

Supplemental Materials for: “Solving Random Quadratic Systems of Equations Is Nearly as Easy as Solving Linear Systems”

Yuxin Chen ^{*} Emmanuel J. Candès ^{*†}

1 Main Theorems

For convenience of presentation, we repeat the main results as follows. To begin with, the noiseless model, the general additive noise model, and the Poisson noise model are given respectively as follows.

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

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

$$\text{(Poisson noise:)} \quad y_i \stackrel{\text{ind.}}{\sim} \text{Poisson}(|\langle \mathbf{a}_i, \mathbf{x} \rangle|^2), \quad i = 1, \dots, m. \quad (3)$$

Theorem 1 (Exact recovery). *Consider the noiseless case (1) with an arbitrary signal $\mathbf{x} \in \mathbb{R}^n$. Suppose that the step size μ_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) obey*

$$\text{dist}(\mathbf{z}^{(t)}, \mathbf{x}) \leq \nu(1 - \rho)^t \|\mathbf{x}\|, \quad \forall t \in \mathbb{N}, \quad (4)$$

provided that

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

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

$$m \geq c_0 n, \quad \mu \leq \mu_0, \quad \text{and} \quad \|\boldsymbol{\eta}\|_\infty \leq c_1 \|\mathbf{x}\|^2, \quad (5)$$

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

$$\text{dist}(\mathbf{z}^{(t)}, \mathbf{x}) \lesssim \frac{\|\boldsymbol{\eta}\|}{\sqrt{m} \|\mathbf{x}\|} + (1 - \rho)^t \|\mathbf{x}\|, \quad \forall t \in \mathbb{N} \quad (6)$$

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.

In particular, under the Poisson noise model, there exists an event of probability at least $1 - c_2 \exp(-c_3 m)$ on which

$$\mathbb{P}\left\{\text{dist}(\mathbf{z}^{(t)}, \mathbf{x}) \lesssim 1 + (1 - \rho)^t \|\mathbf{x}\|, \quad \forall t \in \mathbb{N} \mid \{\mathbf{a}_i\}_{1 \leq i \leq m}\right\} \rightarrow 1. \quad (7)$$

holds for all $\mathbf{x} \in \mathbb{R}^n$ satisfying $\|\mathbf{x}\| \geq \log^{1.5} m$. In what follows, we prove the above two theorems for a broader range of algorithmic parameters, as summarized in Table 1.

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.

^{*}Department of Statistics, Stanford University, Stanford, CA 94305, U.S.A.

[†]Department of Mathematics, Stanford University, Stanford, CA 94305, U.S.A.

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

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

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

$$0 < \alpha_z^{\text{lb}} \leq 0.1, \quad \alpha_z^{\text{ub}} \geq 5, \quad \alpha_h \geq 6, \quad \alpha_y \geq 3, \quad \text{and} \quad \alpha_p \geq 5. \quad (9)$$

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

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 κ independent of n , and n is sufficiently large. For any $K \geq \log^{1.5} m$, define¹*

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

With probability approaching one, the minimax risk under the Poisson model (3) obeys

$$\inf_{\hat{\mathbf{x}}} \sup_{\mathbf{x} \in \Upsilon(K)} \mathbb{E}[\text{dist}(\hat{\mathbf{x}}, \mathbf{x}) \mid \{\mathbf{a}_i\}_{1 \leq i \leq m}] \geq \frac{\varepsilon_1}{\sqrt{\kappa}}, \quad (10)$$

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

2 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 follows similar argument as in [1, Section 7.8], 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 large constant.

The remaining section then boils down to establishing convergence for the gradient flow stage. To this end, we recall a (local) regularity condition given in [1], which has been shown to be 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 $\text{RC}(\mu, \lambda, \epsilon)$, if

$$\left\langle \mathbf{h}, -\frac{1}{m} \nabla \ell_{\text{tr}}(\mathbf{z}) \right\rangle \geq \frac{\mu}{2} \left\| \frac{1}{m} \nabla \ell_{\text{tr}}(\mathbf{z}) \right\|^2 + \frac{\lambda}{2} \|\mathbf{h}\|^2 \quad (11)$$

holds for all \mathbf{z} obeying $\|\mathbf{z} - \mathbf{x}\| \leq \epsilon \|\mathbf{x}\|$, where $0 < \epsilon < 1$ is some constant. Such an ϵ -ball around \mathbf{x} is

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

sometimes referred to as a *basin of attraction*. Formally, under $\text{RC}(\mu, \lambda, \epsilon)$, a little algebra gives

$$\begin{aligned}
\text{dist}^2\left(\mathbf{z} + \frac{\mu}{m} \nabla \ell_{\text{tr}}(\mathbf{z}), \mathbf{x}\right) &\leq \left\| \mathbf{z} + \frac{\mu}{m} \nabla \ell_{\text{tr}}(\mathbf{z}) - \mathbf{x} \right\|^2 \\
&= \|\mathbf{h}\|^2 + \left\| \frac{\mu}{m} \nabla \ell_{\text{tr}}(\mathbf{z}) \right\|^2 + 2\mu \left\langle \mathbf{h}, \frac{1}{m} \nabla \ell_{\text{tr}}(\mathbf{z}) \right\rangle \\
&\leq \|\mathbf{h}\|^2 + \left\| \frac{\mu}{m} \nabla \ell_{\text{tr}}(\mathbf{z}) \right\|^2 - \mu^2 \left\| \frac{1}{m} \nabla \ell_{\text{tr}}(\mathbf{z}) \right\|^2 - \mu\lambda \|\mathbf{h}\|^2 \\
&= (1 - \mu\lambda) \text{dist}^2(\mathbf{z}, \mathbf{x})
\end{aligned} \tag{12}$$

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 $\text{RC}(\mu, \lambda, \epsilon)$ holds for some nonzero μ and λ . This is stated in the following proposition.

Proposition 1 (Local error contraction). *Consider the noiseless case (1). Under the condition (8), there exist some universal constants $0 < \rho_0 < 1$ and $c_0, c_1, c_2 > 0$ such that with probability exceeding $1 - c_1 \exp(-c_2 m)$,*

$$\text{dist}^2\left(\mathbf{z} + \frac{\mu}{m} \nabla \ell_{\text{tr}}(\mathbf{z}), \mathbf{x}\right) \leq (1 - \rho_0) \text{dist}^2(\mathbf{z}, \mathbf{x}) \tag{13}$$

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

$$\frac{\text{dist}(\mathbf{z}, \mathbf{x})}{\|\mathbf{z}\|} \leq \min \left\{ \frac{1}{11}, \frac{\alpha_z^{\text{lb}}}{3\alpha_h}, \frac{\alpha_z^{\text{lb}}}{6}, \frac{5.7(\alpha_z^{\text{lb}})^2}{2\alpha_z^{\text{ub}} + \alpha_z^{\text{lb}}} \right\}, \tag{14}$$

provided that $m \geq c_0 n$ and that μ is some constant obeying $0 < \mu \leq \mu_0 := \frac{0.994 - \zeta_1 - \zeta_2 - \sqrt{2/(9\pi)\alpha_h^{-1}}}{2(1.02 + 0.665/\alpha_h)}$.

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 before, under the hypothesis $\text{RC}(\mu, \lambda, \epsilon)$ one can conclude

$$\text{dist}^2\left(\mathbf{z} + \frac{\mu}{m} \nabla \ell_{\text{tr}}(\mathbf{z}), \mathbf{x}\right) \leq (1 - \mu\lambda) \text{dist}^2(\mathbf{z}, \mathbf{x}), \quad \forall (\mathbf{z}, \mathbf{x}) \text{ with } \text{dist}(\mathbf{z}, \mathbf{x}) \leq \epsilon. \tag{15}$$

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.

2.1 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 :

$$\mathcal{E}_1^i(\mathbf{z}) := \left\{ \alpha_z^{\text{lb}} \leq \frac{|\mathbf{a}_i^\top \mathbf{z}|}{\|\mathbf{z}\|} \leq \alpha_z^{\text{ub}} \right\}; \tag{16}$$

$$\mathcal{E}_2^i(\mathbf{z}) := \left\{ |y_i - |\mathbf{a}_i^\top \mathbf{z}|^2| \leq \frac{\alpha_h}{m} \|\mathbf{y} - \mathcal{A}(\mathbf{z}\mathbf{z}^\top)\|_1 \frac{|\mathbf{a}_i^\top \mathbf{z}|}{\|\mathbf{z}\|} \right\}, \tag{17}$$

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 $1 - C \exp(-c_1 \zeta^2 m)$,*

$$0.9(1 - \zeta) \|\mathbf{M}\|_{\text{F}} \leq \frac{1}{m} \|\mathcal{A}(\mathbf{M})\|_1 \leq (1 + \zeta) \sqrt{2} \|\mathbf{M}\|_{\text{F}} \tag{18}$$

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.

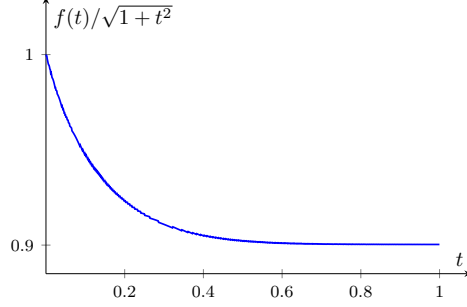


Figure 1: $\frac{f(t)}{\sqrt{1+t^2}}$ as a function of t .

Proof. Since [2, 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 [2, Lemma 3.2] that with high probability,

$$\frac{1}{m} \|\mathcal{A}(\mathbf{M})\|_1 \geq (1 - \zeta) f(t), \quad (19)$$

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. 1. The case where $t \in [-1, 0]$ is an immediate consequence from [2, Lemma 3.1]. \square

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

$$\sqrt{2-4\delta} \|\mathbf{z} - \mathbf{x}\| \|\mathbf{z}\| \leq \|\mathbf{x}\mathbf{x}^\top - \mathbf{z}\mathbf{z}^\top\|_F \leq (2+\delta) \|\mathbf{z} - \mathbf{x}\| \|\mathbf{z}\|. \quad (20)$$

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

$$\begin{aligned} \|\mathbf{x}\mathbf{x}^\top - \mathbf{z}\mathbf{z}^\top\|_F^2 &= \|\mathbf{h}\mathbf{z}^\top + \mathbf{z}\mathbf{h}^\top + \mathbf{h}\mathbf{h}^\top\|_F^2 \\ &= \|\mathbf{h}\mathbf{z}^\top + \mathbf{z}\mathbf{h}^\top\|_F^2 + \|\mathbf{h}\|^4 - 2\langle \mathbf{h}\mathbf{z}^\top + \mathbf{z}\mathbf{h}^\top, \mathbf{h}\mathbf{h}^\top \rangle \\ &= 2\|\mathbf{z}\|^2 \|\mathbf{h}\|^2 + 2|\mathbf{h}^\top \mathbf{z}|^2 + \|\mathbf{h}\|^4 - 2\|\mathbf{h}\|^2 (\mathbf{h}^\top \mathbf{z} + \mathbf{z}^\top \mathbf{h}). \end{aligned}$$

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

$$2\|\mathbf{z}\|^2 \|\mathbf{h}\|^2 - 4\|\mathbf{z}\| \|\mathbf{h}\|^3 \leq \|\mathbf{x}\mathbf{x}^\top - \mathbf{z}\mathbf{z}^\top\|_F^2 \leq 4\|\mathbf{z}\|^2 \|\mathbf{h}\|^2 + 4\|\mathbf{h}\|^3 \|\mathbf{z}\| + \|\mathbf{h}\|^4, \quad (21)$$

$$\Rightarrow \sqrt{(2\|\mathbf{z}\| - 4\|\mathbf{h}\|) \|\mathbf{z}\|} \cdot \|\mathbf{h}\| \leq \|\mathbf{x}\mathbf{x}^\top - \mathbf{z}\mathbf{z}^\top\|_F \leq (2\|\mathbf{z}\| + \|\mathbf{h}\|) \cdot \|\mathbf{h}\| \quad (22)$$

as claimed. \square

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

$$1.15 \|\mathbf{z} - \mathbf{x}\| \|\mathbf{z}\| \leq \frac{1}{m} \|\mathcal{A}(\mathbf{x}\mathbf{x}^\top - \mathbf{z}\mathbf{z}^\top)\|_1 \leq 3 \|\mathbf{z} - \mathbf{x}\| \|\mathbf{z}\| \quad (23)$$

holds simultaneously for all \mathbf{z} and \mathbf{x} satisfying $\|\mathbf{h}\| \leq \frac{1}{11} \|\mathbf{z}\|$. Conditional on (23), the inclusion

$$\mathcal{E}_3^i \subseteq \mathcal{E}_2^i \subseteq \mathcal{E}_4^i \quad (24)$$

holds with respect to the following events

$$\mathcal{E}_3^i : = \{ \left| |\mathbf{a}_i^\top \mathbf{x}|^2 - |\mathbf{a}_i^\top \mathbf{z}|^2 \right| \leq 1.15 \alpha_h \|\mathbf{h}\| \cdot |\mathbf{a}_i^\top \mathbf{z}| \}, \quad (25)$$

$$\mathcal{E}_4^i : = \{ \left| |\mathbf{a}_i^\top \mathbf{x}|^2 - |\mathbf{a}_i^\top \mathbf{z}|^2 \right| \leq 3 \alpha_h \|\mathbf{h}\| \cdot |\mathbf{a}_i^\top \mathbf{z}| \}. \quad (26)$$

The point of introducing these new events is that \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*

$$\mathcal{D}_\gamma^i := \{ |\mathbf{a}_i^\top \mathbf{x}|^2 - |\mathbf{a}_i^\top \mathbf{z}|^2 \leq \gamma \|\mathbf{h}\| |\mathbf{a}_i^\top \mathbf{z}| \}, \quad (27)$$

$$\mathcal{D}_\gamma^{i,1} := \left\{ \frac{|\mathbf{a}_i^\top \mathbf{h}|}{\|\mathbf{h}\|} \leq \gamma \right\}, \quad (28)$$

$$\text{and } \mathcal{D}_\gamma^{i,2} := \left\{ \left| \frac{\mathbf{a}_i^\top \mathbf{h}}{\|\mathbf{h}\|} - \frac{2\mathbf{a}_i^\top \mathbf{z}}{\|\mathbf{h}\|} \right| \leq \gamma \right\}. \quad (29)$$

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 $\frac{\|\mathbf{h}\|}{\|\mathbf{z}\|} \leq \frac{\alpha_z^{\text{lb}}}{\gamma}$, then the following inclusion holds

$$\left(\mathcal{D}_\gamma^{i,1} \cap \mathcal{E}_1^i \right) \cup \left(\mathcal{D}_\gamma^{i,2} \cap \mathcal{E}_1^i \right) \subseteq \mathcal{D}_\gamma^i \cap \mathcal{E}_1^i \subseteq \left(\mathcal{D}_\gamma^{i,1} \cap \mathcal{E}_1^i \right) \cup \left(\mathcal{D}_\gamma^{i,2} \cap \mathcal{E}_1^i \right). \quad (30)$$

2.2 Proof of the regularity condition

By definition, one step towards proving the regularity condition (11) is to control the norm of the truncated 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 truncation rule \mathcal{E}_1^i as well as the inclusion property (24) that

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

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

$$\frac{1}{m} \|\nabla \ell_{\text{tr}}(\mathbf{z})\| = \frac{1}{m} \|\mathbf{A}^\top \mathbf{v}\| \leq \frac{1}{m} \|\mathbf{A}\| \cdot \|\mathbf{v}\| \lesssim \|\mathbf{h}\|. \quad (31)$$

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

$$-\left\langle \mathbf{h}, \frac{1}{m} \nabla \ell_{\text{tr}}(\mathbf{z}) \right\rangle \gtrsim \|\mathbf{h}\|^2, \quad (32)$$

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 (8), if $m > c_1 n$, then with probability exceeding $1 - C \exp(-c_0 m)$,*

$$-\left\langle \mathbf{h}, \frac{1}{m} \nabla \ell_{\text{tr}}(\mathbf{z}) \right\rangle \geq 2 \left\{ 1.99 - 2(\zeta_1 + \zeta_2) - \sqrt{8/(9\pi)} \alpha_h^{-1} - \epsilon \right\} \|\mathbf{h}\|^2 \quad (33)$$

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

$$\frac{\|\mathbf{h}\|}{\|\mathbf{z}\|} \leq \min \left\{ \frac{1}{11}, \frac{\alpha_z^{\text{lb}}}{3\alpha_h}, \frac{\alpha_z^{\text{lb}}}{6}, \frac{5.7 (\alpha_z^{\text{lb}})^2}{2\alpha_z^{\text{ub}} + \alpha_z^{\text{lb}}} \right\}. \quad (34)$$

Here, $c_0, c_1, C > 0$ are some universal constants, and ζ_1 and ζ_2 are defined in (8).

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

$$\begin{aligned} -\frac{1}{2m} \nabla \ell_{\text{tr}}(\mathbf{z}) &= \frac{1}{m} \sum_{i=1}^m \frac{(\mathbf{a}_i^\top \mathbf{z})^2 - (\mathbf{a}_i^\top \mathbf{x})^2}{\mathbf{a}_i^\top \mathbf{z}} \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i} \\ &= \frac{1}{m} \sum_{i=1}^m 2(\mathbf{a}_i^\top \mathbf{h}) \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i} - \frac{1}{m} \sum_{i=1}^m \frac{(\mathbf{a}_i^\top \mathbf{h})^2}{\mathbf{a}_i^\top \mathbf{z}} \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i}. \end{aligned} \quad (35)$$

One would expect the contribution of the second term of (35) (which is a second-order quantity) to be small as $\|\mathbf{h}\| / \|\mathbf{z}\|$ decreases.

To facilitate analysis, we rewrite (35) in terms of the more convenient events $\mathcal{D}_{\gamma}^{i,1}$ and $\mathcal{D}_{\gamma}^{i,2}$. Specifically, the inclusion property (24) together with Lemma 3 reveals that

$$\mathcal{D}_{\gamma_3}^{i,1} \cap \mathcal{E}_1^i \subseteq \mathcal{E}_3^i \cap \mathcal{E}_1^i \subseteq \mathcal{E}_2^i \cap \mathcal{E}_1^i \subseteq \mathcal{E}_4^i \cap \mathcal{E}_1^i \subseteq (\mathcal{D}_{\gamma_4}^{i,1} \cup \mathcal{D}_{\gamma_4}^{i,2}) \cap \mathcal{E}_1^i, \quad (36)$$

where the parameters γ_3, γ_4 are given by

$$\gamma_3 := 0.476\alpha_h, \quad \text{and} \quad \gamma_4 := 3\alpha_h. \quad (37)$$

This taken collectively with the identity (35) leads to a lower estimate

$$-\left\langle \frac{1}{2m} \nabla \ell_{\text{tr}}(\mathbf{z}), \mathbf{h} \right\rangle \geq \frac{2}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h})^2 \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{D}_{\gamma_3}^{i,1}} - \frac{1}{m} \sum_{i=1}^m \frac{|\mathbf{a}_i^\top \mathbf{h}|^3}{|\mathbf{a}_i^\top \mathbf{z}|} \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{D}_{\gamma_4}^{i,1}} - \frac{1}{m} \sum_{i=1}^m \frac{|\mathbf{a}_i^\top \mathbf{h}|^3}{|\mathbf{a}_i^\top \mathbf{z}|} \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{D}_{\gamma_4}^{i,2}}, \quad (38)$$

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 truncation step does not discard too many samples, 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 (38) 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 $\mathcal{O}(\|\mathbf{h}\|^3/\|\mathbf{z}\|)$, 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 (38) 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 noticable 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 (16) and (28), respectively. Set

$$\zeta_1 := 1 - \min \left\{ \mathbb{E} \left[\xi^2 \mathbf{1}_{\{\sqrt{1.01}\alpha_z^{\text{lb}} \leq |\xi| \leq \sqrt{0.99}\alpha_z^{\text{ub}}\}} \right], \mathbb{E} \left[\mathbf{1}_{\{\sqrt{1.01}\alpha_z^{\text{lb}} \leq |\xi| \leq \sqrt{0.99}\alpha_z^{\text{ub}}\}} \right] \right\} \quad (39)$$

$$\text{and } \zeta_2 := \mathbb{E} \left[\xi^2 \mathbf{1}_{\{|\xi| > \sqrt{0.99}\gamma\}} \right], \quad (40)$$

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 $1 - C \exp(-c_0 \epsilon^2 m)$,

$$\frac{1}{m} \sum_{i=1}^m |\mathbf{a}_i^\top \mathbf{h}|^2 \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{D}_{\gamma}^{i,1}} \geq (1 - \zeta_1 - \zeta_2 - \epsilon) \|\mathbf{h}\|^2 \quad (41)$$

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 (38). For any fixed $\gamma > 0$, the definition of \mathcal{E}_1^i gives rise to an upper estimate

$$\frac{1}{m} \sum_{i=1}^m \frac{|\mathbf{a}_i^\top \mathbf{h}|^3}{|\mathbf{a}_i^\top \mathbf{z}|} \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{D}_{\gamma}^{i,1}} \leq \frac{1}{\alpha_z^{\text{lb}} \|\mathbf{z}\|} \cdot \frac{1}{m} \sum_{i=1}^m |\mathbf{a}_i^\top \mathbf{h}|^3 \mathbf{1}_{\mathcal{D}_{\gamma}^{i,1}} \leq \frac{(1 + \epsilon) \sqrt{8/\pi} \|\mathbf{h}\|^3}{\alpha_z^{\text{lb}} \|\mathbf{z}\|}, \quad (42)$$

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*

$$\frac{1}{m} \sum_{i=1}^m |\mathbf{a}_i^\top \mathbf{h}|^3 \mathbf{1}_{\mathcal{D}_\gamma^{i,1}} \leq (1 + \epsilon) \sqrt{8/\pi} \|\mathbf{h}\|^3, \quad \forall \mathbf{h} \in \mathbb{R}^n \quad (43)$$

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 (38). 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 γ (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*

$$\frac{1}{m} \sum_{i=1}^m \mathbf{1}_{\{|\mathbf{a}_i^\top \mathbf{h}| \geq \gamma \|\mathbf{h}\|\}} \leq \frac{1}{0.49\gamma} \exp(-0.485\gamma^2) + \frac{\epsilon}{\gamma^2}, \quad \forall \mathbf{h} \in \mathbb{R}^n \setminus \{\mathbf{0}\} \text{ and } \gamma \geq 2 \quad (44)$$

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 (38), we recognize that when $\gamma \leq \frac{\alpha_z^{\text{lb}} \|\mathbf{z}\|}{\|\mathbf{h}\|}$, one has

$$\mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{D}_\gamma^{i,2}} \leq \mathbf{1}_{\{|\mathbf{a}_i^\top \mathbf{h}| \geq \alpha_z^{\text{lb}} \|\mathbf{z}\|\}}. \quad (45)$$

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

$$\frac{|\mathbf{a}_i^\top \mathbf{h}|}{\|\mathbf{h}\|} \geq \frac{2|\mathbf{a}_i^\top \mathbf{z}|}{\|\mathbf{h}\|} - \gamma \geq \frac{2\alpha_z^{\text{lb}} \|\mathbf{z}\|}{\|\mathbf{h}\|} - \gamma \geq \frac{\alpha_z^{\text{lb}} \|\mathbf{z}\|}{\|\mathbf{h}\|}, \quad (46)$$

where the last inequality comes from our assumption on γ . With Lemma 6 in place, (45) immediately gives

$$\begin{aligned} \sum_{i=1}^m \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{D}_\gamma^{i,2}} &\leq \frac{\|\mathbf{h}\|}{0.49\alpha_z^{\text{lb}} \|\mathbf{z}\|} \exp\left(-0.485 \left(\frac{\alpha_z^{\text{lb}} \|\mathbf{z}\|}{\|\mathbf{h}\|}\right)^2\right) + \frac{\epsilon \|\mathbf{h}\|^2}{(\alpha_z^{\text{lb}})^2 \|\mathbf{z}\|^2} \\ &\leq \frac{1}{9800} \left(\frac{\|\mathbf{h}\|}{\alpha_z^{\text{lb}} \|\mathbf{z}\|}\right)^4 + \frac{\epsilon}{(\alpha_z^{\text{lb}})^2} \left(\frac{\|\mathbf{h}\|}{\|\mathbf{z}\|}\right)^2 \end{aligned} \quad (47)$$

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

$$\frac{|\mathbf{a}_i^\top \mathbf{h}|^3}{|\mathbf{a}_i^\top \mathbf{z}|} \leq \frac{|2\mathbf{a}_i^\top \mathbf{z}| + \gamma \|\mathbf{h}\|}{|\mathbf{a}_i^\top \mathbf{z}|} (2\alpha_z^{\text{ub}} \|\mathbf{z}\| + \gamma \|\mathbf{h}\|)^2 \quad (48)$$

$$\leq \left(2 + \frac{\gamma \|\mathbf{h}\|}{\alpha_z^{\text{lb}} \|\mathbf{z}\|}\right) \left(2\alpha_z^{\text{ub}} + \gamma \frac{\|\mathbf{h}\|}{\|\mathbf{z}\|}\right)^2 \|\mathbf{z}\|^2, \quad (49)$$

where both inequalities are immediate consequences from the definitions of $\mathcal{D}_\gamma^{i,2}$ and \mathcal{E}_1^i (see (29) and (16)). Taking this together with the cardinality bound (47) and picking ϵ appropriately, we get

$$\frac{1}{m} \sum_{i=1}^m \frac{|\mathbf{a}_i^\top \mathbf{h}|^3}{|\mathbf{a}_i^\top \mathbf{z}|} \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{D}_\gamma^{i,2}} \leq \left\{ \underbrace{\frac{\left(2 + \frac{\gamma \|\mathbf{h}\|}{\alpha_z^{\text{lb}} \|\mathbf{z}\|}\right) \left(2\alpha_z^{\text{ub}} + \gamma \frac{\|\mathbf{h}\|}{\|\mathbf{z}\|}\right)^2}{9800 (\alpha_z^{\text{lb}})^4}}_{\vartheta_1} \frac{\|\mathbf{h}\|^2}{\|\mathbf{z}\|^2} + \epsilon \right\} \|\mathbf{h}\|^2. \quad (50)$$

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 (50) by observing that $\vartheta_1 \leq \frac{1}{100}$, which results in

$$\frac{1}{m} \sum_{i=1}^m \frac{|\mathbf{a}_i^\top \mathbf{h}|^3}{|\mathbf{a}_i^\top \mathbf{z}|} \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{D}_\gamma^{i,2}} \leq \left(\frac{1}{100} + \epsilon \right) \|\mathbf{h}\|^2. \quad (51)$$

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

$$\begin{aligned} -\left\langle \mathbf{h}, \frac{1}{2m} \nabla \ell_{\text{tr}}(\mathbf{z}) \right\rangle &\geq \left\{ 1.99 - 2(\zeta_1 + \zeta_2) - \sqrt{8/\pi} \frac{\|\mathbf{h}\|}{\alpha_z^{\text{lb}} \|\mathbf{z}\|} - 3\epsilon \right\} \|\mathbf{h}\|^2 \\ &\geq \left\{ 1.99 - 2(\zeta_1 + \zeta_2) - \sqrt{8/\pi} (3\alpha_h)^{-1} - 3\epsilon \right\} \|\mathbf{h}\|^2 \end{aligned} \quad (52)$$

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

$$\frac{\|\mathbf{h}\|}{\|\mathbf{z}\|} \leq \min \left\{ \frac{\alpha_z^{\text{lb}}}{3\alpha_h}, \frac{\alpha_z^{\text{lb}}}{6}, \frac{\sqrt{98/3} (\alpha_z^{\text{lb}})^2}{2\alpha_z^{\text{ub}} + \alpha_z^{\text{lb}}}, \frac{1}{11} \right\} \quad (53)$$

as claimed in Proposition 2.

To conclude this section, we provide a tighter estimate about the norm of the truncated 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 $1 - C \exp(-cm)$,*

$$\frac{1}{m} \|\nabla \ell_{\text{tr}}(\mathbf{z})\| \leq (1 + \delta) \cdot 4\sqrt{1.02 + 0.665/\alpha_h} \|\mathbf{h}\| \quad (54)$$

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

$$-\left\langle \mathbf{h}, \frac{1}{m} \nabla \ell_{\text{tr}}(\mathbf{z}) \right\rangle \geq \frac{2 \left\{ 1.99 - 2(\zeta_1 + \zeta_2) - \sqrt{8/(9\pi)} \alpha_h^{-1} - \epsilon \right\}}{(1 + \delta)^2 \cdot 16 (1.02 + 0.665/\alpha_h)} \left\| \frac{1}{m} \nabla \ell_{\text{tr}}(\mathbf{z}) \right\|^2. \quad (55)$$

Taking ϵ and δ to be sufficiently small we arrive at a feasible range (cf. Definition (11))

$$\mu \leq \frac{0.994 - \zeta_1 - \zeta_2 - \sqrt{2/(9\pi)} \alpha_h^{-1}}{2 (1.02 + 0.665/\alpha_h)} := \mu_0. \quad (56)$$

This establishes Proposition 1 and in turn Theorem 1 when μ_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 (56), which we defer to Appendix D.

3 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

$$\frac{1}{m} \|\mathbf{y} - \mathcal{A}(\mathbf{z}\mathbf{z}^\top)\|_1 \leq \frac{1}{m} \|\mathcal{A}(\mathbf{x}\mathbf{x}^\top - \mathbf{z}\mathbf{z}^\top)\|_1 + \frac{1}{m} \|\boldsymbol{\eta}\|_1 \lesssim \|\mathbf{h}\| \|\mathbf{z}\| + \frac{1}{m} \|\boldsymbol{\eta}\|_1 \lesssim \|\mathbf{h}\| \|\mathbf{z}\| + \frac{1}{\sqrt{m}} \|\boldsymbol{\eta}\|, \quad (57)$$

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 truncation rules \mathcal{E}_1^i and \mathcal{E}_2^i give

$$\begin{aligned} |v_i| &\lesssim \|\mathbf{h}\| + \frac{\|\boldsymbol{\eta}\|}{\sqrt{m}\|\mathbf{z}\|} \\ \implies \left\| \frac{1}{m} \nabla \ell_{\text{tr}}(\mathbf{z}) \right\| &= \frac{1}{m} \|\mathbf{A}^\top \mathbf{v}\| \leq \left\| \frac{1}{\sqrt{m}} \mathbf{A} \right\| \frac{1}{\sqrt{m}} \|\mathbf{v}\| \stackrel{(i)}{\lesssim} \frac{1}{\sqrt{m}} \|\mathbf{v}\| \lesssim \|\mathbf{h}\| + \frac{\|\boldsymbol{\eta}\|}{\sqrt{m}\|\mathbf{z}\|}, \end{aligned} \quad (58)$$

where (i) arises from [3, Corollary 5.35].

As discussed before, the estimation error is contractive if $-\frac{1}{m} \nabla \ell_{\text{tr}}(\mathbf{z})$ satisfies the regularity condition. With (58) in place, RC reduces to

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

Unfortunately, (59) 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 (59) holds for all \mathbf{h} obeying

$$c_3 \frac{\|\boldsymbol{\eta}\|/\sqrt{m}}{\|\mathbf{z}\|} \leq \|\mathbf{h}\| \leq c_4 \|\mathbf{x}\| \quad (60)$$

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

$$\text{dist}\left(\mathbf{z} + \frac{\mu}{m} \nabla \ell_{\text{tr}}(\mathbf{z}), \mathbf{x}\right) \leq (1 - \rho) \text{dist}(\mathbf{z}, \mathbf{x}) \quad (61)$$

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

$$\|\mathbf{h}\| \leq \frac{c_3 \|\boldsymbol{\eta}\|}{\sqrt{m}\|\mathbf{z}\|} \quad (62)$$

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 $\mathcal{O}(\frac{\|\boldsymbol{\eta}\|}{\sqrt{m}\|\mathbf{z}\|})$, compare (58). 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,

$$\text{dist}\left(\mathbf{z} + \frac{\mu}{m} \nabla \ell_{\text{tr}}(\mathbf{z}), \mathbf{x}\right) \leq c_5 \frac{\|\boldsymbol{\eta}\|}{\sqrt{m}\|\mathbf{x}\|} \quad (63)$$

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 (63) and (61), the estimation errors will never exceed the order of $\frac{\|\boldsymbol{\eta}\|}{\sqrt{m}\|\mathbf{x}\|}$ from then on. Putting these together establishes (6), namely, the first part of the theorem.

Below we justify the condition (59) for Regime 1, for which we start by gathering additional properties of the truncation rules. By Cauchy-Schwartz, $\frac{1}{m} \|\boldsymbol{\eta}\|_1 \leq \frac{1}{\sqrt{m}} \|\boldsymbol{\eta}\| \leq \frac{1}{c_3} \|\mathbf{h}\| \|\mathbf{z}\|$. When c_3 is sufficiently large, applying Lemmas 1 and 2 gives

$$\begin{aligned} \frac{1}{m} \sum_{l=1}^m \left| y_l - |\mathbf{a}_l^\top \mathbf{z}|^2 \right| &\leq \frac{1}{m} \|\mathcal{A}(\mathbf{x}\mathbf{x}^\top - \mathbf{z}\mathbf{z}^\top)\|_1 + \frac{1}{m} \|\boldsymbol{\eta}\|_1 \leq 2.98 \|\mathbf{h}\| \|\mathbf{z}\|; \\ \frac{1}{m} \sum_{l=1}^m \left| y_l - |\mathbf{a}_l^\top \mathbf{z}|^2 \right| &\geq \frac{1}{m} \|\mathcal{A}(\mathbf{x}\mathbf{x}^\top - \mathbf{z}\mathbf{z}^\top)\|_1 - \frac{1}{m} \|\boldsymbol{\eta}\|_1 \geq 1.151 \|\mathbf{h}\| \|\mathbf{z}\|. \end{aligned} \quad (64)$$

From now on, we shall denote $\tilde{\mathcal{E}}_2^i := \left\{ |\mathbf{a}_i^\top \mathbf{x}|^2 - |\mathbf{a}_i^\top \mathbf{z}|^2 \leq \frac{\alpha_h}{m} \|\mathbf{y} - \mathcal{A}(\mathbf{z}\mathbf{z}^\top)\|_1 \frac{|\mathbf{a}_i^\top \mathbf{z}|}{\|\mathbf{z}\|} \right\}$ 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_ϵ must be bounded as n scales, since

$$\|\boldsymbol{\eta}\|^2 \geq \sum_{i \notin \mathcal{G}} \eta_i^2 \geq (m - |\mathcal{G}|) \cdot C_\epsilon^2 \|\boldsymbol{\eta}\|^2 / m \geq \epsilon C_\epsilon^2 \|\boldsymbol{\eta}\|^2 \Rightarrow C_\epsilon \leq 1/\sqrt{\epsilon}. \quad (65)$$

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

$$\begin{aligned} \nabla_{\text{tr}} \ell(\mathbf{z}) &= \underbrace{2 \sum_{i \in \mathcal{G}} \frac{|\mathbf{a}_i^\top \mathbf{x}|^2 - |\mathbf{a}_i^\top \mathbf{z}|^2}{\mathbf{a}_i^\top \mathbf{z}} \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i} + 2 \sum_{i \notin \mathcal{G}} \frac{|\mathbf{a}_i^\top \mathbf{x}|^2 - |\mathbf{a}_i^\top \mathbf{z}|^2}{\mathbf{a}_i^\top \mathbf{z}} \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i \cap \tilde{\mathcal{E}}_2^i}}_{:= \nabla_{\text{tr}}^{\text{clean}} \ell(\mathbf{z})} \\ &\quad + \underbrace{2 \sum_{i \in \mathcal{G}} \frac{\eta_i}{\mathbf{a}_i^\top \mathbf{z}} \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i}}_{:= \nabla_{\text{tr}}^{\text{noise}} \ell(\mathbf{z})} + \underbrace{2 \sum_{i \notin \mathcal{G}} \left(\frac{y_i - |\mathbf{a}_i^\top \mathbf{z}|^2}{\mathbf{a}_i^\top \mathbf{z}} \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i} - \frac{|\mathbf{a}_i^\top \mathbf{x}|^2 - |\mathbf{a}_i^\top \mathbf{z}|^2}{\mathbf{a}_i^\top \mathbf{z}} \mathbf{1}_{\mathcal{E}_1^i \cap \tilde{\mathcal{E}}_2^i} \right) \mathbf{a}_i}_{:= \nabla_{\text{tr}}^{\text{extra}} \ell(\mathbf{z})}. \quad (66) \end{aligned}$$

- For each index $i \in \mathcal{G}$, the inclusion property (24) (i.e. $\mathcal{E}_3^i \subseteq \mathcal{E}_2^i \subseteq \mathcal{E}_4^i$) holds. To see this, observe that

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

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 (24) immediately from (64). As a result, all the proof arguments for Proposition 2 carry over to $\nabla_{\text{tr}}^{\text{clean}} \ell(\mathbf{z})$, suggesting that

$$-\left\langle \mathbf{h}, \frac{1}{m} \nabla_{\text{tr}}^{\text{clean}} \ell(\mathbf{z}) \right\rangle \geq 2 \left\{ 1.99 - 2(\zeta_1 + \zeta_2) - \sqrt{8/(9\pi)} \alpha_h^{-1} - \epsilon \right\} \|\mathbf{h}\|^2. \quad (67)$$

- 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

$$\left\| \frac{1}{m} \nabla_{\text{tr}}^{\text{noise}} \ell(\mathbf{z}) \right\| = \left\| \frac{1}{m} \mathbf{A}^\top \mathbf{w} \right\| \leq \left\| \frac{1}{\sqrt{m}} \mathbf{A} \right\| \left\| \frac{1}{\sqrt{m}} \mathbf{w} \right\| \stackrel{\text{(ii)}}{\leq} \frac{1 + \delta}{\sqrt{m}} \|\mathbf{w}\| \leq (1 + \delta) \frac{2\|\boldsymbol{\eta}\|/\sqrt{m}}{\alpha_z^{\text{lb}} \|\mathbf{z}\|}, \quad (68)$$

when m/n is sufficiently large. Here, (ii) arises from [3, Corollary 5.35], and the last inequality is a consequence of the upper estimate

$$\|\mathbf{w}\|^2 \leq 4 \sum_{i=1}^m \frac{|\eta_i|^2}{(\mathbf{a}_i^\top \mathbf{z})^2} \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i} \leq 4 \sum_{i=1}^m \frac{|\eta_i|^2}{(\alpha_z^{\text{lb}} \|\mathbf{z}\|)^2} = \frac{4 \|\boldsymbol{\eta}\|^2}{(\alpha_z^{\text{lb}} \|\mathbf{z}\|)^2}. \quad (69)$$

In turn, this immediately gives

$$\left| \left\langle \mathbf{h}, \frac{1}{m} \nabla_{\text{tr}}^{\text{noise}} \ell(\mathbf{z}) \right\rangle \right| \leq \|\mathbf{h}\| \left\| \frac{1}{m} \nabla_{\text{tr}}^{\text{noise}} \ell(\mathbf{z}) \right\| \leq \frac{2(1 + \delta)}{\alpha_z^{\text{lb}}} \frac{\|\boldsymbol{\eta}\|}{\sqrt{m} \|\mathbf{z}\|} \|\mathbf{h}\|. \quad (70)$$

- 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 (64), the weight $q_i := 2 \left(\frac{y_i - |\mathbf{a}_i^\top \mathbf{z}|^2}{\mathbf{a}_i^\top \mathbf{z}} \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{E}_2^i} - \frac{|\mathbf{a}_i^\top \mathbf{x}|^2 - |\mathbf{a}_i^\top \mathbf{z}|^2}{\mathbf{a}_i^\top \mathbf{z}} \mathbf{1}_{\mathcal{E}_1^i \cap \tilde{\mathcal{E}}_2^i} \right) \mathbf{1}_{\{i \notin \mathcal{G}\}}$ is bounded in magnitude by $6\|\mathbf{h}\|$. This gives

$$\begin{aligned} \|\mathbf{q}\| &\leq \sqrt{m - |\mathcal{G}|} \cdot 6\|\mathbf{h}\| \leq 6\sqrt{\epsilon m} \|\mathbf{h}\|, \\ \Rightarrow \left| \left\langle \frac{1}{m} \nabla_{\text{tr}}^{\text{extra}} \ell(\mathbf{z}), \mathbf{h} \right\rangle \right| &\leq \|\mathbf{h}\| \cdot \left\| \frac{1}{m} \nabla_{\text{tr}}^{\text{extra}} \ell(\mathbf{z}) \right\| = \frac{1}{m} \|\mathbf{h}\| \cdot \|\mathbf{A}^\top \mathbf{q}\| \leq 6(1 + \delta) \sqrt{\epsilon} \|\mathbf{h}\|^2. \quad (71) \end{aligned}$$

Taking the above bounds together yields

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

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

$$-\frac{1}{m} \langle \nabla \ell_{\text{tr}}(\mathbf{z}), \mathbf{h} \rangle \geq 2 \left\{ 1.95 - 2(\zeta_1 + \zeta_2) - \sqrt{8/(9\pi)} \alpha_h^{-1} \right\} \|\mathbf{h}\|^2 \quad (72)$$

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

$$\frac{1}{m} \|\boldsymbol{\eta}\|^2 \approx \frac{1}{m} \sum_{i=1}^m \mathbb{E} [\eta_i^2] = \frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{x})^2 \approx \|\mathbf{x}\|^2. \quad (73)$$

Substitution into (6) completes the proof.

4 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:

$$\text{KL}(P||Q) := \int \log \left(\frac{dP}{dQ} \right) dP, \quad (74)$$

The basic idea is to adopt the general reduction scheme discussed in [4, 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 $\mathbf{w}^{(l)}, \mathbf{w}^{(j)} \in \mathcal{M}$,

$$1/\sqrt{8} - (2n)^{-1/2} \leq \|\mathbf{w}^{(l)} - \mathbf{w}^{(j)}\| \leq 3/2 + n^{-1/2}; \quad (75)$$

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

$$\frac{|\mathbf{a}_i^\top (\mathbf{w} - \mathbf{x})|^2}{|\mathbf{a}_i^\top \mathbf{x}|^2} \leq \frac{\|\mathbf{w} - \mathbf{x}\|^2}{\|\mathbf{x}\|^2} \{2 + 17 \log^3 m\}, \quad 1 \leq i \leq m. \quad (76)$$

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

$$\frac{|\mathbf{a}_i^\top (\mathbf{w} - \mathbf{x})|^2}{|\mathbf{a}_i^\top \mathbf{x}|^2} \lesssim \frac{\|\mathbf{w} - \mathbf{x}\|^2}{\|\mathbf{x}\|^2} \log^3 m \lesssim \frac{1}{\log^3 m} \log^3 m \lesssim 1. \quad (77)$$

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 $\mathbf{r} \in \mathbb{R}^n$,*

$$\text{KL}(\mathbb{P}(\mathbf{y} | \mathbf{w} + \mathbf{r}) \parallel \mathbb{P}(\mathbf{y} | \mathbf{w})) \leq \sum_{i=1}^m |\mathbf{a}_i^\top \mathbf{r}|^2 \left(8 + \frac{2|\mathbf{a}_i^\top \mathbf{r}|^2}{|\mathbf{a}_i^\top \mathbf{w}|^2} \right). \quad (78)$$

Lemma 9 and (77) 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

$$\text{KL}(\mathbb{P}(\mathbf{y} | \mathbf{w}) \parallel \mathbb{P}(\mathbf{y} | \mathbf{x})) \leq c_3 \sum_{i=1}^m |\mathbf{a}_i^\top (\mathbf{w} - \mathbf{x})|^2 \leq 2c_3 m \|\mathbf{w} - \mathbf{x}\|^2, \quad \forall \mathbf{w} \in \mathcal{M}; \quad (79)$$

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

$$\|\mathbf{w} - \mathbf{x}\| \asymp \delta, \quad \text{and} \quad \|\mathbf{w} - \tilde{\mathbf{w}}\| \asymp \delta, \quad \forall \mathbf{w}, \tilde{\mathbf{w}} \in \mathcal{M} \text{ with } \mathbf{w} \neq \tilde{\mathbf{w}}. \quad (80)$$

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 (79). Indeed, suppose that

$$\frac{1}{M-1} \sum_{\mathbf{w} \in \mathcal{M} \setminus \{\mathbf{x}\}} \text{KL}(\mathbb{P}(\mathbf{y} | \mathbf{w}) \parallel \mathbb{P}(\mathbf{y} | \mathbf{x})) \leq \frac{1}{10} \log(M-1) \quad (81)$$

holds, then the Fano-type minimax lower bound [4, Theorem 2.7] asserts that

$$\inf_{\hat{\mathbf{x}}} \sup_{\mathbf{x} \in \mathcal{M}} \mathbb{E}[\|\hat{\mathbf{x}} - \mathbf{x}\| \mid \{\mathbf{a}_i\}] \gtrsim \min_{\mathbf{w}, \tilde{\mathbf{w}} \in \mathcal{M}, \mathbf{w} \neq \tilde{\mathbf{w}}} \|\mathbf{w} - \tilde{\mathbf{w}}\|. \quad (82)$$

Since $M = \exp(n/30)$, (81) would follow from

$$2c_3 \|\mathbf{w} - \mathbf{x}\|^2 \leq n/(300m), \quad \mathbf{w} \in \mathcal{M}. \quad (83)$$

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

$$\inf_{\hat{\mathbf{x}}} \sup_{\mathbf{x} \in \mathcal{M}} \mathbb{E}[\|\hat{\mathbf{x}} - \mathbf{x}\| \mid \{\mathbf{a}_i\}] \gtrsim \min_{\mathbf{w}, \tilde{\mathbf{w}} \in \mathcal{M}, \mathbf{w} \neq \tilde{\mathbf{w}}} \|\mathbf{w} - \tilde{\mathbf{w}}\| \gtrsim \sqrt{n/m}. \quad (84)$$

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 δ 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.

A Proofs for Section 2

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

$$(\mathbf{a}_i^\top \mathbf{z})^2 \geq \alpha_z^{\text{lb}} \|\mathbf{z}\| \cdot |\mathbf{a}_i^\top \mathbf{z}| \geq \gamma \|\mathbf{h}\| |\mathbf{a}_i^\top \mathbf{z}|. \quad (85)$$

Solving the quadratic inequality that specifies \mathcal{D}_γ^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

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

Assuming (85) 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

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

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:

$$\chi_z(\tau) := \begin{cases} 1, & \text{if } |\tau| \in [\sqrt{1.01}\alpha_z^{\text{lb}}, \sqrt{0.99}\alpha_z^{\text{ub}}]; \\ -100(\alpha_z^{\text{ub}})^{-2}\tau^2 + 100, & \text{if } |\tau| \in [\sqrt{0.99}\alpha_z^{\text{ub}}, \alpha_z^{\text{ub}}]; \\ 100(\alpha_z^{\text{lb}})^{-2}\tau^2 - 100, & \text{if } |\tau| \in [\alpha_z^{\text{lb}}, \sqrt{1.01}\alpha_z^{\text{lb}}]; \\ 0, & \text{else.} \end{cases} \quad (86)$$

$$\chi_h(\tau) := \begin{cases} 1, & \text{if } |\tau| \in [0, \sqrt{0.99}\gamma]; \\ -\frac{100}{\gamma^2}\tau^2 + 100, & \text{if } |\tau| \in [\sqrt{0.99}\gamma, \gamma]; \\ 0, & \text{else.} \end{cases} \quad (87)$$

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

$$0 \leq \chi_z(\mathbf{a}_i^\top \mathbf{z}) \leq \mathbf{1}_{\mathcal{E}_1^i}, \quad \text{and} \quad 0 \leq \chi_h(\mathbf{a}_i^\top \mathbf{h}) \leq \mathbf{1}_{\mathcal{D}_\gamma^{i,1}} \quad (88)$$

and thus,

$$\frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h})^2 \mathbf{1}_{\mathcal{E}_1^i \cap \mathcal{D}_\gamma^{i,1}} \geq \frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h})^2 \chi_z(\mathbf{a}_i^\top \mathbf{z}) \chi_h(\mathbf{a}_i^\top \mathbf{h}). \quad (89)$$

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

$$\mathbf{U}_z = \begin{bmatrix} \mathbf{z}^\top / \|\mathbf{z}\| \\ \vdots \end{bmatrix} \quad (90)$$

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

$$\tilde{\mathbf{h}} := \mathbf{U}_z \mathbf{h}, \quad \text{and} \quad \tilde{\mathbf{a}}_i := \mathbf{U}_z \mathbf{a}_i. \quad (91)$$

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

$$\begin{aligned} \mathbb{E} \left[(\mathbf{a}_i^\top \mathbf{h})^2 \chi_z(\mathbf{a}_i^\top \mathbf{z}) \chi_h(\mathbf{a}_i^\top \mathbf{h}) \right] &\geq \mathbb{E} \left[(\mathbf{a}_i^\top \mathbf{h})^2 \chi_z(\mathbf{a}_i^\top \mathbf{z}) \right] - \mathbb{E} \left[(\mathbf{a}_i^\top \mathbf{h})^2 (1 - \chi_h(\mathbf{a}_i^\top \mathbf{h})) \right] \\ &\geq \mathbb{E} \left[(\tilde{a}_{i,1} \tilde{h}_1)^2 \chi_z(\mathbf{a}_i^\top \mathbf{z}) \right] + \mathbb{E} \left[(\tilde{\mathbf{a}}_{i,\setminus 1}^\top \tilde{\mathbf{h}}_{\setminus 1})^2 \right] \mathbb{E} [\chi_z(\mathbf{a}_i^\top \mathbf{z})] - \|\mathbf{h}\|^2 \mathbb{E} \left[\xi^2 \mathbf{1}_{\{|\xi| > \sqrt{0.99}\gamma\}} \right] \\ &\geq |\tilde{h}_1|^2 (1 - \zeta_1) + \|\tilde{\mathbf{h}}_{\setminus 1}\|^2 (1 - \zeta_1) - \zeta_2 \|\mathbf{h}\|^2 \\ &\geq (1 - \zeta_1 - \zeta_2) \|\mathbf{h}\|^2, \end{aligned} \quad (92)$$

where the identity (92) arises from (39) and (40). 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 $\mathcal{O}(\gamma^2 \|\mathbf{h}\|^2)$. Apply the Hoeffding-type inequality [3, Proposition 5.10] to deduce that for any $\epsilon > 0$,

$$\frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h})^2 \chi_z(\mathbf{a}_i^\top \mathbf{z}) \chi_h(\mathbf{a}_i^\top \mathbf{h}) \geq \mathbb{E} \left[(\mathbf{a}_i^\top \mathbf{h})^2 \chi_z(\mathbf{a}_i^\top \mathbf{z}) \chi_h(\mathbf{a}_i^\top \mathbf{h}) \right] - \epsilon \|\mathbf{h}\|^2 \quad (93)$$

$$\geq (1 - \zeta_1 - \zeta_2 - \epsilon) \|\mathbf{h}\|^2 \quad (94)$$

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 ϵ -net argument. Specifically, we construct an ϵ -net \mathcal{N}_ϵ with cardinality $|\mathcal{N}_\epsilon| \leq (1 + 2/\epsilon)^{2n}$ (cf. [3]) 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

$$\frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h}_0)^2 \chi_z(\mathbf{a}_i^\top \mathbf{z}_0) \chi_h(\mathbf{a}_i^\top \mathbf{h}_0) \geq (1 - \zeta_1 - \zeta_2 - \epsilon) \|\mathbf{h}_0\|^2, \quad \forall \mathbf{h}_0, \mathbf{z}_0 \in \mathcal{N}_\epsilon \quad (95)$$

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 $\mathcal{O}(1)$. This guarantees that for each *unit* vector pair \mathbf{h} and \mathbf{z} ,

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

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

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

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

With the assertion (95) in place, we see that with high probability,

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

for all unit vectors \mathbf{h} and \mathbf{z} . Since ϵ can be arbitrary, putting this and (89) 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 τ with Lipschitz norm $\mathcal{O}(\gamma)$. 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 [3, Proposition 5.16] gives

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

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

From [3, Lemma 5.2], there exists an ϵ -net \mathcal{N}_ϵ 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 χ_2 implies

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

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

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

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

$$\chi_3(\tau, \gamma) := \begin{cases} 1, & \text{if } \sqrt{\tau} \geq \psi_{\text{lb}}(\gamma); \\ \frac{100\tau}{\psi_{\text{lb}}^2(\gamma)} - 99, & \text{if } \sqrt{\tau} \in [\sqrt{0.99}\psi_{\text{lb}}(\gamma), \psi_{\text{lb}}(\gamma)]; \\ 0, & \text{else.} \end{cases} \quad (96)$$

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 $0 < \delta \leq 0.005$,

$$\mathbf{1}_{\{|\mathbf{a}_i^\top \mathbf{h}| \geq \gamma\}} \leq \chi_3(|\mathbf{a}_i^\top \mathbf{h}|^2, \gamma) \leq \mathbf{1}_{\{|\mathbf{a}_i^\top \mathbf{h}| \geq \sqrt{0.99}\psi_{\text{lb}}(\gamma)\}} \leq \mathbf{1}_{\{|\mathbf{a}_i^\top \mathbf{h}| \geq 0.99\gamma\}}. \quad (97)$$

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

$$\begin{aligned} \mathbb{P}\left\{\chi_3(|\mathbf{a}_i^\top \mathbf{h}|^2, \gamma) \geq t\right\} &\leq \mathbb{P}\left\{\mathbf{1}_{\{|\mathbf{a}_i^\top \mathbf{h}| \geq 0.99\gamma\}} \geq t\right\} = \mathbb{P}\left\{\mathbf{1}_{\{|\mathbf{a}_i^\top \mathbf{h}| \geq 0.99\gamma\}} = 1\right\} \\ &= 2 \int_{0.99\gamma}^{\infty} \phi(x) dx \leq \frac{2}{0.99\gamma} \phi(0.99\gamma), \end{aligned} \quad (98)$$

where $\phi(x)$ is the density of a standard normal, and (98) 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 $\mathcal{O}(\gamma^{-2})$ (cf. [3, Definition 5.13]). We apply the Bernstein-type inequality for the sum of sub-exponential random variables [3, Corollary 5.17], which indicates that for any fixed \mathbf{h} and γ 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 ϵ -net \mathcal{N}_ϵ 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 τ with the Lipschitz constant bounded above by $\frac{100}{\psi_{\text{lb}}^2(\gamma_0)}$. With this in mind, for any (\mathbf{h}, γ) 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

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

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

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

with probability exceeding $1 - \frac{\log m}{\log(1+\delta)} (1 + \frac{2}{\epsilon})^n \exp(-c\epsilon^2 m)$. This establishes (44) 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 [2, 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 (35) that

$$\left\| \frac{1}{m} \nabla_{\text{tr}} \ell(\mathbf{z}) \right\| = \left\| \frac{1}{m} \mathbf{A}^\top \mathbf{v} \right\| \leq \frac{1}{m} \|\mathbf{A}\| \cdot \|\mathbf{v}\| \leq (1+\delta) \frac{\|\mathbf{v}\|}{\sqrt{m}} \quad (99)$$

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 [3, Corollary 5.35].

Everything then comes down to controlling $\|\mathbf{v}\|$. To this end, making use of the inclusion (36) 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 [2, 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

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

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 (42) and (51).

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

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

whenever $\frac{\|\mathbf{h}\|}{\|\mathbf{z}\|} \leq \min \left\{ \frac{\alpha_z^{\text{lb}}}{3\alpha_h}, \frac{\alpha_z^{\text{lb}}}{6}, \frac{\sqrt{98/3}(\alpha_z^{\text{lb}})^2}{2\alpha_z^{\text{ub}} + \alpha_z^{\text{lb}}}, \frac{1}{11} \right\}$. Substituting this into (99) completes the proof.

B Proofs for Section 4

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 ,

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

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

$$\min_{1 \leq i \leq m} |\mathbf{a}_i^\top \mathbf{x}| \geq \frac{1}{m \log m} \|\mathbf{x}\|; \quad (100)$$

$$\sum_{i=1}^m \mathbf{1}_{\{|\mathbf{a}_i^\top \mathbf{x}| \leq \frac{\|\mathbf{x}\|}{5 \log m}\}} \leq \frac{m}{4 \log m} := k. \quad (101)$$

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 (76). 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 [3, Corollary 5.35] gives

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

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

$$\begin{aligned} 0.5 - n^{-1/2} &\leq \|\mathbf{w}^{(l)} - \mathbf{w}^{(j)}\| \leq 1.5 + n^{-1/2}, \quad \forall l \neq j \\ 1/\sqrt{8} - (2n)^{-1/2} &\leq \|\mathbf{w}^{(l)} - \mathbf{x}\| \leq \sqrt{9/8} + (2n)^{-1/2}, \quad 1 \leq l \leq M_1 \end{aligned} \quad (102)$$

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

The next step is to show that many vectors in \mathcal{M}_1 satisfy (76). 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

$$\frac{|\mathbf{a}_i^\top \mathbf{r}|^2}{|\mathbf{a}_i^\top \mathbf{x}|^2} \leq \frac{2|a_{i,1}r_\parallel|^2 + 2|\mathbf{a}_{i,\perp}^\top \mathbf{r}_\perp|^2}{|a_{i,1}|^2 \|\mathbf{x}\|^2} \leq \frac{2|r_\parallel|^2}{\|\mathbf{x}\|^2} + \frac{2|\mathbf{a}_{i,\perp}^\top \mathbf{r}_\perp|^2}{|a_{i,1}|^2 \|\mathbf{x}\|^2} \leq \frac{2\|\mathbf{r}\|^2}{\|\mathbf{x}\|^2} + \frac{2|\mathbf{a}_{i,\perp}^\top \mathbf{r}_\perp|^2}{|a_{i,1}|^2 \|\mathbf{x}\|^2}. \quad (103)$$

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

$$\xi_i^l := \begin{cases} \mathbf{1}\{|\mathbf{a}_{i,\perp}^\top \mathbf{r}_\perp^{(l)}| \leq \frac{1}{m} \sqrt{\frac{n-1}{2n}}\}, & 1 \leq i \leq k, \\ \mathbf{1}\{|\mathbf{a}_{i,\perp}^\top \mathbf{r}_\perp^{(l)}| \leq \sqrt{\frac{2(n-1) \log n}{n}}\}, & i > k, \end{cases} \quad (104)$$

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 (76) when n is sufficiently large. To see this, note that for the first group of indices, $\xi_i^l = 1$ requires

$$|\mathbf{a}_{i,\perp}^\top \mathbf{r}_\perp^{(l)}| \leq \frac{1}{m} \sqrt{\frac{n-1}{2n}} \leq \frac{2}{m} \frac{\sqrt{n-1}}{\sqrt{n}-2} \|\mathbf{r}^{(l)}\| \leq \frac{3}{m} \|\mathbf{r}^{(l)}\|, \quad 1 \leq i \leq k, \quad (105)$$

where the second inequality follows from (102). This taken collectively with (100) and (103) 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

$$|\mathbf{a}_{i,\perp}^\top \mathbf{r}_\perp^{(l)}| \leq \sqrt{\frac{2(n-1) \log n}{n}} \leq \sqrt{17 \log n} \|\mathbf{r}^{(l)}\|, \quad i = k+1, \dots, m, \quad (106)$$

where the last inequality again follows from (102). Plugging (106) and (101) into (103) 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, (76) 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 $\mathcal{O}(1/m)$ 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

$$\sum_{l=1}^{M_1} \prod_{i=1}^k \xi_i^l \geq \frac{1}{2} \frac{M_1}{(2\pi)^{k/2} (1 + 4\sqrt{k/n})^{k/2}} \left(\frac{1}{\sqrt{2\pi m}} \right)^k := \widetilde{M}_1 \quad (107)$$

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 (107). 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,

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

For sufficiently large m , one has $k = \frac{m}{4 \log m} \leq \frac{1}{4}n$. Using [3, Corollary 5.35] we get

$$\left\| \frac{1}{n} \mathbf{B}^\top \mathbf{B} - \mathbf{I} \right\| \leq 4\sqrt{k/n} \quad (108)$$

with probability $1 - \exp(-\Omega(k))$. Thus, for any constant $0 < \epsilon < \frac{1}{2}$, conditional on $\{\mathbf{b}_i\}$ and (108) we obtain

$$\begin{aligned} \mathbb{P} \left\{ \bigcap_{i=1}^k \left\{ |\mathbf{b}_i^\top \mathbf{u}| \leq \frac{1}{m} \right\} \right\} &\geq (2\pi)^{-\frac{k}{2}} \det^{-\frac{1}{2}} \left(\frac{1}{n} \mathbf{B}^\top \mathbf{B} \right) \int_{\mathbf{b}_\mathbf{u} \in \Upsilon} \exp \left(-\frac{1}{2} \mathbf{b}_\mathbf{u}^\top \left(\frac{1}{n} \mathbf{B}^\top \mathbf{B} \right)^{-1} \mathbf{b}_\mathbf{u} \right) d\mathbf{b}_\mathbf{u} \\ &\geq (2\pi)^{-\frac{k}{2}} \left(1 + 4\sqrt{k/n} \right)^{-\frac{k}{2}} \int_{\mathbf{b}_\mathbf{u} \in \Upsilon} \exp \left(-\frac{1}{2} \left(1 - 4\sqrt{k/n} \right)^{-1} \sum_{i=1}^k b_{i,\mathbf{u}}^2 \right) d\mathbf{b}_\mathbf{u} \quad (109) \\ &\geq (2\pi)^{-\frac{k}{2}} \left(1 + 4\sqrt{k/n} \right)^{-\frac{k}{2}} (\sqrt{2\pi}m)^{-k}, \quad (110) \end{aligned}$$

where $\Upsilon := \{\tilde{\mathbf{b}} \mid |\tilde{b}_i| \leq m^{-1}, 1 \leq i \leq k\}$ and (109) is a direct consequence from (108).

When it comes to our quantity of interest, the above lower bound (110) indicates that on an event (defined via $\{\mathbf{a}_i\}$) of probability approaching 1, we have

$$\mathbb{E}\left[\sum_{l=1}^{M_1} \prod_{i=1}^k \xi_i^l\right] \geq M_1 (2\pi)^{-\frac{k}{2}} \left(1 + 4\sqrt{k/n}\right)^{-\frac{k}{2}} (\sqrt{2\pi m})^{-k}. \quad (111)$$

Since conditional on $\{\mathbf{a}_i\}$, $\prod_{i=1}^k \xi_i^l$ are independent across l , applying the Chernoff-type bound [5, 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 χ^2 -divergence between two probability measures P and Q as

$$\chi^2(P\|Q) := \int \left(\frac{dP}{dQ}\right)^2 dQ - 1. \quad (112)$$

It is well known (e.g. [4, Lemma 2.7]) that

$$\text{KL}(P\|Q) \leq \log(1 + \chi^2(P\|Q)), \quad (113)$$

and hence it suffices to develop an upper bound on the χ^2 divergence.

Under independence, for any $\mathbf{w}_0, \mathbf{w}_1 \in \mathbb{R}^n$, the decoupling identity of the χ^2 divergence [4, Page 96] gives

$$\begin{aligned} \chi^2(\mathbb{P}(\mathbf{y} | \mathbf{w}_1) \parallel \mathbb{P}(\mathbf{y} | \mathbf{w}_0)) &= \prod_{i=1}^m (1 + \chi^2(\mathbb{P}(y_i | \mathbf{w}_1) \parallel \mathbb{P}(y_i | \mathbf{w}_0))) - 1 \\ &= \exp\left(\sum_{i=1}^m \frac{(|\mathbf{a}_i^\top \mathbf{w}_1|^2 - |\mathbf{a}_i^\top \mathbf{w}_0|^2)^2}{|\mathbf{a}_i^\top \mathbf{w}_0|^2}\right) - 1. \end{aligned} \quad (114)$$

The preceding identity (114) arises from the following computation: by definition of $\chi^2(\cdot\|\cdot)$,

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

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

$$\text{KL}(\mathbb{P}(\mathbf{y} | \mathbf{w}_1) \parallel \mathbb{P}(\mathbf{y} | \mathbf{w}_0)) \leq \sum_{i=1}^m \frac{(|\mathbf{a}_i^\top \mathbf{w}_1|^2 - |\mathbf{a}_i^\top \mathbf{w}_0|^2)^2}{|\mathbf{a}_i^\top \mathbf{w}_0|^2} \quad (115)$$

$$\begin{aligned} &\leq \sum_{i=1}^m \frac{|\mathbf{a}_i^\top \mathbf{r}|^2 (2|\mathbf{a}_i^\top \mathbf{w}_0| + |\mathbf{a}_i^\top \mathbf{r}|)^2}{|\mathbf{a}_i^\top \mathbf{w}_0|^2} \\ &= \sum_{i=1}^m |\mathbf{a}_i^\top \mathbf{r}|^2 \left(\frac{8|\mathbf{a}_i^\top \mathbf{w}_0|^2 + 2|\mathbf{a}_i^\top \mathbf{r}|^2}{|\mathbf{a}_i^\top \mathbf{w}_0|^2} \right). \end{aligned} \quad (116)$$

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 $y_i = \mathbf{a}_i^\top \mathbf{x} + \eta_i$ and $\mathbf{a}_i \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I})$. Suppose that $\|\boldsymbol{\eta}\|_\infty \leq \varepsilon \|\mathbf{x}\|^2$ 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

$$\text{dist}(\mathbf{z}^{(0)}, \mathbf{x}) \leq \delta \|\mathbf{x}\|, \quad (117)$$

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 [2, Lemma 3.1] that $\frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{x})^2 \in [1 \pm \varepsilon] \|\mathbf{x}\|^2$. Under the hypothesis $\|\boldsymbol{\eta}\|_\infty \leq \varepsilon \|\mathbf{x}\|^2$, one has $\frac{1}{m} \|\boldsymbol{\eta}\|_1 \leq \varepsilon \|\mathbf{x}\|^2$, which yields

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

with probability $1 - \exp(-\Omega(m))$. This in turn implies that

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

and, hence,

$$\underbrace{\frac{1}{m} \sum_{i=1}^m \mathbf{a}_i \mathbf{a}_i^\top (\mathbf{a}_i^\top \mathbf{x})^2 \mathbf{1}_{\{ |\mathbf{a}_i^\top \mathbf{x}| \leq \sqrt{(1-2\varepsilon)\alpha_y^2 - \varepsilon} \}}}_{:= \mathbf{Y}_2} \preceq \mathbf{Y} \preceq \underbrace{\frac{1}{m} \sum_{i=1}^m \mathbf{a}_i \mathbf{a}_i^\top (\mathbf{a}_i^\top \mathbf{x})^2 \mathbf{1}_{\{ |\mathbf{a}_i^\top \mathbf{x}| \leq \sqrt{(1+2\varepsilon)\alpha_y^2 + \varepsilon} \}}}_{:= \mathbf{Y}_1}. \quad (118)$$

Letting $\xi \sim \mathcal{N}(0, 1)$, one can compute

$$\mathbb{E}[\mathbf{Y}_1] = \beta_1 \mathbf{x} \mathbf{x}^\top + \beta_2 \mathbf{I}, \quad \text{and} \quad \mathbb{E}[\mathbf{Y}_2] = \beta_3 \mathbf{x} \mathbf{x}^\top + \beta_4 \mathbf{I}, \quad (119)$$

where $\beta_1 := \mathbb{E}[\xi^4 \mathbf{1}_{\{|\xi| \leq \sqrt{(1+2\varepsilon)\alpha_y^2 + \varepsilon}\}}] - \mathbb{E}[\xi^2 \mathbf{1}_{\{|\xi| \leq \sqrt{(1+2\varepsilon)\alpha_y^2 + \varepsilon}\}}]$, $\beta_2 := \mathbb{E}[\xi^2 \mathbf{1}_{\{|\xi| \leq \sqrt{(1+2\varepsilon)\alpha_y^2 + \varepsilon}\}}]$, $\beta_3 := \mathbb{E}[\xi^4 \mathbf{1}_{\{|\xi| \leq \sqrt{(1-2\varepsilon)\alpha_y^2 - \varepsilon}\}}] - \mathbb{E}[\xi^2 \mathbf{1}_{\{|\xi| \leq \sqrt{(1-2\varepsilon)\alpha_y^2 - \varepsilon}\}}]$ and $\beta_4 := \mathbb{E}[\xi^2 \mathbf{1}_{\{|\xi| \leq \sqrt{(1-2\varepsilon)\alpha_y^2 - \varepsilon}\}}]$. 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 [3, Equation (5.26)] to deduce

$$\|\mathbf{Y}_1 - \mathbb{E}[\mathbf{Y}_1]\| \leq \delta, \quad \|\mathbf{Y}_2 - \mathbb{E}[\mathbf{Y}_2]\| \leq \delta \quad (120)$$

with probability $1 - \exp(-\Omega(m))$, provided that m/n exceeds some large constant. Besides, when ε is sufficiently small, one further has $\|\mathbb{E}[\mathbf{Y}_1] - \mathbb{E}[\mathbf{Y}_2]\| \leq \delta$. These taken together with (118) give

$$\|\mathbf{Y} - \beta_1 \mathbf{x} \mathbf{x}^\top - \beta_2 \mathbf{I}\| \leq 3\delta. \quad (121)$$

Fix $\tilde{\delta} > 0$. With (121) in place, repeating the same proof arguments as in [1, Section 7.8] (which we omit in the current paper) and taking δ, ε to be sufficiently small, we obtain

$$\text{dist}(\mathbf{z}^{(0)}, \mathbf{x}) \leq \tilde{\delta} \quad (122)$$

as long as m/n is sufficiently large, as claimed. \square

We now justify that the Poisson model (3) satisfies the condition $\|\boldsymbol{\eta}\| \leq \varepsilon \|\mathbf{x}\|^2$ whenever $\|\mathbf{x}\| \geq \log^{1.5} m$. 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

$$\mathbb{P}(y_i - \mu_i \geq \tau) \leq \frac{\mathbb{E}[e^{t y_i}]}{\exp(t(\mu_i + \tau))} = \frac{\exp(\mu_i(e^t - 1))}{\exp(t(\mu_i + \tau))} = \exp(\mu_i(e^t - t - 1) - t\tau), \quad \forall t \geq 0.$$

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

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

where (i) follows since $e^t \leq 1 + t + t^2$ ($0 \leq t \leq 1$). Letting $\kappa_i = \mu_i/\|\mathbf{x}\|^2$ and setting $\tilde{\varepsilon} = \varepsilon/2\kappa_i$, we obtain

$$\mathbb{P}(y_i - \mu_i \geq \varepsilon\|\mathbf{x}\|^2) = \mathbb{P}(y_i - \mu_i \geq 2\tilde{\varepsilon}\mu_i) \leq \exp(-\kappa_i\|\mathbf{x}\|^2\tilde{\varepsilon}^2) = \exp\left(-\frac{\varepsilon^2\|\mathbf{x}\|^2}{4\kappa_i}\right).$$

In addition, standard results on Gaussian measures indicate that $\max_{1 \leq i \leq m} \kappa_i \lesssim \log n$. As a consequence, if $\|\mathbf{x}\|^2 \gtrsim \log^3 m$, then $\frac{\|\mathbf{x}\|^2}{\kappa_i} \gtrsim \log^2 m$ ($1 \leq i \leq m$), which further gives

$$\mathbb{P}(\forall i : \eta_i \geq \varepsilon\|\mathbf{x}\|^2) = \mathbb{P}(\forall i : y_i - \mu_i \geq \varepsilon\|\mathbf{x}\|^2) \leq m \exp(-\Omega(\varepsilon^2 \log^2 m))$$

from the union bound. Similarly, applying the same argument on $-y_i$ we get $\eta_i \geq -\varepsilon\|\mathbf{x}\|^2$ for all i , which together with (123) establish that

$$\|\boldsymbol{\eta}\|_\infty \leq \varepsilon\|\mathbf{x}\|^2 \quad (123)$$

with high probability. In conclusion, the claim (117) applies to the Poisson model.

D Local error contraction with backtracking line search

In this paper, we also consider a backtracking line search with truncated objective to determine the learning rate. 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 truncated version of the objective function. Specifically, put

$$\widehat{\ell}(\mathbf{z}) := \sum_{i \in \widehat{\mathcal{T}}(\mathbf{z})} \{y_i \log(|\mathbf{a}_i^\top \mathbf{z}|^2) - |\mathbf{a}_i^\top \mathbf{z}|^2\}, \quad (124)$$

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

1. Start with $\tau = 1$;
2. Repeat $\tau \leftarrow \beta\tau$ until

$$\frac{1}{m} \widehat{\ell}(\mathbf{z}^{(t)} + \tau \mathbf{p}^{(t)}) \geq \frac{1}{m} \widehat{\ell}(\mathbf{z}^{(t)}) + \frac{1}{2} \tau \|\mathbf{p}^{(t)}\|^2, \quad (125)$$

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

3. Set $\mu_t = \tau$.

By definition (124), 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 $\mathcal{O}(\log 1/\beta)$ different τ 's, and hence the total cost amounts to computing $\mathbf{A}\mathbf{z}^{(t)}$, $\mathbf{A}\mathbf{p}^{(t)}$ as well as $\mathcal{O}(m \log 1/\beta)$ 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 $\mathcal{O}(m \log 1/\beta)$ flops, is negligible compared to that of computing the gradient even once.

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*

$$\|\mathbf{h}\|/\|\mathbf{z}\| \leq \epsilon_{\text{tr}} \quad (126)$$

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 μ_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

$$\|\mathbf{p}\|/\|\mathbf{z}\| \leq \epsilon \quad (127)$$

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 (52) that in the regime under study, one has

$$\|\mathbf{p}\| \geq 2 \left\{ 1.99 - 2(\zeta_1 + \zeta_2) - \sqrt{8/\pi}(3\alpha_h)^{-1} - o(1) \right\} \|\mathbf{h}\| \geq 3.64 \|\mathbf{h}\|. \quad (128)$$

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

$$\begin{aligned} (b+h)^2 \log \frac{(b+\delta)^2}{b^2} - (b+\delta)^2 + b^2 &= (b+h)^2 \log(1+b_\delta) - b^2 b_\delta \\ &\stackrel{(i)}{\leq} (b+h)^2 \{b_\delta - 0.4875b_\delta^2\} - b^2 b_\delta = ((b+h)^2 - b^2)b_\delta - 0.4875(b+h)^2 b_\delta^2 \\ &= h\delta(2+h/b)(2+\delta/b) - 0.4875(1+h/b)^2 |\delta(2+\delta/b)|^2 \\ &= 4h\delta + \frac{2h^2\delta}{b} + \frac{2h\delta^2}{b} + \frac{h^2\delta^2}{b^2} - 0.4875\delta^2 \left(1 + \frac{h}{b}\right)^2 \left(2 + \frac{\delta}{b}\right)^2, \end{aligned} \quad (129)$$

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

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

Plugging these two identities into (129) yields

$$\begin{aligned} (129) &\leq 4h\delta + \frac{2h^2\delta}{b} - \left(0.95 \left(1 + \frac{h}{b}\right)^2 + 1\right) \delta^2 - 0.4875\delta^2 \left(1 + \frac{h}{b}\right)^2 \frac{4\delta}{b} \\ &\leq 4h\delta - 1.95\delta^2 + \frac{2h^2|\delta|}{|b|} + \frac{1.9|h|}{|b|} \delta^2 + \frac{1.95|\delta^3|}{|b|} \left(1 + \frac{h}{b}\right)^2. \end{aligned}$$

Replacing respectively b , δ , 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 [2, Lemma 3.1] combined with the truncation rule. For the first term, applying the AM-GM inequality we get

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

Secondly, it follows from Lemma 4 that

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

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

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

Finally, the last term is bounded by

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

Under the hypothesis (128), 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

$$\begin{aligned} \frac{1}{m} \sum_{i=1}^m \{\ell_i(\mathbf{z} + \tau \mathbf{p}) - \ell_i(\mathbf{z})\} \mathbf{1}_{\mathcal{E}_3^i} &\leq \frac{1}{m} \sum_{i=1}^m (I_{1,i} + I_{2,i} + I_{3,i} + I_{4,i} + I_{5,i}) \mathbf{1}_{\mathcal{E}_3^i} \\ &\leq \tau(1.1 + \delta) \|\mathbf{p}\|^2 - 1.95(1 - \tilde{\zeta}_1 - \tilde{\zeta}_2) \tau^2 \|\mathbf{p}\|^2 + \tau \tilde{\epsilon} \|\mathbf{p}\|^2 \\ &= \left\{ \tau(1.1 + \delta) - 1.95(1 - \tilde{\zeta}_1 - \tilde{\zeta}_2) \tau^2 + \tau \tilde{\epsilon} \right\} \|\mathbf{p}\|^2 \end{aligned} \tag{130}$$

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

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 (130), 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 δ 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 τ is bounded away from 0. We omit this part for conciseness.

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] 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.
- [3] R. Vershynin. Introduction to the non-asymptotic analysis of random matrices. *Compressed Sensing, Theory and Applications*, pages 210 – 268, 2012.
- [4] A. B. Tsybakov and V. Zaiats. *Introduction to nonparametric estimation*, volume 11. Springer, 2009.
- [5] M. Mitzenmacher and E. Upfal. *Probability and computing*. Cambridge University Press, 2005.