

Supplementary Materials for “Solving Random Systems of Quadratic Equations via Truncated Generalized Gradient Flow”

G. Wang^{*,†} and G. B. Giannakis^{*}

^{*} ECE Dept. and Digital Tech. Center
University of Minnesota
Minneapolis, MN 55455, USA

[†] School of Automation
Beijing Institute of Technology
Beijing 100081, China
`{gangwang, georgios}@umn.edu`

October 24, 2016

1 Algorithm and main theorem

For convenience of presentation, we begin with repeating our model and main assumptions: Consider the noise-free real Gaussian model

$$\psi_i = |\mathbf{a}_i^\top \mathbf{x}|, \quad i \in [m] := \{1, 2, \dots, m\} \quad (1)$$

where $\mathbf{x} \in \mathbb{R}^n$ is the wanted unknown, and $\{\mathbf{a}_i\}_{i=1}^m$ are drawn independently and identically from the n -dimensional real Gaussian distribution, i.e., $\mathbf{a}_i \in \mathbb{R}^n \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$. For notational brevity, define $\mathbf{A} := [\mathbf{a}_1 \cdots \mathbf{a}_m]^\top$, $\boldsymbol{\psi} := [\psi_1 \cdots \psi_m]^\top$, and $\mathbf{y} := [y_1 \cdots y_m]^\top$, where $y_i := \psi_i^2$ stands for the squared magnitudes or intensity. Assume that the quadratic system (1) admits a unique solution, which indeed holds true as long as $m \geq 2n - 1$ generic measurements are taken [1]. Throughout the subsequent analysis, fix \mathbf{x} to be any solution of the given system in (1). Note that if \mathbf{x} satisfies the system in (1), so does $-\mathbf{x}$; i.e., the solution set for real-valued models becomes $\{\mathbf{x}, -\mathbf{x}\}$. Our theoretical analysis focuses on \mathbf{x} , rather than $-\mathbf{x}$. Introduce the notion of Euclidean distance of any estimate \mathbf{z} to the solution set: $\text{dist}(\mathbf{z}, \mathbf{x}) := \min \|\mathbf{z} \pm \mathbf{x}\|$ for real signals; and $\text{dist}(\mathbf{z}, \mathbf{x}) := \min_{\phi \in [0, 2\pi)} \|\mathbf{z} - \mathbf{x}e^{i\phi}\|$ for complex ones [2]. For concreteness, our results focus on real-valued models, define also the unrecoverable global phase factor for real-valued signals [2]

$$\phi(\mathbf{z}) := \begin{cases} 0, & \|\mathbf{z} - \mathbf{x}\| \leq \|\mathbf{z} + \mathbf{x}\| \\ \pi, & \text{otherwise.} \end{cases} \quad (2)$$

Henceforth, we always presume $\phi(\mathbf{z}) = 0$; otherwise, \mathbf{z} is replaced by $e^{-j\phi(\mathbf{z})}\mathbf{z}$, but for simplicity of presentation, the constant phase adaptation term is dropped whenever it is clear from the context [3].

Algorithm 1 and Theorem 1 are repeated next.

Algorithm 1 Truncated generalized gradient flow (TGGF) solver

- 1: **Input:** Data $\{\psi_i\}_{i=1}^m$ and features $\{\mathbf{a}_i\}_{i=1}^m$; maximum iterations T ; by default, set constant step size $\mu = 0.6/1$ for real/complex-valued models, thresholds $|\bar{\mathcal{I}}_0| = \lceil \frac{1}{6}m \rceil$,¹ and $\gamma = 0.7$.
- 2: **Evaluate** $\psi_i/\|\mathbf{a}_i\|$, $i = 1, \dots, m$, and find $\bar{\mathcal{I}}_0$ comprising indices associated with the $|\bar{\mathcal{I}}_0|$ largest $(\psi_i/\|\mathbf{a}_i\|)$'s.
- 3: **Initialize** \mathbf{z}_0 to $\sqrt{\sum_{i=1}^m \psi_i^2/m} \tilde{\mathbf{z}}_0$, where $\tilde{\mathbf{z}}_0$ is the unit leading eigenvector of $\bar{\mathbf{Y}}_0 := \frac{1}{|\bar{\mathcal{I}}_0|} \sum_{i \in \bar{\mathcal{I}}_0} \frac{\mathbf{a}_i \mathbf{a}_i^\top}{\|\mathbf{a}_i\|^2}$.
- 4: **Loop:** for $t = 0$ to $T - 1$

$$\mathbf{z}_{t+1} = \mathbf{z}_t - \frac{\mu}{m} \sum_{i \in \mathcal{I}_{t+1}} \left(\mathbf{a}_i^\top \mathbf{z}_t - \psi_i \frac{\mathbf{a}_i^\top \mathbf{z}_t}{|\mathbf{a}_i^\top \mathbf{z}_t|} \right) \mathbf{a}_i$$

where $\mathcal{I}_{t+1} = \left\{ 1 \leq i \leq m \mid |\mathbf{a}_i^\top \mathbf{z}_t| \geq \frac{1}{1+\gamma} \psi_i \right\}$.

- 5: **Output:** \mathbf{z}_T
-

Theorem 1. [4] Let $\mathbf{x} \in \mathbb{R}^n$ be an arbitrary signal, and consider (noise-free) measurements $\psi_i = |\mathbf{a}_i^\top \mathbf{x}|$, in which $\mathbf{a}_i \stackrel{i.i.d.}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$, $1 \leq i \leq m$. Then with probability at least $1 - (m+5)e^{-n/2} - e^{-c_0 m} - 1/n^2$ for some universal constant $c_0 > 0$, the initialization \mathbf{z}_0 returned by the orthogonality-promoting method in Algorithm 1 satisfies

$$\text{dist}(\mathbf{z}_0, \mathbf{x}) \leq \rho \|\mathbf{x}\| \quad (3)$$

with $\rho = 1/10$ (or any sufficiently small positive constant), provided that $m \geq c_1 |\bar{\mathcal{I}}_0| \geq c_2 n$ for some numerical constants $c_1, c_2 > 0$ and sufficiently large n . Further, choosing a constant step size $\mu \leq \mu_0$ along with a fixed truncation level $\gamma \geq 1/2$, and starting from any guess \mathbf{z}_0 satisfying (5), successive estimates of the TGGF algorithm (tabulated in Algorithm 1) obey

$$\text{dist}(\mathbf{z}_t, \mathbf{x}) \leq \rho (1 - \nu)^t \|\mathbf{x}\|, \quad \forall t = 1, 2, \dots \quad (4)$$

for some $0 < \nu < 1$, which holds with probability exceeding $1 - (m+5)e^{-n/2} - 8e^{-c_0 m} - 1/n^2$. Typical parameter values are $\mu = 0.6$, and $\gamma = 0.7$.

2 Proofs

To prove Theorem 1, this section establishes a few lemmas and the main ideas, while technical details are deferred to the Appendix. Relative to WF and TWF, our objective function involves nonsmoothness and nonconvexity, rendering the proof of exact recovery of TGGF nontrivial. In addition, our initialization method starts from a rather different perspective than the spectral alternatives, so the thoughts and tools involved in proving performance of our initialization deviate from those of the spectral methods [5, 2, 3]. Part of the proof is adapted from [2, 3] and [6].

The proof of Theorem 1 consists of two parts: Section 2.1 justifies the performance of the proposed orthogonality-promoting initialization, which essentially achieves any given constant relative error as soon as the number of equations is on the order of the number of unknowns, namely, $m \asymp n$.²

²The notations $\phi(n) = \mathcal{O}(g(n))$ or $\phi(n) \gtrsim g(n)$ (respectively, $\phi(n) \lesssim g(n)$) means there exists a numerical constant $c > 0$ such that $\phi(n) \leq cg(n)$, while $\phi(n) \asymp g(n)$ means $\phi(n)$ and $g(n)$ are orderwise equivalent.

Section 2.2 demonstrates theoretical convergence of TGGF to the solution of the quadratic system in (1) at a geometric rate provided that the initial estimate has a sufficiently small constant relative error as in (3). The two stages of TGGF can be performed independently, meaning that other better initialization methods, if available, could be adopted to initialize our truncated generalized gradient iterations; likewise, our initialization method can also be applied to initialize other iterative optimization algorithms.

2.1 Constant Relative Error by Orthogonality-promoting Initialization

This section concentrates on proving guaranteed performance of the proposed orthogonality-promoting initialization method, as asserted in the following proposition.

Proposition 1. *Fix $\mathbf{x} \in \mathbb{R}^n$ arbitrarily, and consider the noiseless case $\psi_i = |\mathbf{a}_i^\top \mathbf{x}|$, where $\mathbf{a}_i \stackrel{i.i.d.}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$, $1 \leq i \leq m$. Then with probability at least $1 - (m + 5)e^{-n/2} - e^{-c_0 m} - 1/n^2$ for some universal constant $c_0 > 0$, the initialization \mathbf{z}_0 returned by the orthogonality-promoting method satisfies*

$$\text{dist}(\mathbf{z}_0, \mathbf{x}) \leq \rho \|\mathbf{x}\| \quad (5)$$

for $\rho = 1/10$ or any positive constant, with the proviso that $m \geq c_1 |\bar{\mathcal{I}}_0| \geq c_2 n$ for some numerical constants $c_1, c_2 > 0$ and sufficiently large n .

Due to homogeneity in (5), it suffices to work with the case where $\|\mathbf{x}\| = 1$. Assume for the moment that $\|\mathbf{x}\| = 1$ is known and \mathbf{z}_0 has been scaled such that $\|\mathbf{z}_0\| = 1$. Subsequently, the error between the employed \mathbf{x} 's norm estimate $\sqrt{\frac{1}{m} \sum_{i=1}^m y_i}$ and the unknown norm $\|\mathbf{x}\| = 1$ will be accounted for at the end of this Section. Instrumental in proving Proposition 1 is the following result, whose proof is deferred to Appendix A.1.

Lemma 1. *Consider the noiseless data $\psi_i = |\mathbf{a}_i^\top \mathbf{x}|$, where $\mathbf{a}_i \stackrel{i.i.d.}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$, $1 \leq i \leq m$. For any unit vector $\mathbf{x} \in \mathbb{R}^n$, there exists a vector $\mathbf{u} \in \mathbb{R}^n$ with $\mathbf{u}^\top \mathbf{x} = 0$ and $\|\mathbf{u}\| = 1$ such that*

$$\frac{1}{2} \|\mathbf{x}\mathbf{x}^\top - \mathbf{z}_0 \mathbf{z}_0^\top\|_F^2 \leq \frac{\|\bar{\mathbf{S}}_0 \mathbf{u}\|^2}{\|\bar{\mathbf{S}}_0 \mathbf{x}\|^2} \quad (6)$$

for $\mathbf{z}_0 = \tilde{\mathbf{z}}_0$, where the unit vector $\tilde{\mathbf{z}}_0$ is given by

$$\tilde{\mathbf{z}}_0 := \arg \max_{\|\mathbf{z}\|=1} \mathbf{z}^\top \bar{\mathbf{Y}}_0 \mathbf{z} \quad (7)$$

with $\bar{\mathbf{Y}}_0 := \frac{1}{|\bar{\mathcal{I}}_0|} \bar{\mathbf{S}}_0^\top \bar{\mathbf{S}}_0$, and $\bar{\mathbf{S}}_0$ is formed by removing the rows of $\mathbf{S} := [\mathbf{a}_1 / \|\mathbf{a}_1\| \ \cdots \ \mathbf{a}_m / \|\mathbf{a}_m\|]^\top \in \mathbb{R}^{m \times n}$, if their indices do not belong to the set $\bar{\mathcal{I}}_0$ specified in Algorithm 1.

We now turn to prove Proposition 1. The first step consists in upper-bounding the term on the right-hand-side of (6). Specifically, its numerator term will be upper bounded, and the denominator term lower bounded, which are summarized in Lemma 2 and Lemma 3, whose proofs can be found in Appendix A.2 and Appendix A.3, respectively.

Lemma 2. *In the setup of Lemma 1, if $|\bar{\mathcal{I}}_0| \geq c'_1 n$, then the next*

$$\|\bar{\mathbf{S}}_0 \mathbf{u}\|^2 \leq 1.01 |\bar{\mathcal{I}}_0| / n \quad (8)$$

holds with probability at least $1 - 2e^{-c_K n}$, where c'_2 and c_K are some universal constants.

Lemma 3. *In the setup of Lemma 1, the following holds with probability at least $1 - (m+1)e^{-n/2} - e^{-c_0 m} - 1/n^2$,*

$$\|\bar{\mathbf{S}}_0 \mathbf{x}\|^2 \geq \frac{0.99|\bar{\mathcal{I}}_0|}{2.3n} \left[1 + \log(m/|\bar{\mathcal{I}}_0|) \right] \quad (9)$$

provided that $|\bar{\mathcal{I}}_0| \geq c'_1 n$, $m \geq c'_2 |\bar{\mathcal{I}}_0|$, and $m \geq c'_3 n$ for some absolute constants $c'_1, c'_2, c'_3 > 0$, and sufficiently large n .

Therefore, putting the upper and lower bounds in (8) and (9) together, one arrives at

$$\frac{\|\bar{\mathbf{S}}_0 \mathbf{u}\|^2}{\|\bar{\mathbf{S}}_0 \mathbf{x}\|^2} \leq \frac{2.4}{1 + \log(m/|\bar{\mathcal{I}}_0|)} \triangleq \kappa \quad (10)$$

which holds with probability at least $1 - (m+3)e^{-n/2} - e^{-c_0 m} - 1/n^2$, with the proviso that $m \geq c'_1 |\bar{\mathcal{I}}_0|$, and $m \geq c'_2 n$, $|\bar{\mathcal{I}}_0| \geq c'_3 n$ for some absolute constants $c'_1, c'_2, c'_3 > 0$, and sufficiently large n .

Apparently, the bound κ in (10) is meaningful only when the ratio $\log(m/|\bar{\mathcal{I}}_0|) > 1.4$, i.e., $m/|\bar{\mathcal{I}}_0| > 4$, because the left hand side expressible in terms of $\sin^2 \theta$ enjoys a trivial upper bound 1. Henceforth, we will work with the case where $m/|\bar{\mathcal{I}}_0| > 4$. Empirically, $\lfloor m/|\bar{\mathcal{I}}_0| \rfloor = 6$ or equivalently $|\bar{\mathcal{I}}_0| = \lceil \frac{1}{6} m \rceil$ in Algorithm 1 works well when m/n is relatively small. Note further that the bound κ can be made arbitrarily small by letting $m/|\bar{\mathcal{I}}_0|$ be large enough. Without any loss of generality, let us take $\kappa := 0.001$. An additional step leads to the wanted bound on the distance between $\tilde{\mathbf{z}}_0$ and \mathbf{x} ; similar arguments can be found in [2, Section 7.8]. Recall that

$$|\mathbf{x}^\top \tilde{\mathbf{z}}_0|^2 = \cos^2 \theta = 1 - \sin^2 \theta \geq 1 - \kappa, \quad (11)$$

so one has

$$\begin{aligned} \text{dist}^2(\tilde{\mathbf{z}}_0, \mathbf{x}) &\leq \|\tilde{\mathbf{z}}_0\|^2 + \|\mathbf{x}\|^2 - 2|\mathbf{x}^\top \tilde{\mathbf{z}}_0| \\ &\leq (2 - 2\sqrt{1 - \kappa}) \|\mathbf{x}\|^2 \\ &\approx \kappa \|\mathbf{x}\|^2. \end{aligned} \quad (12)$$

Coming back to the case in which $\|\mathbf{x}\|$ is unknown stated prior to Lemma 1, the unit eigenvector $\tilde{\mathbf{z}}_0$ is scaled by the estimate of $\|\mathbf{x}\|$ to yield the initial guess $\mathbf{z}_0 = \sqrt{\frac{1}{m} \sum_{i=1}^m y_i} \tilde{\mathbf{z}}_0$. Using the results in Lemma 7.8 in [2], the following holds with high probability

$$\|\mathbf{z}_0 - \tilde{\mathbf{z}}_0\| = \|\|\mathbf{z}_0\| - 1\| \leq (1/20) \|\mathbf{x}\|. \quad (13)$$

Summarizing the two inequalities, we conclude that

$$\text{dist}(\mathbf{z}_0, \mathbf{x}) \leq \|\mathbf{z}_0 - \tilde{\mathbf{z}}_0\| + \text{dist}(\tilde{\mathbf{z}}_0, \mathbf{x}) \leq (1/10) \|\mathbf{x}\|. \quad (14)$$

The initialization thus obeys $\text{dist}(\mathbf{z}_0, \mathbf{x})/\|\mathbf{x}\| \leq 1/10$ for any $\mathbf{x} \in \mathbb{R}^n$ with high probability provided that $m \geq c_1 |\bar{\mathcal{I}}_0| \geq c_2 n$ holds for some universal constants $c_1, c_2 > 0$ and sufficiently large n .

2.2 Exact Recovery from Noiseless Data

We now prove that with accurate enough initial estimates, TGGF converges at a geometric rate to \mathbf{x} with high probability (i.e., the second part of Theorem 1). To be specific, with initialization

obeying (5) in Proposition 1, TGGF reconstructs the solution \mathbf{x} exactly in linear time. In this direction, it suffices to demonstrate that the TGGF's update rule (i.e., Step 4 in Algorithm 1) is locally contractive within a sufficiently small neighborhood of \mathbf{x} , as asserted in the following proposition.

Proposition 2 (Local error contraction). *Consider the noise-free measurements $\psi_i = |\mathbf{a}_i^\top \mathbf{x}|$ with i.i.d. Gaussian design vectors $\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$, $1 \leq i \leq m$, and fix any $\gamma \geq 1/2$. Then there exist universal constants $c_0, c_1 > 0$ and $0 < \nu < 1$ such that with probability at least $1 - 7e^{-c_0 m}$, the following holds*

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

for all $\mathbf{x}, \mathbf{z} \in \mathbb{R}^n$ obeying the condition (5) for sufficiently small $\rho > 0$ with the proviso that $m \geq c_1 n$ and that the constant step size μ satisfying $0 < \mu \leq \mu_0$ for some $\mu_0 > 0$.

Proposition 2 demonstrates that the distance of TGGF's successive iterates to \mathbf{x} is monotonically decreasing once the algorithm enters a small-size neighborhood around \mathbf{x} . This neighborhood is commonly referred to as the *basin of attraction*; see further discussions in [2, 7, 3]. In other words, as soon as one lands within the basin of attraction, TGGF's iterates remain in this region and will be attracted to \mathbf{x} exponentially fast. To substantiate Proposition 2, recall the concept of the *local regularity condition*, which was first developed in [2] and plays a fundamental role in establishing linear convergence to global optimum of nonconvex optimization approaches such as WF/TWF [2, 7, 3]. Now consider the update rule of TGGF

$$\mathbf{z}_{t+1} = \mathbf{z}_t - \frac{\mu}{m} \nabla \ell_{\text{tr}}(\mathbf{z}_t), \quad \forall t \geq 0, \quad (16)$$

where the truncated gradient $\nabla \ell_{\text{tr}}(\mathbf{z}_t)$ (See Remark 1 in [8] for more discussion) evaluated at some point $\mathbf{z}_t \in \mathbb{R}^n$ is given by

$$\frac{1}{m} \nabla \ell_{\text{tr}}(\mathbf{z}_t) \triangleq \frac{1}{m} \sum_{i \in \mathcal{I}} \left(\mathbf{a}_i^\top \mathbf{z}_t - \psi_i \frac{\mathbf{a}_i^\top \mathbf{z}_t}{|\mathbf{a}_i^\top \mathbf{z}_t|} \right) \mathbf{a}_i.$$

The truncated gradient $\nabla \ell_{\text{tr}}(\mathbf{z})$ is said to satisfy the local regularity condition, or LRC(μ, λ, ϵ) for some constant $\lambda > 0$, provided that

$$\left\langle \frac{1}{m} \nabla \ell_{\text{tr}}(\mathbf{z}), \mathbf{h} \right\rangle \geq \frac{\mu}{2} \left\| \frac{1}{m} \nabla \ell_{\text{tr}}(\mathbf{z}) \right\|^2 + \frac{\lambda}{2} \|\mathbf{h}\|^2 \quad (17)$$

holds for all $\mathbf{z} \in \mathbb{R}^n$ such that $\|\mathbf{h}\| = \|\mathbf{z} - \mathbf{x}\| \leq \epsilon \|\mathbf{x}\|$ for some constant $0 < \epsilon < 1$, where the ball $\|\mathbf{z} - \mathbf{x}\| \leq \epsilon \|\mathbf{x}\|$ is the so-called *basin of attraction*. Simple linear algebra along with the regularity condition in (17) leads to

$$\begin{aligned} \text{dist}^2 \left(\mathbf{z} - \frac{\mu}{m} \nabla \ell_{\text{tr}}(\mathbf{z}), \mathbf{x} \right) &= \left\| \mathbf{z} - \frac{\mu}{m} \nabla \ell_{\text{tr}}(\mathbf{z}) - \mathbf{x} \right\|^2 \\ &= \|\mathbf{h}\|^2 - 2\mu \left\langle \mathbf{h}, \frac{1}{m} \nabla \ell_{\text{tr}}(\mathbf{z}) \right\rangle + \left\| \frac{\mu}{m} \nabla \ell_{\text{tr}}(\mathbf{z}) \right\|^2 \end{aligned} \quad (18)$$

$$\begin{aligned} &\leq \|\mathbf{h}\|^2 - 2\mu \left(\frac{\mu}{2} \left\| \frac{1}{m} \nabla \ell_{\text{tr}}(\mathbf{z}) \right\|^2 + \frac{\lambda}{2} \|\mathbf{h}\|^2 \right) + \left\| \frac{\mu}{m} \nabla \ell_{\text{tr}}(\mathbf{z}) \right\|^2 \\ &= (1 - \lambda\mu) \|\mathbf{h}\|^2 = (1 - \lambda\mu) \text{dist}^2(\mathbf{z}, \mathbf{x}) \end{aligned} \quad (19)$$

for all \mathbf{z} obeying $\|\mathbf{h}\| \leq \epsilon \|\mathbf{x}\|$. Clearly, if the LRC(μ, λ, ϵ) is proved for TGGF, our goal (15) follows upon letting $\nu := \lambda\mu$.

2.2.1 Proof of the local regularity condition in (17)

By definition, justifying the local regularity condition in (17) entails controlling the norm of the truncated gradient $\frac{1}{m}\nabla\ell_{\text{tr}}(\mathbf{z})$, i.e., bounding the last term in (18). Recall that

$$\frac{1}{m}\nabla\ell_{\text{tr}}(\mathbf{z}) = \frac{1}{m}\sum_{i\in\mathcal{I}}\left(\mathbf{a}_i^{\mathcal{T}}\mathbf{z} - \psi_i\frac{\mathbf{a}_i^{\mathcal{T}}\mathbf{z}}{|\mathbf{a}_i^{\mathcal{T}}\mathbf{z}|}\right)\mathbf{a}_i \triangleq \frac{1}{m}\mathbf{A}\mathbf{v} \quad (20)$$

where $\mathcal{I} := \{1 \leq i \leq m \mid |\mathbf{a}_i^{\mathcal{T}}\mathbf{z}| \geq |\mathbf{a}_i^{\mathcal{T}}\mathbf{x}|/(1+\gamma)\}$, and $\mathbf{v} := [v_1 \ \dots \ v_m]^{\mathcal{T}} \in \mathbb{R}^m$ with $v_i := \frac{\mathbf{a}_i^{\mathcal{T}}\mathbf{z}}{|\mathbf{a}_i^{\mathcal{T}}\mathbf{z}|} (|\mathbf{a}_i^{\mathcal{T}}\mathbf{z}| - \psi_i) \mathbb{1}_{\{|\mathbf{a}_i^{\mathcal{T}}\mathbf{z}| \geq |\mathbf{a}_i^{\mathcal{T}}\mathbf{x}|/(1+\gamma)\}}$. Now, consider

$$|v_i|^2 = \left| (|\mathbf{a}_i^{\mathcal{T}}\mathbf{z}| - |\mathbf{a}_i^{\mathcal{T}}\mathbf{x}|) \mathbb{1}_{\{|\mathbf{a}_i^{\mathcal{T}}\mathbf{z}| \geq |\mathbf{a}_i^{\mathcal{T}}\mathbf{x}|/(1+\gamma)\}} \right|^2 \leq \left| |\mathbf{a}_i^{\mathcal{T}}\mathbf{z}| - |\mathbf{a}_i^{\mathcal{T}}\mathbf{x}| \right|^2 \leq |\mathbf{a}_i^{\mathcal{T}}\mathbf{h}|^2 = a_{i,1}^2 \|\mathbf{h}\|^2, \quad (21)$$

where $\mathbf{h} = \mathbf{z} - \mathbf{x}$. Observe that $a_{i,1}^2$ obeys the *Chi-square* distribution with $k = 1$ degrees of freedom; yet due to our working assumption $\|\mathbf{a}_i\| \leq \sqrt{2.3n}$, it has mean $\mathbb{E}[a_{i,1}^2] \leq k = 1$. So fixing any $0 < \delta' < 1$ and applying the one-sided Bernstein-type inequality, the following holds with probability at least $1 - e^{-m\delta'^2/2}$ [9, Proposition 5.16]

$$\|\mathbf{v}\|^2 = \sum_{i=1}^m v_i^2 \leq \sum_{i=1}^m a_{i,1}^2 \|\mathbf{h}\|^2 \leq (1 + \delta')m\|\mathbf{h}\|^2. \quad (22)$$

On the other hand, standard matrix concentration results confirm that the largest singular value of $\mathbf{A} = [\mathbf{a}_1 \ \dots \ \mathbf{a}_m]^{\mathcal{T}}$ with i.i.d. Gaussian $\{\mathbf{a}_i\}$ satisfies $\sigma_1 := \|\mathbf{A}\| \leq (1 + \delta'')\sqrt{m}$ for some $\delta'' > 0$ with probability exceeding $1 - 2e^{-c_0 m}$ as soon as $m \geq c_1 n$ for sufficiently large $c_1 > 0$, where $c_1 > 0$ is a universal constant depending on δ'' [9, Remark 5.25]. Putting together (20), (21), and (22) yields

$$\left\| \frac{1}{m}\nabla\ell_{\text{tr}}(\mathbf{z}) \right\| \leq \frac{1}{m} \|\mathbf{A}\| \cdot \|\mathbf{v}\| \leq (1 + \delta')(1 + \delta'')\|\mathbf{h}\| \leq (1 + \delta)^2 \|\mathbf{h}\|, \quad \delta := \max\{\delta', \delta''\} \quad (23)$$

which holds with high probability. This condition essentially asserts that the truncated gradient of the objective function $\ell(\mathbf{z})$ or the search direction is well behaved (the function value does not vary too much).

Notice that to prove the LRC, it suffices to show that the truncated gradient $\frac{1}{m}\nabla\ell_{\text{tr}}(\mathbf{z})$ ensures sufficient descent, i.e., it obeys a uniform lower bound along the search direction \mathbf{h} taking the form

$$\left\langle \frac{1}{m}\nabla\ell_{\text{tr}}(\mathbf{z}), \mathbf{h} \right\rangle \gtrsim \|\mathbf{h}\|^2 \quad (24)$$

which occupies the remaining of this section. Formally, this can be stated as follows.

Proposition 3. *Consider the noiseless measurements $\psi_i = |\mathbf{a}_i^{\mathcal{T}}\mathbf{x}|$ and fix any sufficiently small constant $\epsilon > 0$. There exist universal constants $c_0, c_1 > 0$ such that if $m > c_1 n$, then the following holds with probability exceeding $1 - 4e^{-c_0 m}$*

$$\left\langle \mathbf{h}, \frac{1}{m}\nabla\ell_{\text{tr}}(\mathbf{z}) \right\rangle \geq 2(1 - \zeta_1 - \zeta_2 - 2\epsilon) \|\mathbf{h}\|^2 \quad (25)$$

for all $\mathbf{x}, \mathbf{z} \in \mathbb{R}^n$ such that $\|\mathbf{h}\| / \|\mathbf{x}\| \leq \rho$ for $0 < \rho \leq 1/10$ and any fixed $\gamma \geq 1/2$.

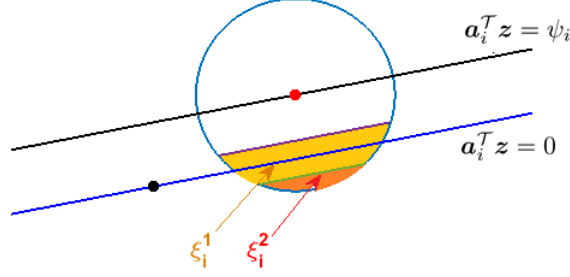


Figure 1: The geometric understanding of the proposed truncation rule on the i -th gradient component involving $\mathbf{a}_i^T \mathbf{x} = \psi_i$, where the red dot denotes the solution \mathbf{x} and the black one is the origin. Hyperplanes $\mathbf{a}_i^T \mathbf{z} = \psi_i$ and $\mathbf{a}_i^T \mathbf{z} = 0$ (of $\mathbf{z} \in \mathbb{R}^n$) passing through points $\mathbf{z} = \mathbf{x}$ and $\mathbf{z} = \mathbf{0}$, respectively, are shown.

Lemma 4. Fix any $\gamma > 0$. For each $i \in [m]$, define the following events

$$\mathcal{E}_i := \left\{ \frac{|\mathbf{a}_i^T \mathbf{z}|}{|\mathbf{a}_i^T \mathbf{x}|} \geq \frac{1}{1+\gamma} \right\}, \quad (26)$$

$$\mathcal{D}_i := \left\{ \frac{|\mathbf{a}_i^T \mathbf{h}|}{|\mathbf{a}_i^T \mathbf{x}|} \geq \frac{2+\gamma}{1+\gamma} \right\}, \quad (27)$$

$$\text{and } \mathcal{K}_i := \left\{ \frac{\mathbf{a}_i^T \mathbf{z}}{|\mathbf{a}_i^T \mathbf{z}|} \neq \frac{\mathbf{a}_i^T \mathbf{x}}{|\mathbf{a}_i^T \mathbf{x}|} \right\} \quad (28)$$

where $\mathbf{h} = \mathbf{z} - \mathbf{x}$. Under the condition $\|\mathbf{h}\| / \|\mathbf{x}\| \leq \rho$, the following inclusion holds

$$\mathcal{E}_i \cap \mathcal{K}_i \subseteq \mathcal{D}_i. \quad (29)$$

Proof. From Fig. 1, it is clear that if $\mathbf{z} \in \xi_i^2$, then the sign of $\mathbf{a}_i^T \mathbf{z}$ will be different than that of $\mathbf{a}_i^T \mathbf{x}$. The region ξ_i^2 , however, can be specified by the conditions that $\frac{\mathbf{a}_i^T \mathbf{z}}{|\mathbf{a}_i^T \mathbf{z}|} \neq \frac{\mathbf{a}_i^T \mathbf{x}}{|\mathbf{a}_i^T \mathbf{x}|}$ and $\frac{|\mathbf{a}_i^T \mathbf{h}|}{|\mathbf{a}_i^T \mathbf{x}|} \geq 1 + \frac{1}{1+\gamma} = \frac{2+\gamma}{1+\gamma}$. Under our initialization condition $\|\mathbf{h}\| / \|\mathbf{x}\| \leq \rho$, it is self-evident that \mathcal{D}_i describes two spherical caps that contain ξ_i^2 . Hence, it holds that $\mathcal{E}_i \cap \mathcal{K}_i = \xi_i^2 \subseteq \mathcal{D}_i$. \square

Rewrite the truncated gradient as

$$\begin{aligned} \frac{1}{2m} \nabla \ell_{\text{tr}}(\mathbf{z}) &= \frac{1}{m} \sum_{i=1}^m \left(\mathbf{a}_i^T \mathbf{z} - |\mathbf{a}_i^T \mathbf{x}| \frac{\mathbf{a}_i^T \mathbf{z}}{|\mathbf{a}_i^T \mathbf{z}|} \right) \mathbf{a}_i \mathbb{1}_{\mathcal{E}_i} \\ &= \frac{1}{m} \sum_{i=1}^m \mathbf{a}_i \mathbf{a}_i^T \mathbf{h} \mathbb{1}_{\mathcal{E}_i} - \frac{1}{m} \sum_{i=1}^m \left(\frac{\mathbf{a}_i^T \mathbf{z}}{|\mathbf{a}_i^T \mathbf{z}|} - \frac{\mathbf{a}_i^T \mathbf{x}}{|\mathbf{a}_i^T \mathbf{x}|} \right) |\mathbf{a}_i^T \mathbf{x}| \mathbf{a}_i \mathbb{1}_{\mathcal{E}_i}. \end{aligned} \quad (30)$$

Using the definitions and properties in Lemma 4, one further arrives at

$$\begin{aligned}
\left\langle \frac{1}{2m} \nabla \ell_{\text{tr}}(\mathbf{z}), \mathbf{h} \right\rangle &\geq \frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h})^2 \mathbb{1}_{\mathcal{E}_i} - \frac{2}{m} \sum_{i=1}^m |\mathbf{a}_i^\top \mathbf{x}| |\mathbf{a}_i^\top \mathbf{h}| \mathbb{1}_{\mathcal{E}_i \cap \mathcal{K}_i} \\
&\geq \frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h})^2 \mathbb{1}_{\mathcal{E}_i} - \frac{2}{m} \sum_{i=1}^m |\mathbf{a}_i^\top \mathbf{x}| |\mathbf{a}_i^\top \mathbf{h}| \mathbb{1}_{\mathcal{D}_i} \\
&\geq \frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h})^2 \mathbb{1}_{\mathcal{E}_i} - \frac{1+\gamma}{2+\gamma} \cdot \frac{2}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h})^2 \mathbb{1}_{\mathcal{D}_i}
\end{aligned} \tag{31}$$

where the last inequality arises from the property $|\mathbf{a}_i^\top \mathbf{x}| \leq \frac{1+\gamma}{2+\gamma} |\mathbf{a}_i^\top \mathbf{h}|$ by the definition of \mathcal{D}_i .

Proving the regularity condition boils down to lower bounding the right-hand side of (31), specifically, to lower bounding the first term and to upper bounding the second one. Apparently, the first term approximately gives $\|\mathbf{h}\|^2$ by the SLLN as long as our truncation procedure does not eliminate too many generalized gradient components (i.e., summands in the first term). Regarding the second, one would expect its contribution to be small under our initialization condition in (5) and as the relative error $\|\mathbf{h}\| / \|\mathbf{x}\|$ decreases. Specifically, under our initialization, \mathcal{D}_i is provably a rare event, thus eliminating the possibility of the second term exerting a noticeable influence on the first term. Rigorous analyses concerning the two terms are elaborated in Lemma 5 and Lemma 6, whose proofs can be found in Appendix A.4 and Appendix A.5, respectively.

Lemma 5. Fix $\gamma \geq 1/2$ and $\rho \leq 1/10$, and let \mathcal{E}_i be defined in (26). For independent random variables $W \sim \mathcal{N}(0, 1)$ and $Z \sim \mathcal{N}(0, 1)$, set

$$\zeta_1 := 1 - \min \left\{ \mathbb{E} \left[\mathbb{1}_{\left\{ \left| \frac{1-\rho}{\rho} + \frac{W}{Z} \right| \geq \frac{\sqrt{1.01}}{\rho(1+\gamma)} \right\}} \right], \mathbb{E} \left[Z^2 \mathbb{1}_{\left\{ \left| \frac{1-\rho}{\rho} + \frac{W}{Z} \right| \geq \frac{\sqrt{1.01}}{\rho(1+\gamma)} \right\}} \right] \right\}. \tag{32}$$

Then for any $\epsilon > 0$ and any vector \mathbf{h} obeying $\|\mathbf{h}\| / \|\mathbf{x}\| \leq \rho$, the following holds with probability exceeding $1 - 2e^{-c_5 \epsilon^2 m}$

$$\frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h})^2 \mathbb{1}_{\mathcal{E}_i} \geq (1 - \zeta_1 - \epsilon) \|\mathbf{h}\|^2, \tag{33}$$

provided that $m > (c_6 \cdot \epsilon^{-2} \log \epsilon^{-1})n$ for some universal constants $c_5, c_6 > 0$.

To have a sense of how large the quantities involved in (5) are, when $\gamma = 0.7$ and $\rho = 1/10$, it holds $\mathbb{E} \left[\mathbb{1}_{\left\{ \left| \frac{1-\rho}{\rho} + \frac{W}{Z} \right| \geq \frac{\sqrt{1.01}}{\rho(1+\gamma)} \right\}} \right] \approx 0.92$, and $\mathbb{E} \left[Z^2 \mathbb{1}_{\left\{ \left| \frac{1-\rho}{\rho} + \frac{W}{Z} \right| \geq \frac{\sqrt{1.01}}{\rho(1+\gamma)} \right\}} \right] \approx 0.99$, hence leading to $\zeta_1 \approx 0.08$.

Having derived a lower bound for the first term in the right-hand side of (31), it remains to deal with the second one.

Lemma 6. Fix $\gamma > 0$ and $\rho \leq 1/10$, and let \mathcal{D}_i be defined in (27). For any constant $\epsilon > 0$, there exists some universal constants $c_5, c_6 > 0$ such that

$$\frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h})^2 \mathbb{1}_{\mathcal{D}_i} \leq (\zeta_2' + \epsilon) \|\mathbf{h}\|^2 \tag{34}$$

holds with probability at least $1 - 2e^{-c_5 \epsilon^2 m}$ provided that $m/n > (c_6 \cdot \epsilon^{-2} \log \epsilon^{-1})$ for some universal constants $c_5, c_6 > 0$, where $\zeta_2' = 1.3785 \sqrt{\rho \tau / (0.99 \tau^2 - \rho^2)}$ with $\tau := (2 + \gamma)/(1 + \gamma)$.

With our TGGF default parameters $\rho = 1/10$ and $\gamma = 0.7$, we have $\zeta'_2 \approx 0.3483$. Taking results in (31), (33), and (34) together, choosing m/n exceeding some sufficiently large constant such that $c_0 \leq c_5 \epsilon^2$, and denoting $\zeta_2 := \zeta'_2(1 + \gamma)/(2 + \gamma)$, then with probability exceeding $1 - 4e^{-c_0 m}$, the following

$$\left\langle \mathbf{h}, \frac{1}{2m} \nabla \ell_{\text{tr}}(\mathbf{z}) \right\rangle \geq (1 - \zeta_1 - \zeta_2 - 2\epsilon) \|\mathbf{h}\|^2 \quad (35)$$

holds for all \mathbf{x} and \mathbf{z} such that $\|\mathbf{h}\| / \|\mathbf{x}\| \leq \rho$ for $0 < \rho \leq 1/10$ and any fixed $\gamma \geq 1/2$. This combining with (17) and (19) proves Proposition 2 for appropriately chosen $\mu > 0$ and $\lambda > 0$.

To conclude this section, an estimate for the working step size is provided next. To be specific, plugging the results in (23) and (25) into (18) suggests that

$$\text{dist}^2 \left(\mathbf{z} - \frac{\mu}{m} \nabla \ell_{\text{tr}}(\mathbf{z}), \mathbf{x} \right) = \|\mathbf{h}\|^2 - 2\mu \left\langle \mathbf{h}, \frac{1}{m} \nabla \ell_{\text{tr}}(\mathbf{z}) \right\rangle + \left\| \frac{\mu}{m} \nabla \ell_{\text{tr}}(\mathbf{z}) \right\|^2 \quad (36)$$

$$\begin{aligned} &\leq \{1 - \mu [4(1 - \xi_1 - \xi_2 - 2\epsilon) - \mu(1 + \delta)^4]\} \|\mathbf{h}\|^2 \\ &\triangleq (1 - \nu) \|\mathbf{h}\|^2. \end{aligned} \quad (37)$$

Taking ϵ and δ to be sufficiently small, one obtains the feasible range of the step size for TGGF

$$\mu \leq \frac{4(0.99 - \xi_1 - \xi_2)}{1.02} \triangleq \mu_0, \quad (38)$$

thus concluding the proof of Theorem 1.

A Proof details

By homogeneity, it suffices to work with the case where $\|\mathbf{x}\| = 1$.

A.1 Proof of Lemma 1

It is easy to check that

$$\begin{aligned} \frac{1}{2} \|\mathbf{x}\mathbf{x}^\mathcal{T} - \tilde{\mathbf{z}}_0\tilde{\mathbf{z}}_0^\mathcal{T}\|_F^2 &= \frac{1}{2} \|\mathbf{x}\|^4 + \frac{1}{2} \|\tilde{\mathbf{z}}_0\|^4 - |\mathbf{x}^\mathcal{T} \tilde{\mathbf{z}}_0|^2 \\ &= 1 - |\mathbf{x}^\mathcal{T} \tilde{\mathbf{z}}_0|^2 \\ &= 1 - \cos^2 \theta \end{aligned} \quad (39)$$

where $0 \leq \theta \leq \pi$ is the angle between the spaces spanned by \mathbf{x} and $\tilde{\mathbf{z}}_0$. Then one can write

$$\mathbf{x} = \cos \theta \tilde{\mathbf{z}}_0 + \sin \theta \tilde{\mathbf{z}}_0^\perp, \quad (40)$$

where $\tilde{\mathbf{z}}_0^\perp \in \mathbb{R}^n$ is a unit vector that is orthogonal to $\tilde{\mathbf{z}}_0$ and has a nonnegative inner product with \mathbf{x} . Likewise, one can express

$$\mathbf{x}^\perp := -\sin \theta \tilde{\mathbf{z}}_0 + \cos \theta \tilde{\mathbf{z}}_0^\perp, \quad (41)$$

in which $\mathbf{x}^\perp \in \mathbb{R}^n$ is a unit vector orthogonal to \mathbf{x} .

Since $\tilde{\mathbf{z}}_0$ is the solution to the maximum eigenvalue problem

$$\tilde{\mathbf{z}}_0 := \arg \max_{\|\mathbf{z}\|=1} \mathbf{z}^\mathcal{T} \bar{\mathbf{Y}}_0 \mathbf{z} \quad (42)$$

for $\bar{\mathbf{Y}}_0 := \frac{1}{|\bar{\mathcal{I}}_0|} \bar{\mathbf{S}}_0^T \bar{\mathbf{S}}_0$, it is the leading eigenvector of $\bar{\mathbf{Y}}_0$, i.e., $\bar{\mathbf{Y}}_0 \tilde{\mathbf{z}}_0 = \lambda_1 \tilde{\mathbf{z}}_0$, where $\lambda_1 > 0$ is the largest eigenvalue of $\bar{\mathbf{Y}}_0$. Premultiplying (40) and (41) by $\bar{\mathbf{S}}_0$ yields

$$\bar{\mathbf{S}}_0 \mathbf{x} = \cos \theta \bar{\mathbf{S}}_0 \tilde{\mathbf{z}}_0 + \sin \theta \bar{\mathbf{S}}_0 \tilde{\mathbf{z}}_0^\perp, \quad (43a)$$

$$\bar{\mathbf{S}}_0 \mathbf{x}^\perp = -\sin \theta \bar{\mathbf{S}}_0 \tilde{\mathbf{z}}_0 + \cos \theta \bar{\mathbf{S}}_0 \tilde{\mathbf{z}}_0^\perp. \quad (43b)$$

Pythagoras' relationship now gives

$$\|\bar{\mathbf{S}}_0 \mathbf{x}\|^2 = \cos^2 \theta \|\bar{\mathbf{S}}_0 \tilde{\mathbf{z}}_0\|^2 + \sin^2 \theta \|\bar{\mathbf{S}}_0 \tilde{\mathbf{z}}_0^\perp\|^2, \quad (44a)$$

$$\|\bar{\mathbf{S}}_0 \mathbf{x}^\perp\|^2 = \sin^2 \theta \|\bar{\mathbf{S}}_0 \tilde{\mathbf{z}}_0\|^2 + \cos^2 \theta \|\bar{\mathbf{S}}_0 \tilde{\mathbf{z}}_0^\perp\|^2, \quad (44b)$$

where the cross-terms vanish because $\tilde{\mathbf{z}}_0^T \bar{\mathbf{S}}_0^T \bar{\mathbf{S}}_0 \tilde{\mathbf{z}}_0^\perp = |\bar{\mathcal{I}}_0| \tilde{\mathbf{z}}_0^T \bar{\mathbf{Y}}_0 \tilde{\mathbf{z}}_0^\perp = \lambda_1 |\bar{\mathcal{I}}_0| \tilde{\mathbf{z}}_0^T \tilde{\mathbf{z}}_0^\perp = 0$ following from the definition of $\tilde{\mathbf{z}}_0^\perp$.

We next construct the following expression

$$\begin{aligned} & \sin^2 \theta \|\bar{\mathbf{S}}_0 \mathbf{x}\|^2 - \|\bar{\mathbf{S}}_0 \mathbf{x}^\perp\|^2 \\ &= \sin^2 \theta \cos^2 \theta \|\bar{\mathbf{S}}_0 \tilde{\mathbf{z}}_0\|^2 + \sin^4 \theta \|\bar{\mathbf{S}}_0 \tilde{\mathbf{z}}_0^\perp\|^2 - \sin^2 \theta \|\bar{\mathbf{S}}_0 \tilde{\mathbf{z}}_0\|^2 - \cos^2 \theta \|\bar{\mathbf{S}}_0 \tilde{\mathbf{z}}_0^\perp\|^2 \\ &= \sin^2 \theta \left(\cos^2 \theta \|\bar{\mathbf{S}}_0 \tilde{\mathbf{z}}_0\|^2 - \|\bar{\mathbf{S}}_0 \tilde{\mathbf{z}}_0\|^2 + \sin^2 \theta \|\bar{\mathbf{S}}_0 \tilde{\mathbf{z}}_0^\perp\|^2 \right) - \cos^2 \theta \|\bar{\mathbf{S}}_0 \tilde{\mathbf{z}}_0^\perp\|^2 \\ &= \sin^4 \theta (\|\bar{\mathbf{S}}_0 \tilde{\mathbf{z}}_0^\perp\|^2 - \|\bar{\mathbf{S}}_0 \tilde{\mathbf{z}}_0\|^2) - \cos^2 \theta \|\bar{\mathbf{S}}_0 \tilde{\mathbf{z}}_0^\perp\|^2 \\ &\leq 0 \end{aligned}$$

where $\bar{\mathbf{S}}_0^T \bar{\mathbf{S}}_0 \succeq \mathbf{0}$, so $\|\bar{\mathbf{S}}_0 \tilde{\mathbf{z}}_0^\perp\|^2 - \|\bar{\mathbf{S}}_0 \tilde{\mathbf{z}}_0\|^2 \leq 0$ holds for any unit vector $\tilde{\mathbf{z}}_0^\perp \in \mathbb{R}^n$ arising from the fact that $\tilde{\mathbf{z}}_0$ maximizes the term in (7), hence yielding

$$\sin^2 \theta = 1 - \cos^2 \theta \leq \frac{\|\bar{\mathbf{S}}_0 \mathbf{x}^\perp\|^2}{\|\bar{\mathbf{S}}_0 \mathbf{x}\|^2}. \quad (45)$$

Upon letting $\mathbf{u} = \mathbf{x}^\perp$, the last inequality taken together with (39) concludes the proof of (6).

A.2 Proof of Lemma 2

Recall that rows in $\bar{\mathbf{S}}_0 \in \mathbb{R}^{|\bar{\mathcal{I}}_0| \times n}$, hereafter denoted by $\mathbf{s}_i^T \in \mathbb{R}^{1 \times n}$, $\forall i \in [|\bar{\mathcal{I}}_0|]$, are drawn uniformly on the unit sphere. The uniformly spherical distribution is rotationally invariant, so it suffices to prove the results in the case where $\mathbf{x} = \mathbf{e}_1$ with \mathbf{e}_1 being the first canonical vector in \mathbb{R}^n . Indeed, any unit vector \mathbf{x} can be expressed as $\mathbf{x} = \mathbf{U} \mathbf{e}_1$ for some orthogonal transformation $\mathbf{U} \in \mathbb{R}^{n \times n}$. To see this, consider the following [10]

$$|\langle \mathbf{s}_i, \mathbf{x} \rangle|^2 = |\langle \mathbf{s}_i, \mathbf{U} \mathbf{e}_1 \rangle|^2 = |\langle \mathbf{U}^T \mathbf{s}_i, \mathbf{e}_1 \rangle|^2 \stackrel{d}{=} |\langle \mathbf{s}_i, \mathbf{e}_1 \rangle|^2, \quad (46)$$

where $\stackrel{d}{=}$ means terms involved on both sides of the equality have the same distribution. Thus, the problem of finding any unit-normed \mathbf{x} is equivalent to that of finding \mathbf{e}_1 . Henceforth, we assume without any loss of generality that $\mathbf{x} = \mathbf{e}_1$.

Considering a unit vector \mathbf{x}^\perp such that $\mathbf{x}^T \mathbf{x}^\perp = \mathbf{e}_1^T \mathbf{x}^\perp = 0$, there exists a unit vector $\mathbf{d} \in \mathbb{R}^{n-1}$ such that $\mathbf{x}^\perp = [0 \ \mathbf{d}^T]^T$. So it holds that

$$\|\bar{\mathbf{S}}_0 \mathbf{x}^\perp\|^2 = \left\| \bar{\mathbf{S}}_0 [0 \ \mathbf{d}^T]^T \right\|^2 = \|\mathbf{F} \mathbf{d}\|^2, \quad (47)$$

where $\mathbf{F} \in \mathbb{R}^{|\bar{\mathcal{I}}_0| \times (n-1)}$ is obtained through deleting the first column in $\bar{\mathbf{S}}_0$, denoted by $\bar{\mathbf{S}}_{0,1}$, i.e., $\bar{\mathbf{S}}_0 = [\bar{\mathbf{S}}_{0,1} \ \mathbf{F}]$. Letting $\mathbf{F} := [\mathbf{f}_1 \ \cdots \ \mathbf{f}_{|\bar{\mathcal{I}}_0|}]^\top$, one can readily write $\mathbf{s}_i = [s_{i,1} \ \mathbf{f}_i^\top]^\top$, $\forall i \in [|\bar{\mathcal{I}}_0|] := \{1, \dots, |\bar{\mathcal{I}}_0|\}$. Uniformly spherically distributed $\mathbf{s}_i \in \mathbb{R}^n$ has statistics $\mathbb{E}[\mathbf{s}_i] = \mathbf{0}$, and $\mathbb{E}[\mathbf{s}_i \mathbf{s}_i^\top] = \frac{1}{n} \mathbf{I}_n$ [11]. Leveraging the linearity of expectation operator, one arrives at

$$\mathbb{E}[\mathbf{s}_i] = \mathbb{E} \begin{bmatrix} s_{i,1} \\ \mathbf{f}_i \end{bmatrix} = \begin{bmatrix} \mathbb{E}[s_{i,1}] \\ \mathbb{E}[\mathbf{f}_i] \end{bmatrix} = \mathbf{0}, \ \forall i \quad (48)$$

to yield

$$\mathbb{E}[\mathbf{f}_i] = \mathbf{0}, \ \forall i. \quad (49)$$

A similar argument holds for the second-order moment

$$\mathbb{E}[\mathbf{s}_i \mathbf{s}_i^\top] = \begin{bmatrix} \mathbb{E}[s_{i,1}^2] & \mathbb{E}[s_{i,1} \mathbf{f}_i^\top] \\ \mathbb{E}[s_{i,1} \mathbf{f}_i] & \mathbb{E}[\mathbf{f}_i \mathbf{f}_i^\top] \end{bmatrix} = \frac{1}{n} \mathbf{I}_n, \ \forall i \quad (50)$$

leading to

$$\mathbb{E}[\mathbf{f}_i \mathbf{f}_i^\top] = \frac{1}{n} \mathbf{I}_{n-1}, \ \forall i. \quad (51)$$

Recall that a random vector $\mathbf{z} \in \mathbb{R}^n$ is said to be *isotropic* if it has zero-mean and identity covariance matrix [9, Definition 5.19]. Then recognize, from (49) and (51), that a proper scaling of \mathbf{f}_i renders $\sqrt{n} \mathbf{f}_i$ isotropic. Further, it is known that a spherical random vector is subgaussian, and its subgaussian norm is bounded by an absolute constant [9]. Indeed, this comes from the following geometric argument: using rotational invariance of the uniform spherical distribution \mathcal{S}^{n-1} in \mathbb{R}^n , it holds that, given any $\epsilon \geq 0$, the spherical cap $\{\mathbf{s}_i \in \mathcal{S}^{n-1} : s_{i,1} > \epsilon\}$ consists of at most $e^{-\epsilon^2 n/2}$ proportion of the total area on the sphere. A similar argument carries over to \mathbf{f}_i , and thus, \mathbf{f}_i is subgaussian as well.

Standard concentration inequalities on the sum of random positive semi-definite matrices composed of independent isotropic subgaussian rows [9, Remark 5.40] confirm that

$$\left\| \frac{1}{|\bar{\mathcal{I}}_0|} (\sqrt{n} \mathbf{F})^\top (\sqrt{n} \mathbf{F}) - \mathbf{I}_{n-1} \right\| \leq \sigma \|\mathbf{I}_{n-1}\| \quad (52)$$

holds with probability at least $1 - 2e^{-c_K n}$ as long as $|\bar{\mathcal{I}}_0|/n$ is sufficiently large, where σ is a numerical constant that can take arbitrarily small values and $c_K > 0$ is a universal constant. Without loss of generality, let us work with $\sigma := 0.01$ in (52), so for any unit vector $\mathbf{d} \in \mathbb{R}^{n-1}$, the following inequality holds with probability at least $1 - 2e^{-c_K n}$,

$$\left| \frac{n}{|\bar{\mathcal{I}}_0|} \mathbf{d}^\top \mathbf{F}^\top \mathbf{F} \mathbf{d} - \mathbf{d}^\top \mathbf{d} \right| \leq 0.01 \mathbf{d}^\top \mathbf{d}, \quad (53)$$

or equivalently,

$$\|\mathbf{F} \mathbf{d}\|^2 = |\mathbf{d}^\top \mathbf{F}^\top \mathbf{F} \mathbf{d}| \leq 1.01 |\bar{\mathcal{I}}_0|/n. \quad (54)$$

Combining the last with (47), one readily concludes that

$$\|\bar{\mathbf{S}}_0 \mathbf{x}^\perp\|^2 \leq 1.01 |\bar{\mathcal{I}}_0|/n \quad (55)$$

holds with probability at least $1 - 2e^{-c_K n}$, provided that $|\bar{\mathcal{I}}_0|/n$ exceeds some constant. Note that c_K depends on the maximum subgaussian norm of the rows of $\sqrt{n} \mathbf{F}$, and we assume without loss of generality $c_K \geq 1/2$. Hence, $\|\bar{\mathbf{S}}_0 \mathbf{u}\|^2$ in (6) is upper bounded simply by letting $\mathbf{u} = \mathbf{x}^\perp$ in (55).

A.3 Proof of Lemma 3

We next pursue a meaningful lower bound for $\|\bar{\mathbf{S}}_0 \mathbf{x}\|^2$ in (9). When $\mathbf{x} = \mathbf{e}_1$, one has $\|\bar{\mathbf{S}}_0 \mathbf{x}\|^2 = \|\bar{\mathbf{S}}_0 \mathbf{e}_1\|^2 = \sum_{i=1}^{|\bar{\mathcal{I}}_0|} s_{i,1}^2$. It is further worth mentioning that all squared entries of any spherical random vector \mathbf{s}_i obey the *Beta* distribution with parameters $\alpha = \frac{1}{2}$, and $\beta = \frac{n-1}{2}$, i.e., $s_{i,j}^2 \sim \text{Beta}(\frac{1}{2}, \frac{n-1}{2})$, $\forall i, j$, [11, Lemma 2]. Although they have closed-form probability density functions (pdfs) that may facilitate deriving a wanted lower bound, we shall take another easier route detailed as follows. A simple yet useful inequality is established first.

Lemma 7. *Given m fractions obeying $1 > \frac{p_1}{q_1} \geq \frac{p_2}{q_2} \geq \dots \geq \frac{p_m}{q_m} > 0$, in which $p_i, q_i > 0, \forall i \in [m]$, the following holds for all $1 \leq k \leq m$*

$$\sum_{i=1}^k \frac{p_i}{q_i} \geq \sum_{i=1}^k \frac{p_{[i]}}{q_{[1]}} \quad (56)$$

where $p_{[i]}$ denotes the i -th largest one among $\{p_i\}_{i=1}^m$, and hence, $q_{[1]}$ is the maximum in $\{q_i\}_{i=1}^m$.

Proof. For any $k \in [m]$, according to the definition of $q_{[i]}$, it holds that $p_{[1]} \geq p_{[2]} \geq \dots \geq p_{[k]}$, so $\frac{p_{[1]}}{q_{[1]}} \geq \frac{p_{[2]}}{q_{[1]}} \geq \dots \geq \frac{p_{[k]}}{q_{[1]}}$. Considering $q_{[1]} \geq q_i, \forall i \in [m]$, and letting $j_i \in [m]$ be the index such that $p_{j_i} = p_{[i]}$, then $\frac{p_{j_i}}{q_{j_i}} = \frac{p_{[i]}}{q_{j_i}} \geq \frac{p_{[i]}}{q_{[1]}}$ holds for any $i \in [k]$. Therefore, $\sum_{i=1}^k \frac{p_{j_i}}{q_{j_i}} = \sum_{i=1}^k \frac{p_{[i]}}{q_{j_i}} \geq \sum_{i=1}^k \frac{p_{[i]}}{q_{[1]}}$. Note that $\left\{ \frac{p_{[i]}}{q_{j_i}} \right\}_{i=1}^k$ comprise a subset of terms in $\left\{ \frac{p_i}{q_i} \right\}_{i=1}^m$. On the other hand, according to our assumption, $\sum_{i=1}^k \frac{p_i}{q_i}$ is the largest among all sums of k summands; hence, $\sum_{i=1}^k \frac{p_i}{q_i} \geq \sum_{i=1}^k \frac{p_{[i]}}{q_{j_i}}$ yields $\sum_{i=1}^k \frac{p_i}{q_i} \geq \sum_{i=1}^k \frac{p_{[i]}}{q_{[1]}}$ concluding the proof. \square

Without loss of generality and for simplicity of exposition, let us assume that indices of \mathbf{a}_i 's have been re-ordered such that

$$\frac{a_{1,1}^2}{\|\mathbf{a}_1\|^2} \geq \frac{a_{2,1}^2}{\|\mathbf{a}_2\|^2} \geq \dots \geq \frac{a_{m,1}^2}{\|\mathbf{a}_m\|^2}, \quad (57)$$

where $a_{i,1}$ denotes the first element of \mathbf{a}_i . Therefore, writing $\|\bar{\mathbf{S}}_0 \mathbf{e}_1\|^2 = \sum_{i=1}^{|\bar{\mathcal{I}}_0|} a_{i,1}^2 / \|\mathbf{a}_i\|^2$, the next task amounts to finding the sum of the $|\bar{\mathcal{I}}_0|$ largest out of all m entities in (57). Applying the result (56) in Lemma 7 gives

$$\sum_{i=1}^{|\bar{\mathcal{I}}_0|} \frac{a_{i,1}^2}{\|\mathbf{a}_i\|^2} \geq \sum_{i=1}^{|\bar{\mathcal{I}}_0|} \frac{a_{[i],1}^2}{\max_{i \in [m]} \|\mathbf{a}_i\|^2}, \quad (58)$$

in which $a_{[i],1}^2$ stands for the i -th largest entity in $\{a_{i,1}^2\}_{i=1}^m$.

Observe that for i.i.d. random vectors $\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$, the property $\mathbb{P}(\|\mathbf{a}_i\|^2 \geq 2.3n) \leq e^{-n/2}$ holds for large enough n (e.g., $n \geq 20$), which can be understood upon substituting $\xi := n/2$ into the following standard result [12, Lemma 1]

$$\mathbb{P}\left(\|\mathbf{a}_i\|^2 - n \geq 2\sqrt{\xi} + 2\xi\right) \leq e^{-\xi}. \quad (59)$$

In addition, one readily concludes that $\mathbb{P}\left(\max_{i \in [m]} \|\mathbf{a}_i\| \leq \sqrt{2.3n}\right) \geq 1 - me^{-n/2}$. We will henceforth build our subsequent proofs on this event without stating this explicitly each time encountering it.

Therefore, (58) can be lower bounded by

$$\|\bar{\mathbf{S}}\mathbf{x}\|^2 = \sum_{i=1}^{|\bar{\mathcal{I}}_0|} \frac{a_{i,1}^2}{\|\mathbf{a}_i\|^2} \geq \sum_{i=1}^{|\bar{\mathcal{I}}_0|} \frac{a_{[i],1}^2}{\max_{i \in [m]} \|\mathbf{a}_i\|^2} \geq \frac{1}{2.3n} \sum_{i=1}^{|\bar{\mathcal{I}}_0|} |a_{[i],1}|^2, \quad (60)$$

which holds with probability at least $1 - me^{-n/2}$. The task left for bounding $\|\bar{\mathbf{S}}\mathbf{x}\|^2$ is to derive a meaningful lower bound for $\sum_{i=1}^{|\bar{\mathcal{I}}_0|} a_{[i],1}^2$. Roughly speaking, because the ratio $|\bar{\mathcal{I}}_0|/m$ is small, e.g., $|\bar{\mathcal{I}}_0|/m \leq 1/5$, a trivial result consists of bounding $(1/|\bar{\mathcal{I}}_0|) \sum_{i=1}^{|\bar{\mathcal{I}}_0|} a_{[i],1}^2$ by its sample average $(1/m) \sum_{i=1}^m a_{[i],1}^2$. The latter can be bounded using its ensemble mean, i.e., $\mathbb{E}[a_{i,1}^2] = 1$, $\forall i \in [\bar{\mathcal{I}}_0]$, to yield $(1/m) \sum_{i=1}^m a_{[i],1}^2 \geq (1 - \epsilon) \mathbb{E}[a_{i,1}^2] = 1 - \epsilon$, which holds with high probability for some numerical constant $\epsilon > 0$ [10, Lemma 3.1]. Therefore, one has a candidate lower bound $\sum_{i=1}^{|\bar{\mathcal{I}}_0|} a_{[i],1}^2 \geq (1 - \epsilon) |\bar{\mathcal{I}}_0|$. Nonetheless, this lower bound is in general too loose, and it contributes to a relatively large upper bound on the wanted term in (6).

To obtain an alternative bound, let us examine first the typical size of the maximum in $\{a_{i,1}^2\}_{i=1}^m$. Observe obviously that the modulus $|a_{i,1}|$ follows the half-normal distribution having the pdf $p(r) = \sqrt{2/\pi} \cdot e^{-r^2/2}$, $r > 0$, and it is easy to verify that

$$\mathbb{E}[|a_{i,1}|] = \sqrt{2/\pi}. \quad (61)$$

Then integrating the pdf from 0 to $+\infty$ yields the corresponding accumulative distribution function (cdf) expressible in terms of the error function $\mathbb{P}(|a_{i,1}| > \xi) = 1 - \text{erf}(\xi/2)$, i.e., $\text{erf}(\xi) := 2/\sqrt{\pi} \cdot \int_0^\xi e^{-r^2} dr$. Appealing to a lower bound on the complimentary error function $\text{erfc}(\xi) := 1 - \text{erf}(\xi)$ from [13, Theorem 2], one establishes that $\mathbb{P}(|a_{i,1}| > \xi) = 1 - \text{erf}(\xi/2) \geq (3/5)e^{-\xi^2/2}$. Additionally, direct application of probability theory and Taylor expansion confirms that

$$\begin{aligned} \mathbb{P}\left(\max_{i \in [m]} |a_{i,1}| \geq \xi\right) &= 1 - [\mathbb{P}(|a_{i,1}| \leq \xi)]^m \\ &\geq 1 - \left(1 - 0.6e^{-\xi^2/2}\right)^m \\ &\geq 1 - e^{-0.6me^{-\xi^2/2}}. \end{aligned} \quad (62)$$

Choosing now $\xi := \sqrt{2 \log n}$ leads to

$$\mathbb{P}\left(\max_{i \in [m]} |a_{i,1}| \geq \sqrt{2 \log n}\right) \geq 1 - e^{-0.6m/n} \geq 1 - o(1) \quad (63)$$

which holds with the proviso that m/n is large enough, and the symbol $o(1)$ represents a small constant probability. Thus, provided that m/n exceeds some large constant, the event $\max_{i \in [m]} a_{i,1}^2 \geq 2 \log n$ occurs with high probability. Hence, one may expect a tighter lower bound than $(1 - \epsilon_0) |\bar{\mathcal{I}}_0|$, which is on the same order of m under the assumption that $|\bar{\mathcal{I}}_0|/m$ is about a constant.

Although $a_{i,1}^2$ obeys the *Chi-square* distribution with $k = 1$ degrees of freedom, its cdf is rather complicated and does not admit a nice closed-form expression. A small trick is hence taken in the sequel. Postulate without loss of generality that both m and $|\bar{\mathcal{I}}_0|$ are even. Grouping two consecutive $a_{[i],1}^2$'s together, introduce a new variable $\vartheta[i] := a_{[2k-1],1}^2 + a_{[2k],1}^2$, $\forall k \in [m/2]$, hence yielding a sequence of ordered numbers, i.e., $\vartheta[1] \geq \vartheta[2] \geq \dots \geq \vartheta[m/2] > 0$. Then, one can equivalently write the wanted sum as

$$\sum_{i=1}^{|\bar{\mathcal{I}}_0|} a_{[i],1}^2 = \sum_{i=1}^{|\bar{\mathcal{I}}_0|/2} \vartheta[i]. \quad (64)$$

On the other hand, for i.i.d. standard normal random variables $\{a_{i,1}\}_{i=1}^m$, let us consider grouping randomly two of them and denote the corresponding sum of their squares by $\chi_k := a_{k_i,1}^2 + a_{k_j,1}^2$, where $k_i \neq k_j \in [m]$, and $k \in [m/2]$. It is self-evident that the χ_k 's are identically distributed obeying the *Chi-square* distribution with $k = 2$ degrees of freedom, having the pdf

$$p(r) = \frac{1}{2}e^{-\frac{r}{2}}, \quad r \geq 0, \quad (65)$$

and the following complementary cdf (ccdf)

$$\mathbb{P}(\chi_k \geq \xi) := \int_{\xi}^{\infty} \frac{1}{2}e^{-\frac{r}{2}}dr = e^{-\frac{\xi}{2}}, \quad \forall \xi \geq 0. \quad (66)$$

Ordering all χ_k 's, summing the $|\bar{\mathcal{I}}_0|/2$ largest ones, and comparing the resultant sum with the one in (64) confirm that

$$\sum_{i=1}^{|\bar{\mathcal{I}}_0|/2} \chi_{[i]} \leq \sum_{i=1}^{|\bar{\mathcal{I}}_0|/2} \vartheta_{[i]} = \sum_{i=1}^{|\bar{\mathcal{I}}_0|} a_{[i],1}^2, \quad \forall |\bar{\mathcal{I}}_0| \in [m]. \quad (67)$$

Upon setting $\mathbb{P}(\chi_k \geq \xi) = |\bar{\mathcal{I}}_0|/m$, one obtains an estimate of $\chi_{|\bar{\mathcal{I}}_0|/2}$, the $(|\bar{\mathcal{I}}_0|/2)$ -th largest value in $\{\chi_k\}_{k=1}^{m/2}$ as follows

$$\hat{\chi}_{|\bar{\mathcal{I}}_0|/2} := 2 \log(m/|\bar{\mathcal{I}}_0|). \quad (68)$$

Furthermore, applying the Hoeffding-type inequality [9, Proposition 5.10] and leveraging the convexity of the ccdf in (66), one readily establishes that

$$\mathbb{P}(\hat{\chi}_{|\bar{\mathcal{I}}_0|/2} - \chi_{|\bar{\mathcal{I}}_0|/2} > \xi) \leq e^{-\frac{1}{4}m\xi^2 e^{-\xi} (|\bar{\mathcal{I}}_0|/m)^2}, \quad \forall \xi > 0. \quad (69)$$

Taking without loss of generality $\xi := 0.05\hat{\chi}_{|\bar{\mathcal{I}}_0|/2} = 0.1 \log(m/|\bar{\mathcal{I}}_0|)$ gives

$$\mathbb{P}(\chi_{|\bar{\mathcal{I}}_0|/2} < 0.95\hat{\chi}_{|\bar{\mathcal{I}}_0|/2}) \leq e^{-c_0 m} \quad (70)$$

for some universal constants $c_0, c_{\chi} > 0$, and sufficiently large n such that $|\bar{\mathcal{I}}_0|/m \gtrsim c_{\chi} > 0$. The remaining part in this section assumes that this event occurs.

Choosing $\xi := 4 \log n$ and substituting this into the ccdf in (66) leads to

$$\mathbb{P}(\chi \leq 4 \log n) = 1 - 1/n^2. \quad (71)$$

Notice that each summand in $\sum_{i=1}^{|\bar{\mathcal{I}}_0|/2} \chi_{[i]} \geq \sum_{i=1}^{m/2} \chi_i \mathbb{1}_{\bar{\mathcal{E}}_i}$ is Chi-square distributed, and hence could be unbounded, so we choose to work with the truncation $\sum_{i=1}^{m/2} \chi_i \mathbb{1}_{\bar{\mathcal{E}}_i}$, where the $\mathbb{1}_{\bar{\mathcal{E}}_i}$'s are independent copies of $\mathbb{1}_{\bar{\mathcal{E}}}$, and $\mathbb{1}_{\bar{\mathcal{E}}}$ denotes the indicator function for the ensuing events

$$\bar{\mathcal{E}} := \left\{ \chi \geq \hat{\chi}_{|\bar{\mathcal{I}}_0|/2} \right\} \cap \{ \chi \leq 4 \log n \}. \quad (72)$$

Apparently, it holds that $\sum_{i=1}^{|\bar{\mathcal{I}}_0|/2} \chi_{[i]} \geq \sum_{i=1}^{m/2} \chi_i \mathbb{1}_{\bar{\mathcal{E}}_i}$. One further establishes that

$$\begin{aligned} \mathbb{E}[\chi_i \mathbb{1}_{\bar{\mathcal{E}}_i}] &:= \int_{\hat{\chi}_{|\bar{\mathcal{I}}_0|/2}}^{4 \log n} \frac{1}{2} r e^{-r/2} dr \\ &= (\hat{\chi}_{|\bar{\mathcal{I}}_0|/2} + 2) e^{-\hat{\chi}_{|\bar{\mathcal{I}}_0|/2}/2} - (4 \log n + 2) e^{-2 \log n} \\ &= \frac{2|\bar{\mathcal{I}}_0|}{m} \left[1 + \log(m/|\bar{\mathcal{I}}_0|) \right] - \frac{(4 \log n + 2)}{n^2}. \end{aligned} \quad (73)$$

The task of bounding $\sum_{i=1}^{|\bar{\mathcal{I}}_0|} a_{[i],1}^2$ in (67) now boils down to bounding $\sum_{i=1}^{m/2} \chi_i \mathbb{1}_{\bar{\mathcal{E}}_i}$ from its expectation in (73). A convenient way to accomplish this is using the Bernstein inequality [9, Proposition 5.16], that deals with bounded random variables. That also justifies the reason of introducing the upper-bound truncation on χ in (72). Specifically, let us define

$$\vartheta_i := \chi_i \mathbb{1}_{\bar{\mathcal{E}}_i} - \mathbb{E} [\chi_i \mathbb{1}_{\bar{\mathcal{E}}_i}], \quad \forall i \in [m/2]. \quad (74)$$

Thus, $\{\vartheta_i\}_{i=1}^{m/2}$ are i.i.d. centered and bounded random variables following from the mean-subtraction and the upper-bound truncation. Further, according to the cdf (66) and the definition of sub-exponential random variables [9, Definition 5.13], the terms $\{\vartheta_i\}_{i=1}^{m/2}$ are sub-exponential. Then, the following

$$\left| \sum_{i=1}^{m/2} \vartheta_i \right| \geq \tau \quad (75)$$

holds with probability at least $1 - 2e^{-c_s \min(\tau/K_s, \tau^2/K_s^2)}$, in which $c_s > 0$ is a universal constant, and $K_s := \max_{i \in [m/2]} \|\vartheta_i\|_{\psi_1}$ represents the maximum subexponential norm of the ϑ_i 's. Indeed, K_s can be found as follows [9, Definition 5.13]

$$\begin{aligned} K_s &:= \sup_{p \geq 1} p^{-1} (\mathbb{E} [|\vartheta_i|^p])^{1/p} \\ &\leq \left(4 \log n - 2 \log (m/|\bar{\mathcal{I}}_0|) \right) \left[|\bar{\mathcal{I}}_0|/m - 1/n^2 \right] \\ &\leq \frac{2|\bar{\mathcal{I}}_0|}{m} \log (n^2 |\bar{\mathcal{I}}_0|/m) \\ &\leq \frac{4|\bar{\mathcal{I}}_0|}{m} \log n. \end{aligned} \quad (76)$$

Choosing $\tau := 8|\bar{\mathcal{I}}_0|/(c_s m) \cdot \log^2 n$ in (75) yields

$$\begin{aligned} \sum_{i=1}^{m/2} \chi_i \mathbb{1}_{\bar{\mathcal{E}}_i} &\geq |\bar{\mathcal{I}}_0| \left[1 + \log (m/|\bar{\mathcal{I}}_0|) \right] - 8|\bar{\mathcal{I}}_0|/(c_s m) \cdot \log^2 n - m(2 \log n + 1)/n^2 \\ &\geq (1 - \epsilon_s) |\bar{\mathcal{I}}_0| \left[1 + \log (m/|\bar{\mathcal{I}}_0|) \right] \end{aligned} \quad (77)$$

for some small constant $\epsilon_s > 0$, which holds with probability at least $1 - me^{-n/2} - e^{-c_0 m} - 3/n^2$ as long as m/n exceeds some numerical constant and n is sufficiently large. Therefore, combining (60), (67), and (77), one concludes that the following holds with high probability

$$\|\bar{\mathbf{S}}_0 \mathbf{x}\|^2 = \sum_{i=1}^{|\bar{\mathcal{I}}_0|} \frac{a_{i,1}^2}{\|\mathbf{a}_i\|^2} \geq (1 - \epsilon_s) \frac{|\bar{\mathcal{I}}_0|}{2.3n} \left[1 + \log (m/|\bar{\mathcal{I}}_0|) \right]. \quad (78)$$

Taking $\epsilon_s := 0.01$ without loss of generality concludes the proof of Lemma 3.

A.4 Proof of Lemma 5

Let us first prove the argument for a fixed pair \mathbf{h} and \mathbf{x} , so \mathbf{h} and \mathbf{z} are independent of $\{\mathbf{a}_i\}_{i=1}^m$, and then apply a covering argument. To start, introduce a Lipschitz-continuous counterpart for the

discontinuous indicator function [3, A.2]

$$\chi_E(\theta) := \begin{cases} 1, & |\theta| \geq \frac{\sqrt{1.01}}{1+\gamma}, \\ 100(1+\gamma)^2\theta^2 - 100, & \frac{1}{1+\gamma} \leq |\theta| < \frac{\sqrt{1.01}}{1+\gamma}, \\ 0, & |\theta| < \frac{1}{1+\gamma} \end{cases} \quad (79)$$

with Lipschitz constant $\mathcal{O}(1)$. Recall that $\mathcal{E}_i := \left\{ \left| \frac{\mathbf{a}_i^\top \mathbf{z}}{\mathbf{a}_i^\top \mathbf{x}} \right| \geq \frac{1}{1+\gamma} \right\}$, so it holds that $0 \leq \chi_E \left(\left| \frac{\mathbf{a}_i^\top \mathbf{z}}{\mathbf{a}_i^\top \mathbf{x}} \right| \right) \leq \mathbb{1}_{\mathcal{E}_i}$ for any $\mathbf{x} \in \mathbb{R}^n$ and $\mathbf{h} \in \mathbb{R}^n$, thus yielding

$$\frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h})^2 \mathbb{1}_{\mathcal{E}_i} \geq \frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h})^2 \chi_E \left(\left| \frac{\mathbf{a}_i^\top \mathbf{z}}{\mathbf{a}_i^\top \mathbf{x}} \right| \right) = \frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h})^2 \chi_E \left(\left| 1 + \frac{\mathbf{a}_i^\top \mathbf{h}}{\mathbf{a}_i^\top \mathbf{x}} \right| \right). \quad (80)$$

By homogeneity and rotational invariance property of normal distributions, it suffices to prove the case where $\mathbf{x} = \mathbf{e}_1$ and $\|\mathbf{h}\|/\|\mathbf{x}\| = \|\mathbf{h}\| \leq \rho$. According to (80), lower bounding the first term in (31) can be achieved by lower bounding $\sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h})^2 \chi_E \left(\left| 1 + \frac{\mathbf{a}_i^\top \mathbf{h}}{\mathbf{a}_i^\top \mathbf{x}} \right| \right)$ instead. To that end, let us find the mean of $(\mathbf{a}_i^\top \mathbf{h})^2 \chi_E \left(\left| 1 + \frac{\mathbf{a}_i^\top \mathbf{h}}{\mathbf{a}_i^\top \mathbf{x}} \right| \right)$. Note that $(\mathbf{a}_i^\top \mathbf{h})^2$ and $\chi_E \left(\left| 1 + \frac{\mathbf{a}_i^\top \mathbf{h}}{\mathbf{a}_i^\top \mathbf{x}} \right| \right)$ are dependent. Introduce an orthonormal matrix \mathbf{U}_h that contains $\mathbf{h}^\top/\|\mathbf{h}\|$ as its first row, i.e.,

$$\mathbf{U}_h := \begin{bmatrix} \mathbf{h}^\top/\|\mathbf{h}\| \\ \tilde{\mathbf{U}}_h \end{bmatrix} \quad (81)$$

for some orthogonal matrix $\tilde{\mathbf{U}}_h \in \mathbb{R}^{(n-1) \times n}$ such that \mathbf{U}_h is orthonormal. Moreover, define $\tilde{\mathbf{h}} := \mathbf{U}_h \mathbf{h}$, and $\tilde{\mathbf{a}}_i := \mathbf{U}_h \mathbf{a}_i$; and let $\tilde{a}_{i,1}$ and $\tilde{\mathbf{a}}_{i,\setminus 1}$ denote the first entry and the remaining entries in vector $\tilde{\mathbf{a}}_i$; likewise for vector $\tilde{\mathbf{h}}$. Then, for any \mathbf{h} such that $\|\mathbf{h}\| \leq \rho$, the next holds

$$\begin{aligned} \mathbb{E} \left[(\mathbf{a}_i^\top \mathbf{h})^2 \chi_E \left(\left| 1 + \frac{\mathbf{a}_i^\top \mathbf{h}}{\mathbf{a}_i^\top \mathbf{x}} \right| \right) \right] &= \mathbb{E} \left[(\tilde{a}_{i,1} \tilde{h}_1)^2 \chi_E \left(\left| 1 + \frac{\mathbf{a}_i^\top \mathbf{h}}{\mathbf{a}_i^\top \mathbf{x}} \right| \right) \right] + \mathbb{E} \left[(\tilde{\mathbf{a}}_{i,\setminus 1}^\top \tilde{\mathbf{h}}_{\setminus 1})^2 \chi_E \left(\left| 1 + \frac{\mathbf{a}_i^\top \mathbf{h}}{\mathbf{a}_i^\top \mathbf{x}} \right| \right) \right] \\ &= \tilde{h}_1^2 \mathbb{E} \left[\tilde{a}_{i,1}^2 \chi_E \left(\left| 1 + \frac{\mathbf{a}_i^\top \mathbf{h}}{\mathbf{a}_i^\top \mathbf{x}} \right| \right) \right] + \mathbb{E} \left[(\tilde{\mathbf{a}}_{i,\setminus 1}^\top \tilde{\mathbf{h}}_{\setminus 1})^2 \right] \mathbb{E} \left[\chi_E \left(\left| 1 + \frac{\mathbf{a}_i^\top \mathbf{h}}{\mathbf{a}_i^\top \mathbf{x}} \right| \right) \right] \\ &= \tilde{h}_1^2 \mathbb{E} \left[\tilde{a}_{i,1}^2 \chi_E \left(\left| 1 + \frac{\mathbf{a}_i^\top \mathbf{h}}{\mathbf{a}_i^\top \mathbf{x}} \right| \right) \right] + \|\tilde{\mathbf{h}}_{\setminus 1}\|^2 \mathbb{E} \left[\chi_E \left(\left| 1 + \frac{\mathbf{a}_i^\top \mathbf{h}}{\mathbf{a}_i^\top \mathbf{x}} \right| \right) \right] \\ &\geq (\tilde{h}_1^2 + \|\tilde{\mathbf{h}}_{\setminus 1}\|^2) \min \left\{ \mathbb{E} \left[\tilde{a}_{i,1}^2 \chi_E \left(\left| 1 + h_1 + \frac{\mathbf{a}_{i,\setminus 1}^\top \mathbf{h}_{\setminus 1}}{a_{i,1}} \right| \right) \right], \right. \\ &\quad \left. \mathbb{E} \left[\chi_E \left(\left| 1 + h_1 + \frac{\mathbf{a}_{i,\setminus 1}^\top \mathbf{h}_{\setminus 1}}{a_{i,1}} \right| \right) \right] \right\} \\ &\geq \|\mathbf{h}\|^2 \min \left\{ \mathbb{E} \left[\tilde{a}_{i,1}^2 \chi_E \left(\left| 1 - \rho + \frac{a_{i,2}}{a_{i,1}} \rho \right| \right) \right], \mathbb{E} \left[\chi_E \left(1 - \rho + \frac{a_{i,2}}{a_{i,1}} \rho \right) \right] \right\} \\ &= (1 - \zeta_1) \|\mathbf{h}\|^2 \end{aligned} \quad (82)$$

where the second equality follows from the independence between $\tilde{\mathbf{a}}_{i,\setminus 1}^\top \tilde{\mathbf{h}}_{\setminus 1}$ and $\mathbf{a}_i^\top \mathbf{h}$, the second inequality holds for $\rho \leq 1/10$ and $\gamma > 1/2$, and the last equality comes from the definition of ζ_1 in (74). Notice that $\varrho := (\mathbf{a}_i^\top \mathbf{h})^2 \chi_E \left(\left| 1 + \frac{\mathbf{a}_i^\top \mathbf{h}}{\mathbf{a}_i^\top \mathbf{x}} \right| \right) \leq (\mathbf{a}_i^\top \mathbf{h})^2 \stackrel{d}{=} \|\mathbf{h}\|^2 a_{i,1}^2$ is a subexponential variable, and thus its subexponential norm $\|\varrho\|_{\psi_1} := \sup_{p \geq 1} [\mathbb{E}(|\varrho|^p)]^{1/p}$ is finite.

Direct application of the Bernstein-type inequality [9, Proposition 5.16] confirms that for any $\epsilon > 0$, the following

$$\begin{aligned} \frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h})^2 \chi_E \left(\left| 1 + \frac{\mathbf{a}_i^\top \mathbf{h}}{\mathbf{a}_i^\top \mathbf{x}} \right| \right) &\geq \mathbb{E} \left[(\mathbf{a}_i^\top \mathbf{h})^2 \chi_E \left(\left| 1 + \frac{\mathbf{a}_i^\top \mathbf{h}}{\mathbf{a}_i^\top \mathbf{x}} \right| \right) \right] - \epsilon \|\mathbf{h}\|^2 \\ &\geq (1 - \zeta_1 - \epsilon) \|\mathbf{h}\|^2 \end{aligned} \quad (83)$$

holds with probability at least $1 - e^{-c_5 m \epsilon^2}$ for some numerical constant $c_5 > 0$ provided that $\epsilon \leq \|\varrho\|_{\psi_1}$ by assumption.

To obtain uniform control over all vectors \mathbf{z} and \mathbf{x} such that $\|\mathbf{z} - \mathbf{x}\| \leq \rho$, the net covering argument is applied [9, Definition 5.1]. Let \mathcal{S}_ϵ be an ϵ -net of the unit sphere, \mathcal{L}_ϵ be an ϵ -net of $[0, \rho]$, and define

$$\mathcal{N}_\epsilon := \{(\mathbf{z}, \mathbf{h}, t) : (\mathbf{z}_0, \mathbf{h}_0, t_0) \in \mathcal{S}_\epsilon \times \mathcal{S}_\epsilon \times \mathcal{L}_\epsilon\}. \quad (84)$$

Since the cardinality $|\mathcal{S}_\epsilon| \leq (1 + 2/\epsilon)^n$ [9, Lemma 5.2], then

$$|\mathcal{N}_\epsilon| \leq (1 + 2/\epsilon)^{2n} \rho/\epsilon \leq (1 + 2/\epsilon)^{2n+1} \quad (85)$$

due to the fact that $\rho/\epsilon < 2/\epsilon < 1 + 2/\epsilon$ for $0 \leq \rho < 1$.

Consider now any $(\mathbf{z}, \mathbf{h}, t)$ obeying $\|\mathbf{h}\| = t \leq \rho$. There exists a pair $(\mathbf{z}_0, \mathbf{h}_0, t_0) \in \mathcal{N}_\epsilon$ such that $\|\mathbf{z} - \mathbf{z}_0\|$, $\|\mathbf{h} - \mathbf{h}_0\|$, and $|t - t_0|$ are each at most ϵ . Taking the union bound yields

$$\begin{aligned} \frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h}_0)^2 \chi_E \left(\left| 1 + \frac{\mathbf{a}_i^\top \mathbf{h}_0}{\mathbf{a}_i^\top \mathbf{x}} \right| \right) &\geq \frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h}_0)^2 \chi_E \left(\left| 1 - t_0 + \frac{a_{i,2}}{a_{i,1}} t_0 \right| \right) \\ &\geq (1 - \zeta_1 - \epsilon) \|\mathbf{h}_0\|^2, \quad \forall (\mathbf{z}_0, \mathbf{h}_0, t_0) \in \mathcal{N}_\epsilon \end{aligned} \quad (86)$$

with probability at least $1 - (1 + 2/\epsilon)^{2n+1} e^{-c_5 \epsilon^2 m} \geq 1 - e^{-c_0 m}$, which follows by choosing m such that $m \geq (c_6 \cdot \epsilon^{-2} \log \epsilon^{-1}) n$ for some constant $c_6 > 0$.

Recall that $\chi_E(\tau)$ is Lipschitz continuous, thus

$$\begin{aligned} &\left| \frac{1}{m} \sum_{i=1}^m \left\{ (\mathbf{a}_i^\top \mathbf{h})^2 \chi_E \left(\left| 1 + \frac{\mathbf{a}_i^\top \mathbf{h}}{\mathbf{a}_i^\top \mathbf{x}} \right| \right) - (\mathbf{a}_i^\top \mathbf{h}_0)^2 \chi_E \left(\left| 1 + \frac{\mathbf{a}_i^\top \mathbf{h}_0}{\mathbf{a}_i^\top \mathbf{x}} \right| \right) \right\} \right| \\ &\lesssim \frac{1}{m} \sum_{i=1}^m \left| (\mathbf{a}_i^\top \mathbf{h})^2 - (\mathbf{a}_i^\top \mathbf{h}_0)^2 \right| \\ &= \frac{1}{m} \sum_{i=1}^m |\mathbf{a}_i^\top (\mathbf{h} \mathbf{h}^\top - \mathbf{h}_0 \mathbf{h}_0^\top) \mathbf{a}_i| \\ &\lesssim c_7 \sum_{i=1}^m |\mathbf{h} \mathbf{h}^\top - \mathbf{h}_0 \mathbf{h}_0^\top| \\ &\leq 2.5 c_7 \|\mathbf{h} - \mathbf{h}_0\| \|\mathbf{h}\| \\ &\leq 2.5 c_7 \rho \epsilon \end{aligned} \quad (87)$$

for some numerical constant c_7 and provided that $\epsilon < 1/2$ and $m \geq (c_6 \cdot \epsilon^{-2} \log \epsilon^{-1}) n$, where the first inequality arises from the Lipschitz property of $\chi_E(\tau)$, the second uses the results in Lemma 1 in [3], and the third from Lemma 2 in [3].

Putting all results together confirms that with probability exceeding $1 - 2e^{-c_0 m}$, the following holds

$$\frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h})^2 \chi_E \left(\left| 1 + \frac{\mathbf{a}_i^\top \mathbf{h}}{\mathbf{a}_i^\top \mathbf{x}} \right| \right) \geq [1 - \zeta_1 - (1 + 2.5c_7\rho)\epsilon] \|\mathbf{h}\|^2 \quad (88)$$

for all vectors $\|\mathbf{h}\| / \|\mathbf{x}\| \leq \rho$, concluding the proof.

A.5 Proof of Lemma 6

Similar to the proof in Section A.4, it is convenient to work with the following auxiliary function instead of the discontinuous indicator function

$$\chi_D(\theta) := \begin{cases} 1, & |\theta| \geq \frac{2+\gamma}{1+\gamma} \\ -100 \left(\frac{1+\gamma}{2+\gamma} \right)^2 \theta^2 + 100, & \sqrt{0.99} \cdot \frac{2+\gamma}{1+\gamma} \leq |\theta| < \frac{2+\gamma}{1+\gamma} \\ 0, & |\theta| < \sqrt{0.99} \cdot \frac{2+\gamma}{1+\gamma} \end{cases} \quad (89)$$

which is Lipschitz continuous in θ with Lipschitz constant $\mathcal{O}(1)$. For $\mathcal{D}_i = \left\{ \left| \frac{\mathbf{a}_i^\top \mathbf{h}}{\mathbf{a}_i^\top \mathbf{x}} \right| \geq \frac{2+\gamma}{1+\gamma} \right\}$, it holds that $0 \leq \mathbb{1}_{\mathcal{D}_i} \leq \chi_D \left(\left| \frac{\mathbf{a}_i^\top \mathbf{h}}{\mathbf{a}_i^\top \mathbf{x}} \right| \right)$ for any $\mathbf{x} \in \mathbb{R}^n$ and $\mathbf{h} \in \mathbb{R}^n$. Assume without loss of generality $\mathbf{x} = \mathbf{e}_1$. Then for $\gamma > 0$ and $\rho \leq 1/10$, it holds that

$$\begin{aligned} \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\left\{ \left| \frac{\mathbf{a}_i^\top \mathbf{h}}{\mathbf{a}_i^\top \mathbf{x}} \right| \geq \frac{2+\gamma}{1+\gamma} \right\}} &\leq \frac{1}{m} \sum_{i=1}^m \chi_D \left(\left| \frac{\mathbf{a}_i^\top \mathbf{h}}{\mathbf{a}_i^\top \mathbf{x}} \right| \right) = \frac{1}{m} \sum_{i=1}^m \chi_D \left(\left| \frac{\mathbf{a}_i^\top \mathbf{h}}{a_{i,1}} \right| \right) \\ &= \frac{1}{m} \sum_{i=1}^m \chi_D \left(\left| h_1 + \frac{\mathbf{a}_{i,\setminus 1}^\top \mathbf{h}_{\setminus 1}}{a_{i,1}} \right| \right) \\ &= \frac{1}{m} \sum_{i=1}^m \chi_D \left(\left| h_1 + \frac{a_{i,2}}{a_{i,1}} \|\mathbf{h}_{\setminus 1}\| \right| \right) \\ &\stackrel{(i)}{\leq} \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\left\{ \left| h_1 + \frac{a_{i,2}}{a_{i,1}} \|\mathbf{h}_{\setminus 1}\| \right| \geq \sqrt{0.99} \cdot \frac{2+\gamma}{1+\gamma} \right\}} \end{aligned} \quad (90)$$

where the last inequality arises from the definition of χ_D . Noting that $a_{i,2}/a_{i,1}$ obeys the standard Cauchy distribution, i.e., $a_{i,2}/a_{i,1} \sim \text{Cauchy}(0, 1)$ [14], and particularly, transformation properties of Cauchy distributions assert that $h_1 + \frac{a_{i,2}}{a_{i,1}} \|\mathbf{h}_{\setminus 1}\| \sim \text{Cauchy}(h_1, \|\mathbf{h}_{\setminus 1}\|)$ [15]. Recall that the cdf of a Cauchy distributed random variable $w \sim \text{Cauchy}(\mu_0, \alpha)$ is given by [14]

$$F(w; \mu_0, \alpha) = \frac{1}{\pi} \arctan \left(\frac{w - \mu_0}{\alpha} \right) + \frac{1}{2}. \quad (91)$$

It is easy to check that when $\|\mathbf{h}_{\setminus 1}\| = 0$, the indicator function $\mathbb{1}_{\mathcal{D}_i} = 0$ due to $|h_1| \leq \rho < \sqrt{0.99}(2 + \gamma)/(1 + \gamma)$. Consider only $\|\mathbf{h}_{\setminus 1}\| \neq 0$ next. Define for notational brevity $w := a_{i,2}/a_{i,1}$,

$\alpha := \|\mathbf{h}_{\setminus 1}\|$, as well as $\mu_0 := h_1/\alpha$ and $w_0 := \sqrt{0.99} \frac{2+\gamma}{\alpha(1+\gamma)}$ to yield

$$\begin{aligned}
\mathbb{E} [\mathbb{1}_{\{\mu_0 + w \geq w_0\}}] &= 1 - [F(w_0; \mu_0, 1) - F(-w_0; \mu_0, 1)] \\
&= 1 - \frac{1}{\pi} [\arctan(w_0 - \mu_0) - \arctan(-w_0 - \mu_0)] \\
&\stackrel{(i)}{=} \frac{1}{\pi} \arctan\left(\frac{2w_0}{w_0^2 - \mu_0^2 - 1}\right) \\
&\stackrel{(ii)}{\leq} \frac{1}{\pi} \cdot \frac{2w_0}{w_0^2 - \mu_0^2 - 1} \\
&\stackrel{(iii)}{\leq} \frac{1}{\pi} \cdot \frac{2\sqrt{0.99}\rho(2+\gamma)/(1+\gamma)}{0.99(2+\gamma)^2/(1+\gamma)^2 - \rho^2} \\
&\leq 0.0646
\end{aligned} \tag{92}$$

for all $\gamma > 0$ and $\rho \leq 1/10$. In deriving (i), the property $\arctan(u) + \arctan(v) = \arctan\left(\frac{u+v}{1-uv}\right) \pmod{\pi}$ for any $uv \neq 1$. Concerning (ii), the inequality $\arctan(x) \leq x$ for $x \geq 0$ is employed. Plugging given parameter values and using $\|\mathbf{h}_{\setminus 1}\| \leq \|\mathbf{h}\| \leq \rho$ confirms (iii). Apparently, $\mathbb{1}_{\{\mu_0 + w \geq w_0\}}$ is bounded; and it is known that all bounded random variables are subexponential. Thus, upon applying the Bernstein-type inequality [9, Corollary 5.17], the next holds with probability at least $1 - e^{-c_5 m \epsilon^2}$ for some numerical constant $c_5 > 0$ and any sufficiently small $\epsilon > 0$

$$\begin{aligned}
\frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\left\{\frac{|\mathbf{a}_i^T \mathbf{h}|}{|\mathbf{a}_i^T \mathbf{w}|} \geq \frac{2+\gamma}{1+\gamma}\right\}} &\leq \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\left\{|h_1 + \frac{a_{i,2}}{a_{i,1}} \|\mathbf{h}_{\setminus 1}\| \geq \sqrt{0.99} \frac{2+\gamma}{1+\gamma}\right\}} \leq (1+\epsilon) \mathbb{E} \left[\mathbb{1}_{\left\{|h_1 + \frac{a_{i,2}}{a_{i,1}} \|\mathbf{h}_{\setminus 1}\| \geq \sqrt{0.99} \frac{2+\gamma}{1+\gamma}\right\}} \right] \\
&\leq \frac{1+\epsilon}{\pi} \cdot \frac{2\sqrt{0.99}\rho(2+\gamma)/(1+\gamma)}{0.99(2+\gamma)^2/(1+\gamma)^2 - \rho^2}.
\end{aligned} \tag{93}$$

On the other hand, one can easily establish that the following holds true for all \mathbf{h}

$$\mathbb{E} \left[(\mathbf{a}_i^T \mathbf{h})^4 \right] = \mathbb{E} [a_{i,1}^4] \|\mathbf{h}\|^4 = 3 \|\mathbf{h}\|^4 \tag{94}$$

which has also been established in Lemma 1 [3] and Lemma 6.1 [16]. Further recalling our working assumption $\|\mathbf{a}_i\| \leq 2.3n$, then random variables $(\mathbf{a}_i^T \mathbf{h})^4$ are bounded, and thus they are subexponential [9]. Appealing again to the Bernstein-type inequality for subexponential random variables and provided that $m/n > c_6 \cdot \epsilon^{-2} \log \epsilon^{-1}$ for some numerical constant $c_6 > 0$, then

$$\frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^T \mathbf{h})^4 \leq 3(1+\epsilon) \|\mathbf{h}\|^4 \tag{95}$$

which holds with probability exceeding $1 - e^{-c_5 m \epsilon^2}$ for some universal constant $c_5 > 0$ and any sufficiently small $\epsilon > 0$.

Collecting together results in (93) and (95) and leveraging the Cauchy-Schwartz inequality, one

establishes that for any $\rho \leq 1/10$ and $\gamma > 0$, the following

$$\begin{aligned}
\frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h})^2 \mathbb{1}_{\mathcal{D}_i} &\leq \sqrt{\frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{h})^4} \sqrt{\frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\left\{ \frac{|\mathbf{a}_i^\top \mathbf{h}|}{|\mathbf{a}_i^\top \mathbf{x}|} \geq \frac{2+\gamma}{1+\gamma} \right\}}} \\
&\leq \sqrt{3(1+\epsilon) \|\mathbf{h}\|^4} \sqrt{\frac{1+\epsilon}{\pi} \cdot \frac{2\sqrt{0.99}\rho(2+\gamma)/(1+\gamma)}{0.99(2+\gamma)^2/(1+\gamma)^2 - \rho^2}} \\
&\triangleq (\zeta'_2 + \epsilon') \|\mathbf{h}\|^2
\end{aligned} \tag{96}$$

where $\zeta'_2 := 1.3785\sqrt{\rho\tau/(0.99\tau^2 - \rho^2)}$ with $\tau := (2+\gamma)/(1+\gamma)$, which holds with probability at least $1 - 2e^{-c_0 m}$. The latter arises if choosing $c_0 \leq c_5 \epsilon^2$ in $1 - 2e^{-c_5 m \epsilon^2}$, which can be accomplished by taking m/n sufficiently large.

References

- [1] R. Balan, P. Casazza, and D. Edidin, “On signal reconstruction without phase,” *Appl. Comput. Harmon. Anal.*, vol. 20, no. 3, pp. 345–356, May 2006.
- [2] E. J. Candès, X. Li, and M. Soltanolkotabi, “Phase retrieval via Wirtinger flow: Theory and algorithms,” *IEEE Trans. Inf. Theory*, vol. 61, no. 4, pp. 1985–2007, Apr. 2015.
- [3] Y. Chen and E. J. Candès, “Solving random quadratic systems of equations is nearly as easy as solving linear systems,” *Comm. Pure Appl. Math.*, 2016 (to appear).
- [4] G. Wang and G. B. Giannakis, “Solving random systems of quadratic equations via truncated generalized gradient flow,” in *Adv. Neural Inf. Process. Syst.*, Barcelona, Spain, 2016.
- [5] P. Netrapalli, P. Jain, and S. Sanghavi, “Phase retrieval using alternating minimization,” *IEEE Trans. Signal Process.*, vol. 63, no. 18, pp. 4814–4826, Sept. 2015.
- [6] P. Chen, A. Fannjiang, and G.-R. Liu, “Phase retrieval with one or two diffraction patterns by alternating projections of the null vector,” *arXiv:1510.07379v2*, 2015.
- [7] M. Soltanolkotabi, “Algorithms and theory for clustering and nonconvex quadratic programming,” Ph.D. dissertation, Stanford University, 2014.
- [8] G. Wang, G. B. Giannakis, and Y. C. Eldar, “Solving systems of random quadratic equations via truncated amplitude flow,” *arXiv:1605.08285*, 2016.
- [9] R. Vershynin, “Introduction to the non-asymptotic analysis of random matrices,” *arXiv:1011.3027*, 2010.
- [10] E. J. Candès, T. Strohmer, and V. Voroninski, “PhaseLift: Exact and stable signal recovery from magnitude measurements via convex programming,” *Appl. Comput. Harmon. Anal.*, vol. 66, no. 8, pp. 1241–1274, Nov. 2013.
- [11] S. Cambanis, S. Huang, and G. Simons, “On the theory of elliptically contoured distributions,” *J. Multivar. Anal.*, vol. 11, no. 3, pp. 368–385, Sep. 1981.
- [12] B. Laurent and P. Massart, “Adaptive estimation of a quadratic functional by model selection,” *Ann. Stat.*, vol. 28, no. 5, pp. 1302–1338, 2000.
- [13] S.-H. Chang, P. C. Cosman, and L. B. Milstein, “Chernoff-type bounds for the Gaussian error function,” *IEEE Trans. Commun.*, vol. 59, no. 11, pp. 2939–2944, July 2011.
- [14] T. S. Ferguson, “A representation of the symmetric bivariate Cauchy distribution,” *Ann. Math. Stat.*, vol. 33, no. 4, pp. 1256–1266, 1962.
- [15] H. Y. Lee, G. J. Parka, and H. M. Kim, “A clarification of the Cauchy distribution,” *Commun. Stat. Appl. Methods*, vol. 21, no. 2, pp. 183–191, Mar. 2014.
- [16] J. Sun, Q. Qu, and J. Wright, “A geometric analysis of phase retrieval,” *arXiv:1602.06664*, 2016.