathbf{z}|}} \right]}_{:=I_2}. \end{aligned}$$

Assuming (112) and making use of the observations

$$\begin{aligned} \frac{\gamma \|\mathbf{h}\| |\mathbf{a}_i^\top \mathbf{z}|}{\mathbf{a}_i^\top \mathbf{z} + \sqrt{(\mathbf{a}_i^\top \mathbf{z})^2 - \gamma \|\mathbf{h}\| |\mathbf{a}_i^\top \mathbf{z}|}} &\leq \frac{\gamma \|\mathbf{h}\| |\mathbf{a}_i^\top \mathbf{z}|}{\mathbf{a}_i^\top \mathbf{z}} = \gamma \|\mathbf{h}\| \\ \text{and } \frac{\gamma \|\mathbf{h}\| |\mathbf{a}_i^\top \mathbf{z}|}{\mathbf{a}_i^\top \mathbf{z} + \sqrt{(\mathbf{a}_i^\top \mathbf{z})^2 + \gamma \|\mathbf{h}\| |\mathbf{a}_i^\top \mathbf{z}|}} &\geq \frac{\gamma \|\mathbf{h}\| |\mathbf{a}_i^\top \mathbf{z}|}{(1 + \sqrt{2}) |\mathbf{a}_i^\top \mathbf{z}|} = \frac{\gamma}{1 + \sqrt{2}} \|\mathbf{h}\|, \end{aligned}$$

we obtain the inner and outer bounds

$$\left[ \pm(1 + \sqrt{2})^{-1} \gamma \|\mathbf{h}\| \right] \subseteq I_1, I_2 \subseteq [\pm \gamma \|\mathbf{h}\|].$$

Setting  $\gamma_1 := \frac{\gamma}{1 + \sqrt{2}}$  gives

$$(\mathcal{D}_{\gamma_1}^{i,1} \cap \mathcal{E}_{i,1}) \cup (\mathcal{D}_{\gamma_1}^{i,2} \cap \mathcal{E}_{i,1}) \subseteq \mathcal{D}_\gamma \cap \mathcal{E}_{i,1} \subseteq (\mathcal{D}_\gamma^{i,1} \cap \mathcal{E}_{i,1}) \cup (\mathcal{D}_\gamma^{i,2} \cap \mathcal{E}_{i,1}).$$

Proceeding with the same argument, we can derive exactly the same inner and outer bounds in the regime where  $\mathbf{a}_i^\top \mathbf{z} < 0$ , concluding the proof.

### A.2 Proof of Lemma 4

By homogeneity, it suffices to establish the claim for the case where both  $\mathbf{h}$  and  $\mathbf{z}$  are *unit vectors*.

Suppose for the moment that  $\mathbf{h}$  and  $\mathbf{z}$  are *statistically independent* from  $\{\mathbf{a}_i\}$ . We introduce two auxiliary Lipschitz functions approximating indicator functions:

$$\chi_z(\tau) := \begin{cases} 1, & \text{if } |\tau| \in [\sqrt{1.01}\alpha_z^{\text{lb}}, \sqrt{0.99}\alpha_z^{\text{ub}}]; \\ -100(\alpha_z^{\text{ub}})^{-2}\tau^2 + 100, & \text{if } |\tau| \in [\sqrt{0.99}\alpha_z^{\text{ub}}, \alpha_z^{\text{ub}}]; \\ 100(\alpha_z^{\text{lb}})^{-2}\tau^2 - 100, & \text{if } |\tau| \in [\alpha_z^{\text{lb}}, \sqrt{1.01}\alpha_z^{\text{lb}}]; \\ 0, & \text{else.} \end{cases} \quad (113)$$

$$\chi_h(\tau) := \begin{cases} 1, & \text{if } |\tau| \in [0, \sqrt{0.99}\gamma]; \\ -\frac{100}{\gamma^2}\tau^2 + 100, & \text{if } |\tau| \in [\sqrt{0.99}\gamma, \gamma]; \\ 0, & \text{else.} \end{cases} \quad (114)$$

Since  $\mathbf{h}$  and  $\mathbf{z}$  are assumed to be unit vectors, these two functions obey

$$0 \leq \chi_z(\mathbf{a}_i^\top \mathbf{z}) \leq \mathbf{1}_{\mathcal{E}_1^i}, \quad \text{and} \quad 0 \leq \chi_h(\mathbf{a}_i^\top \mathbf{h}) \leq \mathbf{1}_{\mathcal{D}_\gamma^{i,1}} \quad (115)$$

and thus,

$$\frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h})^2 \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{D}_\gamma^{i,1}} \geq \frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h})^2 \chi_z(\mathbf{a}_i^\top \mathbf{z}) \chi_h(\mathbf{a}_i^\top \mathbf{h}). \quad (116)$$

We proceed to lower bound  $\frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h})^2 \chi_z(\mathbf{a}_i^\top \mathbf{z}) \chi_h(\mathbf{a}_i^\top \mathbf{h})$ .

Firstly, to compute the mean of  $(\mathbf{a}_i^\top \mathbf{h})^2 \chi_z(\mathbf{a}_i^\top \mathbf{z}) \chi_h(\mathbf{a}_i^\top \mathbf{h})$ , we introduce an auxiliary orthonormal matrix

$$\mathbf{U}_z = \begin{bmatrix} \mathbf{z}^\top / \|\mathbf{z}\| \\ \vdots \end{bmatrix} \quad (117)$$

whose first row is along the direction of  $\mathbf{z}$ , and set

$$\tilde{\mathbf{h}} := \mathbf{U}_z \mathbf{h}, \quad \text{and} \quad \tilde{\mathbf{a}}_i := \mathbf{U}_z \mathbf{a}_i. \quad (118)$$

Also, denote by  $\tilde{a}_{i,1}$  (resp.  $\tilde{h}_1$ ) the first entry of  $\tilde{\mathbf{a}}_i$  (resp.  $\tilde{\mathbf{h}}$ ), and  $\tilde{\mathbf{a}}_{i,\setminus 1}$  (resp.  $\tilde{\mathbf{h}}_{\setminus 1}$ ) the remaining entries of  $\tilde{\mathbf{a}}_i$  (resp.  $\tilde{\mathbf{h}}$ ), and let  $\xi \sim \mathcal{N}(0, 1)$ . We have

$$\begin{aligned} & \mathbb{E} \left[ (\mathbf{a}_i^\top \mathbf{h})^2 \chi_z(\mathbf{a}_i^\top \mathbf{z}) \chi_h(\mathbf{a}_i^\top \mathbf{h}) \right] \\ & \geq \mathbb{E} \left[ (\mathbf{a}_i^\top \mathbf{h})^2 \chi_z(\mathbf{a}_i^\top \mathbf{z}) \right] - \mathbb{E} \left[ (\mathbf{a}_i^\top \mathbf{h})^2 (1 - \chi_h(\mathbf{a}_i^\top \mathbf{h})) \right] \\ & \geq \mathbb{E} \left[ (\tilde{a}_{i,1} \tilde{h}_1)^2 \chi_z(\mathbf{a}_i^\top \mathbf{z}) \right] + \mathbb{E} \left[ (\tilde{\mathbf{a}}_{i,\setminus 1}^\top \tilde{\mathbf{h}}_{\setminus 1})^2 \right] \mathbb{E} \left[ \chi_z(\mathbf{a}_i^\top \mathbf{z}) \right] - \|\mathbf{h}\|^2 \mathbb{E} \left[ \xi^2 \mathbf{1}_{\{|\xi| > \sqrt{0.99}\gamma\}} \right] \\ & \geq |\tilde{h}_1|^2 (1 - \zeta_1) + \|\tilde{\mathbf{h}}_{\setminus 1}\|^2 (1 - \zeta_1) - \zeta_2 \|\mathbf{h}\|^2 \\ & \geq (1 - \zeta_1 - \zeta_2) \|\mathbf{h}\|^2, \end{aligned} \quad (119)$$

where the identity (119) arises from (66) and (67). Since  $(\mathbf{a}_i^\top \mathbf{h})^2 \chi_z(\mathbf{a}_i^\top \mathbf{z}) \chi_h(\mathbf{a}_i^\top \mathbf{h})$  is bounded in magnitude by  $\gamma^2 \|\mathbf{h}\|^2$ , it is a sub-Gaussian random variable with sub-Gaussian norm  $O(\gamma^2 \|\mathbf{h}\|^2)$ . Apply the Hoeffding-type inequality [40, Proposition 5.10] to deduce that for any  $\epsilon > 0$ ,

$$\frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h})^2 \chi_z(\mathbf{a}_i^\top \mathbf{z}) \chi_h(\mathbf{a}_i^\top \mathbf{h}) \geq \mathbb{E} \left[ (\mathbf{a}_i^\top \mathbf{h})^2 \chi_z(\mathbf{a}_i^\top \mathbf{z}) \chi_h(\mathbf{a}_i^\top \mathbf{h}) \right] - \epsilon \|\mathbf{h}\|^2 \quad (120)$$

$$\geq (1 - \zeta_1 - \zeta_2 - \epsilon) \|\mathbf{h}\|^2 \quad (121)$$

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

The next step is to obtain uniform control over all *unit vectors*, for which we adopt a basic version of an  $\epsilon$ -net argument. Specifically, we construct an  $\epsilon$ -net  $\mathcal{N}_\epsilon$  with cardinality  $|\mathcal{N}_\epsilon| \leq (1 + 2/\epsilon)^{2n}$  (cf. [40]) such that

for any  $(\mathbf{h}, \mathbf{z})$  with  $\|\mathbf{h}\| = \|\mathbf{z}\| = 1$ , there exists a pair  $\mathbf{h}_0, \mathbf{z}_0 \in \mathcal{N}_\epsilon$  satisfying  $\|\mathbf{h} - \mathbf{h}_0\| \leq \epsilon$  and  $\|\mathbf{z} - \mathbf{z}_0\| \leq \epsilon$ . Now that we have discretized the unit spheres using a finite set, taking the union bound gives

$$\frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h}_0)^2 \chi_z(\mathbf{a}_i^\top \mathbf{z}_0) \chi_h(\mathbf{a}_i^\top \mathbf{h}_0) \geq (1 - \zeta_1 - \zeta_2 - \epsilon) \|\mathbf{h}_0\|^2, \quad \forall \mathbf{h}_0, \mathbf{z}_0 \in \mathcal{N}_\epsilon \quad (122)$$

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

Define  $f_1(\cdot)$  and  $f_2(\cdot)$  such that  $f_1(\tau) := \tau \chi_h(\sqrt{\tau})$  and  $f_2(\tau) := \chi_z(\sqrt{\tau})$ , which are both bounded functions with Lipschitz constant  $O(1)$ . This guarantees that for each *unit* vector pair  $\mathbf{h}$  and  $\mathbf{z}$ ,

$$\begin{aligned} & \left| (\mathbf{a}_i^\top \mathbf{h})^2 \chi_z(\mathbf{a}_i^\top \mathbf{z}) \chi_h(\mathbf{a}_i^\top \mathbf{h}) - (\mathbf{a}_i^\top \mathbf{h}_0)^2 \chi_z(\mathbf{a}_i^\top \mathbf{z}_0) \chi_h(\mathbf{a}_i^\top \mathbf{h}_0) \right| \\ & \leq |\chi_h(\mathbf{a}_i^\top \mathbf{z})| \cdot |(\mathbf{a}_i^\top \mathbf{h})^2 \chi_h(\mathbf{a}_i^\top \mathbf{h}) - (\mathbf{a}_i^\top \mathbf{h}_0)^2 \chi_h(\mathbf{a}_i^\top \mathbf{h}_0)| \\ & \quad + |(\mathbf{a}_i^\top \mathbf{h}_0)^2 \chi_h(\mathbf{a}_i^\top \mathbf{h}_0)| \cdot |\chi_h(\mathbf{a}_i^\top \mathbf{z}) - \chi_h(\mathbf{a}_i^\top \mathbf{z}_0)| \\ & \leq |\chi_h(\mathbf{a}_i^\top \mathbf{z})| \cdot |f_1(|\mathbf{a}_i^\top \mathbf{h}|^2) - f_1(|\mathbf{a}_i^\top \mathbf{h}_0|^2)| \\ & \quad + |(\mathbf{a}_i^\top \mathbf{h}_0)^2 \chi_h(\mathbf{a}_i^\top \mathbf{h}_0)| \cdot |f_2(|\mathbf{a}_i^\top \mathbf{z}|^2) - f_2(|\mathbf{a}_i^\top \mathbf{z}_0|^2)| \\ & \lesssim |(\mathbf{a}_i^\top \mathbf{h})^2 - (\mathbf{a}_i^\top \mathbf{h}_0)^2| + |(\mathbf{a}_i^\top \mathbf{z})^2 - (\mathbf{a}_i^\top \mathbf{z}_0)^2|. \end{aligned}$$

Consequently, there exists some universal constant  $c_3 > 0$  such that

$$\begin{aligned} & \left| \frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h})^2 \chi_z(\mathbf{a}_i^\top \mathbf{z}) \chi_h(\mathbf{a}_i^\top \mathbf{h}) - \frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h}_0)^2 \chi_z(\mathbf{a}_i^\top \mathbf{z}_0) \chi_h(\mathbf{a}_i^\top \mathbf{h}_0) \right| \\ & \lesssim \frac{1}{m} \left\| \mathcal{A}(\mathbf{h}\mathbf{h}^\top - \mathbf{h}_0\mathbf{h}_0^\top) \right\|_1 + \frac{1}{m} \left\| \mathcal{A}(\mathbf{z}\mathbf{z}^\top - \mathbf{z}_0\mathbf{z}_0^\top) \right\|_1 \\ & \stackrel{(i)}{\leq} c_3 \left\{ \|\mathbf{h}\mathbf{h}^\top - \mathbf{h}_0\mathbf{h}_0^\top\|_F + \|\mathbf{z}\mathbf{z}^\top - \mathbf{z}_0\mathbf{z}_0^\top\|_F \right\} \\ & \stackrel{(ii)}{\leq} 2.5c_3 \left\{ \|\mathbf{h} - \mathbf{h}_0\| \cdot \|\mathbf{h}\| + \|\mathbf{z} - \mathbf{z}_0\| \cdot \|\mathbf{z}\| \right\} \leq 5c_3\epsilon, \end{aligned}$$

where (i) results from Lemma 1, and (ii) arises from Lemma 2 whenever  $\epsilon < 1/2$ .

With the assertion (122) in place, we see that with high probability,

$$\frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h})^2 \chi_z(\mathbf{a}_i^\top \mathbf{z}) \chi_h(\mathbf{a}_i^\top \mathbf{h}) \geq (1 - \zeta_1 - \zeta_2 - (5c_3 + 1)\epsilon) \|\mathbf{h}\|^2$$

for all unit vectors  $\mathbf{h}$  and  $\mathbf{z}$ . Since  $\epsilon$  can be arbitrary, putting this and (116) together completes the proof.

### A.3 Proof of Lemma 5

The proof makes use of standard concentration of measure and covering arguments, and it suffices to restrict our attention to *unit vectors*  $\mathbf{h}$ . We find it convenient to work with an auxiliary function

$$\chi_2(\tau) = \begin{cases} |\tau|^{\frac{3}{2}}, & \text{if } |\tau| \leq \gamma^2, \\ -\gamma(|\tau| - \gamma^2) + \gamma^3, & \text{if } \gamma^2 < |\tau| \leq 2\gamma^2, \\ 0, & \text{else.} \end{cases}$$

Apparently,  $\chi_2(\tau)$  is a Lipschitz function of  $\tau$  with Lipschitz norm  $O(\gamma)$ . Recalling the definition of  $\mathcal{D}_\gamma^{i,1}$ , we see that each summand is bounded above by

$$|\mathbf{a}_i^\top \mathbf{h}|^3 \mathbf{1}_{\mathcal{D}_\gamma^{i,1}} \leq \chi_2(|\mathbf{a}_i^\top \mathbf{h}|^2).$$

For each fixed  $\mathbf{h}$  and  $\epsilon > 0$ , applying the Bernstein inequality [40, Proposition 5.16] gives

$$\begin{aligned} \frac{1}{m} \sum_{i=1}^m |\mathbf{a}_i^\top \mathbf{h}|^3 \mathbf{1}_{\mathcal{D}_\gamma^{i,1}} & \leq \frac{1}{m} \sum_{i=1}^m \chi_2(|\mathbf{a}_i^\top \mathbf{h}|^2) \leq \mathbb{E} \left[ \chi_2(|\mathbf{a}_i^\top \mathbf{h}|^2) \right] + \epsilon \\ & \leq \mathbb{E} \left[ |\mathbf{a}_i^\top \mathbf{h}|^3 \right] + \epsilon = \sqrt{8/\pi} + \epsilon \end{aligned}$$

with probability exceeding  $1 - \exp(-\Omega(\epsilon^2 m))$ .

From [40, Lemma 5.2], there exists an  $\epsilon$ -net  $\mathcal{N}_\epsilon$  of the unit sphere with cardinality  $|\mathcal{N}_\epsilon| \leq (1 + \frac{2}{\epsilon})^n$ . For each  $\mathbf{h}$ , suppose that  $\|\mathbf{h}_0 - \mathbf{h}\| \leq \epsilon$  for some  $\mathbf{h}_0 \in \mathcal{N}_\epsilon$ . The Lipschitz property of  $\chi_2$  implies

$$\begin{aligned} \frac{1}{m} \sum_{i=1}^m \left\{ \chi_2(|\mathbf{a}_i^\top \mathbf{h}|^2) - \chi_2(|\mathbf{a}_i^\top \mathbf{h}_0|^2) \right\} &\lesssim \frac{1}{m} \sum_{i=1}^m \left| |\mathbf{a}_i^\top \mathbf{h}|^2 - |\mathbf{a}_i^\top \mathbf{h}_0|^2 \right| \\ &\stackrel{(i)}{\asymp} \|\mathbf{h} - \mathbf{h}_0\| \|\mathbf{h}\| \asymp \epsilon, \end{aligned}$$

where (i) arises by combining Lemmas 1 and 2. This demonstrates that with high probability,

$$\frac{1}{m} \sum_{i=1}^m |\mathbf{a}_i^\top \mathbf{h}|^3 \mathbf{1}_{\mathcal{D}_\gamma^{i,1}} \leq \frac{1}{m} \sum_{i=1}^m \chi_2(|\mathbf{a}_i^\top \mathbf{h}|^2) \leq \sqrt{8/\pi} + O(\epsilon)$$

for all unit vectors  $\mathbf{h}$ , as claimed.

#### A.4 Proof of Lemma 6

Without loss of generality, the proof focuses on the case where  $\|\mathbf{h}\| = 1$ . Fix an arbitrary small constant  $\delta > 0$ . One can eliminate the difficulty of handling the discontinuous indicator functions by working with the following auxiliary function

$$\chi_3(\tau, \gamma) := \begin{cases} 1, & \text{if } \sqrt{\tau} \geq \psi_{\text{lb}}(\gamma); \\ \frac{100\tau}{\psi_{\text{lb}}^2(\gamma)} - 99, & \text{if } \sqrt{\tau} \in [\sqrt{0.99}\psi_{\text{lb}}(\gamma), \psi_{\text{lb}}(\gamma)]; \\ 0, & \text{else.} \end{cases} \quad (123)$$

Here,  $\psi_{\text{lb}}(\cdot)$  is a piecewise constant function defined as

$$\psi_{\text{lb}}(\gamma) := (1 + \delta) \lfloor \frac{\log \gamma}{\log(1 + \delta)} \rfloor,$$

which clearly satisfy  $\frac{\gamma}{1 + \delta} \leq \psi_{\text{lb}}(\gamma) \leq \gamma$ . Such a function is useful for our purpose since for any  $0 < \delta \leq 0.005$ ,

$$\mathbf{1}_{\{|\mathbf{a}_i^\top \mathbf{h}| \geq \gamma\}} \leq \chi_3(|\mathbf{a}_i^\top \mathbf{h}|^2, \gamma) \leq \mathbf{1}_{\{|\mathbf{a}_i^\top \mathbf{h}| \geq \sqrt{0.99}\psi_{\text{lb}}(\gamma)\}} \leq \mathbf{1}_{\{|\mathbf{a}_i^\top \mathbf{h}| \geq 0.99\gamma\}}. \quad (124)$$

For any fixed unit vector  $\mathbf{h}$ , the above argument leads to an upper tail estimate: for any  $0 < t \leq 1$ ,

$$\begin{aligned} \mathbb{P}\left\{ \chi_3(|\mathbf{a}_i^\top \mathbf{h}|^2, \gamma) \geq t \right\} &\leq \mathbb{P}\left\{ \mathbf{1}_{\{|\mathbf{a}_i^\top \mathbf{h}| \geq 0.99\gamma\}} \geq t \right\} = \mathbb{P}\left\{ \mathbf{1}_{\{|\mathbf{a}_i^\top \mathbf{h}| \geq 0.99\gamma\}} = 1 \right\} \\ &= 2 \int_{0.99\gamma}^{\infty} \phi(x) dx \leq \frac{2}{0.99\gamma} \phi(0.99\gamma), \end{aligned} \quad (125)$$

where  $\phi(x)$  is the density of a standard normal, and (125) follows from the tail bound  $\int_x^{\infty} \phi(x) dx \leq \frac{1}{x} \phi(x)$  for all  $x > 0$ . This implies that when  $\gamma \geq 2$ , both  $\chi_3(|\mathbf{a}_i^\top \mathbf{h}|^2, \gamma)$  and  $\mathbf{1}_{\{|\mathbf{a}_i^\top \mathbf{h}| \geq 0.99\gamma\}}$  are sub-exponential with sub-exponential norm  $O(\gamma^{-2})$  (cf. [40, Definition 5.13]). We apply the Bernstein-type inequality for the sum of sub-exponential random variables [40, Corollary 5.17], which indicates that for any fixed  $\mathbf{h}$  and  $\gamma$  as well as any sufficiently small  $\epsilon \in (0, 1)$ ,

$$\begin{aligned} \frac{1}{m} \sum_{i=1}^m \chi_3(|\mathbf{a}_i^\top \mathbf{h}|^2, \gamma) &\leq \frac{1}{m} \sum_{i=1}^m \mathbf{1}_{\{|\mathbf{a}_i^\top \mathbf{h}| \geq 0.99\gamma\}} \leq \mathbb{E}\left[ \mathbf{1}_{\{|\mathbf{a}_i^\top \mathbf{h}| \geq 0.99\gamma\}} \right] + \epsilon \frac{1}{\gamma^2} \\ &\leq \frac{2}{0.99\gamma} \exp(-0.49\gamma^2) + \epsilon \frac{1}{\gamma^2} \end{aligned}$$

holds with probability exceeding  $1 - \exp(-\Omega(\epsilon^2 m))$ .



We now proceed to obtain uniform control over all  $\mathbf{h}$  and  $2 \leq \gamma \leq 2^n$ . To begin with, we consider all  $2 \leq \gamma \leq m$  and construct an  $\epsilon$ -net  $\mathcal{N}_\epsilon$  over the unit sphere such that: (i)  $|\mathcal{N}_\epsilon| \leq (1 + \frac{2}{\epsilon})^n$ ; (ii) for any  $\mathbf{h}$  with  $\|\mathbf{h}\| = 1$ , there exists a unit vector  $\mathbf{h}_0 \in \mathcal{N}_\epsilon$  obeying  $\|\mathbf{h} - \mathbf{h}_0\| \leq \epsilon$ . Taking the union bound gives the following: with probability at least  $1 - \frac{\log m}{\log(1+\delta)} (1 + \frac{2}{\epsilon})^n \exp(-\Omega(\epsilon^2 m))$ ,

$$\frac{1}{m} \sum_{i=1}^m \chi_3(|\mathbf{a}_i^\top \mathbf{h}_0|^2, \gamma_0) \leq (0.495\gamma_0)^{-1} \exp(-0.49\gamma_0^2) + \epsilon\gamma_0^{-2}$$

holds simultaneously for all  $\mathbf{h}_0 \in \mathcal{N}_\epsilon$  and  $\gamma_0 \in \left\{ (1 + \delta)^k \mid 1 \leq k \leq \frac{\log m}{\log(1+\delta)} \right\}$ .

Note that  $\chi_3(\tau, \gamma_0)$  is a Lipschitz function in  $\tau$  with the Lipschitz constant bounded above by  $\frac{100}{\psi_{\text{lb}}^2(\gamma_0)}$ . With this in mind, for any  $(\mathbf{h}, \gamma)$  with  $\|\mathbf{h}\| = 1$  and  $\gamma_0 := (1 + \delta)^k \leq \gamma < (1 + \delta)^{k+1}$ , one has

$$\begin{aligned} \left| \chi_3(|\mathbf{a}_i^\top \mathbf{h}_0|^2, \gamma_0) - \chi_3(|\mathbf{a}_i^\top \mathbf{h}|^2, \gamma) \right| &= \left| \chi_3(|\mathbf{a}_i^\top \mathbf{h}_0|^2, \gamma_0) - \chi_3(|\mathbf{a}_i^\top \mathbf{h}|^2, \gamma_0) \right| \\ &\leq \frac{100}{\psi_{\text{lb}}^2(\gamma_0)} \left| |\mathbf{a}_i^\top \mathbf{h}|^2 - |\mathbf{a}_i^\top \mathbf{h}_0|^2 \right|. \end{aligned}$$

It then follows from Lemmas 1-2 that

$$\begin{aligned} \frac{1}{m} \left| \sum_{i=1}^m \chi_3(|\mathbf{a}_i^\top \mathbf{h}_0|^2, \gamma_0) - \sum_{i=1}^m \chi_3(|\mathbf{a}_i^\top \mathbf{h}|^2, \gamma) \right| &\leq \frac{100}{\psi_{\text{lb}}^2(\gamma_0)} \frac{1}{m} \left\| \mathcal{A}(\mathbf{h}\mathbf{h}^\top - \mathbf{h}_0\mathbf{h}_0^\top) \right\|_1 \\ &\leq \frac{250(1+\delta)^2}{\gamma^2} \|\mathbf{h} - \mathbf{h}_0\| \|\mathbf{h}\| \leq \frac{250(1+\delta)^2 \epsilon}{\gamma^2}. \end{aligned}$$

Putting the above results together gives that for all  $2 \leq \gamma \leq (1 + \delta)^{\frac{\log m}{\log(1+\delta)}} = m$ ,

$$\begin{aligned} \frac{1}{m} \sum_{i=1}^m \chi_3(|\mathbf{a}_i^\top \mathbf{h}|^2, \gamma) &\leq \frac{1}{m} \sum_{i=1}^m \chi_3(|\mathbf{a}_i^\top \mathbf{h}_0|^2, \gamma_0) + \frac{250(1+\delta)^2}{\gamma^2} \epsilon \\ &\leq \frac{1}{0.495\gamma_0} \exp(-0.49\gamma_0^2) + 251(1+\delta)^2 \frac{\epsilon}{\gamma^2} \\ &\leq \frac{1}{0.49\gamma} \exp(-0.485\gamma^2) + 251(1+\delta)^2 \frac{\epsilon}{\gamma^2} \end{aligned}$$

with probability exceeding  $1 - \frac{\log m}{\log(1+\delta)} (1 + \frac{2}{\epsilon})^n \exp(-c\epsilon^2 m)$ . This establishes (71) for all  $2 \leq \gamma \leq m$ .

It remains to deal with the case where  $\gamma > m$ . To this end, we rely on the following observation:

$$\frac{1}{m} \sum_{i=1}^m \mathbf{1}_{\{|\mathbf{a}_i^\top \mathbf{h}| \geq m\}} \leq \frac{1}{m} \sum_{i=1}^m \frac{|\mathbf{a}_i^\top \mathbf{h}|^2}{m^2} \stackrel{(i)}{\leq} \frac{1+\delta}{m^2} \|\mathbf{h}\|^2 \ll \frac{1}{m}, \quad \forall \mathbf{h} \text{ with } \|\mathbf{h}\| = 1,$$

where (i) comes from [6, Lemmas 3.1]. This basically tells us that with high probability, none of the indicator variables can be equal to 1. Consequently,  $\frac{1}{m} \sum_{i=1}^m \mathbf{1}_{\{|\mathbf{a}_i^\top \mathbf{h}| \geq m\}} = 0$ , which proves the claim.

## A.5 Proof of Lemma 7

Fix  $\delta > 0$ . Recalling the notation  $v_i := 2 \left\{ 2\mathbf{a}_i^\top \mathbf{h} - \frac{|\mathbf{a}_i^\top \mathbf{h}|^2}{\mathbf{a}_i^\top \mathbf{z}} \right\} \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i}$ , we see from the expansion (62) that

$$\left\| \frac{1}{m} \nabla_{\text{tr}} \ell(\mathbf{z}) \right\| = \left\| \frac{1}{m} \mathbf{A}^\top \mathbf{v} \right\| \leq \frac{1}{m} \|\mathbf{A}\| \cdot \|\mathbf{v}\| \leq (1 + \delta) \frac{\|\mathbf{v}\|}{\sqrt{m}} \quad (126)$$

as soon as  $m \geq c_1 n$  for some sufficiently large  $c_1 > 0$ . Here, the norm estimate  $\|\mathbf{A}\| \leq \sqrt{m}(1 + \delta)$  arises from standard random matrix results [40, Corollary 5.35].

Everything then comes down to controlling  $\|\mathbf{v}\|$ . To this end, making use of the inclusion (63) yields

$$\begin{aligned}
\frac{1}{4m} \|\mathbf{v}\|^2 &= \frac{1}{m} \sum_{i=1}^m \left( 2\mathbf{a}_i^\top \mathbf{h} - \frac{|\mathbf{a}_i^\top \mathbf{h}|^2}{\mathbf{a}_i^\top \mathbf{z}} \right)^2 \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i} \\
&\leq \frac{1}{m} \sum_{i=1}^m \left( 2|\mathbf{a}_i^\top \mathbf{h}| + \frac{|\mathbf{a}_i^\top \mathbf{h}|^2}{|\mathbf{a}_i^\top \mathbf{z}|} \right)^2 \mathbf{1}_{\mathcal{E}_1^i \cap (\mathcal{D}_{\gamma_4}^{i,1} \cup \mathcal{D}_{\gamma_4}^{i,2})} \\
&\leq \frac{1}{m} \sum_{i=1}^m \left\{ 4(\mathbf{a}_i^\top \mathbf{h})^2 + \left( \frac{4|\mathbf{a}_i^\top \mathbf{h}|^3}{|\mathbf{a}_i^\top \mathbf{z}|} + \frac{|\mathbf{a}_i^\top \mathbf{h}|^4}{|\mathbf{a}_i^\top \mathbf{z}|^2} \right) \mathbf{1}_{\mathcal{E}_1^i \cap (\mathcal{D}_{\gamma_4}^{i,1} \cup \mathcal{D}_{\gamma_4}^{i,2})} \right\} \\
&= \frac{1}{m} \sum_{i=1}^m \left\{ 4(\mathbf{a}_i^\top \mathbf{h})^2 + \left( 4 + \frac{|\mathbf{a}_i^\top \mathbf{h}|}{|\mathbf{a}_i^\top \mathbf{z}|} \right) \frac{|\mathbf{a}_i^\top \mathbf{h}|^3}{|\mathbf{a}_i^\top \mathbf{z}|} \left( \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{D}_{\gamma_4}^{i,1}} + \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{D}_{\gamma_4}^{i,2}} \right) \right\}.
\end{aligned}$$

The first term is controlled by [6, Lemma 3.1] in such a way that with probability  $1 - \exp(-\Omega(m))$ ,

$$\frac{1}{m} \sum_{i=1}^m 4(\mathbf{a}_i^\top \mathbf{h})^2 \leq 4(1 + \delta) \|\mathbf{h}\|^2.$$

Turning to the remaining terms, we see from the definition of  $\mathcal{D}_\gamma^{i,1}$  and  $\mathcal{D}_\gamma^{i,2}$  that

$$\frac{|\mathbf{a}_i^\top \mathbf{h}|}{|\mathbf{a}_i^\top \mathbf{z}|} \leq \begin{cases} \frac{\gamma \|\mathbf{h}\|}{\alpha_z^{\text{lb}} \|\mathbf{z}\|}, & \text{on } \mathcal{E}_1^i \cap \mathcal{D}_\gamma^{i,1} \\ 2 + \frac{\gamma \|\mathbf{h}\|}{\alpha_z^{\text{lb}} \|\mathbf{z}\|}, & \text{on } \mathcal{E}_1^i \cap \mathcal{D}_\gamma^{i,2} \end{cases} \leq \begin{cases} 1, & \text{on } \mathcal{E}_1^i \cap \mathcal{D}_\gamma^{i,1} \\ 3, & \text{on } \mathcal{E}_1^i \cap \mathcal{D}_\gamma^{i,2} \end{cases}$$

as long as  $\gamma \leq \frac{\alpha_z^{\text{lb}} \|\mathbf{z}\|}{\|\mathbf{h}\|}$ . Consequently, one can bound

$$\begin{aligned}
&\frac{1}{m} \sum_{i=1}^m \left( 4 + \frac{|\mathbf{a}_i^\top \mathbf{h}|}{|\mathbf{a}_i^\top \mathbf{z}|} \right) \frac{|\mathbf{a}_i^\top \mathbf{h}|^3}{|\mathbf{a}_i^\top \mathbf{z}|} \left( \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{D}_\gamma^{i,1}} + \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{D}_\gamma^{i,2}} \right) \\
&\leq \frac{5}{m} \sum_{i=1}^m \frac{|\mathbf{a}_i^\top \mathbf{h}|^3}{|\mathbf{a}_i^\top \mathbf{z}|} \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{D}_\gamma^{i,1}} + \frac{7}{m} \sum_{i=1}^m \frac{|\mathbf{a}_i^\top \mathbf{h}|^3}{|\mathbf{a}_i^\top \mathbf{z}|} \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{D}_\gamma^{i,2}} \\
&\leq \frac{5(1 + \delta) \sqrt{8/\pi} \|\mathbf{h}\|^3}{\alpha_z^{\text{lb}} \|\mathbf{z}\|} + \frac{7}{100} (1 + \delta) \|\mathbf{h}\|^2,
\end{aligned}$$

where the last inequality follows from (69) and (78).

Recall that  $\gamma_4 = 3\alpha_h$ . Taken together all these bounds lead to the upper bound

$$\begin{aligned}
\frac{1}{4m} \|\mathbf{v}\|^2 &\leq (1 + \delta) \left\{ 4 + \frac{5\sqrt{8/\pi} \|\mathbf{h}\|}{\alpha_z^{\text{lb}} \|\mathbf{z}\|} + \frac{7}{100} \right\} \|\mathbf{h}\|^2 \\
&\leq (1 + \delta) \left\{ 4 + \frac{5\sqrt{8/\pi}}{3\alpha_h} + \frac{7}{100} \right\} \|\mathbf{h}\|^2
\end{aligned}$$

whenever  $\frac{\|\mathbf{h}\|}{\|\mathbf{z}\|} \leq \min \left\{ \frac{\alpha_z^{\text{lb}}}{3\alpha_h}, \frac{\alpha_z^{\text{lb}}}{6}, \frac{\sqrt{98/3}(\alpha_z^{\text{lb}})^2}{2\alpha_z^{\text{lb}} + \alpha_z^{\text{lb}}}, \frac{1}{11} \right\}$ . Substituting this into (126) completes the proof.

## B Proofs for Section 7

### B.1 Proof of Lemma 8

Firstly, we collect a few results on the magnitudes of  $\mathbf{a}_i^\top \mathbf{x}$  ( $1 \leq i \leq m$ ) that will be useful in constructing the hypotheses. Observe that for any given  $\mathbf{x}$  and any sufficiently large  $m$ ,

$$\begin{aligned}
\mathbb{P} \left\{ \min_{1 \leq i \leq m} |\mathbf{a}_i^\top \mathbf{x}| \geq \frac{1}{m \log m} \|\mathbf{x}\| \right\} &= \left( \mathbb{P} \left\{ |\mathbf{a}_i^\top \mathbf{x}| \geq \frac{1}{m \log m} \|\mathbf{x}\| \right\} \right)^m \\
&\geq \left( 1 - \frac{2}{\sqrt{2\pi}} \frac{1}{m \log m} \right)^m \geq 1 - o(1).
\end{aligned}$$

Besides, since  $\mathbb{E} \left[ \mathbf{1}_{\{|\mathbf{a}_i^\top \mathbf{x}| \leq \frac{\|\mathbf{x}\|}{5 \log m}\}} \right] \leq \frac{1}{\sqrt{2\pi}} \frac{2}{5 \log m} \leq \frac{1}{5 \log m}$ , applying Hoeffding's inequality yields

$$\begin{aligned} & \mathbb{P} \left\{ \sum_{i=1}^m \mathbf{1}_{\{|\mathbf{a}_i^\top \mathbf{x}| \leq \frac{\|\mathbf{x}\|}{5 \log m}\}} > \frac{m}{4 \log m} \right\} \\ &= \mathbb{P} \left\{ \frac{1}{m} \sum_{i=1}^m \left( \mathbf{1}_{\{|\mathbf{a}_i^\top \mathbf{x}| \leq \frac{\|\mathbf{x}\|}{5 \log m}\}} - \mathbb{E} \left[ \mathbf{1}_{\{|\mathbf{a}_i^\top \mathbf{x}| \leq \frac{\|\mathbf{x}\|}{5 \log m}\}} \right] \right) > \frac{1}{20 \log m} \right\} \\ &\leq \exp \left( -\Omega \left( \frac{m}{\log^2 m} \right) \right). \end{aligned}$$

To summarize, with probability  $1 - o(1)$ , one has

$$\min_{1 \leq i \leq m} |\mathbf{a}_i^\top \mathbf{x}| \geq \frac{1}{m \log m} \|\mathbf{x}\|; \quad (127)$$

$$\sum_{i=1}^m \mathbf{1}_{\{|\mathbf{a}_i^\top \mathbf{x}| \leq \frac{\|\mathbf{x}\|}{5 \log m}\}} \leq \frac{m}{4 \log m} := k. \quad (128)$$

In the sequel, we will first produce a set  $\mathcal{M}_1$  of exponentially many vectors surrounding  $\mathbf{x}$  in such a way that every pair is separated by about the same distance, and then verify that a non-trivial fraction of  $\mathcal{M}_1$  obeys (103). Without loss of generality, we assume that  $\mathbf{x}$  takes the form  $\mathbf{x} = [b, 0, \dots, 0]^\top$  for some  $b > 0$ .

The construction of  $\mathcal{M}_1$  follows a standard random packing argument. Let  $\mathbf{w} = [w_1, \dots, w_n]^\top$  be a random vector with

$$w_i = x_i + \frac{1}{\sqrt{2n}} z_i, \quad 1 \leq i \leq n,$$

where  $z_i \stackrel{\text{ind.}}{\sim} \mathcal{N}(0, 1)$ . The collection  $\mathcal{M}_1$  is then obtained by generating  $M_1 = \exp\left(\frac{n}{20}\right)$  independent copies  $\mathbf{w}^{(l)}$  ( $1 \leq l < M_1$ ) of  $\mathbf{w}$ . For any  $\mathbf{w}^{(l)}, \mathbf{w}^{(j)} \in \mathcal{M}_1$ , the concentration inequality [40, Corollary 5.35] gives

$$\begin{aligned} \mathbb{P} \{ 0.5\sqrt{n} - 1 \leq \sqrt{n} \|\mathbf{w}^{(l)} - \mathbf{w}^{(j)}\| \leq 1.5\sqrt{n} + 1 \} &\geq 1 - 2 \exp(-n/8); \\ \mathbb{P} \{ 0.5\sqrt{n} - 1 \leq \sqrt{2n} \|\mathbf{w}^{(l)} - \mathbf{x}\| \leq 1.5\sqrt{n} + 1 \} &\geq 1 - 2 \exp(-n/8). \end{aligned}$$

Taking the union bound over all  $\binom{M_1}{2}$  pairs we obtain

$$\begin{aligned} 0.5 - n^{-1/2} \leq \|\mathbf{w}^{(l)} - \mathbf{w}^{(j)}\| &\leq 1.5 + n^{-1/2}, \quad \forall l \neq j \\ 1/\sqrt{8} - (2n)^{-1/2} \leq \|\mathbf{w}^{(l)} - \mathbf{x}\| &\leq \sqrt{9/8} + (2n)^{-1/2}, \quad 1 \leq l \leq M_1 \end{aligned} \quad (129)$$

with probability exceeding  $1 - 2M_1^2 \exp\left(-\frac{n}{8}\right) \geq 1 - 2 \exp\left(-\frac{n}{40}\right)$ .

The next step is to show that many vectors in  $\mathcal{M}_1$  satisfy (103). For any given  $\mathbf{w}$  with  $\mathbf{r} := \mathbf{w} - \mathbf{x}$ , by letting  $\mathbf{a}_{i,\perp} := [a_{i,2}, \dots, a_{i,n}]^\top$ ,  $r_{\parallel} := r_1$ , and  $\mathbf{r}_{\perp} := [r_2, \dots, r_n]^\top$ , we derive

$$\frac{|\mathbf{a}_i^\top \mathbf{r}|^2}{|\mathbf{a}_i^\top \mathbf{x}|^2} \leq \frac{2|a_{i,1} r_{\parallel}|^2 + 2|\mathbf{a}_{i,\perp}^\top \mathbf{r}_{\perp}|^2}{|a_{i,1}|^2 \|\mathbf{x}\|^2} \leq \frac{2|r_{\parallel}|^2}{\|\mathbf{x}\|^2} + \frac{2|\mathbf{a}_{i,\perp}^\top \mathbf{r}_{\perp}|^2}{|a_{i,1}|^2 \|\mathbf{x}\|^2} \leq \frac{2\|\mathbf{r}\|^2}{\|\mathbf{x}\|^2} + \frac{2|\mathbf{a}_{i,\perp}^\top \mathbf{r}_{\perp}|^2}{|a_{i,1}|^2 \|\mathbf{x}\|^2}. \quad (130)$$

It then boils down to developing an upper bound on  $\frac{|\mathbf{a}_{i,\perp}^\top \mathbf{r}_{\perp}|^2}{|a_{i,1}|^2}$ . This ratio is convenient to work with since the numerator and denominator are stochastically independent. To simplify presentation, we reorder  $\{\mathbf{a}_i\}$  in a way that

$$(m \log m)^{-1} \|\mathbf{x}\| \leq |\mathbf{a}_1^\top \mathbf{x}| \leq |\mathbf{a}_2^\top \mathbf{x}| \leq \dots \leq |\mathbf{a}_m^\top \mathbf{x}|;$$

this will not affect our subsequent analysis concerning  $\mathbf{a}_{i,\perp}^\top \mathbf{r}_{\perp}$  since it is independent of  $\mathbf{a}_i^\top \mathbf{x}$ .

To proceed, we let  $\mathbf{r}_{\perp}^{(l)}$  consist of all but the first entry of  $\mathbf{w}^{(l)} - \mathbf{x}$ , and introduce the indicator variables

$$\xi_i^l := \begin{cases} \mathbf{1}_{\{|\mathbf{a}_{i,\perp}^\top \mathbf{r}_{\perp}^{(l)}| \leq \frac{1}{m} \sqrt{\frac{n-1}{2n}}\}}, & 1 \leq i \leq k, \\ \mathbf{1}_{\{|\mathbf{a}_{i,\perp}^\top \mathbf{r}_{\perp}^{(l)}| \leq \sqrt{\frac{2(n-1) \log n}{n}}\}}, & i > k, \end{cases} \quad (131)$$

where  $k = \frac{m}{4 \log m}$  as before. In words, we divide  $\mathbf{a}_{i,\perp}^\top \mathbf{r}_\perp^{(l)}$ ,  $1 \leq i \leq m$  into two groups, with the first group enforcing far more stringent control than the second group. These indicator variables are useful since any  $\mathbf{w}^{(l)}$  obeying  $\prod_{i=1}^m \xi_i^l = 1$  will satisfy (103) when  $n$  is sufficiently large. To see this, note that for the first group of indices,  $\xi_i^l = 1$  requires

$$\left| \mathbf{a}_{i,\perp}^\top \mathbf{r}_\perp^{(l)} \right| \leq \frac{1}{m} \sqrt{\frac{n-1}{2n}} \leq \frac{2}{m} \frac{\sqrt{n-1}}{\sqrt{n}-2} \|\mathbf{r}^{(l)}\| \leq \frac{3}{m} \|\mathbf{r}^{(l)}\|, \quad 1 \leq i \leq k, \quad (132)$$

where the second inequality follows from (129). This taken collectively with (127) and (130) yields

$$\frac{|\mathbf{a}_i^\top \mathbf{r}^{(l)}|^2}{|\mathbf{a}_i^\top \mathbf{x}|^2} \leq \frac{2\|\mathbf{r}^{(l)}\|^2}{\|\mathbf{x}\|^2} + \frac{\frac{9}{m^2} \|\mathbf{r}^{(l)}\|^2}{\frac{1}{m^2 \log^2 m} \|\mathbf{x}\|^2} \leq \frac{(2 + 9 \log^2 m) \|\mathbf{r}^{(l)}\|^2}{\|\mathbf{x}\|^2}, \quad 1 \leq i \leq k.$$

Regarding the second group of indices,  $\xi_i^l = 1$  gives

$$\left| \mathbf{a}_{i,\perp}^\top \mathbf{r}_\perp^{(l)} \right| \leq \sqrt{\frac{2(n-1) \log n}{n}} \leq \sqrt{17 \log n} \|\mathbf{r}^{(l)}\|, \quad i = k+1, \dots, m, \quad (133)$$

where the last inequality again follows from (129). Plugging (133) and (128) into (130) gives

$$\frac{|\mathbf{a}_i^\top \mathbf{r}^{(l)}|^2}{|\mathbf{a}_i^\top \mathbf{x}|^2} \leq \frac{2\|\mathbf{r}^{(l)}\|^2}{\|\mathbf{x}\|^2} + \frac{17 \|\mathbf{r}^{(l)}\|^2 \log n}{\|\mathbf{x}\|^2 / \log^2 m} \leq \frac{(2 + 17 \log^3 m) \|\mathbf{r}^{(l)}\|^2}{\|\mathbf{x}\|^2}, \quad i \geq k+1.$$

Consequently, (103) is satisfied for all  $1 \leq i \leq m$ . It then suffices to guarantee the existence of exponentially many vectors obeying  $\prod_{i=1}^m \xi_i^l = 1$ .

Note that the first group of indicator variables are quite stringent, namely, for each  $i$  only a fraction  $O(1/m)$  of the equations could satisfy  $\xi_i^l = 1$ . Fortunately,  $M_1$  is exponentially large, and hence even  $M_1/m^k$  is exponentially large. Put formally, we claim that the first group satisfies

$$\sum_{l=1}^{M_1} \prod_{i=1}^k \xi_i^l \geq \frac{1}{2} \frac{M_1}{(2\pi)^{k/2} (1 + 4\sqrt{k/n})^{k/2}} \left( \frac{1}{\sqrt{2\pi m}} \right)^k := \widetilde{M}_1 \quad (134)$$

with probability exceeding  $1 - \exp(-\Omega(k)) - \exp(-\widetilde{M}_1/4)$ . With this claim in place (which will be proved later), one has

$$\sum_{l=1}^{M_1} \prod_{i=1}^k \xi_i^l \geq \frac{1}{2} M_1 \frac{1}{(e^2 m)^k} = \frac{1}{2} \exp\left(\left(\frac{1}{20} - \frac{k(2 + \log m)}{n}\right)n\right) \geq \frac{1}{2} \exp\left(\frac{1}{25}n\right)$$

when  $n$  and  $m/n$  are sufficiently large. In light of this, we will let  $\mathcal{M}_2$  be a collection comprising all  $\mathbf{w}^{(l)}$  obeying  $\prod_{i=1}^k \xi_i^l = 1$ , which has size  $M_2 \geq \frac{1}{2} \exp(\frac{1}{25}n)$  based on the preceding argument. For notational simplicity, it will be assumed that the vectors in  $\mathcal{M}_2$  are exactly  $\mathbf{w}^{(j)}$  ( $1 \leq j \leq M_2$ ).

We now move on to the second group by examining how many vectors  $\mathbf{w}^{(j)}$  in  $\mathcal{M}_2$  further satisfy  $\prod_{i=k+1}^m \xi_i^j = 1$ . Notably, the above construction of  $\mathcal{M}_2$  relies only on  $\{\mathbf{a}_i\}_{1 \leq i \leq k}$  and is independent of the remaining vectors  $\{\mathbf{a}_i\}_{i > k}$ . In what follows the argument proceeds conditional on  $\mathcal{M}_2$  and  $\{\mathbf{a}_i\}_{1 \leq i \leq k}$ . Applying the union bound gives

$$\begin{aligned} \mathbb{E} \left[ \sum_{j=1}^{M_2} \left( 1 - \prod_{i=k+1}^m \xi_i^j \right) \right] &= \sum_{j=1}^{M_2} \mathbb{P} \left\{ \exists i (k < i \leq m) : \left| \mathbf{a}_{i,\perp}^\top \mathbf{r}_\perp^{(l)} \right| > \sqrt{\frac{2(n-1) \log n}{n}} \right\} \\ &\leq \sum_{j=1}^{M_2} \sum_{i=k+1}^m \mathbb{P} \left\{ \left| \mathbf{a}_{i,\perp}^\top \mathbf{r}_\perp^{(l)} \right| > \sqrt{\frac{2(n-1) \log n}{n}} \right\} \leq M_2 m \frac{1}{n^2}. \end{aligned}$$

This combined with Markov's inequality gives

$$\sum_{j=1}^{M_2} \left( 1 - \prod_{i=k+1}^m \xi_i^j \right) \leq \frac{m \log m}{n^2} \cdot M_2$$

with probability  $1 - o(1)$ . Putting the above inequalities together suggests that with probability  $1 - o(1)$ , there exist at least

$$\left(1 - \frac{m \log m}{n^2}\right) M_2 \geq \frac{1}{2} \left(1 - \frac{m \log m}{n^2}\right) \exp\left(\frac{1}{25}n\right) \geq \exp\left(\frac{n}{30}\right)$$

vectors in  $\mathcal{M}_2$  satisfying  $\prod_{l=k+1}^m \xi_i^l = 1$ . We then choose  $\mathcal{M}$  to be the set consisting of all these vectors, which forms a valid collection satisfying the properties of Lemma 8.

Finally, the only remaining step is to establish the claim (134). To start with, consider an  $n \times k$  matrix  $\mathbf{B} := [\mathbf{b}_1, \dots, \mathbf{b}_k]$  of i.i.d. standard normal entries, and let  $\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \frac{1}{n} \mathbf{I}_n)$ . Conditional on the  $\{\mathbf{b}_i\}$ 's,

$$\mathbf{b}_u = \begin{bmatrix} b_{1,\mathbf{u}} \\ \vdots \\ b_{k,\mathbf{u}} \end{bmatrix} := \begin{bmatrix} \mathbf{b}_1^\top \mathbf{u} \\ \vdots \\ \mathbf{b}_k^\top \mathbf{u} \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \frac{1}{n} \mathbf{B}^\top \mathbf{B}\right).$$

For sufficiently large  $m$ , one has  $k = \frac{m}{4 \log m} \leq \frac{1}{4}n$ . Using [40, Corollary 5.35] we get

$$\left\| \frac{1}{n} \mathbf{B}^\top \mathbf{B} - \mathbf{I} \right\| \leq 4\sqrt{k/n} \quad (135)$$

with probability  $1 - \exp(-\Omega(k))$ . Thus, for any constant  $0 < \epsilon < \frac{1}{2}$ , conditional on  $\{\mathbf{b}_i\}$  and (135) we obtain

$$\mathbb{P} \left\{ \bigcap_{i=1}^k \left\{ |\mathbf{b}_i^\top \mathbf{u}| \leq \frac{1}{m} \right\} \right\} \quad (136)$$

$$\begin{aligned} &\geq (2\pi)^{-\frac{k}{2}} \det^{-\frac{1}{2}} \left( \frac{1}{n} \mathbf{B}^\top \mathbf{B} \right) \int_{\mathbf{b}_u \in \Upsilon} \exp\left(-\frac{1}{2} \mathbf{b}_u^\top \left( \frac{1}{n} \mathbf{B}^\top \mathbf{B} \right)^{-1} \mathbf{b}_u\right) d\mathbf{b}_u \\ &\geq (2\pi)^{-\frac{k}{2}} \left(1 + 4\sqrt{k/n}\right)^{-\frac{k}{2}} \int_{\mathbf{b}_u \in \Upsilon} \exp\left(-\frac{1}{2} \left(1 - 4\sqrt{k/n}\right)^{-1} \sum_{i=1}^k b_{i,\mathbf{u}}^2\right) d\mathbf{b}_u \end{aligned} \quad (137)$$

$$\geq (2\pi)^{-\frac{k}{2}} \left(1 + 4\sqrt{k/n}\right)^{-\frac{k}{2}} (\sqrt{2\pi m})^{-k}, \quad (138)$$

where  $\Upsilon := \{\tilde{\mathbf{b}} \mid |\tilde{b}_i| \leq m^{-1}, 1 \leq i \leq k\}$  and (137) is a direct consequence from (135).

When it comes to our quantity of interest, the above lower bound (138) indicates that on an event (defined via  $\{\mathbf{a}_i\}$ ) of probability approaching 1, we have

$$\mathbb{E} \left[ \sum_{l=1}^{M_1} \prod_{i=1}^k \xi_i^l \right] \geq M_1 (2\pi)^{-\frac{k}{2}} \left(1 + 4\sqrt{k/n}\right)^{-\frac{k}{2}} (\sqrt{2\pi m})^{-k}. \quad (139)$$

Since conditional on  $\{\mathbf{a}_i\}$ ,  $\prod_{i=1}^k \xi_i^l$  are independent across  $l$ , applying the Chernoff-type bound [55, Theorem 4.5] gives

$$\sum_{l=1}^{M_1} \prod_{i=1}^k \xi_i^l \geq \frac{M_1}{2} (2\pi)^{-\frac{k}{2}} \left(1 + 4\sqrt{k/n}\right)^{-\frac{k}{2}} (\sqrt{2\pi m})^{-k}$$

with probability exceeding  $1 - \exp\left(-\frac{1}{8} \frac{M_1}{(2\pi)^{k/2} (1+4\sqrt{k/n})^{k/2}} \left(\frac{1}{\sqrt{2\pi m}}\right)^k\right)$ . This concludes the proof.

## B.2 Proof of Lemma 9

Before proceeding, we introduce the  $\chi^2$ -divergence between two probability measures  $P$  and  $Q$  as

$$\chi^2(P\|Q) := \int \left(\frac{dP}{dQ}\right)^2 dQ - 1. \quad (140)$$

It is well known (e.g. [41, Lemma 2.7]) that

$$\text{KL}(P\|Q) \leq \log(1 + \chi^2(P\|Q)), \quad (141)$$

and hence it suffices to develop an upper bound on the  $\chi^2$  divergence.

Under independence, for any  $\mathbf{w}_0, \mathbf{w}_1 \in \mathbb{R}^n$ , the decoupling identity of the  $\chi^2$  divergence [41, Page 96] gives

$$\begin{aligned}\chi^2(\mathbb{P}(\mathbf{y} | \mathbf{w}_1) \parallel \mathbb{P}(\mathbf{y} | \mathbf{w}_0)) &= \prod_{i=1}^m (1 + \chi^2(\mathbb{P}(y_i | \mathbf{w}_1) \parallel \mathbb{P}(y_i | \mathbf{w}_0))) - 1 \\ &= \exp\left(\sum_{i=1}^m \frac{(|\mathbf{a}_i^\top \mathbf{w}_1|^2 - |\mathbf{a}_i^\top \mathbf{w}_0|^2)^2}{|\mathbf{a}_i^\top \mathbf{w}_0|^2}\right) - 1.\end{aligned}\quad (142)$$

The preceding identity (142) arises from the following computation: by definition of  $\chi^2(\cdot \parallel \cdot)$ ,

$$\begin{aligned}\chi^2(\text{Poisson}(\lambda_1) \parallel \text{Poisson}(\lambda_0)) &= \left\{ \sum_{k=0}^{\infty} \frac{(\lambda_1^k \exp(-\lambda_1))^2}{\lambda_0^k \exp(-\lambda_0) k!} \right\} - 1 \\ &= \exp\left(\lambda_0 - 2\lambda_1 + \frac{\lambda_1^2}{\lambda_0}\right) \left\{ \sum_{k=0}^{\infty} \frac{(\lambda_1^2/\lambda_0)^k}{k!} \exp\left(-\frac{\lambda_1^2}{\lambda_0}\right) \right\} - 1 = \exp\left(\frac{(\lambda_1 - \lambda_0)^2}{\lambda_0}\right) - 1.\end{aligned}$$

Set  $\mathbf{r} := \mathbf{w}_1 - \mathbf{w}_0$ . To summarize,

$$\text{KL}(\mathbb{P}(\mathbf{y} | \mathbf{w}_1) \parallel \mathbb{P}(\mathbf{y} | \mathbf{w}_0)) \leq \sum_{i=1}^m \frac{(|\mathbf{a}_i^\top \mathbf{w}_1|^2 - |\mathbf{a}_i^\top \mathbf{w}_0|^2)^2}{|\mathbf{a}_i^\top \mathbf{w}_0|^2} \quad (143)$$

$$\begin{aligned}&\leq \sum_{i=1}^m \frac{|\mathbf{a}_i^\top \mathbf{r}|^2 (2|\mathbf{a}_i^\top \mathbf{w}_0| + |\mathbf{a}_i^\top \mathbf{r}|)^2}{|\mathbf{a}_i^\top \mathbf{w}_0|^2} \\ &= \sum_{i=1}^m |\mathbf{a}_i^\top \mathbf{r}|^2 \left( \frac{8|\mathbf{a}_i^\top \mathbf{w}_0|^2 + 2|\mathbf{a}_i^\top \mathbf{r}|^2}{|\mathbf{a}_i^\top \mathbf{w}_0|^2} \right).\end{aligned}\quad (144)$$

## C Initialization via truncated spectral Method

This section demonstrates that the truncated spectral method works when  $m \asymp n$ , as stated in the proposition below.

**Proposition 3.** Fix  $\delta > 0$  and  $\mathbf{x} \in \mathbb{R}^n$ . Consider the model where  $y_i = |\mathbf{a}_i^\top \mathbf{x}|^2 + \eta_i$  and  $\mathbf{a}_i \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I})$ . Suppose that

$$|\eta_i| \leq \varepsilon \max\{\|\mathbf{x}\|^2, |\mathbf{a}_i^\top \mathbf{x}|^2\}, \quad 1 \leq i \leq m \quad (145)$$

for some sufficiently small constant  $\varepsilon > 0$ . With probability exceeding  $1 - \exp(-\Omega(m))$ , the solution  $\mathbf{z}^{(0)}$  returned by the truncated spectral method obeys

$$\text{dist}(\mathbf{z}^{(0)}, \mathbf{x}) \leq \delta \|\mathbf{x}\|, \quad (146)$$

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

*Proof.* By homogeneity, it suffices to consider the case where  $\|\mathbf{x}\| = 1$ . Recall from [6, Lemma 3.1] that  $\frac{1}{m} \sum_{i=1}^m |\mathbf{a}_i^\top \mathbf{x}|^2 \in [1 \pm \varepsilon] \|\mathbf{x}\|^2$  holds with probability  $1 - \exp(-\Omega(m))$ . Under the hypothesis (145), one has

$$\frac{1}{m} \|\boldsymbol{\eta}\|_1 \leq \frac{1}{m} \sum_{i=1}^m \varepsilon (\|\mathbf{x}\|^2 + |\mathbf{a}_i^\top \mathbf{x}|^2) \leq \varepsilon \|\mathbf{x}\|^2 + \varepsilon(1 + \varepsilon) \|\mathbf{x}\|^2 \leq 3\varepsilon \|\mathbf{x}\|^2,$$

which yields

$$\lambda_0^2 := \frac{1}{m} \sum_{l=1}^m y_l = \frac{1}{m} \sum_{l=1}^m |\mathbf{a}_l^\top \mathbf{x}|^2 + \frac{1}{m} \sum_{l=1}^m \eta_l \in [1 \pm 4\varepsilon] \|\mathbf{x}\|^2 \quad (147)$$

with probability  $1 - \exp(-\Omega(m))$ .

Consequently, when  $|\eta_i| \leq \varepsilon \|\mathbf{x}\|^2$ , one has

$$\begin{aligned} \mathbf{1}_{\{ |(\mathbf{a}_i^\top \mathbf{x})^2 + \eta_i| \leq \alpha_y^2 (\frac{1}{m} \sum_l y_l) \}} &\leq \mathbf{1}_{\{ |\mathbf{a}_i^\top \mathbf{x}|^2 \leq \alpha_y^2 (\frac{1}{m} \sum_l y_l) + |\eta_i| \}} \leq \mathbf{1}_{\{ |\mathbf{a}_i^\top \mathbf{x}|^2 \leq (1+4\varepsilon)\alpha_y^2 + \varepsilon \}}; \\ \mathbf{1}_{\{ |(\mathbf{a}_i^\top \mathbf{x})^2 + \eta_i| \leq \alpha_y^2 (\frac{1}{m} \sum_l y_l) \}} &\geq \mathbf{1}_{\{ |\mathbf{a}_i^\top \mathbf{x}|^2 \leq \alpha_y^2 (\frac{1}{m} \sum_l y_l) - |\eta_i| \}} \geq \mathbf{1}_{\{ |\mathbf{a}_i^\top \mathbf{x}|^2 \leq (1-4\varepsilon)\alpha_y^2 - \varepsilon \}}. \end{aligned} \quad (148)$$

Besides, in the case where  $|\eta_i| \leq \varepsilon |\mathbf{a}_i^\top \mathbf{x}|^2$ ,

$$\begin{aligned} \mathbf{1}_{\{ |(\mathbf{a}_i^\top \mathbf{x})^2 + \eta_i| \leq \alpha_y^2 (\frac{1}{m} \sum_l y_l) \}} &\leq \mathbf{1}_{\{ (1-\varepsilon)|\mathbf{a}_i^\top \mathbf{x}|^2 \leq \alpha_y^2 (\frac{1}{m} \sum_l y_l) \}} \leq \mathbf{1}_{\{ |\mathbf{a}_i^\top \mathbf{x}|^2 \leq \frac{1+4\varepsilon}{1-\varepsilon} \alpha_y^2 \}}; \\ \mathbf{1}_{\{ |(\mathbf{a}_i^\top \mathbf{x})^2 + \eta_i| \leq \alpha_y^2 (\frac{1}{m} \sum_l y_l) \}} &\geq \mathbf{1}_{\{ (1+\varepsilon)|\mathbf{a}_i^\top \mathbf{x}|^2 \leq \alpha_y^2 (\frac{1}{m} \sum_l y_l) \}} \geq \mathbf{1}_{\{ |\mathbf{a}_i^\top \mathbf{x}|^2 \leq \frac{1-4\varepsilon}{1+\varepsilon} \alpha_y^2 \}}. \end{aligned} \quad (149)$$

Taken collectively, these inequalities imply that

$$\underbrace{\frac{1}{m} \sum_{i=1}^m \mathbf{a}_i \mathbf{a}_i^\top (\mathbf{a}_i^\top \mathbf{x})^2 \mathbf{1}_{\{ |\mathbf{a}_i^\top \mathbf{x}| \leq \gamma_2 \}}}_{:= \mathbf{Y}_2} \preceq \mathbf{Y} \preceq \underbrace{\frac{1}{m} \sum_{i=1}^m \mathbf{a}_i \mathbf{a}_i^\top (\mathbf{a}_i^\top \mathbf{x})^2 \mathbf{1}_{\{ |\mathbf{a}_i^\top \mathbf{x}| \leq \gamma_1 \}}}_{:= \mathbf{Y}_1}. \quad (150)$$

where  $\gamma_1 := \max\{\sqrt{(1+4\varepsilon)\alpha_y^2 + \varepsilon}, \sqrt{\frac{1+4\varepsilon}{1-\varepsilon}\alpha_y}\}$  and  $\gamma_2 := \min\{\sqrt{(1-4\varepsilon)\alpha_y^2 - \varepsilon}, \sqrt{\frac{1-4\varepsilon}{1+\varepsilon}\alpha_y}\}$ . Letting  $\xi \sim \mathcal{N}(0, 1)$ , one can compute

$$\mathbb{E}[\mathbf{Y}_1] = \beta_1 \mathbf{x} \mathbf{x}^\top + \beta_2 \mathbf{I}, \quad \text{and} \quad \mathbb{E}[\mathbf{Y}_2] = \beta_3 \mathbf{x} \mathbf{x}^\top + \beta_4 \mathbf{I}, \quad (151)$$

where  $\beta_1 := \mathbb{E}[\xi^4 \mathbf{1}_{\{|\xi| \leq \gamma_1\}}] - \mathbb{E}[\xi^2 \mathbf{1}_{\{|\xi| \leq \gamma_1\}}]$ ,  $\beta_2 := \mathbb{E}[\xi^2 \mathbf{1}_{\{|\xi| \leq \gamma_1\}}]$ ,  $\beta_3 := \mathbb{E}[\xi^4 \mathbf{1}_{\{|\xi| \leq \gamma_2\}}] - \mathbb{E}[\xi^2 \mathbf{1}_{\{|\xi| \leq \gamma_2\}}]$  and  $\beta_4 := \mathbb{E}[\xi^2 \mathbf{1}_{\{|\xi| \leq \gamma_2\}}]$ . Recognizing that  $\mathbf{a}_i \mathbf{a}_i^\top (\mathbf{a}_i^\top \mathbf{x})^2 \mathbf{1}_{\{ |\mathbf{a}_i^\top \mathbf{x}| \leq c \}}$  can be rewritten as  $\mathbf{b}_i \mathbf{b}_i^\top$  for some sub-Gaussian vector  $\mathbf{b}_i := \mathbf{a}_i (\mathbf{a}_i^\top \mathbf{x}) \mathbf{1}_{\{ |\mathbf{a}_i^\top \mathbf{x}| \leq c \}}$ , we apply standard results on random matrices with non-isotropic sub-Gaussian rows [40, Equation (5.26)] to deduce

$$\|\mathbf{Y}_1 - \mathbb{E}[\mathbf{Y}_1]\| \leq \delta, \quad \|\mathbf{Y}_2 - \mathbb{E}[\mathbf{Y}_2]\| \leq \delta \quad (152)$$

with probability  $1 - \exp(-\Omega(m))$ , provided that  $m/n$  exceeds some large constant. Besides, when  $\varepsilon$  is sufficiently small, one further has  $\|\mathbb{E}[\mathbf{Y}_1] - \mathbb{E}[\mathbf{Y}_2]\| \leq \delta$ . These taken collectively with (150) give

$$\|\mathbf{Y} - \beta_1 \mathbf{x} \mathbf{x}^\top - \beta_2 \mathbf{I}\| \leq 3\delta. \quad (153)$$

Fix some small  $\tilde{\delta} > 0$ . With (153) in place, applying the Davis-Kahan sin  $\Theta$  theorem [56] and taking  $\delta, \varepsilon$  to be sufficiently small, we obtain

$$\text{dist}(\tilde{\mathbf{z}}, \mathbf{x}) \leq \tilde{\delta},$$

where  $\tilde{\mathbf{z}}$  is the leading eigenvector of  $\mathbf{Y}$ . Since  $\mathbf{z}^{(0)} := \lambda_0 \tilde{\mathbf{z}}$ , one derives

$$\begin{aligned} \text{dist}(\mathbf{z}^{(0)}, \mathbf{x}) &\leq \text{dist}(\lambda_0 \tilde{\mathbf{z}}, \tilde{\mathbf{z}}) + \text{dist}(\tilde{\mathbf{z}}, \mathbf{x}) \\ &\leq |\lambda_0 - 1| + \tilde{\delta} \leq \max\{\sqrt{1+2\varepsilon} - 1, 1 - \sqrt{1-2\varepsilon}\} + \tilde{\delta} \end{aligned} \quad (154)$$

as long as  $m/n$  is sufficiently large, where the last inequality follows from (147). Picking  $\tilde{\delta}$  and  $\varepsilon$  to be sufficiently small, we establish the claim.  $\square$

We now justify that the Poisson model (4) satisfies the condition (145). Suppose that  $\mu_i = (\mathbf{a}_i^\top \mathbf{x})^2$  and hence  $y_i \sim \text{Poisson}(\mu_i)$ . It follows from the Chernoff bound that for all  $t \geq 0$ ,

$$\mathbb{P}(y_i - \mu_i \geq \tau) \leq \frac{\mathbb{E}[e^{t y_i}]}{\exp(t(\mu_i + \tau))} = \frac{\exp(\mu_i (e^t - 1))}{\exp(t(\mu_i + \tau))} = \exp(\mu_i (e^t - t - 1) - t\tau),$$

Taking  $\tau = 2\tilde{\varepsilon}\mu_i$  and  $t = \tilde{\varepsilon}$  for any  $0 \leq \tilde{\varepsilon} \leq 1$  gives

$$\mathbb{P}(y_i - \mu_i \geq 2\tilde{\varepsilon}\mu_i) \leq \exp(\mu_i (e^{\tilde{\varepsilon}} - \tilde{\varepsilon} - 1 - 2\tilde{\varepsilon})) \stackrel{(i)}{\leq} \exp(\mu_i (t^2 - 2\tilde{\varepsilon}t)) = \exp(-\mu_i \tilde{\varepsilon}^2),$$

where (i) follows since  $e^t \leq 1 + t + t^2$  ( $0 \leq t \leq 1$ ). In addition, for any  $\tilde{\varepsilon} > 1$ , taking  $t = 1$  we get

$$\mathbb{P}(y_i - \mu_i \geq 2\tilde{\varepsilon}\mu_i) \leq \exp(\mu_i(e^t - t - 1 - 2\tilde{\varepsilon}t)) \leq \exp(-0.5\mu_i\tilde{\varepsilon}).$$

Suppose that  $\|\mathbf{x}\| \gtrsim \log m$ . When  $\mu_i \gtrsim \|\mathbf{x}\|^2$ , setting  $\tilde{\varepsilon} = 0.5\varepsilon < 1$  yields

$$\mathbb{P}(y_i - \mu_i \geq \varepsilon\mu_i) \leq \exp(-\mu_i\varepsilon^2/4) \leq \exp(-\Omega(\varepsilon^2\|\mathbf{x}\|^2)) = \exp(-\Omega(\varepsilon^2 \log^2 m)).$$

When  $\mu_i \lesssim \|\mathbf{x}\|^2$ , letting  $\kappa_i = \mu_i/\|\mathbf{x}\|^2$  and setting  $\tilde{\varepsilon} = \varepsilon/2\kappa_i \gtrsim \varepsilon$ , we obtain

$$\begin{aligned} \mathbb{P}(y_i - \mu_i \geq \varepsilon\|\mathbf{x}\|^2) &= \mathbb{P}(y_i - \mu_i \geq 2\tilde{\varepsilon}\mu_i) \leq \exp(-\min\{\tilde{\varepsilon}, 0.5\} \cdot \tilde{\varepsilon}\mu_i) \\ &= \exp\left(-0.5 \min\{\tilde{\varepsilon}, 0.5\} \varepsilon \|\mathbf{x}\|^2\right) = \exp\left(-\Omega(\varepsilon^2 \log^2 m)\right). \end{aligned}$$

In view of the union bound,

$$\mathbb{P}(\exists i : \eta_i \geq \varepsilon \max\{\|\mathbf{x}\|^2, |\mathbf{a}_i^\top \mathbf{x}|^2\}) \leq m \exp(-\Omega(\varepsilon^2 \log^2 m)) \rightarrow 0. \quad (155)$$

Similarly, applying the same argument on  $-y_i$  we get  $\eta_i \geq -\varepsilon \max\{\|\mathbf{x}\|^2, |\mathbf{a}_i^\top \mathbf{x}|^2\}$  for all  $i$ , which together with (155) establishes the condition (145) with high probability. In conclusion, the claim (146) applies to the Poisson model.

## D Local error contraction with backtracking line search

In this section, we verify the effectiveness of a backtracking line search strategy by showing local error contraction. To keep it concise, we only sketch the proof for the noiseless case, but the proof extends to the noisy case without much difficulty. Also we do not strive to obtain an optimized constant. For concreteness, we prove the following proposition.

**Proposition 4.** *The claim in Proposition 1 continues to hold if  $\alpha_h \geq 6$ ,  $\alpha_z^{\text{ub}} \geq 5$ ,  $\alpha_z^{\text{lb}} \leq 0.1$ ,  $\alpha_p \geq 5$ , and*

$$\|\mathbf{h}\|/\|\mathbf{z}\| \leq \epsilon_{\text{tr}} \quad (156)$$

for some constant  $\epsilon_{\text{tr}} > 0$  independent of  $n$  and  $m$ .

Note that if  $\alpha_h \geq 6$ ,  $\alpha_z^{\text{ub}} \geq 5$  and  $\alpha_z^{\text{lb}} \leq 0.1$ , then the boundary step size  $\mu_0$  given in Proposition 1 satisfies

$$\frac{0.994 - \zeta_1 - \zeta_2 - \sqrt{2/(9\pi)}\alpha_h^{-1}}{2(1.02 + 0.665\alpha_h^{-1})} \geq 0.384.$$

Thus, it suffices to show that the step size obtained by a backtracking line search lies within  $(0, 0.384)$ . For notational convenience, we will set

$$\mathbf{p} := m^{-1} \nabla \ell_{\text{tr}}(\mathbf{z}) \quad \text{and} \quad \mathcal{E}_3^i := \{|\mathbf{a}_i^\top \mathbf{z}| \geq \alpha_z^{\text{lb}} \|\mathbf{z}\| \text{ and } |\mathbf{a}_i^\top \mathbf{p}| \leq \alpha_p \|\mathbf{p}\|\}$$

throughout the rest of the proof. We also impose the assumption

$$\|\mathbf{p}\|/\|\mathbf{z}\| \leq \epsilon \quad (157)$$

for some sufficiently small constant  $\epsilon > 0$ , so that  $|\mathbf{a}_i^\top \mathbf{p}|/|\mathbf{a}_i^\top \mathbf{z}|$  is small for all non-truncated terms. It is self-evident from (79) that in the regime under study, one has

$$\|\mathbf{p}\| \geq 2 \left\{ 1.99 - 2(\zeta_1 + \zeta_2) - \sqrt{8/\pi}(3\alpha_h)^{-1} - o(1) \right\} \|\mathbf{h}\| \geq 3.64 \|\mathbf{h}\|. \quad (158)$$

To start with, consider three scalars  $h$ ,  $b$ , and  $\delta$ . Setting  $b_\delta := \frac{(b+\delta)^2 - b^2}{b^2}$ , we get

$$\begin{aligned} (b+h)^2 \log \frac{(b+\delta)^2}{b^2} - (b+\delta)^2 + b^2 &= (b+h)^2 \log(1+b_\delta) - b^2 b_\delta \\ &\stackrel{(i)}{\leq} (b+h)^2 \{b_\delta - 0.4875b_\delta^2\} - b^2 b_\delta = ((b+h)^2 - b^2)b_\delta - 0.4875(b+h)^2 b_\delta^2 \\ &= h\delta(2+h/b)(2+\delta/b) - 0.4875(1+h/b)^2 |\delta(2+\delta/b)|^2 \\ &= 4h\delta + \frac{2h^2\delta}{b} + \frac{2h\delta^2}{b} + \frac{h^2\delta^2}{b^2} - 0.4875\delta^2 \left(1 + \frac{h}{b}\right)^2 \left(2 + \frac{\delta}{b}\right)^2, \end{aligned} \quad (159)$$



where (i) follows from the inequality  $\log(1+x) \leq x - 0.4875x^2$  for sufficiently small  $x$ . To further simplify the bound, observe that

$$\begin{aligned} \delta^2 \left(1 + \frac{h}{b}\right)^2 \left(2 + \frac{\delta}{b}\right)^2 &\geq 4\delta^2 \left(1 + \frac{h}{b}\right)^2 + \delta^2 \left(1 + \frac{h}{b}\right)^2 \frac{4\delta}{b} \\ \text{and } \frac{2h\delta^2}{b} + \frac{h^2\delta^2}{b^2} &= \left(\left(1 + \frac{h}{b}\right)^2 - 1\right) \delta^2. \end{aligned}$$

Plugging these two identities into (159) yields

$$\begin{aligned} (159) &\leq 4h\delta + \frac{2h^2\delta}{b} - \left(0.95 \left(1 + \frac{h}{b}\right)^2 + 1\right) \delta^2 - 0.4875\delta^2 \left(1 + \frac{h}{b}\right)^2 \frac{4\delta}{b} \\ &\leq 4h\delta - 1.95\delta^2 + \frac{2h^2|\delta|}{|b|} + \frac{1.9|h|}{|b|}\delta^2 + \frac{1.95|\delta^3|}{|b|} \left(1 + \frac{h}{b}\right)^2. \end{aligned}$$

Replacing respectively  $b$ ,  $\delta$ , and  $h$  with  $\mathbf{a}_i^\top \mathbf{z}$ ,  $\tau \mathbf{a}_i^\top \mathbf{p}$ , and  $-\mathbf{a}_i^\top \mathbf{h}$ , one sees that the log-likelihood  $\ell_i(\mathbf{z}) = y_i \log(|\mathbf{a}_i^\top \mathbf{z}|^2) - |\mathbf{a}_i^\top \mathbf{z}|^2$  obeys

$$\begin{aligned} \ell_i(\mathbf{z} + \tau \mathbf{p}) - \ell_i(\mathbf{z}) &= y_i \log \frac{|\mathbf{a}_i^\top (\mathbf{z} + \tau \mathbf{p})|^2}{|\mathbf{a}_i^\top \mathbf{z}|^2} - |\mathbf{a}_i^\top (\mathbf{z} + \tau \mathbf{p})|^2 + |\mathbf{a}_i^\top \mathbf{z}|^2 \\ &\leq \underbrace{-4\tau (\mathbf{a}_i^\top \mathbf{h}) (\mathbf{a}_i^\top \mathbf{p})}_{:=I_{1,i}} - \underbrace{1.95\tau^2 (\mathbf{a}_i^\top \mathbf{p})^2}_{:=I_{2,i}} + \underbrace{\frac{2\tau (\mathbf{a}_i^\top \mathbf{h})^2 |\mathbf{a}_i^\top \mathbf{p}|}{|\mathbf{a}_i^\top \mathbf{z}|}}_{:=I_{3,i}} + \underbrace{\frac{1.9\tau^2 |\mathbf{a}_i^\top \mathbf{h}|}{|\mathbf{a}_i^\top \mathbf{z}|} (\mathbf{a}_i^\top \mathbf{p})^2}_{:=I_{4,i}} \\ &\quad + \underbrace{\frac{1.95\tau^3 |\mathbf{a}_i^\top \mathbf{p}|^3}{|\mathbf{a}_i^\top \mathbf{z}|} \left(1 - \frac{\mathbf{a}_i^\top \mathbf{h}}{\mathbf{a}_i^\top \mathbf{z}}\right)^2}_{:=I_{5,i}}. \end{aligned}$$

The next step is then to bound each of these terms separately. Most of the following bounds are straightforward consequences from [6, Lemma 3.1] combined with the truncation rule. For the first term, applying the AM-GM inequality we get

$$\begin{aligned} \frac{1}{m} \sum_{i=1}^m I_{1,i} \mathbf{1}_{\mathcal{E}_3^i} &\leq \frac{4\tau}{3.64m} \sum_{i=1}^m \left\{ \frac{3.64^2}{2} (\mathbf{a}_i^\top \mathbf{h})^2 + \frac{1}{2} (\mathbf{a}_i^\top \mathbf{p})^2 \right\} \\ &\leq \frac{4\tau(1+\delta)}{3.64} \left\{ \frac{3.64^2}{2} \|\mathbf{h}\|^2 + \frac{1}{2} \|\mathbf{p}\|^2 \right\}. \end{aligned}$$

Secondly, it follows from Lemma 4 that

$$\frac{1}{m} \sum_{i=1}^m I_{2,i} \mathbf{1}_{\mathcal{E}_3^i} = -1.95\tau^2 \frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{p})^2 \mathbf{1}_{\mathcal{E}_3^i} \leq -1.95 \left(1 - \tilde{\zeta}_1 - \tilde{\zeta}_2\right) \tau^2 \|\mathbf{p}\|^2,$$

where  $\tilde{\zeta}_1 := \max\{\mathbb{E}[\xi^2 \mathbf{1}_{\{|\xi| \leq \sqrt{1.01}\alpha_z^{\text{lb}}\}}], \mathbb{E}[\mathbf{1}_{\{|\xi| \leq \sqrt{1.01}\alpha_z^{\text{lb}}\}}]\}$  and  $\tilde{\zeta}_2 := \mathbb{E}[\xi^2 \mathbf{1}_{\{|\xi| > \sqrt{0.99}\alpha_h\}}]$ . The third term is controlled by

$$\frac{1}{m} \sum_{i=1}^m I_{3,i} \mathbf{1}_{\mathcal{E}_3^i} \leq 2\tau \frac{\alpha_p \|\mathbf{p}\|}{\alpha_z^{\text{lb}} \|\mathbf{z}\|} \left\{ \frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h})^2 \right\} \lesssim \tau \epsilon \|\mathbf{h}\|^2.$$

Fourthly, it arises from the AM-GM inequality that

$$\begin{aligned} \frac{1}{m} \sum_{i=1}^m I_{4,i} \mathbf{1}_{\mathcal{E}_3^i} &\leq \frac{1.9\tau^2 \alpha_p \|\mathbf{p}\|}{\alpha_z^{\text{lb}} \|\mathbf{z}\|} \frac{1}{m} \sum_{i=1}^m |\mathbf{a}_i^\top \mathbf{h}| |\mathbf{a}_i^\top \mathbf{p}| \lesssim \epsilon \tau^2 \frac{1}{m} \sum_{i=1}^m \left\{ 2 |\mathbf{a}_i^\top \mathbf{h}|^2 + \frac{1}{8} |\mathbf{a}_i^\top \mathbf{p}|^2 \right\} \\ &\lesssim \epsilon \tau^2 \|\mathbf{p}\|^2. \end{aligned}$$

Finally, the last term is bounded by

$$\begin{aligned} \frac{1}{m} \sum_{i=1}^m I_{5,i} \mathbf{1}_{\mathcal{E}_3^i} &\leq \frac{1}{m} \sum_{i=1}^m \frac{1.95\tau^3 |\mathbf{a}_i^\top \mathbf{p}|^3}{|\mathbf{a}_i^\top \mathbf{z}|} \left( \frac{\mathbf{a}_i^\top \mathbf{x}}{\mathbf{a}_i^\top \mathbf{z}} \right)^2 \leq \frac{1.95\tau^3 \alpha_p^3 \|\mathbf{p}\|^3}{(\alpha_z^{\text{lb}})^3 \|\mathbf{z}\|^3} \frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{x})^2 \\ &\lesssim \tau^3 \epsilon \frac{\|\mathbf{x}\|^2}{\|\mathbf{z}\|^2} \|\mathbf{p}\|^2. \end{aligned}$$

Under the hypothesis (158), we can further derive  $\frac{1}{m} \sum_{i=1}^m I_{1,i} \mathbf{1}_{\mathcal{E}_3^i} \leq \tau (1.1 + \delta) \|\mathbf{p}\|^2$ . Putting all the above bounds together yields that the truncated objective function is majorized by

$$\begin{aligned} \frac{1}{m} \sum_{i=1}^m \{\ell_i(\mathbf{z} + \tau\mathbf{p}) - \ell_i(\mathbf{z})\} \mathbf{1}_{\mathcal{E}_3^i} &\leq \frac{1}{m} \sum_{i=1}^m (I_{1,i} + I_{2,i} + I_{3,i} + I_{4,i} + I_{5,i}) \mathbf{1}_{\mathcal{E}_3^i} \\ &\leq \tau (1.1 + \delta) \|\mathbf{p}\|^2 - 1.95 \left(1 - \tilde{\zeta}_1 - \tilde{\zeta}_2\right) \tau^2 \|\mathbf{p}\|^2 + \tau \tilde{\epsilon} \|\mathbf{p}\|^2 \\ &= \left\{ \tau (1.1 + \delta) - 1.95 \left(1 - \tilde{\zeta}_1 - \tilde{\zeta}_2\right) \tau^2 + \tau \tilde{\epsilon} \right\} \|\mathbf{p}\|^2 \end{aligned} \quad (160)$$

for some constant  $\tilde{\epsilon} > 0$  that is linear in  $\epsilon$ .

Note that the backtracking line search seeks a point satisfying  $\frac{1}{m} \sum_{i=1}^m \{\ell_i(\mathbf{z} + \tau\mathbf{p}) - \ell_i(\mathbf{z})\} \mathbf{1}_{\mathcal{E}_3^i} \geq \frac{1}{2} \tau \|\mathbf{p}\|^2$ . Given the above majorization (160), this search criterion is satisfied only if

$$\tau/2 \leq \tau (1.1 + \delta) - 1.95(1 - \tilde{\zeta}_1 - \tilde{\zeta}_2) \tau^2 + \tau \tilde{\epsilon}$$

or, equivalently,

$$\tau \leq \frac{0.6 + \delta + \tilde{\epsilon}}{1.95(1 - \tilde{\zeta}_1 - \tilde{\zeta}_2)} := \tau_{\text{ub}}.$$

Taking  $\delta$  and  $\tilde{\epsilon}$  to be sufficiently small, we see that  $\tau \leq \tau_{\text{ub}} \leq 0.384$ , provided that  $\alpha_z^{\text{lb}} \leq 0.1$ ,  $\alpha_z^{\text{ub}} \geq 5$ ,  $\alpha_h \geq 6$ , and  $\alpha_p \geq 5$ .

Using very similar arguments, one can also show that  $\frac{1}{m} \sum_{i=1}^m \{\ell_i(\mathbf{z} + \tau\mathbf{p}) - \ell_i(\mathbf{z})\} \mathbf{1}_{\mathcal{E}_3^i}$  is minorized by a similar quadratic function, which combined with the stopping criterion  $\frac{1}{m} \sum_{i=1}^m \{\ell_i(\mathbf{z} + \tau\mathbf{p}) - \ell_i(\mathbf{z})\} \mathbf{1}_{\mathcal{E}_3^i} \geq \frac{1}{2} \tau \|\mathbf{p}\|^2$  suggests that  $\tau$  is bounded away from 0. We omit this part for conciseness.

## Acknowledgements

E. C. is partially supported by NSF under grant CCF-0963835 and by the Math + X Award from the Simons Foundation. Y. C. is supported by the same NSF grant. We thank Carlos Sing-Long and Weijie Su for their constructive comments about an early version of the manuscript. E. C. is grateful to Xiaodong Li and Mahdi Soltanolkotabi for many discussions about Wirtinger flows. We thank the anonymous reviewers for helpful comments.

## References

- [1] 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, April 2015.
- [2] A. Ben-Tal and A. Nemirovski. *Lectures on modern convex optimization*, volume 2. SIAM, 2001.
- [3] R. W. Gerchberg. A practical algorithm for the determination of phase from image and diffraction plane pictures. *Optik*, 35:237, 1972.
- [4] J. R. Fienup. Phase retrieval algorithms: a comparison. *Applied optics*, 21:2758–2769, 1982.
- [5] Y. Shechtman, Y. C. Eldar, O. Cohen, H.N. Chapman, J. Miao, and M. Segev. Phase retrieval with application to optical imaging. *IEEE Signal Processing Magazine*, 32(3):87–109, May 2015.

- [6] 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):1017–1026, 2013.
- [7] 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.
- [8] I. Waldspurger, A. d’Aspremont, and S. Mallat. Phase recovery, maxcut and complex semidefinite programming. *Mathematical Programming*, 149(1-2):47–81, 2015.
- [9] Y. Chen, Y. Chi, and A. J. Goldsmith. Exact and stable covariance estimation from quadratic sampling via convex programming. *IEEE Transactions on Information Theory*, 61(7):4034–4059, 2015.
- [10] Y. Shechtman, Y. C. Eldar, A. Szameit, and M. Segev. Sparsity based sub-wavelength imaging with partially incoherent light via quadratic compressed sensing. *Optics express*, 19(16), 2011.
- [11] H. Ohlsson, A. Yang, R. Dong, and S. S. Sastry. Compressive phase retrieval from squared output measurements via semidefinite programming. *Neural Information Processing Systems (NIPS)*, 2011.
- [12] L. Demanet and P. Hand. Stable optimizationless recovery from phaseless linear measurements. *Journal of Fourier Analysis and Applications*, 20(1):199–221, 2014.
- [13] X. Li and V. Voroninski. Sparse signal recovery from quadratic measurements via convex programming. *SIAM Journal on Mathematical Analysis*, 45(5):3019–3033, 2013.
- [14] T. Cai and A. Zhang. ROP: Matrix recovery via rank-one projections. *The Annals of Statistics*, 43(1):102–138, 2015.
- [15] K. Jaganathan, S. Oymak, and B. Hassibi. Recovery of sparse 1-D signals from the magnitudes of their Fourier transform. In *IEEE ISIT*, pages 1473–1477, 2012.
- [16] Y. Chen, X. Yi, and C. Caramanis. A convex formulation for mixed regression with two components: Minimax optimal rates. In *Conf. on Learning Theory*, 2014.
- [17] R. Kueng, H. Rauhut, and U. Terstiege. Low rank matrix recovery from rank one measurements. *arXiv preprint arXiv:1410.6913*, 2014.
- [18] S. Bahmani and J. Romberg. Efficient compressive phase retrieval with constrained sensing vectors. *Advances in Neural Information Processing Systems*, pages 523–531, 2015.
- [19] D. Gross, F. Krahermer, and R. Kueng. A partial derandomization of phaselift using spherical designs. *Journal of Fourier Analysis and Applications*, 21(2):229–266, 2015.
- [20] D. Gross, F. Krahermer, and R. Kueng. Improved recovery guarantees for phase retrieval from coded diffraction patterns. *arXiv preprint arXiv:1402.6286*, 2014.
- [21] P. Netrapalli, P. Jain, and S. Sanghavi. Phase retrieval using alternating minimization. *Advances in Neural Information Processing Systems (NIPS)*, 2013.
- [22] V. Elser. Phase retrieval by iterated projections. *JOSA. A*, 20(1):40–55, 2003.
- [23] P. Schniter and S. Rangan. Compressive phase retrieval via generalized approximate message passing. *IEEE Transactions on Signal Processing*, 63(4):1043–1055, Feb 2015.
- [24] Q. Zheng and J. Lafferty. A convergent gradient descent algorithm for rank minimization and semidefinite programming from random linear measurements. *Advances in Neural Information Processing Systems*, 2015.
- [25] S. Marchesin, Y. Tu, and H. Wu. Alternating projection, ptychographic imaging and phase synchronization. *arXiv:1402.0550*, 2014.

- [26] A. Repetti, E. Chouzenoux, and J-C Pesquet. A nonconvex regularized approach for phase retrieval. *International Conference on Image Processing*, pages 1753–1757, 2014.
- [27] Y. Shechtman, A. Beck, and Y. C. Eldar. GESPAR: Efficient phase retrieval of sparse signals. *IEEE Transactions on Signal Processing*, 62(4):928–938, 2014.
- [28] C. D. White, R. Ward, and S. Sanghavi. The local convexity of solving quadratic equations. *arXiv preprint arXiv:1506.07868*, 2015.
- [29] M. Soltanolkotabi. *Algorithms and Theory for Clustering and Nonconvex Quadratic Programming*. PhD thesis, Stanford University, 2014.
- [30] E. J. Candès, X. Li, and M. Soltanolkotabi. <http://www-bcf.usc.edu/~soltanol/PRWF.html>, 2014.
- [31] J. Nocedal and S. J. Wright. *Numerical Optimization (2nd edition)*. Springer, 2006.
- [32] L. N. Trefethen and D. Bau III. *Numerical linear algebra*, volume 50. SIAM, 1997.
- [33] E. J. Candès, X. Li, and M. Soltanolkotabi. Phase retrieval from coded diffraction patterns. *to appear in Applied and Computational Harmonic Analysis*, 2014.
- [34] Y. Chen and E. J. Candès. Supplemental materials for: “solving random quadratic systems of equations is nearly as easy as solving linear systems”. May 2015.
- [35] K. Wei. Phase retrieval via Kaczmarz methods. *arXiv preprint arXiv:1502.01822*, 2015.
- [36] B. Alexeev, A. S. Bandeira, M. Fickus, and D. G. Mixon. Phase retrieval with polarization. *SIAM Journal on Imaging Sciences*, 7(1):35–66, 2014.
- [37] R. Balan, B. Bodmann, P. Casazza, and D. Edidin. Painless reconstruction from magnitudes of frame coefficients. *Journal of Fourier Analysis and Applications*, 15(4):488–501, 2009.
- [38] Y. C. Eldar and S. Mendelson. Phase retrieval: Stability and recovery guarantees. *Applied and Computational Harmonic Analysis*, 36(3):473–494, 2014.
- [39] A. S. Bandeira, J. Cahill, D. G. Mixon, and A. A. Nelson. Saving phase: Injectivity and stability for phase retrieval. *Applied and Computational Harmonic Analysis*, 37(1):106–125, 2014.
- [40] R. Vershynin. Introduction to the non-asymptotic analysis of random matrices. *Compressed Sensing, Theory and Applications*, pages 210 – 268, 2012.
- [41] A. B. Tsybakov and V. Zaiats. *Introduction to nonparametric estimation*, volume 11. Springer, 2009.
- [42] 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.
- [43] P. Jain, P. Netrapalli, and S. Sanghavi. Low-rank matrix completion using alternating minimization. In *ACM symposium on Theory of computing*, pages 665–674, 2013.
- [44] M. Hardt and M. Wootters. Fast matrix completion without the condition number. *Conference on Learning Theory*, pages 638 – 678, 2014.
- [45] R. Sun and Z. Luo. Guaranteed matrix completion via non-convex factorization. *arXiv:1411.8003*, 2014.
- [46] P. Netrapalli, U. Niranjan, S. Sanghavi, A. Anandkumar, and P. Jain. Non-convex robust PCA. In *Advances in Neural Information Processing Systems*, pages 1107–1115, 2014.
- [47] J. Sun, Q. Qu, and J. Wright. Complete dictionary recovery using nonconvex optimization. *International Conference on Machine Learning*, pages 2351–2360, 2015.
- [48] J. Sun, Q. Qu, and J. Wright. When are nonconvex problems not scary? *arXiv preprint arXiv:1510.06096*, 2015.

- [49] S. Arora, R. Ge, T. Ma, and A. Moitra. Simple, efficient, and neural algorithms for sparse coding. *arXiv:1503.00778*, 2015.
- [50] X. Yi, C. Caramanis, and S. Sanghavi. Alternating minimization for mixed linear regression. *International Conference on Machine Learning*, June 2014.
- [51] S. Balakrishnan, M. J. Wainwright, and B. Yu. Statistical guarantees for the EM algorithm: From population to sample-based analysis. *arXiv preprint arXiv:1408.2156*, 2014.
- [52] A. S. Bandeira, M. Charikar, A. Singer, and A. Zhu. Multireference alignment using semidefinite programming. In *Conference on Innovations in theoretical computer science*, pages 459–470, 2014.
- [53] Q. Huang and L. Guibas. Consistent shape maps via semidefinite programming. *Computer Graphics Forum*, 32(5):177–186, 2013.
- [54] Y. Chen, L. J. Guibas, and Q. Huang. Near-optimal joint optimal matching via convex relaxation. *International Conference on Machine Learning (ICML)*, pages 100 – 108, June 2014.
- [55] M. Mitzenmacher and E. Upfal. *Probability and computing*. Cambridge University Press, 2005.
- [56] Chandler Davis and William Morton Kahan. The rotation of eigenvectors by a perturbation. iii. *SIAM Journal on Numerical Analysis*, 7(1):1–46, 1970.