### FAST COMMUNICATION

# OPTIMAL SAMPLE COMPLEXITY OF SUBGRADIENT DESCENT FOR AMPLITUDE FLOW VIA NON-LIPSCHITZ MATRIX CONCENTRATION\*

PAUL HAND<sup>†</sup>, OSCAR LEONG<sup>‡</sup>, AND VLADISLAV VORONINSKI<sup>§</sup>

Abstract. We consider the problem of recovering a real-valued n-dimensional signal from m phaseless, linear measurements and analyze the amplitude-based non-smooth least squares objective. We establish local convergence of subgradient descent with optimal sample complexity based on the uniform concentration of a random, discontinuous matrix-valued operator arising from the objective's gradient dynamics. While common techniques to establish uniform concentration of random functions exploit Lipschitz continuity, we prove that the discontinuous matrix-valued operator satisfies a uniform matrix concentration inequality when the measurement vectors are Gaussian as soon as  $m = \Omega(n)$  with high probability. We then show that satisfaction of this inequality is sufficient for subgradient descent with proper initialization to converge linearly to the true solution up to the global sign ambiguity. As a consequence, this guarantees local convergence for Gaussian measurements at optimal sample complexity. The concentration methods in the present work have previously been used to establish recovery guarantees for a variety of inverse problems under generative neural network priors. This paper demonstrates the applicability of these techniques to more traditional inverse problems and serves as a pedagogical introduction to those results.

 $\textbf{Keywords.} \ \ \textbf{Phase retrieval; Subgradient descent; Concentration inequality; Non-convex optimization.}$ 

AMS subject classifications. 90C26; 94A12; 94A15.

#### 1. Introduction

Consider the problem of recovering a signal  $x_* \in \mathbb{R}^n$  from m phaseless measurements of the form

$$y := |Ax_*| + \eta$$

where  $A \in \mathbb{R}^{m \times n}$  is a measurement matrix,  $|\cdot|$  acts entrywise, and  $\eta \in \mathbb{R}^m$  denotes noise. This problem is known as phase retrieval as, in practice, the phase of the signal is lost in the forward measurement process due to the underlying physics of the measurement system. We consider the case when the entries of A are i.i.d. Gaussian, which we will refer to as the generic measurement regime. In this work, we aim to recover  $x_*$  by solving the following non-smooth least squares problem

$$\min_{x \in \mathbb{R}^n} f(x) := \frac{1}{2} \||Ax| - y\|^2. \tag{1.1}$$

This objective function is known as Amplitude Flow. For generic measurements, previous works have shown that with proper initialization, gradient descent both with [12] and without [29] truncated gradients can recover the signal with the optimal sample complexity of  $m = \Omega(n)$ .

<sup>\*</sup>Received: November 01, 2020; Accepted (in revised form): May 18, 2021. Communicated by Lexing Ying.

<sup>&</sup>lt;sup>†</sup>Department of Mathematics and College of Computer and Information Science, Northeastern University, Boston, MA 02115, USA (p.hand@northeastern.edu).

<sup>&</sup>lt;sup>‡</sup>Department of Computational and Applied Mathematics, Rice University, Houston, TX 77005-1827, USA (oscar.f.leong@rice.edu).

<sup>§</sup>Helm.ai, Menlo Park, CA 94025, USA (vlad@helm.ai).

Existing proof techniques of convergence guarantees for (sub)gradient descent of (1.1) follow a two-step process: (1) establish that spectral initialization or some variant thereof guarantees an initializer with relative error bounded by a small absolute constant and then (2) show that the objective satisfies a property akin to convexity near the minimizer to guarantee convergence. This latter property is called the *local regularity* condition  $RC(\mu,\lambda,\varepsilon)^1$ . Showing that this condition holds is crucial in establishing local convergence for Amplitude Flow [12,29] and its variants [23].

This proof technique is not unique to Amplitude Flow as it was initially introduced to guarantee convergence for the intensity-based formulation Wirtinger Flow, which aims to solve

$$\min_{x \in \mathbb{R}^n} \frac{1}{2} \||Ax|^2 - y^2\|^2. \tag{1.2}$$

In the original work [3], the aforementioned two-step procedure established exact recovery with sample complexity  $m = \Omega(n \log n)$ . A follow-up variant using truncation [5] improved the sample complexity to  $m = \Omega(n)$  and also employed the  $\mathsf{RC}(\mu, \lambda, \varepsilon)$  to show convergence post-initialization. Recently, [28] established the sufficiency of a deterministic condition for local convergence when solving (1.2) in the lifted domain and a relationship between the condition's accuracy and convergence rate of gradient descent was shown. This deterministic condition is a uniform matrix concentration inequality that is proven to hold for generic measurements when  $m = \Omega(n \log n)$ .

Other works include the global landscape analysis in [25] which showed that (1.2) exhibits benign geometry given a sufficient number of measurements  $m = \Omega(n\log^3 n)$  by carefully analyzing the gradient and Hessian in partitioned regions of space. In [11], the authors considered the robust  $\ell_1$  loss with intensity-based measurements and established local convergence of a prox-linear algorithm using composite optimization theory. Convex approaches based on lifting [2,4,10,15] utilize dual certificates to assert correctness of the minimizers of semidefinite programs. Linear programming approaches have also been studied, whose proof techniques range from using tools in statistical learning theory [1] and geometric probability theory [14], along with elementary approaches using standard concentration estimates of the singular values of random matrices [19]. For a more comprehensive overview of prior work for phase retrieval, we refer the reader to [13].

In this paper, we present a proof technique for solving (1.1) with  $m = \Omega(n)$  based on uniform concentration of random matrix-valued functions that are discontinuous in space. Consider a subgradient descent algorithm with iterates  $\{x_t\}_{t\geqslant 0}$  of the form  $x_{t+1} = x_t - \alpha v_{x_t,x_*}$  where  $\alpha > 0$ ,  $v_{x_t,x_*} \in \partial f(x_t)$ , and  $\partial f(x)$  is the Clarke subdifferential at x (defined in Section 3). Let  $\operatorname{dist}(x,x_*) := \min(\|x-x_*\|, \|x+x_*\|)$ . We first state our main local convergence result in the Gaussian measurement regime.

THEOREM 1.1. There exists positive absolute constants C,  $c_1$ ,  $c_2$ ,  $\rho_1$ , and  $\rho_2$  such that the following holds. Suppose  $A \in \mathbb{R}^{m \times n}$  has i.i.d.  $\mathcal{N}(0,1/m)$  entries and the noise is bounded  $\|\eta\| \leq \rho_1 \|x_*\|$ . Assume the initial iterate  $x_0$  satisfies  $\operatorname{dist}(x_0, x_*) \leq \rho_2 \|x_*\|$  and the step size satisfies  $0 < \alpha \leq 1$ . If  $m \geq Cn$ , then with probability at least  $1 - 3\exp(-c_1 m) - m\exp(-c_2 n)$ , we have that for all  $t \geq 1$ ,

$$\operatorname{dist}(x_t, x_*) \leq \left(1 - \frac{\alpha}{2}\right)^t \operatorname{dist}(x_0, x_*) + 4\|\eta\|.$$

<sup>&</sup>lt;sup>1</sup>A function  $\mathcal{L}$  satisfies  $\mathsf{RC}(\mu, \lambda, \varepsilon)$  at a stationary point y if for all  $x \in \mathbb{R}^n$  such that  $||x - y|| \le \varepsilon ||y||$ ,  $\langle \nabla \mathcal{L}(x), x - y \rangle \ge \frac{\mu}{2} ||x - y||^2 + \frac{\lambda}{2} ||\nabla \mathcal{L}(x)||^2$ .

This result asserts local convergence up to the noise level with optimal sample complexity. In the theorem, note that we require an initializer with relative error less than a sufficiently small constant. There are several schemes to achieve this with  $m = \Omega(n)$  Gaussian measurements, even in the presence of noise [5,11,12,29]. While convergence of subgradient descent without truncation for the Amplitude Flow objective is known [29], the method of proof we present here is novel. In particular, we show that Theorem 1.1 is a consequence of the following two results: (1) a uniform matrix concentration inequality is sufficient to guarantee local convergence with proper initialization and (2) Gaussian matrices satisfy this inequality with high probability when  $m = \Omega(n)$ .

We now detail the high level intuition behind the proof ideas and techniques. Let  $\operatorname{sgn}(z) := z/|z|$  for  $z \neq 0$  and  $\operatorname{sgn}(0) = 0$  act entrywise. For ease of exposition, suppose there is no noise  $\eta = 0$ . The discontinuous, spatially-varying measurement operator  $A_x := \operatorname{diag}(\operatorname{sgn}(Ax))A$  plays a critical role in analyzing subgradient descent as this operator governs the gradient dynamics of f. Specifically, the gradient almost everywhere is given by  $\nabla f(x) = A_x^{\mathrm{T}}(A_x x - A_{x_*} x_*)$ . As will be shown in the next section, the gradient in expectation obeys a property equivalent to the  $\operatorname{RC}(\mu, \lambda, \varepsilon)$  in neighborhoods of the global minimizers  $\pm x_*$ . Hence if we establish concentration of the quantity  $A_x^{\mathrm{T}}A_y$  to its expectation  $\mathbb{E}[A_x^{\mathrm{T}}A_y]$  uniformly in x,y, then this property will also be satisfied by the gradient. This will be shown to guarantee local convergence up to the global sign ambiguity with high probability.

The uniform concentration result we establish is the following: when A has i.i.d.  $\mathcal{N}(0,1/m)$  entries, then for any parameter  $0 < \varepsilon < 1$ , when  $m = \Omega(n)$  we have that with high probability

$$||A_x^{\mathrm{T}} A_y - \Phi_{x,y}|| \leqslant \varepsilon \ \forall x, y \in \mathbb{R}^n$$
(1.3)

where  $\Phi_{x,y} := \mathbb{E}[A_x^{\mathrm{T}}A_y]$  has an analytic expression. As this result holds uniformly in x,y, we have that for any  $x,x_* \in \mathbb{R}^n$ , the gradient  $\nabla f(x) \approx \Phi_{x,x}x - \Phi_{x,x_*}x_*$ . The difficulty of establishing (1.3) uniformly in x,y is due to the fact that  $A_x^{\mathrm{T}}A_y$  is a non-Lipschitz matrix-valued operator. Standard approaches to control these types of quantities exploit Lipschitz continuity by first (1) establishing concentration for fixed x,y, then (2) establishing concentration over all points in a net of the sphere by using a union bound, and finally (3) appealing to Lipschitz continuity to get concentration uniformly in x,y. However, in this case, (3) is not possible as  $A_x^{\mathrm{T}}A_y$  is discontinuous with respect to x,y.

Fortunately, this issue can be solved by concentrating Lipschitz continuous approximations of  $A_x^{\mathrm{T}}A_y$  with respect to x,y. In particular, one can create continuous matrix-valued functions that are upper and lower bounds of  $A_x^{\mathrm{T}}A_y$  with respect to the semidefinite ordering. Then, concentration of these continuous approximations can be established by appealing to the standard arguments outlined above. This, in turn, will be shown to establish concentration of the non-Lipschitz quantity of interest by a squeezing argument. Moreover, using novel tools developed in [9], a more efficient set of coverings of the sphere can be exploited to achieve sample complexity linear in the ambient dimension n. Intuitively, this is achieved by constructing a net of the sphere that does not penalize all directions equally, but instead exploits directions for which the function of interest does not deviate much and penalizes those for which the function exhibits larger change. This intuition is made precise in the proof of Proposition 2.1.

Discussion. While the convergence results presented here have been shown in [29], the contribution of this work lies in the novelty of the analysis. In particular, this work is an illustrative example of using the concentration of non-Lipschitz functions to establish favorable properties of first-order algorithms to solve inverse problems. The

concentration methods of this paper have been used to establish recovery in compressive sensing [20,22], phase retrieval [17,18], and other problems [8,16,21,24] under image priors given by generative neural networks. For example, a similar uniform matrix concentration inequality to the one used in this paper was introduced by the present authors in [17,18] to establish recovery in compressive phase retrieval under a generative prior with information-theoretically optimal sample complexity. This paper demonstrates the applicability of these techniques to more traditional inverse problems and serves as a pedagogical introduction to those results.

### 2. Proof technique

We now establish the sufficiency of a deterministic condition for local convergence in the form of a uniform matrix concentration inequality and show that Gaussian matrices satisfy this condition with high probability when  $m = \Omega(n)$ . A similar condition, known as the Weight Distribution Condition, was first introduced in [20] in the context of compressive sensing under generative neural network priors. The matrix concentration inequality is stated as follows.

DEFINITION 2.1. Fix  $0 < \varepsilon < 1$ . We say that  $A \in \mathbb{R}^{m \times n}$  satisfies the **Measurement** Distribution Condition (MDC) with constant  $\varepsilon$  if

$$||A_x^{\mathrm{T}}A_y - \Phi_{x,y}|| \leq \varepsilon \ \forall \ x, y \in \mathbb{R}^n$$

where

$$\Phi_{x,y} := \begin{cases} \frac{\pi - 2\theta_{x,y}}{\pi} I_n + \frac{2\sin\theta_{x,y}}{\pi} M_{\hat{x} \leftrightarrow \hat{y}} & \text{if } x \neq 0, y \neq 0, \\ 0 & \text{otherwise.} \end{cases}$$
 (2.1)

Here  $\theta_{x,y} := \angle(x,y)$ ,  $\hat{x} := x/\|x\|$ ,  $\hat{y} := y/\|y\|$ ,  $I_n$  is the  $n \times n$  identity matrix, and  $M_{\hat{x} \leftrightarrow \hat{y}}^2$  is the matrix that sends  $\hat{x} \mapsto \hat{y}$ ,  $\hat{y} \mapsto \hat{x}$ , and  $z \mapsto 0$  for any  $z \in span(\{x,y\})^{\perp}$ .

Note that for points x,y with small angle, this condition requires  $A_x^{\rm T}A_y$  to act like an isometry. In the extreme case when  $x=y,\,\Phi_{x,y}$  is the identity. An elementary calculation gives  $\mathbb{E}[A_x^{\rm T}A_y]=\Phi_{x,y}$  for  $x,y\neq 0$  and  $A_{ij}\sim \mathcal{N}(0,1/m)$ .

The first result is that the MDC is sufficient to guarantee the following: a subgradient descent algorithm with proper initialization will converge to the true solution up to the global sign ambiguity. To establish this, we first show that the MDC is sufficient for the objective to satisfy the following regularity condition which states that, within a neighborhood of the true solution, all subgradients point towards the true solution. This result is proven in Section 3.1.

LEMMA 2.1. Fix  $0 < \varepsilon \le 0.001$ . Suppose  $A \in \mathbb{R}^{m \times n}$  satisfies the MDC with constant  $\varepsilon$ . Then for all  $x \in \mathbb{R}^n$  such that  $\operatorname{dist}(x, x_*) \le \varepsilon \|x_*\|$  and any  $v_{x, x_*} \in \partial f(x)$ , we have that  $\|v_{x, x_*} - (x \pm x_*)\| \le \frac{1}{2} \|x \pm x_*\| + 2\|\eta\|$ . Here  $x \pm x_* := x - x_*$  if  $\|x - x_*\| = \operatorname{dist}(x, x_*)$  and  $x + x_*$  otherwise.

Note that the conclusion of this lemma in the noiseless setting is in fact equivalent to

the  $RC(\mu, \lambda, \varepsilon)$  condition<sup>3</sup>. We now show that satisfaction of the MDC implies local convergence.

THEOREM 2.1 (Deterministic local convergence guarantee). Fix  $0 < \varepsilon \le 0.001$ . Suppose  $A \in \mathbb{R}^{m \times n}$  satisfies the MDC with constant  $\varepsilon$ ,  $\|\eta\| \le \frac{\varepsilon}{4} \|x_*\|$ , and  $0 < \alpha \le 1$ . If  $\operatorname{dist}(x_0, x_*) \le \varepsilon \|x_*\|$  then for all  $t \ge 1$ ,

$$\operatorname{dist}(x_t, x_*) \leq \left(1 - \frac{\alpha}{2}\right)^t \operatorname{dist}(x_0, x_*) + 4\|\eta\|.$$

*Proof.* Let  $\mathcal{B}(x_*,r) := \{x \in \mathbb{R}^n : ||x-x_*|| \le r\}$ . Suppose  $x_0 \in \mathcal{B}(x_*,\varepsilon||x_*||)$  as the proof for the case  $x_0 \in \mathcal{B}(-x_*,\varepsilon||x_*||)$  is identical. For  $t \ge 1$ , observe that for any  $v_{x_{t-1},x_*} \in \partial f(x_{t-1})$ , we have

$$||x_{t} - x_{*}|| = ||x_{t-1} - \alpha v_{x_{t-1}, x_{*}} + \alpha (x_{t-1} - x_{*}) - \alpha (x_{t-1} - x_{*}) - x_{*}||$$

$$\leq (1 - \alpha) ||x_{t-1} - x_{*}|| + \alpha ||v_{x_{t-1}, x_{*}} - (x_{t-1} - x_{*})||$$

$$\leq (1 - \alpha) ||x_{t-1} - x_{*}|| + \frac{\alpha}{2} ||x_{t-1} - x_{*}|| + 2\alpha ||\eta||$$

$$= \left(1 - \frac{\alpha}{2}\right) ||x_{t-1} - x_{*}|| + 2\alpha ||\eta||$$

$$(2.2)$$

where in the third line we used Lemma 2.1. We claim that the iterates must stay within a ball of the minimizer. Indeed, if  $x_{t-1} \in \mathcal{B}(x_*, \varepsilon || x_* ||)$ , we have that by Equation (2.2) and our bound on the size of the noise  $||\eta|| \leq \frac{\varepsilon}{4} ||x_*||$  that

$$\|x_t - x_*\| \leqslant \left(1 - \frac{\alpha}{2}\right) \|x_{t-1} - x_*\| + 2\alpha \|\eta\| \leqslant \left(1 - \frac{\alpha}{2}\right) \varepsilon \|x_*\| + \alpha \cdot \frac{\varepsilon}{2} \|x_*\| = \varepsilon \|x_*\|$$

so  $x_t \in \mathcal{B}(x_*, \varepsilon || x_*||)$ . Thus, we can invoke Lemma 2.1 and Equation (2.2) for each  $t \ge 1$ . Letting  $\tau := 1 - \frac{\alpha}{2}$ , starting at t = 1 and repeatedly applying (2.2), we attain

$$||x_t - x_*|| \le \tau^t ||x_0 - x_*|| + 2\alpha(\tau^t + \tau^{t-1} + \dots + 1)||\eta|| \le \tau^t ||x_0 - x_*|| + \frac{2\alpha}{1 - \tau} ||\eta||.$$

Plugging in the definition of  $\tau$  yields the desired inequality.

Finally, using recent tools developed in [9], we show that Gaussian matrices satisfy the MDC with high probability with  $m = \Omega(n)$  sample complexity.

PROPOSITION 2.1. Fix  $0 < \varepsilon < 1$ . Suppose  $A \in \mathbb{R}^{m \times n}$  has i.i.d.  $\mathcal{N}(0,1/m)$  entries. If  $m \ge C_{\varepsilon}n$ , then A satisfies the MDC with constant  $\varepsilon$  with probability at least  $1 - \exp(-cm\varepsilon^2/2) - m\exp(-n/8) - \exp(-m/2)$ . Here  $C_{\varepsilon} = \Omega(\varepsilon^{-2}\log(\varepsilon^{-1}))$  and c is a universal constant.

Combining this result with Theorem 2.1 with  $\varepsilon = 0.001$  proves Theorem 1.1.

Hence this shows that the MDC is sufficient for local convergence of subgradient descent with proper initialization. Moreover, the conclusion holds for generic measurements with high probability as soon as  $m = \Omega(n)$ . We emphasize that the MDC is a global property concerning the measurement matrix. Hence Proposition 2.1 implies one has uniform concentration of subgradients to their expectation with optimal sample complexity. Extending this local convergence result to a result about convergence of subgradient descent with generic initialization is an interesting future direction, as shown in recent works [6, 26].

<sup>&</sup>lt;sup>3</sup>Indeed, note that if the conditions of Lemma 2.1 are satisfied and  $\eta=0$ , then for all  $x\in\mathbb{R}^n$  such that  $\operatorname{dist}(x,x_*)\leqslant \varepsilon \|x_*\|$  and any  $v_{x,x_*}\in\partial f(x)$ ,  $\|v_{x,x_*}-(x\pm x_*)\|\leqslant \frac{1}{2}\|x\pm x_*\|\Longleftrightarrow \langle v_{x,x_*},x\pm x_*\rangle\geqslant \frac{3}{8}\|x\pm x_*\|^2+\frac{1}{2}\|v_{x,x_*}\|^2$ . Thus the MDC is sufficient to guarantee the  $\operatorname{RC}(\mu,\lambda,\varepsilon)$  holds with  $\mu=3/4$ ,  $\lambda=1$ , and our choice of  $\varepsilon$ .

## 3. Proofs

In this section, we prove Lemma 2.1 and Proposition 2.1. We first introduce some notation used in the proofs. Let  $[n] := \{1, ..., n\}$ . Let  $\mathcal{B}(y,r) := \{x \in \mathbb{R}^n : ||x-y|| \le r\}$  and  $\mathcal{B} := \{x \in \mathbb{R}^n : ||x|| \le 1\}$ . For  $x \in \mathbb{R}^n \setminus \{0\}$ , let  $\hat{x} := x/||x||$ . Let  $\mathbb{1}_{\{E\}}$  be the indicator function on the event E. For a random variable X, let X|(E) be the random variable X conditioned on the event E. Let  $I_n$  be the  $n \times n$  identity matrix. Let  $\mathcal{S}^{n-1}$  denote the unit sphere in  $\mathbb{R}^n$ . We write  $\gamma = \Omega(\delta)$  when  $\gamma \ge C\delta$  for some positive constant C. Similarly, we write  $\gamma = O(\delta)$  when  $\gamma \le C\delta$  for some positive constant C.

For a locally Lipschitz function  $f: \mathcal{X} \to \mathbb{R}$  from a Hilbert space  $\mathcal{X}$  to  $\mathbb{R}$ , the Clarke generalized directional derivative [7] of f at  $x \in \mathcal{X}$  in the direction u is defined by

$$f^o(x;u) := \limsup_{y \to x. t \downarrow 0} \frac{f(y+tu) - f(y)}{t}.$$

Then the generalized subdifferential of f at x is defined as

$$\partial f(x) := \{ v \in \mathbb{R}^n : \langle v, u \rangle \leqslant f^o(x; u), \ \forall u \in \mathcal{X} \}.$$

Any  $v_{x,x_*} \in \partial f(x)$  is called a subgradient of f at x. When f is differentiable at x,  $\partial f(x) = {\nabla f(x)}$ . In the proofs, we will make use of the following fact concerning the Clarke subdifferential of the objective function f. Since f is piecewise quadratic, Theorem 9.6 from [7] asserts that for any  $x \in \mathbb{R}^n$ ,  $\partial f(x)$  can be written equivalently as

$$\partial f(x) = \operatorname{conv}(v_1, v_2, \dots, v_s) \tag{3.1}$$

where  $\operatorname{conv}(\cdot)$  denotes the convex hull of  $v_1, \ldots, v_s$ , s is the number of quadratic functions adjoint to x, and  $v_\ell$  is the gradient of the  $\ell$ -th quadratic function of f at x. For each  $v_\ell$ , there exists a  $w_\ell$  and a sufficiently small  $\delta_\ell > 0$  such that f is differentiable at  $x + \delta_\ell w_\ell$  and  $v_\ell = \lim_{\delta_\ell \downarrow 0} \nabla f(x + \delta_\ell w_\ell)$ .

**3.1.** Convexity property of objective. Here we prove Lemma 2.1, the convexity-like property around the minimizer. In essence, it states that when iterates are near the minimizers, all subgradients point towards the true solution.

*Proof.* (**Proof of Lemma 2.1.**) We consider the case  $x \in \mathcal{B}(x_*, \varepsilon || x_* ||)$  as the case  $x \in \mathcal{B}(-x_*, \varepsilon || x_* ||)$  is similar. Suppose f is differentiable at x. First, note the MDC implies that  $||A_x^{\mathrm{T}} A_x - I_n|| \le \varepsilon$ . Moreover, for any  $x, z \in \mathbb{R}^n$ ,  $||A_x z||^2 \le |\langle A_x^{\mathrm{T}} A_x z, z \rangle - ||z||^2| + ||z||^2 \le (1+\varepsilon)||z||^2$ . Hence  $||A_x|| \le 2$  for all  $x \in \mathbb{R}^n$  when  $\varepsilon < 1$ . Thus, we have

$$||v_{x,x_*} - (x - x_*)|| \le ||A_x^{\mathrm{T}}(A_x - A_{x_*})x_*|| + ||A_x^{\mathrm{T}}A_x(x - x_*) - (x - x_*)|| + ||A_x^{\mathrm{T}}\eta||$$

$$\le 2||(A_x - A_{x_*})x_*|| + \varepsilon||x - x_*|| + 2||\eta||.$$
(3.2)

We now show that for sufficiently small  $\varepsilon$ ,  $\|(A_x - A_{x_*})x_*\| \le 1/8\|x - x_*\|$ . Letting  $\{a_i\}_{i=1}^m$  denote the rows of A, observe that

$$\begin{split} \left\| (A_x - A_{x_*}) x_* \right\|^2 &= \sum_{i=1}^m \left( \operatorname{sgn}(\langle a_i, x \rangle) - \operatorname{sgn}(\langle a_i, x_* \rangle) \right)^2 \langle a_i, x_* \rangle^2 \\ &\leqslant \sum_{i=1}^m \left( \operatorname{sgn}(\langle a_i, x \rangle) - \operatorname{sgn}(\langle a_i, x_* \rangle) \right)^2 \langle a_i, (x - x_*) \rangle^2 \\ &= \left\| A_x (x - x_*) \right\|^2 + \left\| A_{x_*} (x - x_*) \right\|^2 - 2 \langle x - x_*, A_x^{\mathrm{T}} A_{x_*} (x - x_*) \rangle. \end{split}$$

Since  $||A_x^{\mathrm{T}}A_x - I_n|| \le \varepsilon$ , we have  $||A_x(x - x_*)||^2 \le (1 + \varepsilon)||x - x_*||^2$ . The same upper bound holds for  $||A_{x_*}(x - x_*)||^2$ . We now bound  $2\langle x - x_*, A_x^{\mathrm{T}}A_{x_*}(x - x_*)\rangle$  from below. By the MDC, we have

$$|\langle x - x_*, (A_x^{\mathsf{T}} A_{x_*} - \Phi_{x,x_*})(x - x_*)\rangle| \le \varepsilon ||x - x_*||^2.$$
 (3.3)

Since  $x \in \mathcal{B}(x_*, \varepsilon ||x_*||)$ , we have that  $|\theta_{x,x_*}| \leq 2\varepsilon$ . Hence  $\Phi_{x,x_*}$  is approximately an isometry since

$$\|\Phi_{x,x_*} - I_n\| \leqslant \frac{2|\theta_{x,x_*}|}{\pi} \|I_n\| + \frac{2|\sin\theta_{x,x_*}|}{\pi} \|M_{\hat{x}\leftrightarrow\hat{x}_*}\| \leqslant \frac{8\varepsilon}{\pi}$$

where we used  $||M_{z\leftrightarrow w}|| \le 1$  for all  $z,w \in \mathcal{S}^{n-1}$ . Combining this with (3.3), we have  $2\langle x-x_*,A_x^{\mathrm{T}}A_{x_*}(x-x_*)\rangle \ge \left(2-\frac{16\varepsilon}{\pi}-2\varepsilon\right)||x-x_*||^2$ . Thus we attain

$$\begin{split} \|(A_x - A_{x_*})x_*\|^2 &\leqslant \|A_x(x - x_*)\|^2 + \|A_{x_*}(x - x_*)\|^2 - 2\langle x - x_*, A_x^{\mathrm{T}}A_{x_*}(x - x_*)\rangle \\ &\leqslant \left(2 + 2\varepsilon - 2 + \frac{16\varepsilon}{\pi} + 2\varepsilon\right)\|x - x_*\|^2 \\ &= \left(4\varepsilon + \frac{16\varepsilon}{\pi}\right)\|x - x_*\|^2. \end{split}$$

Finally, choosing  $\varepsilon$  so that  $\varepsilon \leq 0.001$ , we conclude

$$2\|(A_x - A_{x_*})x_*\| \le 2\sqrt{4\varepsilon + 16\varepsilon/\pi} \|x - x_*\| \le \frac{1}{4} \|x - x_*\|.$$

Combining this inequality,  $\varepsilon < 1/4$ , and (3.2) shows  $||v_{x,x_*} - (x - x_*)|| \le 1/2 ||x - x_*|| + 2||\eta||$ .

Finally, for non-differentiable x, recall that by (3.1) we can write  $v_{x,x_*} = \sum_{\ell=1}^{s} c_{\ell} v_{\ell}$  where  $c_{\ell} \ge 0$ ,  $\sum_{\ell=1}^{s} c_{\ell} = 1$ , and  $v_{\ell} = \lim_{\delta_{\ell} \downarrow 0} \nabla f(x + \delta_{\ell} w_{\ell})$  for some  $w_{\ell} \in \mathbb{R}^{n}$ . Then, using  $\sum_{\ell=1}^{s} c_{\ell} = 1$  and our result for differentiable points, we conclude that for  $x \in \mathcal{B}(x_*, \varepsilon ||x_*||)$ ,

$$\begin{split} \|v_{x,x_*} - (x - x_*)\| &\leqslant \sum_{\ell = 1}^s c_\ell \|v_\ell - (x - x_*)\| \leqslant \sum_{\ell = 1}^s c_\ell \lim_{\delta_\ell \downarrow 0} \|\nabla f(x + \delta_\ell w_\ell) - (x + \delta_\ell w_\ell - x_*)\| \\ &\leqslant \frac{1}{2} \|x - x_*\| + 2 \|\eta\|. \end{split}$$

**3.2.** Gaussian matrices satisfy the MDC. To show that A satisfies the MDC, we will use novel probabilistic tools developed in [9], which improved the sample complexity required for Gaussian matrices to satisfy a related concentration result introduced in [20] known as the Weight Distribution Condition. We first write  $A_x^T A_y$  in a more convenient form. For  $v \in \mathbb{R}^n$ , let  $\operatorname{diag}(v > 0)$  denote the diagonal matrix whose i-th entry is 1 if  $v_i > 0$  and 0 otherwise. Define  $\operatorname{diag}(v < 0)$  analogously. For  $x \in \mathbb{R}^n$ , let  $A_{+,x} := \operatorname{diag}(Ax > 0)A$  and  $A_{-,x} := \operatorname{diag}(Ax < 0)A$ . Since  $\operatorname{sgn}(b) = \mathbbm{1}_{\{b > 0\}} - \mathbbm{1}_{\{b < 0\}}$  for any  $b \in \mathbb{R}$ , observe that

$$A_x^{\rm T} A_y = A_{+,x}^{\rm T} A_{+,y} + A_{-,x}^{\rm T} A_{-,y} - A_{+,x}^{\rm T} A_{-,y} - A_{-,x}^{\rm T} A_{+,y}.$$

We will establish concentration of each term separately. For the first term, [9] recently showed that concentration is possible when  $m = \Omega(n)$ :

LEMMA 3.1 (Theorem 3.2 in [9]). Fix  $\varepsilon > 0$ . If  $A \in \mathbb{R}^{m \times n}$  has i.i.d.  $\mathcal{N}(0, 1/m)$  entries and  $m \ge C\varepsilon^{-2}\log(\varepsilon^{-1})n$ , then with probability at least  $1 - \exp(-cm\varepsilon^2/2) - m\exp(-n/8) - \exp(-m/2)$ , we have

$$||A_{+,x}^{\mathrm{T}}A_{+,y} - Q_{x,y}|| \leq \varepsilon \ \forall \ x,y \in \mathbb{R}^n$$

where  $Q_{x,y} := \frac{\pi - \theta_{x,y}}{2\pi} I_n + \frac{\sin \theta_{x,y}}{2\pi} M_{\hat{x} \leftrightarrow \hat{y}}$  if  $x,y \neq 0$  and  $0_{n \times n}$  otherwise. Here C and c are absolute constants.

An elementary calculation shows  $\mathbb{E}[A_{+,x}^{\mathrm{T}}A_{+,y}] = Q_{x,y}$ . Also by symmetry,  $\mathbb{E}[A_{-,x}^{\mathrm{T}}A_{-,y}] = Q_{x,y}$ . By applying a nearly identical argument as in [9], the analogous result for  $A_{-,x}^{\mathrm{T}}A_{-,y}$  holds.

LEMMA 3.2. Fix  $\varepsilon > 0$ . If  $A \in \mathbb{R}^{m \times n}$  has i.i.d.  $\mathcal{N}(0,1/m)$  entries and  $m \ge C\varepsilon^{-2}\log(\varepsilon^{-1})n$  then with probability at least  $1 - \exp(-cm\varepsilon^2/2) - m\exp(-n/8) - \exp(-m/2)$ , we have

$$||A_{-,x}^{\mathrm{T}}A_{-,y} - Q_{x,y}|| \leq \varepsilon \ \forall \ x,y \in \mathbb{R}^n.$$

Here C and c are absolute constants.

We now extend the argument in [9] for  $A_{+,x}^{\mathrm{T}}A_{-,y}$ . Note that a result for  $A_{-x}^{\mathrm{T}}A_{+,y}$  would be identical. Observe that

$$\mathbb{E}[A_{+,x}^{\mathrm{T}}A_{-,y}] = H_{x,y} := \frac{\theta_{x,y}}{2\pi} I_n - \frac{\sin \theta_{x,y}}{2\pi} M_{\hat{x} \leftrightarrow \hat{y}}.$$

We will prove the following:

LEMMA 3.3. Fix  $\varepsilon > 0$ . If  $A \in \mathbb{R}^{m \times n}$  has i.i.d.  $\mathcal{N}(0,1/m)$  entries and  $m \ge C\varepsilon^{-2}\log(\varepsilon^{-1})n$  then with probability at least  $1 - \exp(-cm\varepsilon^2/2) - m\exp(-n/8) - \exp(-m/2)$ , we have

$$||A_{+,x}^{\mathrm{T}}A_{-,y} - H_{x,y}|| \leq \varepsilon \ \forall \ x,y \in \mathbb{R}^n.$$

Here C and c are absolute constants.

Note that this would complete Proposition 2.1 by observing  $\Phi_{x,y} = 2Q_{x,y} - 2H_{x,y}$  and combining Lemmas 3.1, 3.2, 3.3, and an analogous result for  $A_{-x}^{T}A_{+,y}$ , each satisfied with  $\varepsilon/4$ .

The main probabilistic tool in the proof of Lemma 3.3 is a result concerning concentration of pseudo-Lipschitz functions. Pseudo-Lipschitzness can be considered as a relaxation of standard Lipschitz continuity but with particular attention towards which sets a function is Lipschitz with respect to. When the sets are balls, then the notion of pseudo-Lipschitzness reduces to standard Lipschitzness. Prior to stating the result, we require the following definitions.

DEFINITION 3.1  $((\delta, \gamma)$ -wide system). A set system  $\{B_t \subseteq \mathbb{R}^n : t \in \Theta\}$  is  $(\delta, \gamma)$ -wide if  $B_t = -B_t$ ,  $B_t$  is convex, and  $Vol(B_t \cap \delta \mathcal{B}) \geqslant \gamma Vol(\delta \mathcal{B}) \ \forall \ t \in \Theta$ .

DEFINITION 3.2 (pseudo-Lipschitz function). Suppose there exists a  $(\delta, \gamma)$ -wide system  $\{B_t \subseteq \mathbb{R}^n : t \in \Theta\}$  such that  $|g_t(x) - g_t(y)| \leq \varepsilon$  for any  $t \in \Theta$  and  $x, y \in (\mathbb{R}^n)^d$  with  $x_i - y_i \in B_t$  for all  $i \in [d]$ . Then we say that  $\{g_t\}_{t \in \Theta}$  is  $(\varepsilon, \delta, \gamma)$ -pseudo-Lipschitz.

Note here that a function is pseudo-Lipschitz with respect to a *particular* system of sets. The following theorem establishes favorable concentration for pseudo-Lipschitz functions.

THEOREM 3.1 (Theorem 4.4 in [9]). Let  $\theta$  be a random variable taking values in  $\Theta$ . Let  $\{g_t : (\mathbb{R}^n)^d \to \mathbb{R} : t \in \Theta\}$  be a function family and let  $h : (\mathbb{R}^n)^d \to \mathbb{R}$  be a function. Let  $\varepsilon, \gamma, D > 0$  and  $\delta \in (0,1)$ . Define the spherical shell  $\mathcal{H} := (1+\delta/2)\mathcal{B} \setminus (1-\delta/2)\mathcal{B}$  in  $\mathbb{R}^n$ . Suppose:

- (1) For any fixed  $x \in \mathcal{H}^d$ ,  $\mathbb{P}_{\theta}(g_{\theta}(x) \leq h(x) + \varepsilon) \geq 1 p$ ,
- (2)  $\{g_t\}_{t\in\Theta}$  is  $(\varepsilon,\delta,\gamma)$ -pseudo-Lipschitz,
- (3)  $|h(x)-h(y)| \leq D$  whenever  $x \in (S^{n-1})^d$ ,  $y \in (\mathbb{R}^n)^d$ , and  $||y_i-x_i|| \leq \delta$  for all  $i \in [d]$ . Then

$$\mathbb{P}_{\theta}\left(g_{\theta}(x) \leqslant h(x) + 2\varepsilon + D, \ \forall \ x \in (\mathcal{S}^{n-1})^d\right) \geqslant 1 - \gamma^{-2d} (4/\delta)^{2dn} p.$$

**3.2.1. Proof of Lemma 3.3.** For ease of exposition, assume the entries of A are i.i.d.  $\mathcal{N}(0,1)$ . The main idea is that we will concentrate Lipschitz approximations of  $A_{+,x}^{\mathrm{T}}A_{-,y}$  that are upper and lower bounds with respect to the semidefinite ordering. For  $\varepsilon \in (0,1)$ , define the following continuous relaxations of  $\mathbb{1}_{\{t>0\}}$ :

$$\varphi_{-\varepsilon}^+(t) := \begin{cases} 0 & t \leqslant -\varepsilon \\ 1 + t/\varepsilon & -\varepsilon < t \leqslant 0 \text{ and } \varphi_{\varepsilon}^+(t) := \begin{cases} 0 & t < 0 \\ t/\varepsilon & 0 \leqslant t < \varepsilon \end{cases}.$$

$$1 & t > 0$$

Analogously define the following continuous relaxations for  $\mathbb{1}_{\{t<0\}}$ :

$$\varphi_{-\varepsilon}^-(t) := \begin{cases} 1 & t \leqslant -\varepsilon \\ -t/\varepsilon & -\varepsilon < t \leqslant 0 \text{ and } \varphi_\varepsilon^-(t) := \begin{cases} 1 & t < 0 \\ 1 - t/\varepsilon & 0 \leqslant t < \varepsilon \end{cases}.$$

Then we have that for all  $t \in \mathbb{R}$ ,  $\varphi_{\varepsilon}^+(t) \leqslant \mathbb{1}_{\{t>0\}} \leqslant \varphi_{-\varepsilon}^+(t)$  and  $\varphi_{-\varepsilon}^-(t) \leqslant \mathbb{1}_{\{t<0\}} \leqslant \varphi_{\varepsilon}^-(t)$ . For  $V \in \mathbb{R}^{m \times n}$  with rows  $v_i$  for  $i \in [m]$  and  $x, y \in \mathbb{R}^n$ , define

$$G_{V,\text{up}}(x,y) := \sum_{i=1}^{m} \varphi_{-\varepsilon}^{+}(\langle v_i, x \rangle) \varphi_{\varepsilon}^{-}(\langle v_i, y \rangle) v_i v_i^{\text{T}}$$

and

$$G_{V,\text{low}}(x,y) := \sum_{i=1}^{m} \varphi_{\varepsilon}^{+}(\langle v_{i}, x \rangle) \varphi_{-\varepsilon}^{-}(\langle v_{i}, y \rangle) v_{i} v_{i}^{T}.$$

Note that for any  $x, y \in \mathbb{R}^n$ ,  $G_{A,\text{low}}(x,y) \leq A_{+,x}^T A_{-,y} \leq G_{A,\text{up}}(x,y)$  so it suffices to upper bound  $G_{A,\text{up}}(x,y)$  and lower bound  $G_{A,\text{low}}(x,y)$  uniformly. For the upper bound, we will prove the following:

PROPOSITION 3.1. Fix  $0 < \varepsilon < 1$ . Suppose  $A \in \mathbb{R}^{m \times n}$  has i.i.d.  $\mathcal{N}(0,1)$  entries. Then if  $m \geqslant C\varepsilon^{-2}\log(\varepsilon^{-1})n$ , we have that with probability at least  $1 - \exp(-cm\varepsilon^2/2) - \exp(-n/8) - \exp(m/2)$ ,

$$G_{A,up}(x,y) \leq mH_{x,y} + m\varepsilon I_n \ \forall \ x,y \neq 0.$$

Here C and c are absolute constants.

The central argument can be broken down into three steps and directly follows [9]. We first show that the function  $g_V(x,y) := \frac{1}{m} \langle u, G_{V,\text{up}}(x,y)u \rangle$  is  $(\varepsilon, \delta, \gamma)$ -pseudo-Lipschitz for fixed  $u \in \mathcal{S}^{n-1}$  for appropriate parameters  $\varepsilon, \delta$ , and  $\gamma$ . Second, we use Theorem 3.1 to establish, for fixed u, concentration of  $g_A(x,y)$  uniformly in x,y to  $h(x,y) := \langle u, H_{x,y}u \rangle$ . Finally, we use a standard  $\varepsilon$ -net argument to establish uniform concentration over u, guaranteeing an upper bound on  $G_{A,\text{up}}(x,y)$ . Throughout the proof, we will operate on the set of matrices

$$\Theta := \left\{ V \in \mathbb{R}^{m \times n} : \|V\| \leqslant 3\sqrt{m}, \ \max_{i \in [m]} \|v_i\| \leqslant \sqrt{2n} \right\}.$$

When A is Gaussian, standard results [27] show that  $A \in \Theta$  with high probability.

LEMMA 3.4 ([27]). Suppose  $A \in \mathbb{R}^{m \times n}$  has i.i.d.  $\mathcal{N}(0,1)$  entries. Then with probability at least  $1 - \exp(-m/2) - m \exp(-n/8)$ , we have  $||A|| \leq 3\sqrt{m}$  and  $\max_{i \in [m]} ||a_i|| \leq \sqrt{2n}$ .

Step 1: Establishing pseudo-Lipschitzness. We first establish that  $\{g_V\}_{V\in\Theta}$  is pseudo-Lipschitz with respect to a particular set system.

LEMMA 3.5. Fix  $\varepsilon > 0$  and  $u \in \mathcal{S}^{n-1}$ . For  $V \in \Theta$ , define  $g_V(x,y) := \frac{1}{m} \langle u, G_{V,up}(x,y)u \rangle$ . Then  $\{g_V\}_{V \in \Theta}$  is  $(2\varepsilon, \varepsilon^2/82, 1/2)$ -pseudo-Lipschitz with respect to the set system  $\{B_{M,\varepsilon^2,u}\}_{V \in \Theta}$  where

$$B_{M,\varepsilon^2,u} := \left\{z \in \mathbb{R}^n : \sum_{i=1}^m |\langle v_i,z\rangle| \langle v_i,u\rangle^2 \leqslant \varepsilon^2 m \right\}.$$

*Proof.* We first note that it was shown in Lemma 5.5 of [9] that the set system  $\{B_{M,\varepsilon^2,u}\}_{V\in\Theta}$  is  $(\varepsilon^2/82,1/2)$ -wide. We now show that  $\{g_V\}_{V\in\Theta}$  is  $(2\varepsilon,\varepsilon^2/82,1/2)$ -pseudo-Lipschitz. For  $x,y,\tilde{x},\tilde{y}\in\mathbb{R}^n$ , suppose  $y-\tilde{y}\in B_{M,\varepsilon^2,u}$  and  $x-\tilde{x}\in B_{M,\varepsilon^2,u}$ . Then observe that

$$|g_{V}(x,y) - g_{V}(\tilde{x},\tilde{y})| \leq \frac{1}{m} \sum_{i=1}^{m} [|\varphi_{-\varepsilon}^{+}(\langle v_{i}, x \rangle)\varphi_{\varepsilon}^{-}(\langle v_{i}, y \rangle) - \varphi_{-\varepsilon}^{+}(\langle v_{i}, \tilde{x} \rangle)\varphi_{\varepsilon}^{-}(\langle v_{i}, y \rangle)|$$

$$+ |\varphi_{-\varepsilon}^{+}(\langle v_{i}, \tilde{x} \rangle)\varphi_{\varepsilon}^{-}(\langle v_{i}, y \rangle) - \varphi_{-\varepsilon}^{+}(\langle v_{i}, \tilde{x} \rangle)\varphi_{\varepsilon}^{-}(\langle v_{i}, \tilde{y} \rangle)|]\langle v_{i}, u \rangle^{2}$$

$$\leq \frac{1}{m} \sum_{i=1}^{m} [|\varphi_{-\varepsilon}^{+}(\langle v_{i}, x \rangle) - \varphi_{-\varepsilon}^{+}(\langle v_{i}, \tilde{x} \rangle)|$$

$$+ |\varphi_{\varepsilon}^{-}(\langle v_{i}, y \rangle) - \varphi_{\varepsilon}^{-}(\langle v_{i}, \tilde{y} \rangle)|]\langle v_{i}, u \rangle^{2}$$

$$\leq \frac{1}{m\varepsilon} \sum_{i=1}^{m} [|\langle v_{i}, x - \tilde{x} \rangle| + |\langle v_{i}, y - \tilde{y} \rangle|]\langle v_{i}, u \rangle^{2}$$

$$\leq 2\varepsilon.$$

In the first inequality, we used the triangle inequality. In the second, we used  $|\varphi_{-\varepsilon}^+(t)|, |\varphi_{\varepsilon}^-(t)| \leqslant 1$ . In the third, we used the fact that  $\varphi_{-\varepsilon}^+$  and  $\varphi_{\varepsilon}^-$  are both  $1/\varepsilon$ -Lipschitz. In the last inequality, we used the assumptions  $y - \tilde{y} \in B_{M,\varepsilon^2,u}$  and  $x - \tilde{x} \in B_{M,\varepsilon^2,u}$ .

Step 2: Point-wise concentration. We now show that, for fixed  $u \in S^{n-1}$ ,  $q_V(x,y)$  concentrates around h(x,y) uniformly in x,y by an application of Theorem 3.1.

LEMMA 3.6. Fix  $\varepsilon > 0$  and  $u \in \mathcal{S}^{n-1}$ . Let  $A \in \mathbb{R}^{m \times n}$  have i.i.d.  $\mathcal{N}(0,1)$  entries. Define  $\theta := A|(A \in \Theta)$ . There exist absolute constants c, K, and  $\tilde{C}$  such that

$$\mathbb{P}_{\theta}\left(g_{\theta}(x,y)\leqslant h(x,y)+K\varepsilon\ \forall\ x,y\neq 0\right)\geqslant 1-(\tilde{C}/\varepsilon)^{8n}\exp(-cm\varepsilon^2).$$

*Proof.* We will first bound  $\mathbb{E}[G_{A,\text{up}}(x,y)]$ . Observe that for any  $t \in \mathbb{R}$ , we have  $\varphi_{-\varepsilon}^+(t) \leq \mathbb{1}_{\{t \geq -\varepsilon\}}$  and  $\varphi_{\varepsilon}^-(t) \leq \mathbb{1}_{\{t \leq \varepsilon\}}$ . This implies that for any  $t_1, t_2 \in \mathbb{R}$ ,

$$\varphi_{-\varepsilon}^+(t_1)\varphi_{\varepsilon}^-(t_2) \leqslant \mathbb{1}_{\{t_1\geqslant -\varepsilon\}} \, \mathbb{1}_{\{t_2\leqslant \varepsilon\}} \leqslant \mathbb{1}_{\{t_1>0\}} \, \mathbb{1}_{\{t_2<0\}} + \mathbb{1}_{\{-\varepsilon\leqslant t_1\leqslant 0\}} + \mathbb{1}_{\{0\leqslant t_2\leqslant \varepsilon\}} \, .$$

Thus

$$\mathbb{E}[G_{A,\mathrm{up}}(x,y)] \leq \mathbb{E}\left[\sum_{i=1}^{m} (\mathbb{1}_{\{\langle a_i,x\rangle > 0\}} \mathbb{1}_{\{\langle a_i,y\rangle < 0\}} + \mathbb{1}_{\{-\varepsilon \leqslant \langle a_i,x\rangle \leqslant 0\}} + \mathbb{1}_{\{0 \leqslant \langle a_i,y\rangle \leqslant \varepsilon\}}) a_i a_i^{\mathrm{T}}\right]$$

$$= mH_{x,y} + m\mathbb{E}[\mathbb{1}_{\{-\varepsilon \leqslant \langle a,x\rangle \leqslant 0\}} a a^{\mathrm{T}}] + m\mathbb{E}[\mathbb{1}_{\{0 \leqslant \langle a,y\rangle \leqslant \varepsilon\}} a a^{\mathrm{T}}]$$

where  $a \sim \mathcal{N}(0, I_n)$ . It was shown in Lemma 12 of [20] that  $\mathbb{E}[\mathbb{1}_{\{-\varepsilon \leqslant \langle a, x \rangle \leqslant 0\}} aa^{\mathrm{T}}] \leq \frac{\varepsilon}{2\|x\|} I_n \ \forall \ x \neq 0$ . An analogous bound shows  $\mathbb{E}[\mathbb{1}_{\{0 \leqslant \langle a, y \rangle \leqslant \varepsilon\}} aa^{\mathrm{T}}] \leq \frac{\varepsilon}{2\|y\|} I_n$  for all  $y \neq 0$ . Hence we attain

$$\mathbb{E}[G_{A,\mathrm{up}}(x,y)] \leq mH_{x,y} + m\left(\frac{\varepsilon}{2||x||} + \frac{\varepsilon}{2||y||}\right)I_n \ \forall \ x,y \neq 0.$$

This implies  $\mathbb{E}[g_A(x,y)] \leq h(x,y) + 2\varepsilon$  for fixed  $x,y \in \mathbb{R}^n$  with  $||x||, ||y|| \geq 1/2$ .

Now, we show the probability bound. First consider fixed  $x,y \in \mathbb{R}^n$  with  $\|x\|,\|y\| \geqslant 1/2$ . Observe that  $g_A(x,y) = \frac{1}{m} \sum_{i=1}^m \varphi_{-\varepsilon}^+(\langle a_i,x\rangle) \varphi_{\varepsilon}^-(\langle a_i,y\rangle) \langle a_i,u\rangle^2$  is a sum of subexponential random variables. Hence by Bernstein's inequality, we have for some absolute constant c and any  $\beta > 0$ ,  $\mathbb{P}(g_A(x,y) - \mathbb{E}[g_A(x,y)] > \beta) \leqslant 2\exp(-cm\min(\beta,\beta^2))$ . Taking  $\beta = \varepsilon$  and using  $\mathbb{E}[g_A(x,y)] \leqslant h(x,y) + 2\varepsilon$ , we get  $\mathbb{P}(g_A(x,y) > h(x,y) + 3\varepsilon) \leqslant 2\exp(-cm\varepsilon^2)$ . Since  $\mathbb{P}(A \in \Theta) \geqslant 1/2$ , conditioning on the event  $A \in \Theta$  at most doubles the failure probability so we attain

$$\mathbb{P}(g_{\theta}(x,y) \leqslant h(x,y) + 3\varepsilon) \geqslant 1 - 4\exp(-cm\varepsilon^2). \tag{3.4}$$

To establish uniform concentration in x,y, we note that by a simple modification to Lemma 27 in [17], we have that  $H_{x,y}$  is L-Lipschitz with respect to  $x,y \in \mathcal{S}^{n-1}$  where  $L = 22/\pi$ . Hence  $|h(x,y) - h(\tilde{x},\tilde{y})| \leq L\varepsilon$  if  $||x - \tilde{x}|| \leq \varepsilon$  and  $||y - \tilde{y}|| \leq \varepsilon$ . Thus the result follows by applying Theorem 3.1 to  $\{g_V\}_{V \in \Theta}$  and  $\theta := A|(A \in \Theta)$ . By Lemma 3.5,  $\{g_V\}_{V \in \Theta}$  is  $(2\varepsilon,\varepsilon^2/82,1/2)$ -pseudo-Lipschitz. We can then take  $p = \exp(-cm\varepsilon^2)$  by (3.4) and  $D = 2L\varepsilon$ .

Step 3: Uniform concentration. The last step is to get a uniform bound over all  $u \in \mathcal{S}^{n-1}$ . Augmenting our notation, let  $g_V(x,y,u) := \frac{1}{m} \langle u, G_{V,\text{up}}(x,y)u \rangle$  and  $h(x,y,u) := \langle u, H_{x,y}u \rangle$ .

LEMMA 3.7. Fix  $\varepsilon > 0$ . Suppose  $A \in \mathbb{R}^{m \times n}$  has i.i.d.  $\mathcal{N}(0,1)$  entries. There exist absolute constants c and C such that if  $m \ge C\varepsilon^{-2}\log(\varepsilon^{-1})n$ , then with probability  $1 - \exp(-cm\varepsilon^2/2) - m\exp(-n/8) - \exp(-m/2)$ ,

$$G_{A,up}(x,y) \leq mH_{x,y} + m\varepsilon \ \forall x,y \neq 0.$$

*Proof.* We first show that  $g_V(x,y,u)$  is 18-Lipschitz with respect to  $u \in \mathcal{S}^{n-1}$  when  $V \in \Theta$ . Fix  $x,y \neq 0$ . Observe that for  $u,w \in \mathcal{S}^{n-1}$ ,

$$|g_{V}(x,y,u) - g_{V}(x,y,w)| \leq \frac{1}{m} \sum_{i=1}^{m} |\langle v_{i}, u \rangle^{2} - \langle v_{i}, w \rangle^{2}| \leq \frac{1}{m} \sum_{i=1}^{m} |\langle v_{i}, u - w \rangle| |\langle v_{i}, u + w \rangle|$$

$$\leq \frac{1}{m} ||V(u - w)|| ||V(u + w)||$$

$$\leqslant 18||u - w|| \tag{3.5}$$

where we used  $||V|| \leq 3\sqrt{m}$  along with  $u, w \in \mathcal{S}^{n-1}$  in the last inequality.

Let  $\mathcal{N}_{\varepsilon} \subset \mathcal{S}^{n-1}$  be an  $\varepsilon$ -net of cardinality  $|\mathcal{N}_{\varepsilon}| \leq (3/\varepsilon)^n$ . By Lemma 3.6 and a union bound, it holds with probability at least  $1 - (3/\varepsilon)^n (\tilde{C}/\varepsilon)^{8n} \exp(-cm\varepsilon^2)$  over  $\theta = A|(A \in \Theta)$  that for all  $x, y \in \mathcal{S}^{n-1}$  and  $u \in \mathcal{N}_{\varepsilon}$ ,  $g_{\theta}(x, y, u) \leq h(x, y, u) + K\varepsilon$ . For any  $x, y, u \in \mathcal{S}^{n-1}$ , there exists a  $w \in \mathcal{N}_{\varepsilon}$  with  $||u - w|| \leq \varepsilon$  so using (3.5) we get

$$g_{\theta}(x,y,u) \leqslant g_{\theta}(x,y,w) + 18||u-w|| \leqslant h(x,y,w) + 18\varepsilon.$$

Since  $||H_{x,y}|| \leq 2$ , we have that for  $u, w \in \mathcal{S}^{n-1}$ ,  $|h(x,y,u) - h(x,y,w)| \leq 4||u-w||$ . This further implies  $g_{\theta}(x,y,u) \leq h(x,y,u) + 22\varepsilon$ . We conclude with probability at least  $1 - (3/\varepsilon)^n (C/\varepsilon)^{8n} \exp(-cm\varepsilon^2)$  over  $\theta$ , the desired inequality holds. Using  $\mathbb{P}(A \in \Theta) \geqslant 1 - m\exp(-n/8) - \exp(-m/2)$  and taking  $m \geqslant C\varepsilon^{-2}\log(\varepsilon^{-1})n$  achieves the final result with the desired probability.

This completes the upper bound on  $G_{A,up}$ . The lower bound on  $G_{A,low}$  is identical:

LEMMA 3.8. Fix  $\varepsilon > 0$ . Suppose  $A \in \mathbb{R}^{m \times n}$  has i.i.d.  $\mathcal{N}(0,1)$  entries. There exist absolute constants c and C such that if  $m \ge C\varepsilon^{-2}\log(\varepsilon^{-1})n$ , then with probability  $1 - \exp(-cm\varepsilon^2/2) - m\exp(-n/8) - \exp(-m/2)$ ,

$$G_{A,low}(x,y) \succeq mH_{x,y} - m\varepsilon \ \forall x,y \neq 0.$$

Lemma 3.3 follows by combining Lemma 3.7 and Lemma 3.8.

**Acknowledgements.** PH is supported by NSF Grant DMS-2022205 and NSF CAREER Grant DMS-1848087. OL acknowledges support by the NSF Graduate Research Fellowship under Grant No. DGE-1450681.

### REFERENCES

- [1] S. Bahmani and J. Romberg, Flexible convex relaxation for phase retrieval, Electron. J. Stat., 11(2):5254–5281, 2017. 1
- [2] E.J. Candès and X. Li, Solving quadratic equations via phaselift when there are about as many equations as unknowns, Found. Comput. Math., 14:1017-1026, 2014.
- [3] E.J. Candès, X. Li, and M. Soltanolkotabi, Retrieval via Wirtinger flow: Theory and applications, IEEE Trans. Inf. Theory, 61(4):1985–2007, 2017. 1
- [4] E.J. Candès, T. Strohmer, and V. Voroninski, Phaselift: Exact and stable signal recovery from magnitude measurements via convex programming, Commun. Pure. Appl. Math., 66(8):1241– 1274, 2013.
- [5] Y. Chen and E.J. Candès, Solving random quadratic systems of equations is nearly as easy as solving linear systems, Commun. Pure Appl. Math., 70(5):822–883, 2017. 1, 1
- [6] Y. Chen, Y. Chi, J. Fan, and C. Ma, Gradient descent with random initialization: Fast global convergence for nonconvex phase retrieval, Math. Program., 176(1-2):5-37, 2019.
- [7] C. Clason, Nonsmooth analysis and optimization, arXiv preprint, arXiv:1708.04180, 2017. 3
- [8] J. Cocola, P. Hand, and V. Voroninski, Nonasymptotic guarantees for spiked matrix recovery with generative priors, in H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan and H. Lin (eds.), Adv. Neur. Inf. Process. Syst., 15185-15197, 2020. 1
- [9] C. Daskalakis, D. Rohatgi, and M. Zampetakis, Constant-expansion suffices for compressed sensing with generative priors, in H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan and H. Lin (eds.), Adv. Neur. Inf. Process. Syst., 13917–13926, 2020. 1, 2, 3.2, 3.1, 3.2, 3.2, 3.1, 3.2.1, 3.2.1
- [10] L. Demanet and P. Hand, Stable optimizationless recovery from phaseless linear measurements, J. Fourier Anal. Appl., 20(1):199–221, 2014. 1
- [11] J. Duchi and F. Ruan, Solving (most) of a set of quadratic equalities: Composite optimization for robust phase retrieval, Inf. Inference, 8(3):471-529, 2019. 1, 1
- [12] G. Wang, G.B. Giannakis, and Y.C. Eldar, Solving systems of random quadratic equations via truncated amplitude flow, IEEE Trans. Inf. Theory, 64(2):773-794, 2018. 1, 1

- [13] A. Fannjiang and T. Strohmer, The numerics of phase retrieval, Acta Numer., 29:125–228, 2020.
- [14] T. Goldstein and C. Struder, Phasemax: Convex phase retrieval via basis pursuit, IEEE Trans. Inf. Theory, 64(4):2675–2689, 2018.
- [15] P. Hand, Phaselift is robust to a constant fraction of arbitrary errors, Appl. Comput. Harmonic Anal., 42(3):550–562, 2017. 1
- [16] P. Hand and B. Joshi, Global guarantees for blind demodulation with generative priors, in H. Wallach, H. Larochelle, A. Beygelzimer, F. dAlché-Buc, E. Fox and R. Garnett (eds.), Adv. Neur. Inf. Process. Syst., 11535–11545, 2019.
- [17] P. Hand, O. Leong, and V. Voroninski, Phase retrieval under a generative prior, in S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi and R. Garnett (eds.), Adv. Neur. Inf. Process. Syst., 9154–9164, 2018. 1, 3.2.1
- [18] P. Hand, O. Leong, and V. Voroninski, Compressive phase retrieval: Optimal sample complexity with deep generative priors, arXiv preprint, arXiv:2008.10579, 2020. 1
- [19] P. Hand and V. Voroninski, An elementary proof of convex phase retrieval in the natural parameter space via the linear program PhaseMax, Commun. Math. Sci., 16(7):2047–2051, 2018.
- [20] P. Hand and V. Voroninski, Global guarantees for enforcing deep generative priors by empirical risk, IEEE Trans. Inf. Theory, 66(1):401–418, 2019. 1, 2, 3.2, 3.2.1
- [21] R. Heckel, W. Huang, P. Hand, and V. Voroninski, Rate-optimal denoising with deep neural networks, Inf. Inference, iaaa011, 2020. 1
- [22] W. Huang, P. Hand, R. Heckel, and V. Voroninski, A provably convergent scheme for compressive sensing under random generative priors, J. Fourier Anal. Appl., 27:19, 2021. 1
- [23] Q. Luo, H. Wang and S. Lin, Phase retrieval via smoothed amplitude flow, Signal Process., 177:107719, 2020. 1
- [24] S. Qiu, X. Wei, and Z. Yang, Robust one-bit recovery via ReLu generative networks: Improved statistical rates and global landscape analysis, International Conference on Machine Learning (ICML), 2020. 1
- [25] J. Sun, Q. Qu, and J. Wright, A geometric analysis of phase retrieval, Found. Comput. Math., 18(5):1131–1198, 2018. 1
- [26] Y.S. Tan and R. Vershynin, Online stochastic gradient descent with arbitrary initialization solves non-smooth, non-convex phase retrieval, arXiv preprint, arXiv:1910.12837, 2019.
- [27] R. Vershynin, Introduction to the non-asymptotic analysis of random matrices, in Y.C. Eldar and G. Kutyniok (eds.), Compressed Sensing: Theory and Applications, Cambridge University Press, 210–268, 2012. 3.2.1, 3.4
- [28] B. Yonel and B. Yazici, A deterministic theory for exact non-convex phase retrieval, IEEE Trans. Signal Process., 68:4612–4626, 2020. 1
- [29] H. Zhang, Y. Liang, and Y. Chi, A nonconvex approach for phase retrieval: Reshaped Wirtinger flow and incremental algorithms, J. Mach. Learn. Res., 18(141):1–35, 2017. 1, 1, 1