# The nonsmooth landscape of phase retrieval

DAMEK DAVIS

*School of Operations Research and Information Engineering, Cornell University,*
*Ithaca, NY 14850, USA*
dsd95@cornell.edu

DMITRIY DRUSVYATSKIY

*Department of Mathematics, University of Washington, Seattle, WA 98195, USA*
*Corresponding author: ddrusv@uw.edu

AND

COURTNEY PAQUETTE

*Industrial and Systems Engineering Department, Lehigh University, Bethlehem, PA 18015, US*
cop318@lehigh.edu

We consider a popular nonsmooth formulation of the real phase retrieval problem. We show that under standard statistical assumptions a simple subgradient method converges linearly when initialized within a constant relative distance of an optimal solution. Seeking to understand the distribution of the stationary points of the problem, we complete the paper by proving that as the number of Gaussian measurements increases, the stationary points converge to a codimension two set, at a controlled rate. Experiments on image recovery problems illustrate the developed algorithm and theory.

*Keywords*: phase retrieval; stationary points; subdifferential; variational principle; subgradient method; spectral functions; eigenvalues.

## 1. Introduction

Phase retrieval is a common task in computational science, with numerous applications including imaging, X-ray crystallography and speech processing. In this work, we consider a popular real counterpart of the problem. Given a set of tuples $\{(a_i, b_i)\}_{i=1}^m \subset \mathbb{R}^d \times \mathbb{R}$, the (real) phase retrieval problem seeks to determine a vector $x \in \mathbb{R}^d$ satisfying $(a^T x)^2 = b_i$ for each index $i = 1, \ldots, m$. Due to its combinatorial nature, this problem is known to be NP-hard (Fickus *et al.*, 2014). One can model the real phase retrieval problem in a variety of ways. Here, we consider the following 'robust formulation':

$$\min_x f_S(x) := \frac{1}{m} \sum_{i=1}^m |(a_i^T x)^2 - b_i|.$$

This model of the problem has gained some attention recently with the work of Duchi & Ruan (2017) and Eldar & Mendelson (2014). Indeed, this model exhibits a number of desirable properties, making it amenable to numerical methods. Namely, in contrast to other possible formulations, mild statistical assumptions imply that $f_S$ is both *weakly convex* (Duchi & Ruan, 2017, Corollary 3.2) and

*sharp* (Eldar & Mendelson, 2014, Theorem 2.4), with high probability. That is, there exist numerical constants $\rho, \kappa > 0$ such that

$$\text{the assignment} x \mapsto f_S(x) + \frac{\rho}{2}\|x\|^2 \text{ is a convex function,}$$

and the inequality

$$f_S(x) - \inf f_S \geq \kappa \|x - \bar{x}\| \|x + \bar{x}\| \qquad \text{holds for all } x \in \mathbb{R}^d.$$

Here, $\pm\bar{x}$ are the true signals and $\|\cdot\|$ denotes the $\ell_2$-norm. Weak convexity is a well-studied concept in optimization literature (Federer, 1959; Rockafellar, 1982; Clarke *et al.*, 1995; Poliquin & Rockafellar, 1996), while sharpness and the closely related notion of error bounds (Burke & Ferris, 1993; Luo & Tseng, 1993; Drusvyatskiy & Lewis, 2018) classically underlie rapid local convergence guarantees in nonlinear programming. Building on these observations, Duchi & Ruan (2017) showed that with proper initialization, the so-called *prox-linear algorithm* (Lewis & Wright, 2016; Drusvyatskiy & Lewis, 2018; Drusvyatskiy & Paquette, 2019; Duchi & Ruan, 2018a, 2018b) quadratically converges to $\pm\bar{x}$ (even in presence of outliers). The only limitation of their approach is that the prox-linear method requires, at every iteration, invoking an iterative solver for a convex subproblem. For large-scale instances ($m \gg 1, d \gg 1$), the numerical resolution of such problems is nontrivial. In the current work, we analyse a lower-cost alternative when there are no errors in the measurements.

We will show that the robust phase retrieval objective favourably lends itself to classical subgradient methods. This is somewhat surprising because, until recently, convergence rates of subgradient methods in nonsmooth, nonconvex optimization have remained elusive; see the discussion in Davis & Grimmer (2019). We will prove that under mild statistical assumptions and proper initialization, the standard Polyak subgradient method

$$x_{k+1} = x_k - \left(\frac{f_S(x_k) - \min f_S}{\|g_k\|^2}\right) g_k \qquad \text{with} \qquad g_k \in \partial f_S(x_k)$$

linearly converges to $\pm\bar{x}$, with high probability. We note that high quality initialization, in turn, is straightforward to obtain using a spectral method; see e.g. Duchi & Ruan (2017, Section 3.3) and Wang *et al.* (2017). The argument we present is appealingly simple, relying only on weak convexity and sharpness of the function.

Aside from the current work and that of Duchi & Ruan (2017), we are not aware of other attempts to optimize the robust phase retrieval objective directly. Other works focus on different problem formulations. Notably, the papers (Candès *et al.*, 2013; Candès & Li, 2014) prove exact recovery of the signal for a semidefinite programming (SDP) relaxation called PhaseLift, while Bahmani & Romberg (2017), Hand & Voroninski (2016) and Goldstein & Studer (2018) investigate exact recovery of the linear programming relaxation, PhaseMax. Though the latter linear program is certainly much smaller than the SDP, it can still be quite large. The two papers (Candès *et al.*, 2015; Chen & Candès, 2017) instead optimize the smooth loss $\frac{1}{m}\sum_{i=1}^m (\langle a_i, x \rangle - \sqrt{b_i})^2$ using a specialized gradient method called (truncated) Wirtinger flow, while Sun *et al.* (2017) optimizes the same loss using a second-order trust region method. In a recent paper, Wang *et al.* (2017) instead minimize the highly nonsmooth function $\frac{1}{m}\sum_{i=1}^m (\langle a_i, x \rangle^2 - b_i)^2$ by a gradient descent-like method. Though the truncated Wirtinger flow algorithm in Chen & Candès (2017), the truncated amplitude flow in Wang *et al.* (2017) and the

Polyak subgradient algorithm investigated here all have similar guarantees, one possible advantage of the Polyak subgradient method is that it is completely parameter free. Another closely related recent work is that of Tan & Vershynin (2018). One can interpret their scheme as a stochastic subgradient method on the formulation $\frac{1}{m} \sum_{i=1}^{m} ||\langle a_i, x \rangle| - \sqrt{b_i}|$. Under proper initialization and assuming that $a_i$ are uniformly sampled from a sphere, they prove linear convergence. Their argument relies on sophisticated probabilistic tools. In contrast, we disentangle the probabilistic statements (weak convexity and sharpness) from the deterministic convergence of Algorithm 1. As a proof of concept, we illustrate the proposed subgradient method on synthetic and large-scale real image recovery problems.

It is worth emphasizing that the Polyak subgradient method, which we investigate here, heavily relies on the measurements being exact. This is in contrast to the more involved prox-linear algorithm in Duchi & Ruan (2018a), which is applicable even when the measurements are corrupted by gross outliers. The recent paper (Davis *et al.*, 2018) provides a partial extension of the subgradient method to the noisy setting, based on a geometrically decaying step size. The resulting scheme, however, requires having estimates on the sharpness and weak convexity constants, and therefore is not completely parameter-free.

Weak convexity and sharpness, taken together, imply existence of a small neighbourhood $\mathcal{X}$ of $\{\pm \bar{x}\}$ devoid of extraneous stationary points of $f_S$ (see Lemma 3.1). On the other hand, it is intriguing to determine where the objective function $f_S$ may have stationary points outside of this neighbourhood. We complete the paper by proving that as the number of Gaussian measurements increases, the stationary points of the problem converge to a codimension two set, at a controlled rate. This suggests that there are much larger regions than the neighbourhood $\mathcal{X}$, where the objective function has benign geometry.

We follow an intuitive and transparent strategy, in line with previous work on smooth formulations of phase retrieval (Sun *et al.*, 2018). Setting the groundwork, assume that $a_i$ are i.i.d samples from a normal distribution $N(0, I_{d \times d})$. Hence, the problem $\min f_S$ is an empirical average approximation of the population objective

$$\min_x f_P(x) := \mathbb{E}_a[|(a^T x)^2 - (a^T \bar{x})^2|].$$

Seeking to determine the location of stationary points of $f_S$, we begin by first determining the stationary points of $f_P$. We base our analysis on the elementary observation that $f_P(x)$ depends on $x$ only through the eigenvalues of the rank two matrix $X := xx^T - \bar{x}\bar{x}^T$. More precisely, equality holds:

$$f_P(x) = \frac{4}{\pi} \left[ \text{Tr}(X) \cdot \arctan\left( \sqrt{\left| \frac{\lambda_{\max}(X)}{\lambda_{\min}(X)} \right|} \right) + \sqrt{|\lambda_{\max}(X)\lambda_{\min}(X)|} \right] - \text{Tr}(X).$$

See Fig. 1 for a graphical illustration.

Using basic perturbation properties of eigenvalues, we will show that the stationary points of $f_P$ are precisely

$$\{0\} \cup \{\pm \bar{x}\} \cup \{x \in \bar{x}^\perp : \|x\| = c \cdot \|\bar{x}\|\}, \tag{1.1}$$

where $c \approx 0.4416$ is a numerical constant. Intuitively, this region, excluding $\{\pm \bar{x}\}$, is where numerical methods may stagnate. In particular, $f_P$ has no extraneous stationary points outside of the subspace $\bar{x}^\perp$. Along the way, we prove a number of results in matrix theory, which may be of independent interest. For example, we show that all stationary points of a composition of an orthogonally invariant gauge function with the map $x \mapsto xx^T - \bar{x}\bar{x}^T$ must be either perpendicular or collinear with $\bar{x}$. We note that the smooth formulation of the phase retrieval problem has a similar landscape of stationary points (Sun *et al.*, 2018, Theorem 2.1), though *a priori* there is no reason for this to be the case.
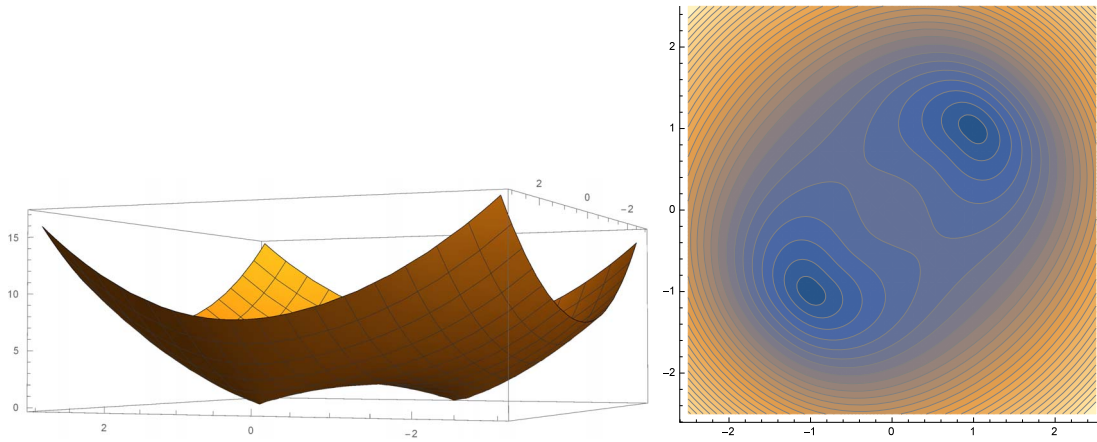
Fig. 1. Depiction of the population objective $f_P$ with $\bar{x} = (1, 1)$: graph (left), contours (right).

Having located the stationary points of the population objective $f_P$, we turn to the stationary points of the subsampled function $f_S$. This is where the techniques commonly used for smooth formulations of the problem, such as those in Sun *et al.* (2018), are no longer applicable; indeed, the subdifferential $\partial f_P(x)$ is usually a very poor approximation of $\partial f_S(x)$. Nonetheless, we show that the *graphs* of the subdifferentials $\partial f_P$ and $\partial f_S$ are close with high probability—a result closely related to the celebrated Attouch *et al.* (1990)'s convergence theorem. The analysis of the stationary points of the subsampled objective flows from there. Namely, we show that there is a constant $C$ such that whenever $m \geq Cd$, all stationary points $x$ of $f_S$ satisfy

$$\frac{\|x\|\|x - \bar{x}\|\|x + \bar{x}\|}{\|\bar{x}\|^3} \lesssim \sqrt[4]{\frac{d}{m}} \quad \text{or} \quad \left\{ \begin{array}{c} \left| \frac{\|x\|}{\|\bar{x}\|} - c \right| \lesssim \sqrt[4]{\frac{d}{m}} \cdot \left(1 + \frac{\|\bar{x}\|}{\|x\|}\right) \\ \frac{|\langle x,\bar{x}\rangle|}{\|x\|\|\bar{x}\|} \lesssim \sqrt[4]{\frac{d}{m}} \cdot \frac{\|\bar{x}\|}{\|x\|} \end{array} \right\},$$

with high probability; compare with (1.1). The argument we present is very general, relying only on weak convexity and concentration of $f_S$ around its mean. Therefore, we believe that the technique may be of independent interest.

We comment in Section B.1 on the structure of stationary points for the variant of the phase retrieval problem, in which the measurements $b$ are corrupted by gross outliers. It is straightforward to obtain a full characterization of the stationary points of the population objective using the techniques developed in earlier sections.

The outline for the paper is as follows. Section 2 summarizes notation and basic results we will need. In Section 3, we analyse the linear convergence of the Polyak subgradient method for a class of nonsmooth, nonconvex functions, which includes the subsampled objective $f_S$. In Section 4, we perform a few proof-of-concept experiments, illustrating the performance of the Polyak subgradient method on synthetic and real large-scale image recovery problems. Section 5 is devoted to characterizing the nonsmooth landscape of the population objective $f_P$. In Section 6, we develop a concentration theorem for the subdifferential graphs of $f_S$ and $f_P$, and briefly comment on robust extensions.

## 2. Notation

Throughout, we mostly follow standard notation. The symbol $\mathbb{R}$ will denote the real line, while $\mathbb{R}_+$ and $\mathbb{R}_{++}$ will denote non-negative and strictly positive real numbers, respectively. We always endow $\mathbb{R}^d$ with the dot product $\langle x, y \rangle = x^T y$ and the induced norm $\|x\| := \sqrt{\langle x, x \rangle}$. The symbol $\mathbb{S}^{d-1}$ will denote the unit sphere in $\mathbb{R}^d$, while $B(x, r) := \{y : \|x - y\| < r\}$ will stand for the open ball around $x$ of radius $r > 0$. For any set $Q \subset \mathbb{R}^d$, the distance function is defined by $\operatorname{dist}(x; Q) := \inf_{y \in Q} \|y - x\|$. The adjoint of a linear map $A \colon \mathbb{R}^d \to \mathbb{R}^m$ will be written as $A^* \colon \mathbb{R}^m \to \mathbb{R}^d$.

Since the main optimization problem we consider is nonsmooth, we will use some basic generalized derivative constructions. For a more detailed discussion, see for example the monographs of Mordukhovich (2006) and Rockafellar & Wets (1998).

Consider a function $f \colon \mathbb{R}^d \to \mathbb{R}$ and a point $\bar{x}$. The *Fréchet subdifferential* of $f$ at $\bar{x}$, denoted $\hat{\partial} f(\bar{x})$, is the set of all vectors $v \in \mathbb{R}^d$ satisfying

$$f(x) \geq f(\bar{x}) + \langle v, x - \bar{x} \rangle + o(\|x - \bar{x}\|) \qquad \text{as } x \to \bar{x}.$$

Thus, $v$ lies in $\hat{\partial} f(\bar{x})$ if and only if the affine function $x \mapsto f(\bar{x}) + \langle v, x - \bar{x} \rangle$ minorizes $f$ near $\bar{x}$ up to first-order. Since the assignment $x \mapsto \hat{\partial} f(x)$ may have poor continuity properties, it is useful to extend the definition slightly. The *limiting subdifferential* of $f$ at $\bar{x}$, denoted $\partial f(\bar{x})$, consists of all vectors $v \in \mathbb{R}^d$ such that there exist sequences $x_i$ and $v_i \in \hat{\partial} f(x_i)$ satisfying $(x_i, f(x_i), v_i) \to (\bar{x}, f(\bar{x}), v)$. We say that $\bar{x}$ is *stationary* for $f$ if the inclusion $0 \in \partial f(\bar{x})$ holds. The *graph* of $\partial f$ is the set

$$\operatorname{gph} \partial f := \{(x, y) \in \mathbb{R}^d \times \mathbb{R}^d : y \in \partial f(x)\}.$$

For essentially all functions that we will encounter, the two subdifferentials, $\hat{\partial} f(\bar{x})$ and $\partial f(\bar{x})$, coincide. This is the case for $C^1$-smooth functions $f$, where $\hat{\partial} f(\bar{x})$ and $\partial f(\bar{x})$ consist only of the gradient $\nabla f(\bar{x})$. Similarly for convex function $f$, both subdifferentials reduce to the subdifferential in the sense of convex analysis:

$$v \in \partial f(\bar{x}) \qquad \Longleftrightarrow \qquad f(x) \geq f(\bar{x}) + \langle v, x - \bar{x} \rangle \qquad \text{for all } x \in \mathbb{R}^d.$$

Most of the nonsmooth functions we will encounter have a simple composite form:

$$F(x) := h(c(x)),$$

where $h \colon \mathbb{R}^m \to \mathbb{R}$ is a finite convex function and $c \colon \mathbb{R}^d \to \mathbb{R}^n$ is a $C^1$-smooth map. For such composite functions, the two subdifferentials coincide, and admit the intuitive chain rule (Rockafellar & Wets, 1998, Theorem 10.6, Corollary 10.9):

$$\partial F(x) = \nabla c(x)^* \partial h(c(x)) \qquad \text{for all } x \in \mathbb{R}^d.$$

Here, $\nabla c(x)^*$ denotes the adjoint of the Jacobian matrix $\nabla c(x)$.

A function $f \colon \mathbb{R}^d \to \mathbb{R}$ is called *$\rho$-weakly convex* if $f + \frac{\rho}{2} \| \cdot \|^2$ is a convex function. It follows immediately from Rockafellar & Wets (1998, Theorem 12.17) that a lower semicontinuous function $f$

is $\rho$-weakly convex if and only if the inequality

$$f(y) \geq f(x) + \langle v, y - x \rangle - \frac{\rho}{2} \|y - x\|^2$$

holds for all points $x, y \in \mathbb{R}^d$ and vectors $v \in \partial f(x)$.

Finally, we will often use implicitly the observation that the Lipschitz constant of any lower semicontinuous function $f$ on a convex open set $U$ coincides with $\sup\{\|\zeta\| : x \in U, \zeta \in \partial f(x)\}$; see e.g. Rockafellar & Wets (1998, Theorem 9.13).

## 3. Subgradient method

In this work, we consider the robust formulation of the (real) phase retrieval problem. Setting the stage, suppose we are given vectors $\{a_i\}_{i=1}^m$ in $\mathbb{R}^d$ and measurements $b := \langle a_i, \bar{x} \rangle^2$, for a fixed but unknown vector $\bar{x}$. The goal of the phase retrieval problem is to recover the vector $\bar{x} \in \mathbb{R}^d$, up to a sign flip. The formulation of the problem we consider in this work is

$$\min_x f_S(x) := \frac{1}{m} \sum_{i=1}^m |\langle a_i, x \rangle^2 - b_i|.$$

The function $f_S$ (in contrast to other possible formulations) has a number of desirable properties, which we will highlight as we continue.

In this section, we show that the landscape of the phase retrieval objective $f_S$ favourably lends itself to classical subgradient methods. Namely, with proper initialization and under appropriate statistical assumptions, the Poljak (1967) subgradient method linearly converges to $\pm x$.

### 3.1 Subgradient method for weakly convex and sharp functions

The linear convergence guarantees that we present are mostly independent of the structure of $f_S$ and instead rely only on a few general regularity properties, which $f_S$ satisfies under mild statistical assumptions. Consequently, it will help the exposition in the current section to abstract away from $f_S$.

Assumption A. Fix a function $g: \mathbb{R}^d \to \mathbb{R}$ such that there exist real $\rho, \mu > 0$ satisfying the following two properties.

1. **Weak convexity.** The function $g + \frac{\rho}{2}\| \cdot \|^2$ is convex.

2. **Sharpness.** The inequality holds:

$$g(x) - \min g \geq \mu \cdot \text{dist}(x; \mathcal{X}) \qquad \text{for all } x \in \mathbb{R}^d,$$

where $\mathcal{X} \neq \emptyset$ is the set of minimizers of $g$.

Duchi & Ruan (2018a), following the work of Eldar & Mendelson (2014), showed that the robust phase retrieval loss $f_S(\cdot)$ satisfies Assumption A, under reasonable statistical assumptions. We will discuss these guarantees in Section 3.2, where we will instantiate the subgradient method on the robust phase retrieval objective. Consider now the standard Polyak subgradient method applied to $g$ (Algorithm 1).

---

**Algorithm 1** Polyak subgradient method

---

**Data**: $x_0 \in \mathbb{R}^d$
**Step** $k$: $(k \geq 1)$
Choose $\zeta_k \in \partial g(x_k)$.

  **if** $\zeta_k \neq 0$ **then**
    |  Set $x_{k+1} = x_k - \frac{g(x_k) - \min g}{\|\zeta_k\|^2} \zeta_k$.
  **else**
    |  Exit algorithm.
  **end**

---

As the first step in the analysis of Algorithm 1, we must ensure that there are no extraneous stationary points of $g$ near $\mathcal{X}$. This is the content of the following lemma.

LEMMA 3.1 (Neighbourhood with no stationary points). Suppose Assumption A holds. Then $g$ has no stationary points $x$ satisfying

$$0 < \mathrm{dist}(x; \mathcal{X}) < \frac{2\mu}{\rho}. \tag{3.1}$$

*Proof.* Consider a stationary point $x$ of $g$, which is outside of $\mathcal{X}$. Let $\bar{x} \in \mathcal{X}$ be a point satisfying $\|x - \bar{x}\| = \mathrm{dist}(x; \mathcal{X})$. Properties 1 and 2 then imply

$$\mu \cdot \mathrm{dist}(x; \mathcal{X}) \leq g(x) - g(\bar{x}) \leq \frac{\rho}{2} \|x - \bar{x}\|^2 = \frac{\rho}{2} \cdot \mathrm{dist}^2(x; \mathcal{X}),$$

where the second inequality follows from weak convexity of $g$ and the inclusion $0 \in \partial g(x)$. Dividing through by $\mathrm{dist}(x; \mathcal{X})$, the result follows. $\qquad \square$

The following Theorem 3.2—the main result of this subsection—shows that when Algorithm 1 is initialized within a certain tube $\mathcal{T}$ of $\mathcal{X}$, the iterates $x_k$ stay within the tube and converge linearly to $\mathcal{X}$. It is interesting to note that the rate of local linear convergence does not depend on the weak convexity constant $\rho$; indeed, the value $\rho$ only dictates the size of the tube $\mathcal{T}$.

THEOREM 3.2 (Linear rate). Suppose Assumption A holds. Fix a real $\gamma \in (0, 1)$ and define the tube

$$\mathcal{T} := \left\{ x \in \mathbb{R}^d : \mathrm{dist}(x; \mathcal{X}) \leq \gamma \cdot \frac{\mu}{\rho} \right\},$$

and the corresponding Lipschitz constant

$$L_g := \sup_{x \in \mathcal{T}, \zeta \in \partial g(x)} \|\zeta\|.$$

Then Algorithm 1 initialized at any point $x_0 \in \mathcal{T}$ produces iterates that converge $Q$-linearly to $\mathcal{X}$ at the rate:

$$\text{dist}^2(x_{k+1}; \mathcal{X}) \leq \left(1 - \frac{(1-\gamma)\mu^2}{L_g^2}\right) \text{dist}^2(x_k; \mathcal{X}). \tag{3.2}$$

*Proof.* We proceed by induction. Suppose that the theorem holds up to iteration $k$. We will prove the inequality (3.2). To this end, let $\bar{x} \in \mathcal{X}$ be a point satisfying $\|x_k - \bar{x}\| = \text{dist}(x_k; \mathcal{X})$. Note that if $x_k$ lies in $\mathcal{X}$, there is nothing to prove. Thus, we may suppose $x_k \notin \mathcal{X}$. Note that the inductive hypothesis implies $\text{dist}(x_k; S) \leq \text{dist}(x_0; S)$ and therefore $x_k$ lies in $\mathcal{T}$. Lemma 3.2 therefore guarantees $\zeta_k \neq 0$. Using Properties 1 and 2, we successively deduce

$$\begin{aligned}
\|x_{k+1} - \bar{x}\|^2 &= \|x_k - \bar{x}\|^2 + 2\langle x_k - \bar{x}, x_{k+1} - x_k \rangle + \|x_{k+1} - x_k\|^2 \\
&= \|x_k - \bar{x}\|^2 + \frac{2(g(x_k) - g(\bar{x}))}{\|\zeta_k\|^2} \cdot \langle \zeta_k, \bar{x} - x_k \rangle + \frac{(g(x_k) - g(\bar{x}))^2}{\|\zeta_k\|^2} \\
&\leq \|x_k - \bar{x}\|^2 + \frac{2(g(x_k) - g(\bar{x}))}{\|\zeta_k\|^2} \left(g(\bar{x}) - g(x_k) + \frac{\rho}{2}\|x_k - \bar{x}\|^2\right) + \frac{(g(x_k) - g(\bar{x}))^2}{\|\zeta_k\|^2} \\
&= \|x_k - \bar{x}\|^2 + \frac{(g(x_k) - g(\bar{x}))}{\|\zeta_k\|^2} \left(\rho\|x_k - \bar{x}\|^2 - (g(x_k) - g(\bar{x}))\right) \\
&\leq \|x_k - \bar{x}\|^2 + \frac{(g(x_k) - g(\bar{x}))}{\|\zeta_k\|^2} \left(\rho\|x_k - \bar{x}\|^2 - \mu\|x_k - \bar{x}\|\right) \\
&= \|x_k - \bar{x}\|^2 + \frac{\rho(g(x_k) - g(\bar{x}))}{\|\zeta_k\|^2} \left(\|x_k - \bar{x}\| - \frac{\mu}{\rho}\right) \|x_k - \bar{x}\|.
\end{aligned}$$

Combining the inclusion $x_k \in \mathcal{T}$ with sharpness (Assumption 2), we therefore deduce

$$\text{dist}^2(x_{k+1}; \mathcal{X}) \leq \|x_{k+1} - \bar{x}\|^2 \leq \left(1 - \frac{(1-\gamma)\mu^2}{\|\zeta_k\|^2}\right) \|x_k - \bar{x}\|^2.$$

The result follows. □

### 3.2 *Convergence for the phase retrieval objective*

We now turn to an application of Theorem 3.2 to the phase retrieval loss $f_S$. In particular, to run the subgradient method, we must only compute a subgradient of $f_S$, which can be easily done using the chain rule:

$$\frac{1}{m}\sum_{i=1}^{m} 2\langle a_i, x \rangle \cdot \text{sign}(\langle a_i, x \rangle^2 - b_i)a_i \in \partial f_S(x).$$

Each iteration of Algorithm 1 thus requires a single pass through the set of measurement vectors. We will see momentarily that under mild statistical assumptions, $\{\pm\bar{x}\}$ are the unique minimizers of $f_S$, as soon as $m > 2d$.

Thus, for a successful application of Theorem 3.2, we must only address the following questions:

(i) Describe the statistical conditions on the data generating mechanism, which insure that Assumption A holds with high probability.

(ii) Estimate the Lipschitz constant of $f_S$ on the union of balls $\mathcal{T} = B\left(\bar{x}, \frac{\gamma\mu}{\rho}\right) \cup B\left(-\bar{x}, \frac{\gamma\mu}{\rho}\right)$.

(iii) Describe a good initialization procedure for producing $x_0 \in \mathcal{T}$.

Essentially all of these points follow from the work of Duchi & Ruan (2017), Eldar & Mendelson (2014) and Wang *et al.* (2016). We summarize them here for the sake of completeness. Henceforth, let us suppose that $a_i \in \mathbb{R}^d$ (for $i = 1, \ldots, m$) are independent realizations of a random vector $a \in \mathbb{R}^d$.

3.2.1 *Sharpness.* In order to ensure sharpness (or rather the stronger 'stability' property, Eldar & Mendelson, 2014), we make the following assumption on the distribution of $a$.

ASSUMPTION B. There exist constants $\kappa_{\mathrm{st}}^*, p_0 > 0$ such that for all $u, v \in \mathbb{S}^{d-1}$, we have

$$\mathbb{P}(|\langle a, v \rangle \langle a, u \rangle| \geq \kappa_{\mathrm{st}}^*) \geq p_0.$$

Roughly speaking, this mild assumption simply says that the random vector $a$ has sufficient support in all directions. In particular, the standard Gaussian $a \sim \mathsf{N}(0, I_d)$ satisfies Assumption B with $\kappa_{\mathrm{st}}^* = 0.365$ and $p_0 = 0.25$; see Duchi & Ruan (2017, Example 1). The following is proved in Duchi & Ruan (2017, Corollary 3.1).

THEOREM 3.3 (Sharpness). Suppose that Assumption B holds. Then there exists a numerical constant $c < \infty$ such that if $mp_0^2 \geq cd$, we have

$$\mathbb{P}\left(f_S(x) - f_S(\bar{x}) \geq \frac{1}{2}\kappa_{\mathrm{st}}^* p_0 \|x - \bar{x}\| \|x + \bar{x}\| \quad \text{for all } x \in \mathbb{R}^d\right) \geq 1 - 2\exp\left(-\frac{mp_0^2}{32}\right).$$

To simplify notation, set $\mathrm{dist}(x; \bar{x}) := \min\{\|x - \bar{x}\|, \|x + \bar{x}\|\}$. Thus, Assumption B implies, with high probability, that $f_S$ is sharp. Indeed, Theorem 3.3 directly implies that with high probability we have

$$\begin{aligned}
f_S(x) - f_S(\bar{x}) &\geq \frac{1}{2}\kappa_{\mathrm{st}}^* p_0 \|x - \bar{x}\| \|x + \bar{x}\| \\
&\geq \frac{1}{2}\kappa_{\mathrm{st}}^* p_0 \cdot \min\{\|x - \bar{x}\|, \|x + \bar{x}\|\} \cdot \max\{\|x - \bar{x}\|, \|x + \bar{x}\|\} \\
&\geq \frac{1}{2}\kappa_{\mathrm{st}}^* p_0 \|\bar{x}\| \cdot \mathrm{dist}(x; \bar{x}).
\end{aligned}$$

Thus, the sharpness condition in Assumption A holds for $g = f_S$ with $\mu = \frac{1}{2}\kappa_{\mathrm{st}}^* p_0 \|\bar{x}\|$.

3.2.2 *Weak convexity.* We next look at weak convexity of the objective $f_S$. We will need the following definition.

DEFINITION 3.4 A random vector $a \in \mathbb{R}^d$ is $\sigma^2$-*sub-Gaussian* if for all unit vectors $v \in \mathbb{S}^{d-1}$, we have

$$\mathbb{E}\left[\exp\left(\frac{\langle a, v \rangle^2}{\sigma^2}\right)\right] \leq e.$$

Assumption C. The random vector $a$ is $\sigma^2$-sub-Gaussian.

The following is a direct consequence of Duchi & Ruan (2017, Corollary 3.2).

THEOREM 3.5 (Weak convexity). Suppose that Assumption C holds. Then there exists a numerical constant $c < \infty$ such that whenever $m \geq cd$, the function $f_S$ is $4\sigma^2$-weakly convex, with probability at least $1 - \exp\left(-\frac{m}{c}\right)$.

*Proof.* This follows almost immediately from Duchi & Ruan (2017, Corollary 3.2). Define the separable function $h(z_1, \ldots, z_m) := \frac{1}{m} \sum_{i=1}^{m} |z_i|$ and the map $F \colon \mathbb{R}^d \to \mathbb{R}^m$ with the $i$th coordinate given by $F_i(x) := (a_i^T x)^2 - b_i$. Observe the equality $f_S(x) = h(F(x))$. Corollary 3.2 in Duchi & Ruan (2017) shows that there exists a numerical constant $c < \infty$ such that whenever $m \geq cd$, with probability at least $1 - \exp\left(-\frac{m}{c}\right)$, we have

$$f_S(y) \geq h(F(x) + \nabla F(x)(y-x)) - 2\sigma^2 \|y-x\|^2 \qquad \text{for all } x, y \in \mathbb{R}^d.$$

Since $h$ is convex, for any vector $v \in \partial h(F(x))$ we have

$$h(F(x) + \nabla F(x)(y-x)) \geq h(F(x)) + \langle v, \nabla F(x)(y-x) \rangle = f_S(x) + \langle \nabla F(x)^* v, y-x \rangle.$$

Taking into account the equality $\partial f_S(x) = \nabla F(x)^* \partial h(F(x))$, we conclude that $f_S$ is $4\sigma^2$-weakly convex. $\qquad\square$

3.2.3 *Lipschitz constant on a ball.* Let us next estimate the Lipschitz constant of $f_S$ on a ball of a fixed radius. To this end, observe the chain of inequalities

$$
\begin{aligned}
|f_S(x) - f_S(y)| &\leq \frac{1}{m} \sum_{i=1}^{m} \left| |\langle a_i, x \rangle^2 - \langle a_i, \bar{x} \rangle^2| - |\langle a_i, y \rangle^2 - \langle a_i, \bar{x} \rangle^2| \right| \\
&\leq \frac{1}{m} \sum_{i=1}^{m} |\langle a_i, x \rangle^2 - \langle a_i, y \rangle^2| \\
&= \|x-y\|\|x+y\| \cdot \frac{1}{m} \sum_{i=1}^{m} |\langle a_i, v \rangle \langle a_i, w \rangle|,
\end{aligned}
\tag{3.3}
$$

where we set $v := \frac{x-y}{\|x-y\|}$ and $w := \frac{x+y}{\|x+y\|}$. Thus, we would like to upper bound the term $\frac{1}{m} \sum_{i=1}^{m} |\langle a_i, v \rangle \langle a_i, w \rangle|$ by a numerical constant, with high probability. Assumptions C quickly yields such a bound.

COROLLARY 3.6 (Lipschitz constant on a ball). Suppose that Assumptions C holds. Then there exist a numerical constant $c > 0$ such that whenever $m \geq cd$, with probability at least $1 - \exp\left(-\frac{m}{c}\right)$, we have

$$|f_S(x) - f_S(y)| \leq 2\sigma^2 \|x-y\|\|x+y\| \qquad \text{for all } x, y \in \mathbb{R}^d, \tag{3.4}$$

and consequently

$$\max_{\zeta \in \partial f_S(x)} \|\zeta\| \leq 4\sigma^2 \|x\| \qquad \text{for all } x \in \mathbb{R}^d. \tag{3.5}$$

*Proof.*   Observe first that for any unit vectors $v, w \in \mathbb{R}^d$, the inequalities hold:

$$\frac{1}{m} \sum_{i=1}^{m} |\langle a_i, v \rangle \langle a_i, w \rangle| \leq \frac{1}{2m} \sum_{i=1}^{m} \langle a_i, v \rangle^2 + \frac{1}{2m} \sum_{i=1}^{m} \langle a_i, w \rangle^2 \leq \left\| \frac{1}{m} \sum_{i=1}^{m} a_i a_i^T \right\|_{op}.$$

By matrix concentration (see e.g. Vershynin, 2010, Thm. 5.39 or Duchi & Ruan, 2017, Lemma 3.1), there exists a numerical constant $c$ such that whenever $m \geq cd$, we have

$$\left\| \frac{1}{m} \sum_{i=1}^{m} a_i a_i^T \right\|_{op} \leq 2\sigma^2$$

with probability at least $1 - \exp\left(-\frac{m}{c}\right)$. Taking into account the estimate (3.3), we deduce

$$|f_S(x) - f_S(y)| \leq 2\sigma^2 \|x - y\| \|x + y\|$$

with high probability, as claimed. Consequently, for any point $z \in \mathbb{R}^d$, the triangle inequality leads to the following bound on the Lipschitz constant of $f_S$:

$$\limsup_{x,y \to z} \frac{|f_S(x) - f_S(y)|}{\|x - y\|} \leq 4\sigma^2 \|z\|.$$

Since the Lipschitz constant of $f_S$ at $x$ coincides with the value $\max_{\zeta \in \partial f(x)} \|\zeta\|$ (see e.g. Rockafellar & Wets, 1998, Theorem 9.13), the estimate (3.5) follows.                                                    □

We now have all the ingredients in place to apply Theorem 3.2 with $\gamma := 1/2$ to the robust phase retrieval objective. Namely, under Assumptions B and C, we may set[1]

$$\rho := 4\sigma^2; \qquad\qquad \mu := \tfrac{1}{2} \kappa_{st}^* p_0 \|\bar{x}\|; \qquad\qquad L_g := 4\sigma^2 \left( 1 + \frac{\kappa_{st}^* p_0}{16\sigma^2} \right) \|\bar{x}\|. \qquad (3.6)$$

Thus, we have proved the following convergence guarantee—the main result of this section. To simplify the formulas, we apply Theorem 3.2 only with $\gamma := 1/2$.

COROLLARY 3.7 (Linear convergence for phase retrieval). Suppose that Assumptions B and C hold. Then there exists a numerical constant $c < \infty$ such that the following is true. Whenever we are in the regime, $\frac{c}{p_0^2} \leq \frac{m}{d}$, and we initialize Algorithm 1 at $x_0$ satisfying

$$\min \left\{ \frac{\|x_0 - \bar{x}\|}{\|\bar{x}\|}, \frac{\|x_0 + \bar{x}\|}{\|\bar{x}\|} \right\} \leq \frac{\kappa_{st}^* p_0}{16\sigma^2}, \qquad (3.7)$$

---

[1] The definition of $L_g$ uses the norm of any point in the tube $\mathcal{T} = B(\bar{x}, \frac{\gamma\mu}{\rho}) \cup B(-\bar{x}, \frac{\gamma\mu}{\rho})$ is clearly upper bounded by $\|\bar{x}\| + \frac{\mu}{2\rho}$.

we can be sure with probability at least

$$1 - 4\exp\left(-m \cdot \min\left\{\frac{p_0^2}{32}, c^{-1}\right\}\right)$$

that the produced iterates $\{x_k\}$ converge to $\{\pm\bar{x}\}$ at the linear rate:

$$\text{dist}^2(x_{k+1}; \bar{x}) \leq \left(1 - \frac{1}{128}\left(\frac{16p_0\kappa_{\text{st}}^*}{16\sigma^2 + \kappa_{\text{st}}^* p_0}\right)^2\right)\text{dist}^2(x_k; \bar{x}). \tag{3.8}$$

Aside from numerical constants, the linear rate depends only on $\kappa_{\text{st}}^*$, $p_0$ and $\sigma$.

Thus, under typical statistical assumptions, the subgradient method converges linearly to $\{\pm\bar{x}\}$, as long as one can initialize the method at a point $x_0$ satisfying the relative error condition $\|x_0 \pm \bar{x}\| \leq R\|\bar{x}\|$, where $R$ is a constant. A number of authors have proposed initialization strategies that can achieve this guarantee using only a constant multiple of $d$ measurements (Candès *et al.*, 2015; Wang *et al.*, 2017; Zhang *et al.*, 2016; Duchi & Ruan, 2017; Tan & Vershynin, 2018). For completeness, we record the strategy that was proposed in Wang *et al.* (2017), and rigorously justified in Duchi & Ruan (2017). To simplify the exposition, we only state the guarantees of the initialization under Gaussian assumptions on the measurement vectors $a_i$.

THEOREM 3.8 (Duchi & Ruan, 2017, Equation (15)). Assume that $a_i \sim \mathsf{N}(0, I_d)$ are i.i.d. standard Gaussian. Define the value $\hat{r}^2 := \frac{1}{m}\sum_{i=1}^m b_i$ and the index set $\mathcal{I}_{\text{sel}} := \left\{i \in [m] \mid b_i \leq \frac{1}{2}\hat{r}^2\right\}$. Set

$$X^{\text{init}} := \sum_{i \in \mathcal{I}_{\text{sel}}} a_i a_i^T \qquad \text{and} \qquad \hat{w} := \underset{w \in \mathbb{S}^{d-1}}{\text{argmin}}\ w^T X^{\text{init}} w.$$

Then as soon as $\frac{m}{d} \gtrsim \varepsilon^{-2}$ the point $x_0 = \hat{r}\hat{w}$ satisfies

$$\min\left\{\frac{\|x_0 - \bar{x}\|}{\|\bar{x}\|}, \frac{\|x_0 + \bar{x}\|}{\|\bar{x}\|}\right\} \lesssim \varepsilon \log\frac{1}{\varepsilon}$$

with probability at least $\geq 1 - 5\exp(-cm\varepsilon^2)$, where $c$ is a numerical constant.

For more details and intuition underlying the initialization procedure, see Duchi & Ruan (2017, Section 3.3).

## 4. Numerical illustration

In this section, as a proof of concept, we apply the subgradient method to medium and large-scale phase retrieval problems. All of our experiments were performed on a standard desktop: Intel(R) Core(TM) i7-4770 CPU3.40 GHz with 8.00 GB RAM.

We begin with simulated data. Set $d = 5000$. We generated a standard Gaussian random matrix $A \in \mathbb{R}^{m \times d}$ for each value $m \in \{12225, 13500, 14750, 16000, 17250, 18500\}$; afterwards, we generated a Gaussian vector $\bar{x} \sim \mathsf{N}(0, I_d)$ and set $b = (A\bar{x})^2$. We then applied the initialization procedure, detailed in Theorem 3.8, followed by the subgradient method. Figure 4 plots the progress of the iterates produced
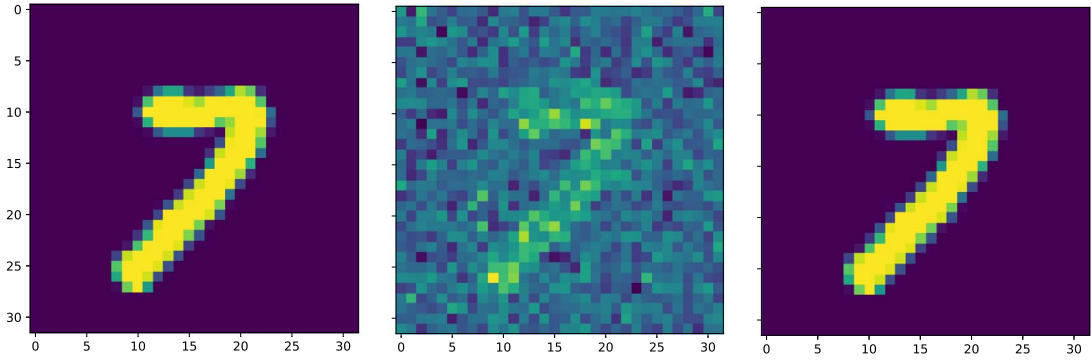
FIG. 2. Convergence plot for the experiment on simulated data (iteration vs. $\|x_k - \bar{x}\| / \|\bar{x}\|$).
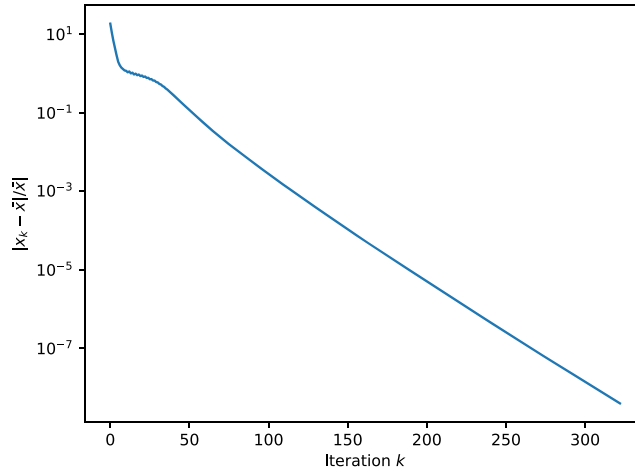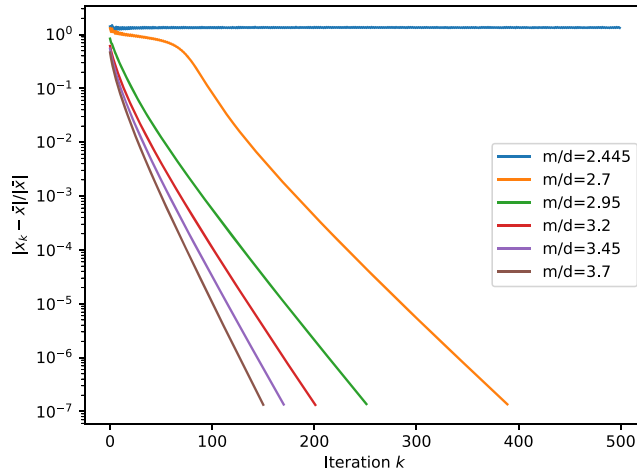


FIG. 3. Digit recovery; left is the true digit, middle is the initial, right is the digit produced by the subgradient method. Dimension of the problem: $(n, d, m) = (32, 3072, 9216)$.

by the subgradient method in each of the seven experiments. The top curve corresponds to $m = 12{,}225$, the bottom curve corresponds to $m = 18\,500$, while the curves for the other values of $m$ interpolate in between. The iterates corresponding to $m = 12\,225$ stagnate; evidently, the number of measurements is too small. Indeed, the iterates do not even converge to a stationary point of the problem; this is in contrast to the prox-linear method in Duchi & Ruan (2017). The iterates for the rest of the experiments converge to the true signal $\pm\bar{x}$ at an impressive linear rate. In our repeated experiments (with $m$ and $d$ fixed), we found that when the subgradient method succeeds, the performance curves, as in Fig. 4, concentrate tightly around their mean.

In our second experiment, we use digit images from the MNIST data set (Lecun *et al.*, 1998); these are relatively small so that the measurement matrices can be stored in memory. We illustrate the generic behaviour of the algorithm on digit seven in Fig. 2. The dimensions of the image we use are $32 \times 32$ (with three RGB channels). Hence, after vectorizing the dimension of the variable is $d = 3072$, while the number of Gaussian measurements is $m = 3d = 9216$. The initialization produced appears to be

FIG. 4. Convergence plot on MNIST digit (iterates vs. $\|x_k - \bar{x}\| / \|\bar{x}\|$).

reasonable; the digit is visually discernible. The true image and the final image produced by the method are essentially identical. The convergence plot appears in Fig. 3.

We next apply the subgradient method for recovering large-scale real images. To allow an easy comparison with previous work, we generate the data using the same process as in Duchi & Ruan (2017, Section 6.3). We first describe how we generate the operator $A$. To this end, let $H \in \{-1, 1\}^{l \times l} / \sqrt{l}$ be a symmetric normalized Hadamard matrix. Consequently, $H$ satisfies the equation $H^2 = I_l$. Note that by the virtue of being Hadamard, matrix vector multiplication $Hv$ requires time $l \log(l)$. For some integer $k$, we then generate $k$ i.i.d. diagonal sign matrices $S_1, \ldots, S_k \in \text{diag}(\{-1, 1\}^l)$ uniformly at random, and define $A = \begin{bmatrix} HS_1 & HS_2 & \ldots & HS_k \end{bmatrix}^T \in \mathbb{R}^{kl \times l}$.

We work with square coloured images, represented as an array $\overline{X} \in \mathbb{R}^{n \times n \times 3}$. The number 3 appears because coloured images have three RGB channels. We then stretch the matrix $\overline{X}$ into a $3n^2$-dimensional vector $\bar{x}$ and set the measurements $b_i := (A(i, \cdot)\bar{x})^2$, where $A(i, \cdot)$ denotes the $i$th row of $A$. Thus, if the image is $n \times n$, the number of variables in the problem formulation is $d := 3n^2$ and the number of measurements is $m := kd = 3kn^2$. We use the initialization procedure proposed in Theorem 3.8, with a standard power method (with a shift) to find the minimal eigenvalue of $X^{\text{init}}$. We complete the experiment by running the subgradient method (Algorithm 1), which requires no parameter tuning.

We perform a large scale experiment on two pictures taken by the Hubble telescope. Figure 5 describes the results of the experiment, while Fig. 6 plots the iterate progress. The image on the left is $1024 \times 1024$ and we use $k = 3$ Hadamard matrices. Hence, the dimensions of the problem are $d \approx 2^{22}$ and $m = 3d \approx 2^{24}$. The image on the right is $2048 \times 2048$ and we use $k = 3$ Hadamard matrices. Hence the dimensions of the problem are $d \approx 2^{24}$ and $m = 3d \approx 2^{25}$. For the image on the left, the entire experiment, including initialization and the subgradient method completed in 3 min. For the image on the right, it completed in 25.6 min. The vast majority of time was taken up by the initialization. Thus, a more careful implementation and/or tuning of the initialization procedure could speed up the experiment.

FIG. 5. Image recovery; top row are the true images, bottom row are the images produced by the subgradient method. We do not record the images produced by the initialization as they were both completely black. Dimensions of the problem: $(n, k, d, m) \approx (1024, 3, 2^{22}, 2^{24})$ (left) and $(n, k, d, m) \approx (2048, 3, 2^{24}, 2^{25})$ (right).

## 5. Nonsmooth landscape of the robust phase retrieval

In this section, we pursue a finer analysis of the stationary points of the robust phase retrieval objective $f_S$. To motivate the discussion, recall that under Assumptions B and C, Lemma 3.1 shows that there are no extraneous stationary points $x$ satisfying

$$\min\left\{\frac{\|x - \bar{x}\|}{\|\bar{x}\|}, \frac{\|x + \bar{x}\|}{\|\bar{x}\|}\right\} < \frac{\kappa_{\mathrm{st}}^* p_0}{4\sigma^2}.$$
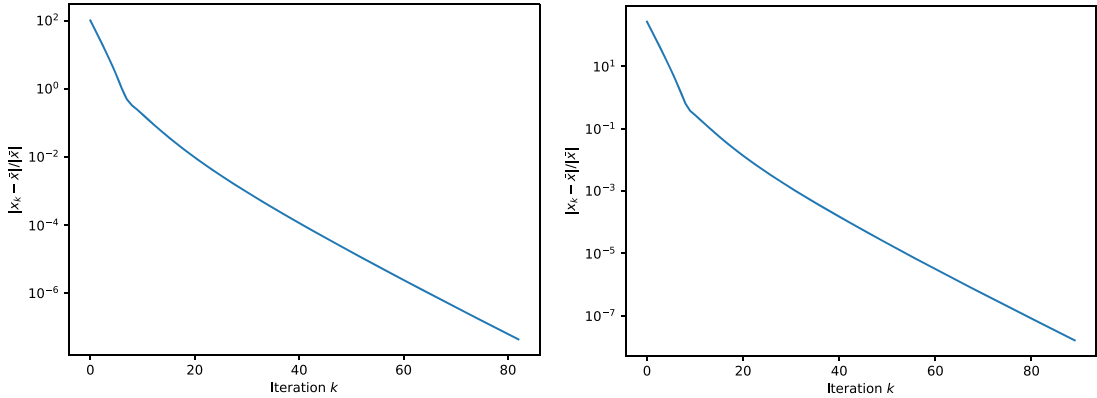
Fig. 6. Convergence plot on the two Hubble images (iterates vs. $\|x_k - \bar{x}\|/\|\bar{x}\|$).

This result is uninformative when $x$ is far away from $\bar{x}$ or when $x$ is close to the origin. Therefore, it is intriguing to determine the location of *all* the stationary points of $f_S$. In this section, we will see that under a Gaussian observation model, the stationary points of $f_S$ cluster around the codimension two set, $\{0, \pm\bar{x}\} \cup (\bar{x}^\perp \cap c \cdot \mathbb{S}^{d-1})$, where $c \approx 0.4416$ is a numerical constant.

We follow a common two-step strategy. First, we analyse the landscape of the population objective

$$f_P(x) := \mathbb{E}_a\left[|\langle a, x\rangle^2 - \langle a, \bar{x}\rangle^2|\right].$$

Indeed, we will see that under Gaussian assumptions, the set of stationary points of $f_P$ is precisely $\{0, \pm\bar{x}\} \cup (\bar{x}^\perp \cap c \cdot \mathbb{S}^{d-1})$. Secondly, we use uniform control on the error $|f_P(x) - f_S(x)|$ (see below) to prove that the subdifferential graphs of $f_P$ and $f_S$ are close with high probability, and therefore, that the stationary points of $f_S$ concentrate around those of $f_P$. This second step deviates from typical landscape arguments, which rely on *pointwise* concentration of gradients. In the nonsmooth setting, pointwise concentration is decidedly false, and we must instead focus on the distance between subdifferential graphs as a whole.

Before continuing, we record the following theorem, which is a special case of Eldar & Mendelson (2014, Theorem 2.8). This result immediately yields uniform control on the error $|f_S(x) - f_P(x)|$, as alluded to above.

THEOREM 5.1 (Concentration). Let $\{a_i\}_{i=1}^m$ be i.i.d. realizations of a normally distributed random vector, $a \sim N(0, I_d)$. Then there exist numerical constants $c_1, c_2, c_3 > 0$ so that with probability at least $1 - 2\exp(-c_2 c_1^2 \min\{m, d^2\})$ the inequality holds:

$$\sup_{v,w \in \mathbb{S}^{d-1}} \left| \frac{1}{m}\sum_{i=1}^m |\langle a_i, v\rangle\langle a_i, w\rangle| - \mathbb{E}_a[|\langle a, v\rangle\langle a, w\rangle|] \right| \leq c_1^3 c_3 \left( \sqrt{\frac{d}{m}} + \frac{d}{m} \right).$$

Rewriting the sampled and population objectives as

$$f_S(x) = \frac{1}{m} \sum_{i=1}^{m} |\langle a_i, x - \bar{x} \rangle \langle a_i, x + \bar{x} \rangle| \qquad \text{and} \qquad f_P(x) = \mathbb{E}_a \left[ |\langle a, x - \bar{x} \rangle \langle a, x + \bar{x} \rangle| \right].$$

and setting $v = \frac{x - \bar{x}}{\|x - \bar{x}\|}$ and $w = \frac{x + \bar{x}}{\|x + \bar{x}\|}$ in Theorem 5.1, directly implies

$$\left| f_S(x) - f_P(x) \right| \leq c_1^3 c_3 \left( \sqrt{\frac{d}{m}} + \frac{d}{m} \right) \|x - \bar{x}\| \|x + \bar{x}\| \qquad \text{for all } x \in \mathbb{R}^d$$

with high probability. Thus, the error $|f_S(x) - f_P(x)|$ is indeed well controlled over all points $x \in \mathbb{R}^d$.

## 5.1 *A matrix analysis interlude*

Before continuing, we introduce some basic matrix notation. We mostly follow Lewis (1996), Lewis (1999) and Drusvyatskiy & Paquette (2018). The symbol $\mathcal{S}^d$ will denote the Euclidean space of real symmetric $d \times d$-matrices with the trace inner product $\langle X, Y \rangle := \text{Tr}(XY)$. A function $f \colon \mathbb{R}^d \to \mathbb{R}$ is called *symmetric* if equality, $f(\sigma x) = f(x)$, holds for all coordinate permutations $\sigma$. For any symmetric function $f \colon \mathbb{R}^d \to \mathbb{R}$, we define the induced function on the symmetric matrices $f_\lambda \colon \mathcal{S}^d \to \mathbb{R}$ as the composition

$$f_\lambda(X) := f(\lambda(X)),$$

where $\lambda \colon \mathcal{S}^d \to \mathbb{R}^d$ assigns to each matrix $X \in \mathcal{S}^d$ its eigenvalues in nonincreasing order

$$\lambda_1(X) \geq \lambda_2(X) \geq \ldots \geq \lambda_n(X).$$

Note that $f$ coincides with the restriction of $f_\lambda$ to diagonal matrices, $f_\lambda(\text{Diag}(x)) = f(x)$. Any function on $\mathcal{S}^d$ that has the form $f_\lambda$ for some symmetric function $f$, is called *spectral*. Equivalently, spectral functions on $\mathcal{S}^d$ are precisely those that are invariant under conjugation by orthogonal matrices. Henceforth, let $\mathbb{O}^d$ be the set of real $d \times d$ orthogonal matrices.

Recall that two matrices $X, V \in \mathcal{S}^d$ commute if, and only if, they can be simultaneously diagonalized. When describing variational properties of convex spectral functions, a stronger notion is needed. We say that $X, V$ admit a *simultaneous ordered spectral decomposition* if there exists a matrix $U \in \mathbb{O}^d$ satisfying

$$UVU^T = \text{Diag}(\lambda(V)) \qquad \text{and} \qquad UXU^T = \text{Diag}(\lambda(X)).$$

Thus, the definition stipulates that $X$ and $V$ admit a simultaneous diagonalization, where the diagonals of the two diagonal matrices are simultaneously ordered.

The following is a foundational theorem in the convex analysis of spectral functions, due to Lewis (1996). An extension to the nonconvex setting was proved in Lewis (1999), while a much simplified argument was recently presented in Drusvyatskiy & Paquette (2018).

THEOREM 5.2 (Spectral convex analysis). Consider a symmetric function $f \colon \mathbb{R}^d \to \mathbb{R} \cup \{+\infty\}$. Then $f$ is convex if and only if $f_\lambda$ is convex. Moreover, if $f$ is convex, then the subdifferential $\partial f_\lambda(X)$ consists
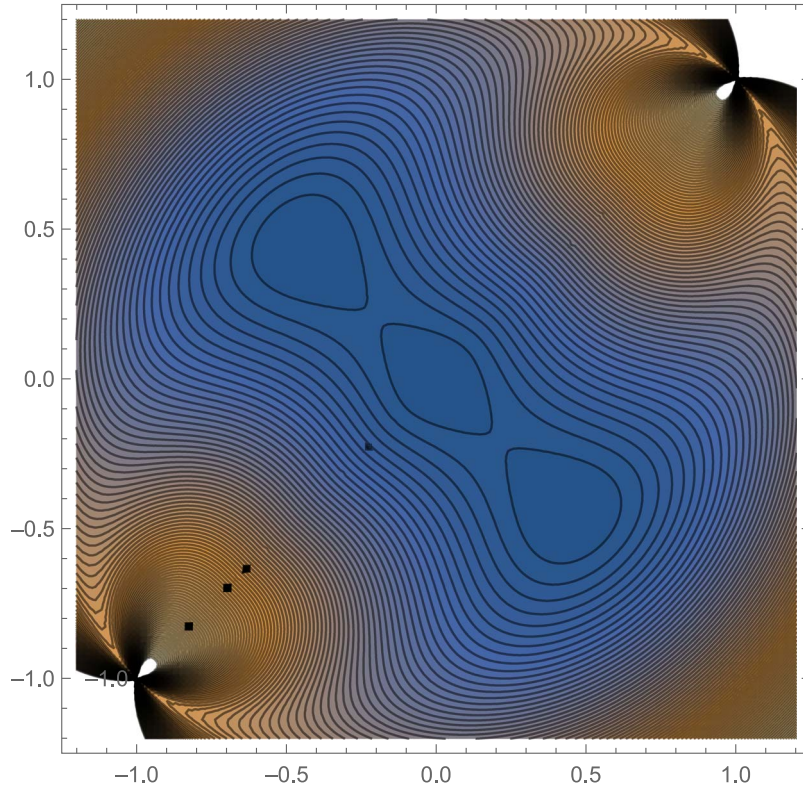
FIG. 7. The contour plot of the function $x \mapsto \|\nabla f_P(x)\|$, where $\bar{x} = (1, 1)$. The global minimizers of $f_P$ are $\pm\bar{x}$, while the three extraneous stationary points are $(0, 0)$ and $\pm c(-1, 1)$, where $c \approx 0.4416$.

of all matrices $V \in \mathcal{S}^d$ satisfying $\lambda(V) \in \partial f(\lambda(X))$ and such that $X$ and $V$ admit a simultaneous ordered spectral decomposition.

### 5.2 *Landscape of the population objective*

Henceforth, we fix a point $0 \neq \bar{x} \in \mathbb{R}^d$ and assume that $a \in \mathbb{R}^d$ is a normally distributed random vector $a \sim \mathsf{N}(0, I_d)$. In this section, we will investigate the population objective of the robust phase retrieval problem:

$$f_P(x) := \mathbb{E}_a \left[ |\langle a, x \rangle^2 - \langle a, \bar{x} \rangle^2| \right].$$

Our aim is to prove the following result; see Fig. 7 for a graphical depiction.

THEOREM 5.3 (Landscape of the population objective). The stationary points of the population objective $f_P$ are precisely

$$\{0\} \cup \{\pm\bar{x}\} \cup \{x \in \bar{x}^{\perp} : \|x\| = c \cdot \|\bar{x}\|\}, \tag{5.1}$$

where $c > 0$ (approx. $c \approx 0.4416$) is the unique solution of the equation $\frac{\pi}{4} = \frac{c}{1+c^2} + \arctan(c)$.

Theorem 5.3 provides an exact characterization of the stationary points of the population objective $f_P$. Looking ahead, when we will pass to the subsampled objective $f_S$ in Section 6, we will show that every stationary point of $f_S$ is *close* to an *approximately* stationary point of $f_P$. Therefore, it will be useful to have an extension of Theorem 5.3 that locates approximately stationary points of $f_P$. This is the content of the following theorem.

THEOREM 5.4 (Location of approximate stationary points). There exists a numerical constant $\gamma > 0$ such that the following holds. For any point $x \in \mathbb{R}^d$ with

$$\varepsilon := \mathrm{dist}(0; \partial f_P(x)) \leq \gamma \|x\|,$$

it must be the case that $\|x\| \lesssim \|\bar{x}\|$ and $x$ satisfies either

$$\|x\| \|x - \bar{x}\| \|x + \bar{x}\| \lesssim \varepsilon \|\bar{x}\|^2 \qquad \text{or} \qquad \left\{ \begin{array}{c} |\|x\| - c\|\bar{x}\|| \lesssim \varepsilon \dfrac{\|\bar{x}\|}{\|x\|} \\[2ex] |\langle x, \bar{x} \rangle| \lesssim \varepsilon \|\bar{x}\| \end{array} \right\},$$

where $c > 0$ is the unique solution of the equation $\frac{\pi}{4} = \frac{c}{1+c^2} + \arctan(c)$.

We present the proofs of Theorem 5.3 in Section 5.3, and defer the proof of Theorem 5.4 to the Appendix (Section B), as the latter requires a much more delicate argument. At their core, the arguments rely on the observation that the population objective $f_P(x)$ depends on the input vector $x$ only through the eigenvalues of the rank two matrix $xx^T - \bar{x}\bar{x}^T$. This observation was already implicitly used by Candès *et al.* (2013). Since this matrix will appear often in the arguments, we will use the symbol $X := xx^T - \bar{x}\bar{x}^T$ throughout. For ease of reference, we record the following simple observation: the matrix $X$ is typically indefinite.

LEMMA 5.5 (Eigenvalues of the rank two matrix). Suppose $x$ and $\bar{x}$ are not collinear. Then $X$ has exactly one strictly positive and one strictly negative eigenvalue.

*Proof.* Suppose the claim is false. Then either $X$ is positive semidefinite or negative semidefinite. Let us dispense with the first case. Observe $X \succeq 0$ if and only if $(x^T v)^2 - (\bar{x}^T v)^2 \geq 0$ for all $v$. Hence, if $X$ were positive semidefinite, we would deduce $x^\perp \subset \bar{x}^\perp$; that is, $x$ and $\bar{x}$ are collinear, a contradiction. The case $X \preceq 0$ is analogous. □

The following lemma, as we alluded to above, shows that $f_P(x)$ depends on $x$ only through the eigenvalues of the rank two matrix $X = xx^T - \bar{x}\bar{x}^T$.

LEMMA 5.6 (Spectral representation of the population objective). For all points $x \in \mathbb{R}^d$, equality holds:

$$f_P(x) = \mathbb{E}_v \left[ \left| \langle \lambda(X), v \rangle \right| \right], \tag{5.2}$$

where $v_i \in \mathbb{R}$ are i.i.d. chi-squared random variables $v_i \sim \chi_1^2$.

*Proof.* Observe the equalities:

$$f_P(x) = \mathbb{E}_a \left[ |\langle a, x \rangle^2 - \langle a, \bar{x} \rangle^2| \right] = \mathbb{E}_a[|\langle a, x - \bar{x} \rangle \langle a, x + \bar{x} \rangle|]$$

$$= \mathbb{E}_a[|(x - \bar{x})^T a a^T (x + \bar{x})|]$$

$$= \mathbb{E}_a \left[ \left| \mathrm{Tr} \left( a^T (x + \bar{x})(x - \bar{x})^T a \right) \right| \right].$$

Thus, in terms of the matrix $M := (x + \bar{x})(x - \bar{x})^T$, we have $f_P(x) = \mathbb{E}_a \left[ |\mathrm{Tr} \left( a^T M a \right)| \right]$. Taking into account the equalities $a^T M a = a^T \left( \frac{M + M^T}{2} \right) a = a^T X a$, we deduce

$$f_P(x) = \mathbb{E}_a \left[ \left| \mathrm{Tr} \left( a^T X a \right) \right| \right].$$

Form now an eigenvalue decomposition $X = U \mathrm{Diag}(\lambda(X)) U^T$, where $U \in \mathbb{R}^{d \times d}$ is an orthogonal matrix. Rotation invariance of the Gaussian distribution then implies

$$\mathbb{E}_a \left[ \left| \mathrm{Tr}(a^T X a) \right| \right] = \mathbb{E}_a \left[ \left| \mathrm{Tr}((Ua)^T X (Ua)) \right| \right] = \mathbb{E}_u \left[ \left| \sum_{i=1}^d \lambda_i(X) u_i^2 \right| \right],$$

where $u_i$ are i.i.d. standard normals. The result follows. $\qquad\square$

Thus, Lemma 5.6 shows that the population objective $f_P$ is a spectral function of $X$. Combined with Lemma 5.5, we deduce that there are two ways to rewrite the population objective in composite form:

$$f_P(x) = \varphi_\lambda(X) \qquad \text{and} \qquad f_P(x) = \zeta(\lambda_1(X), \lambda_d(X)),$$

where

$$\varphi(z) := \mathbb{E}_v \left[ \left| \langle z, v \rangle \right| \right] \qquad \text{and} \qquad \zeta(y_1, y_2) := \mathbb{E}_{v_1, v_2} \left[ |v_1 y_1 + v_2 y_2| \right]. \qquad (5.3)$$

Notice that $\varphi$ and $\zeta$ are norms on $\mathbb{R}^d$ and $\mathbb{R}^2$, respectively. It is instructive to compute $\zeta$ in closed form, yielding the following lemma. Since the proof is a straightforward computation, we have placed it in the appendix.

LEMMA 5.7 (Explicit representation of the outer function). Let $v_1, v_2 \sim \chi_1^2$ be i.i.d. chi-squared. Then for all real $(y_1, y_2) \in \mathbb{R}_+ \times \mathbb{R}_-$, equality holds:

$$\mathbb{E}_{v_1, v_2} \left[ |v_1 y_1 + v_2 y_2| \right] = \frac{4}{\pi} \left[ (y_1 + y_2) \arctan \left( \sqrt{-\frac{y_1}{y_2}} \right) + \sqrt{-y_1 y_2} \right] - (y_1 + y_2).$$

Thus, we have arrived at the following explicit representation of $f_P(x)$. Figure 1 in the introduction depicts the graph and the contours of the population objective.

COROLLARY 5.8 (Explicit representation of the population objective). The explicit representation holds:

$$f_P(x) = \frac{4}{\pi}\left[\text{Tr}(X) \cdot \arctan\left(\sqrt{\left|\frac{\lambda_{\max}(X)}{\lambda_{\min}(X)}\right|}\right) + \sqrt{|\lambda_{\max}(X)\lambda_{\min}(X)|}\right] - \text{Tr}(X).$$

### 5.3 Proof of Theorem 5.3

We next move on to the proof of Theorem 5.3. Let us first dispense with the easy implication, namely that every point in the set (5.1) is indeed stationary for $f_P$; in the process, we will see how the slope $c \approx 0.4416$ arises. Clearly, $\pm\bar{x}$ are minimizers of $f_P$ and are therefore stationary. The chain rule $\partial f_P(x) = 2\partial\varphi_\lambda(X)x$ implies that $x = 0$ is stationary as well. Fix now a point $x \in \bar{x}^\perp \setminus \{0\}$. Observe that the extremal eigenvalues of $X$ are

$$\lambda_1(X) = \|x\|^2 \qquad \text{and} \qquad \lambda_d(X) = -\|\bar{x}\|^2,$$

with corresponding eigenvectors

$$w_1 := \frac{x}{\|x\|} \qquad \text{and} \qquad w_d := \frac{\bar{x}}{\|\bar{x}\|}.$$

Since $\lambda_1(X)$ and $\lambda_d(X)$ each have multiplicity one, the individual eigenvalue functions $\lambda_1(\cdot)$ and $\lambda_d(\cdot)$ are smooth at $X$ with gradients

$$\nabla\lambda_1(X) = w_1 w_1^T \qquad \text{and} \qquad \nabla\lambda_d(X) = w_d w_d^T.$$

See for example Kato (1982, Theorem 5.11). Setting $(y_1, y_2) := (\|x\|^2, -\|\bar{x}\|^2)$ and applying the chain rule to the decomposition $f_P(x) = \zeta(\lambda_1(X), \lambda_d(X))$ shows

$$\nabla f_P(x) = 2\left(\nabla_{y_1}\zeta(y_1, y_2)w_1 w_1^T + \nabla_{y_2}\zeta(y_1, y_2)w_d w_d^T\right)x = 2\nabla_{y_1}\zeta(y_1, y_2)x.$$

Thus, a point $x \in \bar{x}^\perp \setminus \{0\}$ is stationary for $f_P$ if and only if the partial derivative $\nabla_{y_1}\zeta(y_1, y_2)$ vanishes. The points $(y_1, y_2)$ satisfying the equation $0 = \nabla_{y_1}\zeta(y_1, y_2)$ trace out exactly the line depicted in Fig. 8.

LEMMA 5.9 The solutions of the equation $0 = \nabla_{y_1}\zeta(y_1, y_2)$ on $\mathbb{R}_{++} \times \mathbb{R}_{--}$ are precisely the tuples $\{(c^2 y, -y)\}_{y>0}$, where $c > 0$ is the unique solution of the equation

$$\frac{\pi}{4} = \frac{c}{1 + c^2} + \arctan(c).$$

Note $c \approx 0.4416$.

*Proof.* Differentiating shows that $\omega(c) := \frac{c}{1+c^2} + \arctan(c)$ is a continuous strictly increasing function on $[0, +\infty)$ with $\omega(0) = 0$ and $\lim_{c\to+\infty}\omega(c) = \pi/2$. Hence, the equation $\pi/4 = \omega(c)$ has a unique solution in the set $(0, \infty)$. A short computation yields the expression

$$\nabla_{y_1}\zeta(y_1, y_2) = \frac{4}{\pi}\left(\frac{y_1 + y_2}{2\sqrt{-y_1/y_2}(y_1 - y_2)} - \frac{y_2}{2\sqrt{-y_1 y_2}} + \arctan\left(\sqrt{-\frac{y_1}{y_2}}\right)\right) - 1.$$
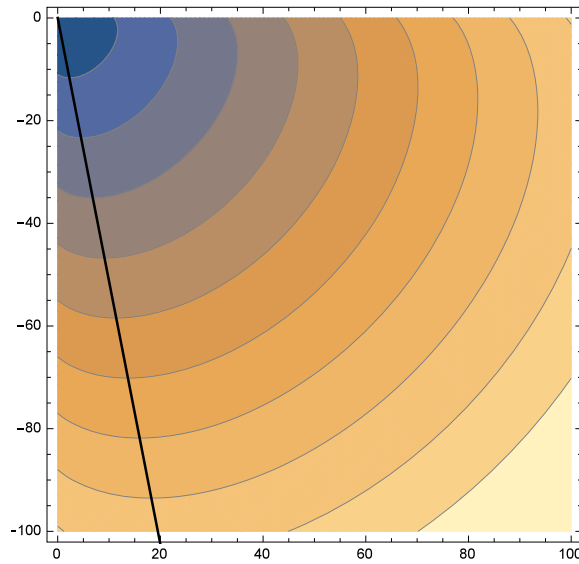
FIG. 8. Contour plot of the function $\zeta(y_1, y_2) := \mathbb{E}_{v_1, v_2}[|v_1 y_1 + v_2 y_2|]$ on $\mathbb{R}_+ \times \mathbb{R}_-$. The black line depicts all points $(y_1, y_2)$ with $\nabla_{y_1} \zeta(y_1, y_2) = 0$; for the explanation of the significance of this line, see Lemma 5.9.

Set $y_1 = -c^2 y_2$ for some $c > 0$ and $y_2 < 0$. Then plugging in this value of $y_1$, equality $0 = \nabla_{y_1} \zeta(y_1, y_2)$ holds if and only if

$$\pi/4 = \left( \frac{c}{1+c^2} + \arctan(c) \right).$$

This equation is independent of $y_1$ and its solution in $c$ is exactly the value satisfying $\pi/4 = \omega(c)$. $\square$

Thus, we have proved the following.

PROPOSITION 5.10 Let $c > 0$ be the unique solution of the equation $\frac{\pi}{4} = \frac{c}{1+c^2} + \arctan(c)$. Then a point $x \in \bar{x}^\perp \setminus \{0\}$ is stationary for $f_P$ if and only if equality $\|x\| = c\|\bar{x}\|$ holds.

In particular, we have proved one implication in Theorem 5.3. To prove the converse, we must show that every stationary point of $f_P$ lies in the set (5.1). Various approaches are possible based either on the decomposition $f_P(x) = \varphi_\lambda(X)$ or $f_P(x) = \zeta(\lambda_1(X), \lambda_d(X))$. We will focus on the former. We will prove a strong result about the location of stationary points of arbitrary convex spectral functions of $X$. Indeed, it will be more convenient to consider the more abstract setting as follows.

Throughout, we fix a symmetric convex function $f : \mathbb{R}^d \to \mathbb{R}$ and a point $0 \neq \bar{x} \in \mathbb{R}^d$, and define the function

$$g(x) := f_\lambda(xx^T - \bar{x}\bar{x}^T).$$

Note that the population objective $f_P$ has this representation with $f = \varphi$. The chain rule directly implies

$$\partial g(x) = 2\partial f_\lambda(X)x.$$

Therefore, using Theorem 5.2 let us also fix a matrix $V \in \partial f_\lambda(X)$ and a matrix $U \in \mathbb{O}^d$ satisfying

$$\lambda(V) \in \partial f(\lambda(X)), \qquad V = U\text{Diag}(\lambda(V))U^T \qquad \text{and} \qquad X = U\text{Diag}(\lambda(X))U^T.$$

The following two elementary lemmas will form the core of the argument.

LEMMA 5.11 (Eigenvalue correlation). The following are true.

1. **Eigenvalues.** We have $\lambda_i(X) = \langle U_i, x \rangle^2 - \langle U_i, \bar{x} \rangle^2$ for $i \in \{1, d\}$, and consequently

$$0 \le \lambda_1(X) \le \langle U_1, x \rangle^2 \le \|x\|^2$$
$$0 \le -\lambda_d(X) \le \langle U_d, \bar{x} \rangle^2 \le \|\bar{x}\|^2. \tag{5.4}$$

2. **Anticorrelation.** Equality holds:

$$\langle U_1, x \rangle \langle U_d, x \rangle = \langle U_1, \bar{x} \rangle \langle U_d, \bar{x} \rangle.$$

3. **Correlation.** Provided $x \notin \{\pm \bar{x}\}$, we have $\text{span}\{x, \bar{x}\} \subset \text{span}\{U_1, U_d\}$ and

$$\langle x, \bar{x} \rangle = \langle U_1, x \rangle \langle U_1, \bar{x} \rangle + \langle U_d, x \rangle \langle U_d, \bar{x} \rangle.$$

*Proof.* From the eigenvalue decomposition, we obtain

$$\lambda_1(X) = U_1^T X U_1 = \langle U_1, x \rangle^2 - \langle U_1, \bar{x} \rangle^2$$
$$\lambda_d(X) = U_d^T X U_d = \langle U_d, x \rangle^2 - \langle U_d, \bar{x} \rangle^2.$$

Taking into account that always $\lambda_1(X) \ge 0$ and $\lambda_1(X) \le 0$ (Lemma 5.5), we conclude $\lambda_1(X) \le \langle U_1, x \rangle^2$ and $\lambda_d(X) \ge -\langle U_d, \bar{x} \rangle^2$. Claim 1 follows. For Claim 2, simply observe

$$0 = U_d^T X U_1 = \langle U_1, x \rangle \langle U_d, x \rangle - \langle U_1, \bar{x} \rangle \langle U_d, \bar{x} \rangle.$$

To see Claim 3, for each $i \in \{1, d\}$ notice

$$\langle U_i, x \rangle x - \langle U_i, \bar{x} \rangle \bar{x} = X U_i = \lambda_i(X) U_i.$$

Suppose $x \notin \{\pm \bar{x}\}$. Then if $x$ and $\bar{x}$ are not collinear, we may divide through by $\lambda_i(X)$ and deduce, $\text{span}\{U_1, U_d\} = \text{span}\{x, \bar{x}\}$. On the other hand, if $x$ and $\bar{x}$ are collinear, then exactly one $\lambda_1$ or $\lambda_d$ is nonzero, and then $x$ lies in the span of the corresponding column of $U$. In either case, we may write $x = \langle U_1, x \rangle U_1 + \langle U_d, x \rangle U_d$ and $\bar{x} = \langle U_1, \bar{x} \rangle U_1 + \langle U_d, \bar{x} \rangle U_d$ in terms of their orthogonal expansions. We deduce

$$\langle x, \bar{x} \rangle = \langle \langle U_1, x \rangle U_1 + \langle U_d, x \rangle U_d, \langle U_1, \bar{x} \rangle U_1 + \langle U_d, \bar{x} \rangle U_d \rangle = \langle U_1, x \rangle \langle U_1, \bar{x} \rangle + \langle U_d, x \rangle \langle U_d, \bar{x} \rangle,$$

as claimed.                                                                                                                       □

LEMMA 5.12 (Spectral subdifferential). The following hold:

$$\max\left\{|\lambda_1(V)\langle U_1,x\rangle|, |\lambda_d(V)\langle U_d,x\rangle|\right\} \leq \|Vx\|, \tag{5.5}$$

and

$$g(x) - g(\bar{x}) \leq \lambda_1(V)\lambda_1(X) + \lambda_d(V)\lambda_d(X). \tag{5.6}$$

*Proof.* To see (5.5), observe that for all unit vectors $z \in \mathbb{S}^{d-1}$, we have $\|Vx\| \geq \langle z, Vx\rangle$. Thus, testing against all $z \in \{\pm U_1, \pm U_d\}$ yields the lower bounds (5.5). To prove the final bound (5.6), we exploit the convexity of $f_\lambda$. The subgradient inequality implies

$$f_\lambda(X) - f_\lambda(0) \leq \langle V, X\rangle = \lambda_1(V)\lambda_1(X) + \lambda_d(V)\lambda_d(X).$$

The result follows. $\qquad \square$

The following corollary follows quickly from the previous two lemmas.

COROLLARY 5.13 (Stationary point inclusion). Suppose that $x$ is stationary for $g$, that is $Vx = 0$. Then one of the following conditions holds:

1. $g(x) \leq g(\bar{x})$

2. $x = 0$

3. $\langle x, \bar{x}\rangle = 0$, $\lambda_1(V) = 0$.

Moreover, if $\bar{x}$ minimizes $g$, then a point $x$ is stationary for $g$ if and only if $x$ satisfies 1, 2 or 3.

*Proof.* Suppose $Vx = 0$ and that the first two conditions fail, that is $x \neq 0$ and $g(x) > g(\bar{x})$. We will show that the third condition holds. To this end, inequalities (5.5) and (5.6), along with Part 3 of Lemma 5.11, directly imply the following:

$$0 < g(x) - g(\bar{x}) \leq \lambda_1(V)\lambda_1(X) + \lambda_d(V)\lambda_d(X), \tag{5.7}$$

$$0 = \lambda_1(V)\langle U_1,x\rangle \tag{5.8}$$

$$0 = \lambda_d(V)\langle U_d,x\rangle, \tag{5.9}$$

$$x = \langle U_1,x\rangle U_1 + \langle U_d,x\rangle U_d. \tag{5.10}$$

Aiming towards a contradiction, suppose $\lambda_1(V) \neq 0$. Then (5.8) and (5.10) imply $\langle U_1,x\rangle = 0$ and $\langle U_d,x\rangle \neq 0$, since otherwise $x$ would be 0. Equation (5.9), in turn, yields $\lambda_d(V) = 0$. Appealing to Lemma 5.11, we moreover deduce

$$0 \leq \lambda_1(X) = \langle U_1,x\rangle^2 - \langle U_1,\bar{x}\rangle^2 \leq 0.$$

Thus, $\lambda_1(X) = 0$ and therefore the right-hand side of (5.7) is zero, a contradiction. We have shown the equality $\lambda_1(V) = 0$, as claimed.

Inequality (5.7) implies $\lambda_d(V) \neq 0$ and $\lambda_d(X) \neq 0$, and hence by inequality (5.9), we have $\langle U_d, x \rangle = 0$. Combining the latter equality with Part 2 of Lemma 5.11, we conclude $0 = \langle U_1, x \rangle \langle U_d, \rangle x = \langle U_1, \bar{x} \rangle \langle U_d, \bar{x} \rangle$. Note $\langle U_d, \bar{x} \rangle \neq 0$, since otherwise we would get $\lambda_d(X) = 0$ by (5.4). We conclude $\langle U_1, \bar{x} \rangle = 0$. Finally, Part 3 of Lemma 5.11 then yields

$$\langle x, \bar{x} \rangle = \langle U_1, x \rangle \langle U_1, \bar{x} \rangle + \langle U_d, x \rangle \langle U_d, \bar{x} \rangle = 0,$$

thereby completing the proof.

Now suppose that $\bar{x}$ minimizes $g$. If $g(x) \leq g(\bar{x})$, then $x$ is a minimizer of $g$ and thus, a stationary point. In addition, 0 is a stationary point of $g$ because $V \cdot 0 = 0$. Thus, it remains to show that all points satisfying 3 are stationary. Thus, suppose $x$ satisfies 3 and $x \neq 0$. Then the eigenvalues of $X$ are precisely $\|x\|^2$ and $-\|\bar{x}\|^2$ with eigenvectors $U_1 = \pm\frac{x}{\|x\|}$ and $U_d = \pm\frac{\bar{x}}{\|\bar{x}\|}$, respectively. Thus, we have $U^T V x = \text{Diag}(\lambda(V)) U^T x = (\lambda_1(V) \langle U_1, x \rangle, 0, \ldots, 0, \lambda_d(V) \langle U_d, x \rangle)^T = 0$. We conclude $Vx = 0$, as required. □

The proof of Theorem 5.3 is now immediate.

*Proof of Theorem* 5.3.   We have already proved that every point in the set (5.1) is stationary for $f_P$ (Proposition 5.10). Thus, we focus on the converse. In light of Proposition 5.10, it is sufficient to show that every stationary point $x$ of $f_P$ lies in the set $\{0, \pm\bar{x}\} \cup x^\perp$. This is immediate from Corollary 5.13 under the identification $f_P(x) = g(x) = \varphi_\lambda(X)$. □

## 6. Concentration and stability

Having determined the stationary points of the population objective $f_P$, we next turn to the stationary points of $f_S$. Our strategy is to show that with high probability, every stationary point of $f_S$ is close to some stationary point of $f_P$. The difficulty is that it is not true that $\partial f_S(x)$ concentrates around $\partial f_P(x)$. Instead, we will see that the graphs of the two subdifferentials $\partial f_S$ and $\partial f_P$ concentrate, which is sufficient for our purposes. Our argument will rely on two basic properties, namely (1) the subsampled objective $f_S$ concentrates well around $f_P$, and (2) the function $f_S$ is weakly convex.

### 6.1   *Concentration of subdifferential graphs*

Armed with the concentration (Theorem 5.1) and the weak convexity (Theorem 3.5) guarantees, we can show that the graphs of $\partial f_P$ and $\partial f_S$ are close. The following theorem will be our main technical tool, and is of interest in its own right. In essence, the result is a quantitative extension of the celebrated Attouch *et al.* (1990)'s convergence theorem in convex analysis. Henceforth, for any function $l: \mathbb{R}^d \to \overline{\mathbb{R}}$ and a point $\bar{x} \in \mathbb{R}^d$, with $l(\bar{x})$ finite, we define the local Lipschitz constant

$$\text{lip}(l; \bar{x}) := \limsup_{x \to \bar{x}} \frac{|l(x) - l(\bar{x})|}{\|x - \bar{x}\|}.$$

THEOREM 6.1 (Comparison). Consider four lsc functions $f, g, l, u \colon \mathbb{R}^d \to \overline{\mathbb{R}}$ and a pair $(x, v) \in \operatorname{gph} \partial g$. Suppose that $l$ is locally Lipschitz continuous and that the following conditions

$$\left\{ \begin{array}{l} l(y) \leq f(y) - g(y) \leq u(y) \\[2mm] g(y) \geq g(x) + \langle v, y - x \rangle - \dfrac{\rho}{2} \|y - x\|^2 \end{array} \right\} \qquad \text{hold for all points } y \in \mathbb{R}^d.$$

Then for any $\gamma > 0$, there exists a point $\hat{x}$, with $l(\hat{x})$ finite, satisfying

$$\|\hat{x} - x\| \leq 2\gamma \qquad \text{and} \qquad \operatorname{dist}(v; \partial f(\hat{x})) \leq 2\rho\gamma + \frac{u(x) - l(x)}{\gamma} + \operatorname{lip}(l; \hat{x}).$$

In particular, if $l(\cdot)$ is constant, we have the estimate

$$\operatorname{dist}\Big((x, v), \operatorname{gph} \partial f\Big) \leq \sqrt{4(\rho + \sqrt{2 + \rho^2})} \cdot \sqrt{u(x) - l(x)}. \tag{6.1}$$

*Proof.* From the two assumptions, for any point $y \in \mathbb{R}^d$ we have

$$f(y) \geq g(y) + l(y) \geq g(x) + l(y) + \langle v, y - x \rangle - \frac{\rho}{2} \|y - x\|^2.$$

Define the function

$$\zeta_x(y) := f(y) - \langle v, y - x \rangle + \frac{\rho}{2} \|y - x\|^2 - l(y).$$

Clearly then we have

$$\zeta_x(x) - \inf_y \zeta_x \leq f(x) - l(x) - g(x) \leq u(x) - l(x). \tag{6.2}$$

Choose now any minimizer

$$\hat{x} \in \underset{y}{\operatorname{argmin}} \left\{ \zeta_x(y) + \frac{u(x) - l(x)}{4\gamma^2} \cdot \|y - x\|^2 \right\}.$$

First-order optimality conditions and the sum rule (Rockafellar & Wets, 1998, Exercise 10.10) immediately imply

$$\frac{u(x) - l(x)}{2\gamma^2} \cdot (x - \hat{x}) \in \partial \zeta_x(\hat{x}) \subset \partial f(\hat{x}) - v + \rho(\hat{x} - x) + \operatorname{lip}(l; \hat{x}) B(0, 1),$$

and hence

$$\operatorname{dist}(v; \partial f(\hat{x})) \leq \frac{u(x) - l(x)}{2\gamma^2} \cdot \|\hat{x} - x\| + \rho \|\hat{x} - x\| + \operatorname{lip}(l; \hat{x}). \tag{6.3}$$

Next, we estimate the distance $\|\hat{x} - x\|$. To this end, observe from the definition of $\hat{x}$, we have

$$\zeta_x(\hat{x}) + \frac{u(x) - l(x)}{4\gamma^2} \cdot \|\hat{x} - x\|^2 \leq \zeta_x(x),$$

and hence

$$\frac{u(x) - l(x)}{4\gamma^2} \cdot \|\hat{x} - x\|^2 \le \zeta_x(x) - \zeta_x(\hat{x}) \le u(x) - l(x), \tag{6.4}$$

where the last inequality follows from (6.2). In the case $u(x) = l(x)$, we deduce $\zeta_x(x) = \zeta_x(\hat{x})$. Thus, we equally well could have set $\hat{x} = x$, and the theorem follows immediately from (6.3). On the other hand, in the setting $u(x) > l(x)$, the inequality (6.4) immediately yields $\|\hat{x} - x\| \le 2\gamma$, as claimed. Combining this inequality with (6.3) then gives the desired guarantee

$$\text{dist}(v; \partial f(\hat{x})) \le 2\rho\gamma + \frac{u(x) - l(x)}{\gamma} + \text{lip}(l; \hat{x}).$$

Supposing $l$ is a constant, we have the estimate

$$\text{dist}\left((x, v), \text{gph } \partial f\right) \le \sqrt{4\gamma^2 + \left(2\rho\gamma + \frac{u(x) - l(x)}{\gamma}\right)^2}.$$

Minimizing the right-hand side in $\gamma$ yields the choice $\gamma = \frac{\sqrt{u(x) - l(x)}}{(8 + 4\rho^2)^{1/4}}$. With this value of $\gamma$, a quick computation yields the claimed guarantee (6.1). □

Let us now specialize the theorem to the setting where the lower and upper bounds $l(\cdot), u(\cdot)$ are functions of the product $\|x - \bar{x}\| \cdot \|x + \bar{x}\|$, as in phase retrieval.

COROLLARY 6.2   Fix two functions $f, g: \mathbb{R}^d \to \mathbb{R}$. Suppose that $g$ is $\rho$-weakly convex and that there is a point $\bar{x}$ and a real $\delta > 0$ such that the inequality

$$|f(x) - g(x)| \le \delta \|x - \bar{x}\| \cdot \|x + \bar{x}\| \qquad \text{holds for all } x \in \mathbb{R}^d.$$

Then for any stationary point $x$ of $g$, there exists a point $\hat{x}$ satisfying

$$\left\{ \begin{array}{c} \|x - \hat{x}\| \le \sqrt{\frac{4\delta}{\rho + 2\delta}} \cdot \sqrt{\|x - \bar{x}\| \|x + \bar{x}\|}, \\[2mm] \text{dist}(0; \partial f(\hat{x})) \le (\delta + 2\sqrt{\delta(\rho + 2\delta)}) \cdot (\|x - \bar{x}\| + \|x + \bar{x}\|) \end{array} \right\}.$$

*Proof.*   Set $u(y) := \delta \|y - \bar{x}\| \cdot \|y + \bar{x}\|$ and $l(y) := -\delta \|y - \bar{x}\| \cdot \|y + \bar{x}\|$ and observe $\text{lip}(l; y) \le \delta(\|y - \bar{x}\| + \|y + \bar{x}\|)$. Applying Theorem 6.1, we deduce that for any $\gamma > 0$, there exists a point $\hat{x}$ satisfying

$$\|\hat{x} - x\| \le 2\gamma \qquad \text{and} \qquad \text{dist}(0; \partial f(\hat{x})) \le 2\rho\gamma + \frac{2\delta \|x - \bar{x}\| \|x + \bar{x}\|}{\gamma} + \delta(\|\hat{x} - \bar{x}\| + \|\hat{x} + \bar{x}\|).$$

The triangle inequality implies

$$\|\hat{x} - \bar{x}\| \le 2\gamma + \|x - \bar{x}\| \qquad \text{and} \qquad \|\hat{x} + \bar{x}\| \le 2\gamma + \|x + \bar{x}\|,$$

and therefore

$$\text{dist}(0; \partial f(\hat{x})) \leq 2(\rho + 2\delta)\gamma + \frac{2\delta \|x - \bar{x}\| \|x + \bar{x}\|}{\gamma} + \delta(\|x - \bar{x}\| + \|x + \bar{x}\|).$$

Minimizing this expression in $\gamma > 0$ yields the choice $\gamma := \sqrt{\frac{\delta \|x-\bar{x}\| \|x+\bar{x}\|}{\rho + 2\delta}}$. Plugging in this value of $\gamma$ and applying the Arithmetic Mean–Geometric Mean (AM–GM) inequality then implies

$$\text{dist}(0; \partial f(\hat{x})) \leq 4\sqrt{\delta(\rho + 2\delta)\|x - \bar{x}\| \|x + \bar{x}\|} + \delta(\|x - \bar{x}\| + \|x + \bar{x}\|)$$

$$\leq (\delta + 2\sqrt{\delta(\rho + 2\delta)})(\|x - \bar{x}\| + \|x + \bar{x}\|).$$

The result follows. □

We now arrive at the main result of the section.

COROLLARY 6.3 (Subsampled stationary points). Consider the robust phase retrieval objective $f_S(\cdot)$ generated from i.i.d. standard Gaussian vectors. There exist numerical constants $c_1, c_2 > 0$ such that whenever $m \geq c_1 d$, then with probability at least $1 - 2\exp(-\min\{m/c_1, c_2 m, d^2\})$, every stationary point $x$ of $f_S$ satisfies $\|x\| \lesssim \|\bar{x}\|$ and one of the two conditions:

$$\frac{\|x\| \|x - \bar{x}\| \|x + \bar{x}\|}{\|\bar{x}\|^3} \lesssim \sqrt[4]{\frac{d}{m}} \qquad \text{or} \qquad \left\{ \begin{array}{l} \left| \frac{\|x\|}{\|\bar{x}\|} - c \right| \lesssim \sqrt[4]{\frac{d}{m}} \cdot \left( 1 + \frac{\|\bar{x}\|}{\|x\|} \right) \\ \\ \frac{|\langle x, \bar{x} \rangle|}{\|x\| \|\bar{x}\|} \lesssim \sqrt[4]{\frac{d}{m}} \cdot \frac{\|\bar{x}\|}{\|x\|} \end{array} \right\},$$

where $c > 0$ is the unique solution of the equation $\frac{\pi}{4} = \frac{c}{1+c^2} + \arctan(c)$.

*Proof.* Theorem 5.1 shows that there exist constants $c_1, c_2 > 0$ such with probability at least $1 - 2\exp(-c_1 \min\{m, d^2\})$, we have

$$\left| f_S(x) - f_P(x) \right| \leq \frac{c_2}{2} \left( \sqrt{\frac{d}{m}} + \frac{d}{m} \right) \|x - \bar{x}\| \|x + \bar{x}\| \qquad \text{for all } x \in \mathbb{R}^d. \tag{6.5}$$

Lemma 3.5, in turn, shows that there exist numerical constants $c_3, \rho > 0$ such that provided $m \geq c_3 d$, the function $f_S$ is $\rho$-weakly convex, with probability at least $1 - \exp\left(-\frac{m}{c_3}\right)$. Let us now try to apply Corollary 6.2. To simplify notation, define $\Delta := \sqrt{\frac{d}{m}}$ and set $\delta := c_2 \Delta$. Notice $\delta \geq \frac{c_2}{2}(\Delta + \Delta^2)$ and hence we may apply Corollary 6.2. We deduce that with high probability, for any stationary point $x$ of $f_S$ there exists a point $\hat{x} \in \mathbb{R}^d$ satisfying

$$\left\{ \begin{array}{l} \|x - \hat{x}\| \leq \sqrt{\frac{4c_2 \Delta}{\rho + 2c_2 \Delta}} \cdot \sqrt{\|x - \bar{x}\| \|x + \bar{x}\|}, \\ \\ \text{dist}(0; \partial f_P(\hat{x})) \leq (c_2 \Delta + 2\sqrt{c_2 \Delta(\rho + 2c_2 \Delta)}) \cdot (\|x - \bar{x}\| + \|x + \bar{x}\|) \end{array} \right\}. \tag{6.6}$$

Notice $\sqrt{\frac{4c_2\Delta}{\rho+2c_2\Delta}} \leq \sqrt{\Delta} \cdot \sqrt{\frac{4c_2}{\rho}} \leq 2C'\sqrt{\Delta}$ and $(c_2\Delta + 2\sqrt{c_2\Delta(\rho+2c_2\Delta)}) \leq C'\sqrt{\Delta}$ for some numerical constant $C'$. For notational convenience, set $D_x := \|x-\bar{x}\| + \|x+\bar{x}\|$. Thus, by the AM–GM inequality, the inclusion $\hat{x} \in B(x, C'\sqrt{\Delta}D_x)$ holds.

**Claim 1.** There exist constants $C'', \tau > 0$ such that with high probability, for all $\Delta < C''$, the inequality $\|x\| \leq \tau\|\bar{x}\|$ holds for any stationary point $x$ of $f_S$.

*Proof.* We may assume that $\|\bar{x}\| \leq \|x\|$ since otherwise the result is trivial. Next, observe that $\|x\|$ and $\|\hat{x}\|$ have comparable norms:

$$\|\hat{x}\| \leq \|x\| + C'\sqrt{\Delta}D_x \leq (1 + 4C'\sqrt{\Delta})\|x\|,$$

$$\|\hat{x}\| \geq \|x\| - C'\sqrt{\Delta}D_x \geq (1 - 4C'\sqrt{\Delta})\|x\|,$$

where we have used the bound $D_x \leq 4\|x\|$ twice. To make the last bound meaningful, we may set $C'' < \left(\frac{1}{8C'}\right)^2$, thereby ensuring $1 - 4C'\sqrt{\Delta} \geq 1/2$. Because the norms are comparable, we deduce

$$\operatorname{dist}(0; \partial f_P(\hat{x})) \leq C'\sqrt{\Delta}D_x \leq 4C'\sqrt{\Delta}\|x\| \leq \frac{4C'\sqrt{\Delta}}{(1-4C'\sqrt{\Delta})}\|\hat{x}\|. \tag{6.7}$$

Let us now decrease $C''$ if necessary to have $C'' < \min\left\{\left(\frac{1}{8C'}\right)^2, \left(\frac{\gamma}{8C'}\right)^2\right\}$, where $\gamma$ is the fixed constant from Theorem 5.4. Then for all $\Delta < C''$, we have $1 - 4C'\sqrt{\Delta} \geq \frac{1}{2}$ and $\frac{4C'\sqrt{\Delta}}{1-4C'\sqrt{\Delta}} \leq 8C'\sqrt{\Delta} \leq \gamma$. Now we can apply Theorem 5.4 to $\hat{x}$, which guarantees that $\|\hat{x}\| \lesssim \|\bar{x}\|$. Thus, because the norms of $\|x\|$ and $\|\hat{x}\|$ are comparable, we obtain the desired result. $\qquad\square$

Let us now examine two cases. First, we suppose that $\|\hat{x}\| \leq \frac{C'\sqrt{\Delta}D_x}{\gamma}$. Then we have

$$\|x\| \leq \|x - \hat{x}\| + \|\hat{x}\| \leq C'\sqrt{\Delta}D_x + \frac{C'\sqrt{\Delta}D_x}{\gamma} \lesssim \sqrt{\Delta}D_x.$$

In addition, Claim 1 implies that

$$\|x - \bar{x}\|\|x + \bar{x}\| \lesssim \|\bar{x}\|^2.$$

Multiplying these two estimates together yields the claimed bound.

Next suppose that $\|\hat{x}\| > \frac{C'\sqrt{\Delta}D_x}{\gamma}$. Then

$$\operatorname{dist}(0; \partial f_P(\hat{x})) \leq \varepsilon := C'\sqrt{\Delta}D_x \leq \gamma\|\hat{x}\|.$$

Applying Theorem 5.4 we find that the point $\hat{x} \in B(x, C'\sqrt{\Delta}D_x)$ satisfies either

$$\|\hat{x}\|\|\hat{x} - \bar{x}\|\|\hat{x} + \bar{x}\| \lesssim \sqrt{\Delta}D_x\|\bar{x}\|^2 \qquad \text{or} \qquad \left\{ \begin{array}{c} \left|\|\hat{x}\| - c\|\bar{x}\|\right| \lesssim \sqrt{\Delta}D_x\dfrac{\|\bar{x}\|}{\|\hat{x}\|} \\[2mm] |\langle\hat{x}, \bar{x}\rangle| \lesssim \sqrt{\Delta}D_x\|\bar{x}\| \end{array} \right\}. \tag{6.8}$$

Applying the triangle inequality and the bound $D_x \leq (2 + 2\tau)\|\bar{x}\|$, the claimed inequalities all follow (see Appendix A for a detailed explanation). □

REMARK 1 The distribution of stationary points of $f_S$, presented in Corollary 6.3, is very similar to the distribution of stationary points for the smooth formulation of phase retrieval (Sun *et al.*, 2018). Indeed, the main result of Sun *et al.* (2018) shows a much stronger property, namely that every stationary point is either a global minimizer, or a point where the Hessian of the objective function has strictly negative curvature. Their argument proceeds by showing that the Hessian of the sampled function concentrates around the Hessian of the population objective. A similar analysis seems out of reach in our setting since the objective function is nonsmooth.

## REFERENCES

ATTOUCH, H., NDOUTOUME, J. L. & THÉRA, M. (1990) Epigraphical convergence of functions and convergence of their derivatives in Banach spaces. *Sém. Anal. Convexe*, **20**, 45.

BAHMANI, S. & ROMBERG, J. (2017) Phrase retrieval meets statistical learning theory: a flexible convex relaxation. *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics, in PMLR*. **54**, 252–260.

BURKE, J. V. & FERRIS, M. C. (1993) Weak sharp minima in mathematical programming. *SIAM J. Control Optim.*, **31**, 1340–1359.

CANDÈS, E. J. & LI, X. (2014) Solving quadratic equations via PhaseLift when there are about as many equations as unknowns. *Found. Comput. Math.*, **14**, 1017–1026.

CANDÈS, E. J., LI, X. & SOLTANOLKOTABI, M. (2015) Phase retrieval via Wirtinger flow: theory and algorithms. *IEEE Trans. Inf. Theory*, **61**, 1985–2007.

CANDÈS, E.J., STROHMER, T. & VORONINSKI, V. (2013) Phaselift: exact and stable signal recovery from magnitude measurements via convex programming. *Commun. Pure Appl. Math.*, **66**, 1241–1274.

CHEN, Y. & CANDÈS, E. J. (2017) Solving random quadratic systems of equations is nearly as easy as solving linear systems. *Commun. Pure Appl. Math.*, **70**, 822–883.

CLARKE, F. H., STERN, R. J. & WOLENSKI, P.R. (1995) Proximal smoothness and the lower-$C^2$ property. *J. Convex Anal.*, **2**, 117–144.

DAVIS, D., DRUSVYATSKIY, D., MACPHEE, K. J. & PAQUETTE, C. (2018) Subgradient methods for sharp weakly convex functions. *J. Optimiz. Theory App.*, **179**, 962–982

DAVIS, D. & GRIMMER, B. (2019) Proximally guided stochastic method for nonsmooth, nonconvex problems. *SIAM J. Optim.*, **29**, 1908–1930.

DRUSVYATSKIY, D. & LEWIS, A. S. (2018) Error bounds, quadratic growth, and linear convergence of proximal methods. *Math. Oper. Res.*, **43**, 919–948.

DRUSVYATSKIY, D. & PAQUETTE, C. (2019) Efficiency of minimizing compositions of convex functions and smooth maps. *Math. Program.,* **178**, 503–558.

DRUSVYATSKIY, D. & PAQUETTE, C. (2018) Variational analysis of spectral functions simplified. *J. Convex Anal.*, **25**.

DUCHI, J. C. & RUAN, F. (2018a) Solving (most) of a set of quadratic equalities: composite optimization for robust phase retrieval. *Information and Inference: A Journal of the IMA*, **8**, 471–529.

DUCHI, J. C. & RUAN, F. (2018b) Stochastic methods for composite and weakly convex optimization problems. *SIAM J. Optim.,* **28**, 3229–3259.

ELDAR, Y.C. & MENDELSON, S. (2014) Phase retrieval: stability and recovery guarantees. *Appl. Comput. Harmon. Anal.*, **36**, 473–494.

FEDERER, H. (1959) Curvature measures. *Trans. Amer. Math. Soc.*, **93**, 418–491.

FICKUS, M., MIXON, D. G., NELSON, A. A. & WANG, Y. (2014) Phase retrieval from very few measurements. *Linear Algebra Appl.*, **449**, 475–499.

GOLDSTEIN, T. & STUDER, C. (2018) PhaseMax: convex phase retrieval via basis pursuit. *IEEE Trans. Inf. Theory*, **64**, 2675–2689.

HAND, P. & VORONINSKI, V. (2016) An elementary proof of convex phase retrieval in the natural parameter space via the linear program phasemax. *Preprint arXiv:1611.03935*.

HORN, R. A. & JOHNSON, C. R. (2013) *Matrix Analysis*, 2nd edn. Cambridge: Cambridge University Press.

KATO, T. (1982) *A Short Introduction to Perturbation Theory for Linear Operators*. New York-Berlin: Springer Science & Business Media.

LECUN, Y., BOTTOU, L., BENGIO, Y. & HAFFNER, P. (1998) Gradient-based learning applied to document recognition. *Proc. IEEE*, **86**, 2278–2324.

LEWIS, A. S. (1996) Convex analysis on the Hermitian matrices. *SIAM J. Optim.*, **6**, 164–177.

LEWIS, A. S. (1999) Nonsmooth analysis of eigenvalues. *Math. Programming*, **84**, 1–24.

LEWIS, A. S. & WRIGHT, S. J. (2016) A proximal method for composite minimization. *Math. Program*, **158**, 501–546.

LUO, Z.-Q. & TSENG, P. (1993) Error bounds and convergence analysis of feasible descent methods: a general approach. *Ann. Oper. Res.*, **46**, 157–178.

MORDUKHOVICH, B. S. (2006) *Variational Analysis and Generalized Differentiation I: Basic Theory*. Grundlehren der mathematischen Wissenschaften, vol. **330**. Berlin: Springer.

POLIQUIN, R. A. & ROCKAFELLAR, R. T. (1996) Prox-regular functions in variational analysis. *Trans. Am. Math. Soc.*, **348**, 1805–1838.

POLJAK, B. T. (1967) A general method for solving extremal problems. *Dokl. Akad. Nauk SSSR*, **174**, 33–36.

ROCKAFELLAR, R. T. (1982) Favorable classes of Lipschitz-continuous functions in subgradient optimization. *Progress in nondifferentiable optimization*, vol. 8. IIASA Collaborative Proceedings Series CP-82, pp. 125–143. Laxenburg: International Institute for Applied Systems Analysis.

ROCKAFELLAR, R. T. & WETS, R. J.-B. (1998) *Variational Analysis*. Grundlehren der mathematischen Wissenschaften, vol. **317**. Berlin: Springer.

SUN, J., QU, Q. & WRIGHT, J. (2018) A geometric analysis of phase retrieval. *Found. Comput. Math.* **18**, 1131–1198.

TAN, Y. S. & VERSHYNIN, R. (2018) Phase retrieval via randomized Kaczmarz: theoretical guarantees. *Information and Inference: A Journal of the IMA*, **8**, 97–123.

VERSHYNIN, R. (2010) *Introduction to the non-asymptotic analysis of random matrices. arXiv preprint:1011.3027*.

WANG, G., GIANNAKIS, G. B. & ELDAR, Y. C. (2017) Solving systems of random quadratic equations via truncated amplitude flow. *IEEE Transactions on Information Theory*, **64**, 773–794.

ZHANG, H., CHI, Y. & LIANG, Y. (2016) Provable non-convex phase retrieval with outliers: median truncated Wirtinger flow. *Proceedings of the 33rd International Conference on International Conference on Machine Learning*, vol. 48. pp. 1022–1031. Proceedings of Machine Learning Research, New York, USA. JMLR.org.

## Appendices

## A. Auxiliary computations

*Proof of Lemma* 5.7. We let $\sigma_1 = y_1$ and $\sigma_2 = -y_2$. We may write

$$
\begin{aligned}
\mathbb{E}_v\left[|\sigma_1 v_1^2 - \sigma_2 v_2^2|\right] &= \frac{1}{2\pi} \int_{\mathbb{R}^2} |\sigma_1 v_1^2 - \sigma_2 v_2^2| \exp\left(-\left(\frac{v_1^2 + v_2^2}{2}\right)\right) dv_1\, dv_2 \\
&= \frac{1}{2\pi} \int_{R_1} \left(\sigma_1 v_1^2 - \sigma_2 v_2^2\right) \exp\left(-\left(\frac{v_1^2 + v_2^2}{2}\right)\right) dv_1\, dv_2 \\
&\quad + \frac{1}{2\pi} \int_{R_2} \left(\sigma_2 v_2^2 - \sigma_1 v_1^2\right) \exp\left(-\left(\frac{v_1^2 + v_2^2}{2}\right)\right) dv_1\, dv_2,
\end{aligned}
$$

where

$$
\begin{aligned}
R_1 &= \left\{(v_1, v_2) : \sqrt{\sigma_1}|v_1| \geq \sqrt{\sigma_2}|v_2|\right\} \\
R_2 &= \left\{(v_1, v_2) : \sqrt{\sigma_2}|v_2| \geq \sqrt{\sigma_1}|v_1|\right\}.
\end{aligned}
$$

Using the convention $\arctan(\theta) \in \left[\frac{-\pi}{2}, \frac{\pi}{2}\right]$, we define the angle $\theta_1 := \arctan\left(\sqrt{\frac{\sigma_1}{\sigma_2}}\right)$. Passing to the polar coordinates, we deduce

$$
\begin{aligned}
\frac{1}{2\pi} \int_{R_1} (\sigma_1 v_1^2 - \sigma_2 v_2^2) \exp\left(-\left(\frac{v_1^2 + v_2^2}{2}\right)\right) dv_1\, dv_2 \\
= \frac{1}{2\pi} \int_{R_1} r^3 (\sigma_1 \cos^2(\theta) - \sigma_2 \sin^2(\theta))\, e^{-r^2/2}\, dr\, d\theta.
\end{aligned}
$$

We break up the region $R_1$ into three wedges corresponding to the angles $[0, \theta_1]$, $[2\pi, 2\pi - \theta_1]$ and $[\pi + \theta_1, \pi - \theta_1]$. We will compute the integral over one of the regions. The rest will follow analogously.

To this end, we successively deduce

$$\frac{1}{2\pi} \int_0^{\theta_1} \int_0^{\infty} r^3 \left( \sigma_1 \cos^2(\theta) - \sigma_2 \sin^2(\theta) \right) e^{-r^2/2} \, dr \, d\theta$$

$$= \frac{1}{2\pi} \int_0^{\theta_1} \sigma_1 \left(1 + \cos(2\theta)\right) - \sigma_2 \left(1 - \cos(2\theta)\right) d\theta$$

$$= \frac{1}{2\pi} \left( (\sigma_1 - \sigma_2)\theta + (\sigma_1 + \sigma_2) \sin(\theta) \cos(\theta) \right) \Big|_0^{\theta_1}$$

$$= \frac{1}{2\pi} \left( (\sigma_1 - \sigma_2)\theta_1 + (\sigma_1 + \sigma_2) \sin(\theta_1) \cos(\theta_1) \right)$$

$$\frac{1}{2\pi} \int_{2\pi-\theta_1}^{2\pi} \int_0^{\infty} r^3 \left( \sigma_1 \cos^2(\theta) - \sigma_2 \sin^2(\theta) \right) e^{-r^2/2} \, dr \, d\theta$$

$$= \frac{1}{2\pi} \left( (\sigma_1 - \sigma_2)\theta_1 + (\sigma_1 + \sigma_2) \sin(\theta_1) \cos(\theta_1) \right)$$

$$\frac{1}{2\pi} \int_{\pi-\theta_1}^{\pi+\theta_1} \int_0^{\infty} r^3 \left( \sigma_1 \cos^2(\theta) - \sigma_2 \sin^2(\theta) \right) e^{-r^2/2} \, dr \, d\theta$$

$$= \frac{1}{2\pi} \left( 2(\sigma_1 - \sigma_2)\theta_1 + 2(\sigma_1 + \sigma_2) \sin(\theta_1) \cos(\theta_1) \right).$$

Similarly, we see that for the region $R_2$, we have

$$\frac{1}{2\pi} \int_{R_2} (\sigma_2 v_2^2 - \sigma_1 v_1^2) \exp\left( -\left( \frac{v_1^2 + v_2^2}{2} \right) \right) dv_1 \, dv_2$$

$$= \frac{1}{2\pi} \int_{R_2} r^3 (\sigma_2 \sin^2(\theta) - \sigma_1 \cos^2(\theta)) \, e^{-r^2/2} \, dr \, d\theta.$$

We break up the region $R_2$ into two wedges where the angles range from $[\theta_1, \pi - \theta_1]$ and $[\pi + \theta_1, 2\pi - \theta_1]$ as we did in $R_1$. We will show the explicit computation for one of these terms and note the rest following

using similar computations:

$$\frac{1}{2\pi} \int_{\theta_1}^{\pi-\theta_1} \int_0^\infty r^3 \left( \sigma_2 \sin^2(\theta) - \sigma_1 \cos^2(\theta) \right) e^{-r^2/2} \, dr \, d\theta$$

$$= \frac{1}{2\pi} \int_{\theta_1}^{\pi-\theta_1} \sigma_2 \left( 1 - \cos(2\theta) \right) - \sigma_1 \left( 1 + \cos(2\theta) \right) d\theta$$

$$= \frac{1}{2\pi} \left( (\sigma_2 - \sigma_1)\theta - (\sigma_1 + \sigma_2) \sin(\theta) \cos(\theta) \right) \Big|_{\theta_1}^{\pi-\theta_1}$$

$$= \frac{1}{2\pi} \left( (\sigma_2 - \sigma_1)(\pi - 2\theta_1) + 2(\sigma_1 + \sigma_2) \sin(\theta_1) \cos(\theta_1) \right)$$

$$\frac{1}{2\pi} \int_{\pi+\theta_1}^{2\pi-\theta_1} \int_0^\infty r^3 \left( \sigma_2 \sin^2(\theta) - \sigma_1 \cos^2(\theta) \right) e^{-r^2/2} \, dr \, d\theta$$

$$= \frac{1}{2\pi} \left( (\sigma_2 - \sigma_1)(\pi - 2\theta_1) + 2(\sigma_1 + \sigma_2) \sin(\theta_1) \cos(\theta_1) \right).$$

By combining the computed integrals, we arrive at the full answer

$$\mathbb{E}_v \left[ |\sigma_1 v_1^2 - \sigma_2 v_2^2| \right] = \frac{4}{\pi} \left[ (\sigma_1 - \sigma_2) \arctan\left( \sqrt{\frac{\sigma_1}{\sigma_2}} \right) + \sqrt{\sigma_1 \sigma_2} \right] - (\sigma_1 - \sigma_2),$$

as claimed. $\square$

*Proof.* showing Equation (6.8) implies Corollary 6.3 We observe that $D_x \leq 2\|x\| + 2\|\bar{x}\|$, which by Claim 6.4 gives $D_x \leq (2\tau + 2)\|\bar{x}\|$. First by applying the triangle inequality with $\|\hat{x} - x\| \leq C'\sqrt{\Delta}D_x$ and (6.8), we obtain

$$\left| \|x\| - c\|\bar{x}\| \right| \leq \left| \|x - \hat{x}\| + \|\hat{x}\| - c\|\bar{x}\| \right| \lesssim C'\sqrt{\Delta}D_x + \sqrt{\Delta}D_x \frac{\|\bar{x}\|}{\|\hat{x}\|}.$$

Using the bound on $D_x$ gives the desired inequality. Next, we conclude

$$|\langle x, \bar{x} \rangle| \leq \|\bar{x}\| \|x - \hat{x}\| + |\langle \hat{x}, \bar{x} \rangle|$$
$$\lesssim C'\sqrt{\Delta}D_x\|\bar{x}\| + \sqrt{\Delta}D_x\|\bar{x}\|.$$

Applying the bound on $D_x$, the result is shown. Lastly, using $\|x\| \leq \tau\|\bar{x}\|$ and $\|\hat{x}\| \lesssim \|\bar{x}\|$, we conclude

$$\|x\| \|x - \bar{x}\| \|x + \bar{x}\| \leq (\|x - \hat{x}\| + \|\hat{x}\|)(\|x - \bar{x}\| \|x + \bar{x}\|)$$
$$\leq D_x^2\|x - \hat{x}\| + \|\hat{x}\| \|x - \bar{x}\| \|x + \bar{x}\|$$
$$\lesssim \|\bar{x}\|^2 D_x\sqrt{\Delta} + \|\hat{x}\|(\|\hat{x} - \bar{x}\| + \|\hat{x} - x\|)\|x + \bar{x}\|$$
$$\lesssim \|\bar{x}\|^3\sqrt{\Delta} + \|\bar{x}\|^2\sqrt{\Delta}D_x + \|\hat{x}\| \|\hat{x} - \bar{x}\| \|x + \bar{x}\|$$
$$\lesssim \|\bar{x}\|^3\sqrt{\Delta} + \|\hat{x}\| \|\hat{x} - \bar{x}\|(\|x - \hat{x}\| + \|\hat{x} + \bar{x}\|)$$
$$\lesssim \|\bar{x}\|^3\sqrt{\Delta} + \|\bar{x}\|^2 D_x\sqrt{\Delta} + \|\hat{x}\| \|\hat{x} - \bar{x}\| \|\hat{x} + \bar{x}\|$$
$$\lesssim \|\bar{x}\|^3\sqrt{\Delta} + \sqrt{\Delta}D_x\|\bar{x}\|^2.$$

Dividing through by $\|\bar{x}\|^3$ finishes the proof. $\square$

## B. Proof of Theorem 5.4

In this section, we will prove Theorem 5.4. Contrasting with Theorem 5.3, the proof of Theorem 5.4 is much more delicate, in large part relying on perturbation bounds on eigenvalues; e.g. Gershgorin theorem (Horn & Johnson, 2013, Corollary 6.1.3). We continue using the notation of Section 5.3. Namely, fix a symmetric convex function $f \colon \mathbb{R}^d \to \mathbb{R}$ and a point $\bar{x} \in \mathbb{R}^d \setminus \{0\}$, and define the function

$$g(x) := f_\lambda(xx^T - \bar{x}\bar{x}^T).$$

The chain rule directly implies

$$\partial g(x) = 2\partial f_\lambda(X)x.$$

Therefore, using Theorem 5.2 let us also fix a matrix $V \in \partial f_\lambda(X)$ and a matrix $U \in \mathbb{O}^d$ satisfying

$$\lambda(V) \in \partial f(\lambda(X)), \qquad V = U\mathrm{Diag}(\lambda(V))U^T \qquad \text{and} \qquad X = U\mathrm{Diag}(\lambda(X))U^T.$$

We begin with two technical lemmas.

LEMMA B.1    Suppose that there exists $\kappa > 0$ such that the inequality

$$g(x) - g(\bar{x}) \geq \kappa \|x - \bar{x}\| \|x + \bar{x}\| \qquad \text{holds for all } x \in \mathbb{R}^d.$$

Then for any $x \notin \{\pm\bar{x}\}$, we have $\max\{|\lambda_1(V)|, |\lambda_d(V)|\} \geq \kappa/2$.

*Proof.*    Using Lemma 5.11, for $i \in \{1, d\}$ we obtain

$$|\lambda_i(X)| = |\langle U_i, x\rangle^2 - \langle U_i, \bar{x}\rangle^2| = |\langle U_i, x - \bar{x}\rangle\langle U_i, x + \bar{x}\rangle| \leq \|x - \bar{x}\| \|x + \bar{x}\|.$$

Taking into account (5.6), yields

$$\kappa\|x - \bar{x}\|\|x + \bar{x}\| \leq g(x) - g(\bar{x}) \leq \lambda_1(V)\lambda_1(X) + \lambda_d(V)\lambda_d(X)$$
$$\leq 2\max\{|\lambda_1(V)|, |\lambda_d(V)|\}\|x - \bar{x}\|\|x + \bar{x}\|,$$

as desired.    □

LEMMA B.2    Suppose that there exists $\kappa > 0$ such that the inequality

$$g(x) - g(\bar{x}) \geq \kappa \|x - \bar{x}\| \|x + \bar{x}\| \qquad \text{holds for all } x \in \mathbb{R}^d. \tag{B.1}$$

Then any point $x \in \mathbb{R}^d \setminus \{0\}$ satisfies

$$2\left(\frac{\kappa\|x - \bar{x}\|\|x + \bar{x}\|}{\|x\|} - \frac{(|\lambda_1(V)| + |\lambda_d(V)|)\|\bar{x}\|^2}{\|x\|}\right) \leq \mathrm{dist}(0; \partial g(x)).$$

*Proof.*    First, note that for $x \in \{\pm\bar{x}\}$, the result holds trivially, so we may assume $x \notin \{\pm\bar{x}\}$. Recall the equality $\partial g(x) = 2\partial f_\lambda(X)x$. Fix now a vector $V \in \partial f_\lambda(X)$ satisfying $\mathrm{dist}(0; \partial g(x)) = 2\|Vx\|$. Using convexity, we deduce

$$g(x) - g(\bar{x}) = f_\lambda(xx^T - \bar{x}\bar{x}^T) - f_\lambda(0) \leq \langle V, xx^T - \bar{x}\bar{x}^T\rangle \leq \frac{\|x\|\mathrm{dist}(0, \partial g(x))}{2} + |\bar{x}^T V\bar{x}|. \tag{B.2}$$

We next upper bound the term $|\bar{x}^T V\bar{x}|$. To this end, fix a matrix $U \in \mathbb{O}^d$ satisfying $V = U\mathrm{Diag}(\lambda(V))U^T$ and $X = U\mathrm{Diag}(\lambda(X))U^T$, and such that the inclusion $\lambda(V) \in \partial f(\lambda(X))$ holds. Taking into account $\bar{x} \in \mathrm{span}\{U_1, U_d\}$ (Lemma 5.11), we deduce

$$|\bar{x}^T V\bar{x}| = |\lambda_1(V)\langle U_1, \bar{x}\rangle^2 + \lambda_d(V)\langle U_d, \bar{x}\rangle^2| \leq (|\lambda_1(V)| + |\lambda_d(V)|)\|\bar{x}\|^2.$$

Combining this estimate with (B.2) and (B.1) completes the proof. □

We next prove a quantitative version of Corollary 5.13. The argument follows a similar outline.

THEOREM B.3 (Quantitative Version of Corollary 5.13). Suppose that there exists a constant $\kappa > 0$ such that the inequality

$$g(y) - g(\bar{x}) \geq \kappa \|y - \bar{x}\| \|y + \bar{x}\| \qquad \text{holds for all} y \in \mathbb{R}^d.$$

Suppose $|\lambda_1(V)|, |\lambda_d(V)|$ are both upper bounded by a numerical constant[2] and set $\varepsilon := \|Vx\|$. Then there exists a numerical constant $\gamma > 0$, such that whenever $\varepsilon \leq \gamma \cdot \|\bar{x}\|$, we have that $\|x\| \lesssim \|\bar{x}\|$ and $x$ satisfies either

$$\|x\| \|x - \bar{x}\| \|x + \bar{x}\| \lesssim \varepsilon \|\bar{x}\|^2 \qquad \text{or} \qquad \begin{cases} |\lambda_1(V)| \lesssim \varepsilon/\|x\| \\ |\langle x, \bar{x} \rangle| \lesssim \varepsilon \|\bar{x}\| \end{cases}.$$

*Proof.* Clearly, we may suppose $x \notin \{0, \pm\bar{x}\}$ and $\varepsilon \neq 0$, since otherwise the theorem would hold vacuously. We will prove the following precise bound, which immediately implies the statement of the theorem; there exists a numerical constant $\gamma > 0$, such that whenever $\varepsilon \leq \gamma \|\bar{x}\|$, the inequalities $\|x\| \leq \delta \|\bar{x}\|$ and

$$\min\left\{ \frac{\|x - \bar{x}\| \|x + \bar{x}\|}{\frac{2}{\kappa} \max\left\{ \left( \frac{\|x\|}{\sqrt{2}} + \frac{\sqrt{2}\|\bar{x}\|^2}{\|x\|} \right), \frac{\|x\|(\kappa\sqrt{2}+2|\lambda_d(V)|)}{\kappa} \right\}}, \right.$$
$$\left. \max\left\{ \frac{|\lambda_1(V)| \|x\|}{\sqrt{2}}, \frac{\kappa|\langle x, \bar{x} \rangle|}{2\sqrt{2}\delta\|x\| + 2\|\bar{x}\|} \right\} \right\} \leq \|Vx\| \qquad \text{(B.3)}$$

hold, where we define the numerical constant

$$\delta := \sqrt{\frac{2(|\lambda_1(V)| + |\lambda_d(V)|)}{\kappa}} + 1.$$
□

As a first step, we show that $\|x\|$ is within a numerical constant of $\|\bar{x}\|$.

Claim B.4 Provided $\gamma < \frac{\kappa(1-1/\delta)^2}{2}$, the inequality, $\|x\| \leq \delta \|\bar{x}\|$, holds.

*Proof.* Assume for sake of contradiction $\frac{\|x\|}{\|\bar{x}\|} > \delta := \sqrt{\frac{2(|\lambda_1(V)|+|\lambda_d(V)|)}{\kappa}} + 1$. Lemma B.1 shows $\max\{|\lambda_1(V)|, |\lambda_d(V)|\} \geq \frac{\kappa}{2}$, and therefore $\delta > 1$. Using the bound $\text{dist}(0; \partial g(x)) \leq 2\|Vx\| = 2\varepsilon$ and Lemma B.2, we deduce

$$\frac{\kappa \|x - \bar{x}\| \|x + \bar{x}\|}{\|x\| \|\bar{x}\|} - \frac{\varepsilon}{\|\bar{x}\|} \leq \frac{(|\lambda_1(V)| + |\lambda_d(V)|)\|\bar{x}\|^2}{\|x\| \|\bar{x}\|}.$$

Clearly, we have

$$\frac{\kappa \|x - \bar{x}\| \|x + \bar{x}\|}{\|x\| \|\bar{x}\|} \geq \frac{\kappa(\|x\| - \|\bar{x}\|)^2}{\|x\| \|\bar{x}\|} \geq \frac{\|x\|}{\|\bar{x}\|} \frac{\kappa(\|x\| - \|\bar{x}\|)^2}{\|x\|^2} \geq \frac{\|x\|}{\|\bar{x}\|} \kappa(1 - 1/\delta)^2,$$

---

[2] This holds whenever $(t, s) \mapsto f(t, s, 0, \ldots, 0)$ is Lipschitz continuous.

where the first inequality follows because $\|x\| > \delta\|\bar{x}\| \geq \|\bar{x}\|$. Let us now choose $\gamma < \frac{\kappa(1-1/\delta)^2}{2}$, thereby guaranteeing $\frac{\varepsilon}{\|\bar{x}\|} \leq \frac{\kappa(1-1/\delta)^2}{2}$. Hence, we obtain

$$\frac{\|x\|}{\|\bar{x}\|} \cdot \frac{\kappa(1-1/\delta)^2}{2(|\lambda_1(V)| + |\lambda_d(V)|)} \leq \frac{1}{|\lambda_1(V)| + |\lambda_d(V)|}\left(\frac{\kappa\|x-\bar{x}\|\|x+\bar{x}\|}{\|x\|\|\bar{x}\|} - \frac{\varepsilon}{\|\bar{x}\|}\right) \leq \frac{\|\bar{x}\|}{\|x\|}.$$

Rearranging yields

$$\frac{\kappa}{2(|\lambda_1(V)| + |\lambda_d(V)|)} \leq \frac{\|\bar{x}\|^2}{\|x\|^2(1-1/\delta)^2} < \frac{1}{(\delta-1)^2},$$

a contradiction. $\qquad\square$

Looking back at the expression, define the values:

$$\rho_1 = \frac{\delta\|\bar{x}\|}{\sqrt{2}} \quad \text{and} \quad \rho_3 = \frac{2}{\kappa}\max\left\{\left(\frac{\|x\|}{\sqrt{2}} + \frac{\sqrt{2}\|\bar{x}\|^2}{\|x\|}\right), \frac{\|x\|(\kappa\sqrt{2} + 2|\lambda_d(V)|)}{\kappa}\right\}.$$

Notice that the inequality, $\varepsilon\rho_3 \geq \|x-\bar{x}\|\|x+\bar{x}\|$, would immediately imply the validity of the theorem. Thus, we assume $\varepsilon\rho_3 < \|x-\bar{x}\|\|x+\bar{x}\|$ throughout. It suffices now to show

$$|\lambda_1(V)| \leq \varepsilon/\rho_1 \quad \text{and} \quad |\langle x,\bar{x}\rangle| \leq \frac{\varepsilon}{\kappa}\left(2\sqrt{2}\delta\|x\| + \|\bar{x}\|\right).$$

We do so in order. We begin by observing that the inequality (5.5) guarantees

$$\max\{|\lambda_1(V)\langle U_1, x\rangle|, |\lambda_d(V)\langle U_d, x\rangle|\} \leq \varepsilon. \tag{B.4}$$

**Claim B.5** The inequality $|\lambda_1(V)| < \varepsilon/\rho_1$ holds.

*Proof.* Let us assume the contrary, $|\lambda_1(V)| \geq \varepsilon/\rho_1$. Inequality (B.4) then implies $|\langle U_1, x\rangle| \leq \rho_1$, while Lemma 5.11 in turn guarantees

$$0 \leq \lambda_1(X) = \langle U_1, x\rangle^2 - \langle U_1, \bar{x}\rangle^2 \leq \rho_1^2.$$

Taking into account $\langle U_1, x\rangle^2 + \langle U_d, x\rangle^2 = \|x\|^2$ (Lemma 5.11, correlation), we deduce $\langle U_d, x\rangle^2 \geq \|x\|^2 - \rho_1^2$. Combining this with (B.4), we deduce

$$|\lambda_d(V)| \leq \frac{\varepsilon}{|\langle U_d, x\rangle|} \leq \frac{\varepsilon}{\sqrt{\|x\|^2 - \rho_1^2}}.$$

Therefore, using the correlation inequality (5.6), we find

$$\varepsilon\rho_3\kappa < \kappa\|x-\bar{x}\|\|x+\bar{x}\| \leq g(x) - g(\bar{x}) \leq \lambda_1(V)\lambda_1(X) + \lambda_d(V)\lambda_d(X)$$

$$\leq |\lambda_1(V)|(\langle U_1, x\rangle^2 - \langle U_1, \bar{x}\rangle^2) + \frac{\varepsilon}{\sqrt{\|x\|^2 - \rho_1^2}}\left(\langle U_d, \bar{x}\rangle^2 - \langle U_d, x\rangle^2\right)$$

$$\leq \varepsilon|\langle U_1, x\rangle| + \frac{\varepsilon\langle U_d, \bar{x}\rangle^2}{\sqrt{\|x\|^2 - \rho_1^2}}$$

$$\leq \varepsilon\left(\rho_1 + \frac{\|\bar{x}\|^2}{\sqrt{\|x\|^2 - \rho_1^2}}\right).$$

Dividing through by $\varepsilon$ and plugging in the value of $\rho_1$ yields

$$\rho_3 \kappa < \frac{\|x\|}{\sqrt{2}} + \frac{\sqrt{2}\|\bar{x}\|^2}{\|x\|},$$

which contradicts the definition of $\rho_3$. □

Let us now decrease $\gamma > 0$ further by ensuring $\gamma < \min\left\{\frac{\kappa(1-1/\delta)^2}{2}, \frac{\kappa}{2\sqrt{2}}\right\}$. Thus, from Claim B.5 and our standing assumption $\|Vx\| \leq \frac{\kappa\|x\|}{2\sqrt{2}}$, we conclude

$$|\lambda_1(V)| < \frac{\sqrt{2}\varepsilon}{\|x\|} < \frac{\kappa}{2}.$$

Lemma B.1 guarantees, $\max\{|\lambda_1(V)|, |\lambda_d(V)|\} \geq \kappa/2$; thus, we deduce $|\lambda_d(V)| \geq \kappa/2$. Applying (B.4), we find that

$$|\langle U_d, x\rangle| \leq \frac{\varepsilon}{|\lambda_d(V)|} \leq \frac{2\varepsilon}{\kappa}. \tag{B.5}$$

Thus, by Lemma 5.11, we have

$$|\langle U_1, \bar{x}\rangle\langle U_d, \bar{x}\rangle| = |\langle U_1, x\rangle\langle U_d, x\rangle| \leq \frac{2\|x\|\varepsilon}{\kappa}. \tag{B.6}$$

**Claim B.6** The inequality $|\langle U_d, \bar{x}\rangle| > |\langle U_1, \bar{x}\rangle|$ holds.

*Proof.* Let us assume the contrary $|\langle U_d, \bar{x}\rangle| \leq |\langle U_1, \bar{x}\rangle|$. Then from (B.6) we obtain[3] $\langle U_d, \bar{x}\rangle^2 < \frac{2\|x\|\varepsilon}{\kappa}$. Hence, from Lemma 5.11, we find that $|\lambda_d(X)| \leq \langle U_d, \bar{x}\rangle^2 \leq \frac{2\|x\|\varepsilon}{\kappa}$. Putting these facts together with the correlation inequality (5.6), we successively deduce

$$\varepsilon\rho_3\kappa < \kappa\|x - \bar{x}\|\|x + \bar{x}\| \leq g(x) - g(\bar{x}) \leq \lambda_1(V)\lambda_1(X) + |\lambda_d(V)| \cdot |\lambda_d(X)|$$

$$\leq \frac{\sqrt{2}\varepsilon}{\|x\|} \cdot \lambda_1(X) + |\lambda_d(V)| \cdot \frac{2\|x\|\varepsilon}{\kappa}$$

$$\leq \frac{\kappa\sqrt{2}\varepsilon\|x\|}{\kappa} + \frac{2|\lambda_d(V)|\varepsilon\|x\|}{\kappa},$$

where the last inequality uses the bound $\lambda_1(X) \leq \|x\|^2$. Therefore, we have reached a contradiction to the definition of $\rho_3$. □

Combining Claim B.6 with the expression $\langle U_1, \bar{x}\rangle^2 + \langle U_d, \bar{x}\rangle^2 = \|\bar{x}\|^2$, we conclude $\langle U_d, \bar{x}\rangle^2 \geq \frac{\|\bar{x}\|^2}{2}$. Therefore, (B.6) and Claim B.4 imply the strong result:

$$|\langle U_1, \bar{x}\rangle| \leq \frac{2\sqrt{2}\varepsilon\|x\|}{\kappa\|\bar{x}\|} \leq \frac{2\sqrt{2}\varepsilon\delta}{\kappa}. \tag{B.7}$$

---

[3] If $ab < \delta$, then $\min\{a, b\}^2 < \delta$.

Thus, combining Claim B.4, Lemma 5.11 and (B.5) we conclude

$$|\langle x, \bar{x}\rangle| = |\langle U_1, x\rangle\langle U_1, \bar{x}\rangle + \langle U_d, x\rangle\langle U_d, \bar{x}\rangle| \leq |\langle U_1, \bar{x}\rangle| \cdot \|x\| + |\langle U_d, x\rangle| \cdot \|\bar{x}\|$$

$$\leq \frac{\varepsilon}{\kappa}\left(2\sqrt{2}\delta\|x\| + 2\|\bar{x}\|\right).$$

The proof is complete.

In order to interpret the conclusion of Theorem B.3 on the phase retrieval objective $f_P$, we must show that the condition

$$\left\{ \begin{array}{c} |\lambda_1(V)| \lesssim \varepsilon/\|x\| \\ |\langle x, \bar{x}\rangle| \lesssim \varepsilon\|\bar{x}\| \end{array} \right\}$$

guarantees that the equation $\|x\| = c \cdot \|\bar{x}\|$ almost holds, where $c$ is defined in Theorem 5.3. This is the content of the following two lemmas. Note that it is easy to verify the equality $\lambda_1(V) = \nabla_{y_1}\zeta(y_1, y_2)$, where we set $(y_1, y_2) := (\lambda_1(X), \lambda_d(X))$.

LEMMA B.7 (Extension of Lemma 5.9). Fix a real constant $0 \leq \varepsilon < 1$. The solutions of the inequality $|\nabla_{y_1}\zeta(y_1, y_2)| \leq \varepsilon$ on $\mathbb{R}_{++} \times \mathbb{R}_{--}$ are precisely the elements of the open cone

$$\{(c^2 y, -y) : 0 < y, 0 < c_1 \leq c \leq c_2\},$$

where $c_1, c_2$ are the unique solutions of the equations

$$\frac{\pi}{4}(1 + \varepsilon) = \frac{c_2}{1 + c_2^2} + \arctan\left(c_2\right),$$

and

$$\frac{\pi}{4}(1 - \varepsilon) = \frac{c_1}{1 + c_1^2} + \arctan\left(c_1\right).$$

Moreover, considering $c_1$ and $c_2$ as functions of $\varepsilon$, we have $c_2(\varepsilon) - c_1(\varepsilon) \leq 5\pi\varepsilon$ whenever $0 < \varepsilon < 1/2$.

*Proof.* The proof is completely analogous to that of Lemma 5.9. We leave the details to the reader. The only point worth commenting is the inequality $c_2(\varepsilon) - c_1(\varepsilon) \leq 5\pi\varepsilon$ whenever $0 < \varepsilon < 1/2$. To get this bound, observe that $0 < c_2(\varepsilon) \leq c_2(0.5) \leq 0.83$ for all $\varepsilon \leq 1/2$ as $c_2$ is a increasing function of $\varepsilon$. Therefore,

$$\frac{\pi}{2}\varepsilon = \frac{c_2}{1 + c_2^2} - \frac{c_1}{1 + c_1^2} + \arctan(c_2) - \arctan(c_1) \geq \frac{c_2}{1 + c_2^2} - \frac{c_1}{1 + c_1^2} = \frac{1 - c_1 c_2}{(1 + c_1^2)(1 + c_2^2)}(c_2 - c_1)$$

(B.8)

$$\geq \frac{1 - c_2^2}{(1 + c_1^2)(1 + c_2^2)}(c_2 - c_1) \geq \frac{1 - c_2^2}{(1 + c_2^2)^2}(c_2 - c_1).$$

Thus, we have

$$c_2(\varepsilon) - c_1(\varepsilon) \leq \frac{\pi\varepsilon}{2}\frac{(1 + c_2^2(\varepsilon))^2}{1 - c_2^2(\varepsilon)} \leq \frac{\pi\varepsilon}{2}\frac{(1 + 0.83^2)^2}{1 - 0.83^2} \leq 5\pi\varepsilon,$$

as claimed.                                                                                       □

LEMMA B.8 Fix a real constant $0 \le \varepsilon < \frac{1}{3}$ and vectors $x, \bar{x} \in \mathbb{R}^d \setminus \{0\}$. Suppose $\lambda_1(X) = -c^2 \lambda_d(X)$ for some real constant $c > 0$ and $|\langle x, \bar{x} \rangle| \le \varepsilon \|\bar{x}\| \|x\|$. Then we have

$$1 - (1 + c^2)(\varepsilon + \varepsilon^2) \le c^2 \frac{\|\bar{x}\|^2}{\|x\|^2} \le 1 + (1 + c^2)(\varepsilon + \varepsilon^2).$$

*Proof.* Fix a decomposition $x = \frac{\langle x, \bar{x} \rangle}{\|\bar{x}\|^2} \bar{x} + v$, where $v \in \bar{x}^{\perp}$. Note inequality $|\langle x, \bar{x} \rangle| \le \varepsilon \|x\| \|\bar{x}\|$ implies that $\bar{x}$ and $x$ are not collinear, and therefore $\|v\| > 0$. Define the constant $\alpha = \frac{\langle x, \bar{x} \rangle}{\|\bar{x}\|^2}$. Then a quick computation shows the following decomposition:

$$X = \begin{bmatrix} \frac{\bar{x}}{\|\bar{x}\|} & \frac{v}{\|v\|} \end{bmatrix} \begin{bmatrix} (\alpha^2 - 1)\|\bar{x}\|^2 & \alpha \|\bar{x}\| \|v\| \\ \alpha \|\bar{x}\| \|v\| & \|v\|^2 \end{bmatrix} \begin{bmatrix} \frac{\bar{x}}{\|\bar{x}\|} & \frac{v}{\|v\|} \end{bmatrix}^T.$$

Notice that the above $2 \times 2$-matrix is invertible, and therefore its eigenvalues must be $\lambda_1(X)$ and $\lambda_d(X)$. By the Gershgorin theorem (Horn & Johnson, 2013, Corollary 6.1.3) applied to the $2 \times 2$ matrix, we know that $\lambda_1(X)$ and $\lambda_d(X)$ must lie in the union of the intervals

$$\bar{D}_1 = \{z : |z - \|v\|^2| \le |\alpha| \|\bar{x}\| \|v\|\} \quad \text{and} \quad \bar{D}_2 = \{z : |z - (\alpha^2 - 1)\|\bar{x}\|^2| \le |\alpha| \|\bar{x}\| \|v\|\}.$$

We next prove the following claim.

Claim B.9 The intervals $\bar{D}_1$ and $\bar{D}_2$ are contained in the following intervals around $\|x\|^2$ and $-\|\bar{x}\|^2$, respectively:

$$\bar{D}_1 \subset D_1 := \{z : |z - \|x\|^2| \le (\varepsilon^2 + \varepsilon)\|x\|^2\},$$
$$\bar{D}_2 \subset D_2 := \{z : |z + \|\bar{x}\|^2| \le (\varepsilon^2 + \varepsilon)\|x\|^2\}.$$

Moreover, we have $D_1 \cap D_2 = \emptyset$ and $D_1 \subset \mathbb{R}_{++}$.

*Proof.* Consider the interval $\bar{D}_1$. A routine computation shows

$$|\alpha| \le \frac{\varepsilon \|x\|}{\|\bar{x}\|}, \quad 0 \le \alpha^2 \le \frac{\varepsilon^2 \|x\|^2}{\|\bar{x}\|^2} \quad \text{and} \quad 0 \le \alpha \langle x, \bar{x} \rangle \le \varepsilon^2 \|x\|^2.$$

Using $\|x\| \ge \|v\|$ and $\|v\|^2 = \|x\|^2 - 2\alpha \langle x, \bar{x} \rangle + \alpha^2 \|\bar{x}\|^2$, we successively deduce for any $z \in \bar{D}_1$, the inequalities

$$
\begin{array}{ccc}
-|\alpha| \|\bar{x}\| \|x\| & \le & z - \|v\|^2 & \le & |\alpha| \|\bar{x}\| \|x\| \\
-\varepsilon \|x\|^2 & \le & z - \|x\|^2 + 2\alpha \langle x, \bar{x} \rangle - \alpha^2 \|\bar{x}\|^2 & \le & \varepsilon \|x\|^2 \\
-\varepsilon \|x\|^2 + \alpha^2 \|\bar{x}\|^2 - 2\alpha \langle x, \bar{x} \rangle & \le & z - \|x\|^2 & \le & \varepsilon \|x\|^2 + \alpha^2 \|\bar{x}\|^2 - 2\alpha \langle x, \bar{x} \rangle \\
-\varepsilon \|x\|^2 - \varepsilon^2 \|x\|^2 & \le & z - \|x\|^2 & \le & \varepsilon \|x\|^2 + \varepsilon^2 \|x\|^2.
\end{array}
$$

Thus, we have shown $\bar{D}_1 \subset D_1$. Similarly, for all $z \in \bar{D}_2$, we compute

$$
\begin{array}{ccc}
-|\alpha| \|\bar{x}\| \|v\| & \le & z - (\alpha^2 - 1)\|\bar{x}\|^2 & \le & |\alpha| \|\bar{x}\| \|v\| \\
-\varepsilon \|x\|^2 + \alpha^2 \|\bar{x}\|^2 & \le & z + \|\bar{x}\|^2 & \le & \varepsilon \|x\|^2 + \alpha^2 \|\bar{x}\|^2 \\
-\varepsilon \|x\|^2 & \le & z + \|\bar{x}\|^2 & \le & \varepsilon \|x\|^2 + \varepsilon^2 \|x\|^2.
\end{array}
$$

We conclude $\bar{D}_2 \subset D_2$. Provided $\|x\| \ne 0$ and $\varepsilon^2 + \varepsilon < 1$, it is clear $D_1 \subset \mathbb{R}_{++}$. It remains to show that $D_2 \cap D_1 = \emptyset$. Clearly, it is sufficient to guarantee that the sum of the radii of $D_2$ and $D_1$ is strictly

smaller than the distance between the centres:

$$(\varepsilon^2 + \varepsilon)\|x\|^2 + (\varepsilon^2 + \varepsilon)\|x\|^2 < \|x\|^2 - (-\|\bar{x}\|^2).$$

Rearranging, we must guarantee $2(\varepsilon^2 + \varepsilon) - 1 < \frac{\|\bar{x}\|^2}{\|x\|^2}$. Clearly, this is the case as soon as $\varepsilon < 1/3$. The result follows. □

Thus, we have proved $D_1 \cap D_2 = \emptyset$ and $D_1 \subset \mathbb{R}_{++}$. Since $\bar{D}_1$ and $\bar{D}_2$, each contains at least one eigenvalue, it must be the case that $\lambda_d(X)$ lies in $\bar{D}_2$ and $\lambda_1(X)$ lies in $\bar{D}_1$. We thus conclude

$$\left| \lambda_1(X) - \|x\|^2 \right| \le (\varepsilon^2 + \varepsilon)\|x\|^2$$

$$\left| \lambda_d(X) + \|\bar{x}\|^2 \right| \le (\varepsilon^2 + \varepsilon)\|x\|^2.$$

Writing $\lambda_1(X) = -c^2\lambda_d(X)$, we obtain

$$\left| -c^2\lambda_d(X) - c^2\|\bar{x}\|^2 + c^2\|\bar{x}\|^2 - \|x\|^2 \right| \le (\varepsilon^2 + \varepsilon)\|x\|^2,$$

and hence

$$\left| \|x\|^2 - c^2\|\bar{x}\|^2 \right| \le (1 + c^2)(\varepsilon^2 + \varepsilon)\|x\|^2.$$

The result follows. □

Combining Lemmas B.7 and B.8, we arrive at the following.

COROLLARY B.10 (Small $\lambda_1(V)$ and near orthogonality). Fix a real constant $0 \le \varepsilon \le \frac{1}{8}$ and consider a point $x \in \mathbb{R}^d \setminus \{0\}$ satisfying $|\nabla_{y_1} \zeta(\lambda_1(X), \lambda_d(X))| \le \varepsilon$ and $|\langle x, \bar{x} \rangle| \le \varepsilon \|x\| \|\bar{x}\|$. Then $x$ satisfies

$$\left| \|x\| - c\|\bar{x}\| \right| \le 26\varepsilon\|\bar{x}\|,$$

where $c$ is the solution of the equation $\frac{\pi}{4} = \frac{c}{1+c^2} + \arctan(c)$.

*Proof.* Define the quantities $c_1(\varepsilon)$ and $c_2(\varepsilon)$ to be the solutions of the equations

$$\frac{\pi}{4}(1 - \varepsilon) = \frac{c_1}{1 + c_1^2} + \arctan(c_1),$$

$$\frac{\pi}{4}(1 + \varepsilon) = \frac{c_2}{1 + c_2^2} + \arctan(c_2),$$

respectively. First, since $c_2(\cdot)$ is an increasing function, it is easy to verify $c_2(\varepsilon) < 1$ whenever $0 < \varepsilon \le \frac{1}{8}$; thus, we have $2\varepsilon(1 + c_2^2(\varepsilon)) < \frac{1}{2}$. By Lemma B.7, we know that whenever $|\nabla_{y_1} \zeta(\lambda_1(X), \lambda_d(X))| \le \varepsilon$, there exists $\hat{c}$ satisfying $\lambda_1(X) = -\hat{c}^2\lambda_d(X)$ and $0 < c_1(\varepsilon) \le \hat{c} \le c_2(\varepsilon)$. Lemma B.8, in turn, implies

$$(1 - 2\varepsilon(1 + \hat{c}^2))\|x\|^2 \le \hat{c}^2\|\bar{x}\|^2 \le (1 + 2\varepsilon(1 + \hat{c}^2))\|x\|^2.$$

Looking at the right-hand side, we deduce

$$c_1^2(\varepsilon)\|\bar{x}\|^2 \le \hat{c}^2\|\bar{x}\|^2 \le \left(1 + 2\varepsilon(1 + c_2^2(\varepsilon))\right)\|x\|^2,$$

while looking at the left-hand side yields

$$(1 - 2\varepsilon(1 + c_2^2(\varepsilon)))\|x\|^2 \le \hat{c}^2\|\bar{x}\|^2 \le c_2^2(\varepsilon)\|\bar{x}\|^2.$$

Isolating $\|x\|^2$ and taking square roots we obtain

$$\frac{c_1(\varepsilon)}{\sqrt{1 + 2\varepsilon(1 + c_2^2(\varepsilon))}} \|\bar{x}\| \leq \|x\| \leq \frac{c_2(\varepsilon)}{\sqrt{1 - 2\varepsilon(1 + c_2^2(\varepsilon))}} \|\bar{x}\|. \tag{B.9}$$

Applying Lemma B.7 and the inequality $c_2(\varepsilon) < 1$, we upper bound the right-hand side:

$$\frac{c_2(\varepsilon)}{\sqrt{1 - 2\varepsilon(1 + c_2^2(\varepsilon))}} \leq \frac{5\pi\varepsilon + c}{\sqrt{1 - 4\varepsilon}}$$

$$= c\left(1 + \frac{5\pi\varepsilon/c + 1 - \sqrt{1 - 4\varepsilon}}{\sqrt{1 - 4\varepsilon}}\right) \leq c\left(1 + \frac{5\pi\varepsilon/c + 4\varepsilon}{\sqrt{1/2}}\right) \leq c(1 + 57\varepsilon).$$

Exactly the same reasoning shows

$$\frac{c_1(\varepsilon)}{\sqrt{1 + 2\varepsilon(1 + c_2^2(\varepsilon))}} \geq c(1 - 57\varepsilon).$$

Thus, the inequality $\big|\|x\| - c\|\bar{x}\|\big| \leq 57c\varepsilon\|\bar{x}\| \leq 26\varepsilon\|\bar{x}\|$ holds, as claimed. □

We are now ready to prove the inexact extension of Theorem 5.1.

*Proof of Theorem* 5.4. We use the decomposition $g = f_P(X)$ and $f = \varphi$. Let us verify that we may apply Theorem B.3. To this end, observe that the population objective satisfies

$$f_P(x) - f_P(\bar{x}) = \mathbb{E}_a\left[\left\langle a, \frac{x - \bar{x}}{\|x - \bar{x}\|}\right\rangle\left\langle a, \frac{x + \bar{x}}{\|x + \bar{x}\|}\right\rangle\right] \|x - \bar{x}\|\|x + \bar{x}\| \geq \kappa\|x - \bar{x}\|\|x + \bar{x}\|$$

for the numerical constant $\kappa$ (Eldar & Mendelson, 2014, Corollary 3.7). Moreover, clearly $\zeta$ is globally Lipschitz (being a norm), and therefore $|\lambda_1(V)|$ and $|\lambda_d(V)|$ are bounded by a numerical constant. Thus, provided $\varepsilon := \|Vx\|$ satisfies $\varepsilon \leq \gamma \cdot \|x\|$ for the numerical constant $\gamma$, we can be sure that $x$ satisfies either

$$\|x\|\|x - \bar{x}\|\|x + \bar{x}\| \lesssim \varepsilon\|\bar{x}\|^2 \qquad \text{or} \qquad \left\{\begin{array}{c} |\lambda_1(V)| \lesssim \varepsilon/\|x\| \\ |\langle x, \bar{x}\rangle| \lesssim \varepsilon\|\bar{x}\| \end{array}\right\}.$$

Now suppose the latter is the case, and let $C$ be a numerical constant satisfying $|\lambda_1(V)| \leq C\varepsilon/\|x\|$ and $|\langle x, \bar{x}\rangle| \leq C\varepsilon\|\bar{x}\|$. We aim to apply Corollary B.10. To do so, we must ensure

$$|\lambda_1(V)| \leq \frac{C\varepsilon}{\|x\|} \leq \frac{1}{8} \quad \text{and} \quad \left|\left\langle \frac{x}{\|x\|}, \frac{\bar{x}}{\|\bar{x}\|}\right\rangle\right| \leq C\varepsilon \cdot \frac{\|\bar{x}\|}{\|x\|\|\bar{x}\|} = \frac{C\varepsilon}{\|x\|} \leq \frac{1}{8}.$$

Adjusting $\gamma$ if necessary, we can be sure that $\varepsilon/\|x\|$ is below $\frac{1}{8C}$. Applying Corollary B.10, with $\frac{C\varepsilon}{\|x\|}$ in place of $\varepsilon$, we conclude $\big|\|x\| - c\|\bar{x}\|\big| \lesssim \varepsilon\frac{\|\bar{x}\|}{\|x\|}$, as claimed. □

## B.1 *Comments on robustness*

We have thus far assumed that the measurement vector $b = (Ax)^2$ has not been corrupted by errant noise. In this section, we record a few straightforward extensions of earlier results, which hold if the measurements $b$ are noisy.

Assumption D.   Let $b_1, \ldots, b_m$ be $m$ i.i.d. copies of

$$\hat{b} = (a^T x)^2 + \delta \cdot \xi,$$

where $\delta \in \{0, 1\}$, $\xi \in \mathbb{R}$, and $a \in \mathbb{R}^d$ are independent random variables satisfying (1) $p_{\text{fail}} := P(\delta \neq 0) < 1$, (2) $\mathbb{E}[|\xi|] < \infty$, and (3) $a \sim \mathsf{N}(0, I_d)$.

Under corruption by $\delta \cdot \xi$, we define new population and subsampled objectives

$$\hat{f}_P(x) := \mathbb{E}_{a, \xi, \delta}[|(a^T x)^2 - (a^T \bar{x})^2 - \delta \cdot \xi|];$$

$$\hat{f}_S(x) := \frac{1}{m} \sum_{i=1}^m |(a_i^T x)^2 - b_i|.$$

Then by following the outline of the proof of Lemma 5.6, we arrive at a similar characterization of $\hat{f}_P$ as a spectral function.

LEMMA B.11  (Spectral representation of the population objective).  For all points $x \in \mathbb{R}^d$, equality holds:

$$\hat{f}_P(x) = \mathbb{E}_{v, \xi, \delta}\left[\left|\langle \lambda(X), v\rangle - \delta \cdot \xi_i\right|\right],$$

where $v_i \in \mathbb{R}$ are i.i.d. chi-squared random variables $v_i \sim \chi_1^2$.

Thus, we may write

$$\hat{f}_P(x) = \hat{\varphi}(\lambda(X)),$$

where $\hat{\varphi}$ is the convex symmetric function

$$\hat{\varphi}(z) := \mathbb{E}_{v, \xi, \delta}\left[\left|\langle z, v\rangle - \delta \cdot \xi_i\right|\right].$$

Moreover, provided that $\bar{x}$ is a minimizer of $\hat{f}_P$, the complete set of stationary points of $\hat{f}_P$ may be determined from Corollary 5.13. We prove this now.

LEMMA B.12   For all $x \in \mathbb{R}^d$, the following inequality holds:

$$\hat{f}_P(x) - \hat{f}_P(\pm \bar{x}) \geq (1 - 2p_{\text{fail}})f_P(x).$$

Consequently, if $p_{\text{fail}} < 1/2$, the points $\pm \bar{x}$ are the only minimizers of $\hat{f}_P$, and there exists a numerical constant $\kappa$ such that

$$\hat{f}_P(x) - \hat{f}_P(\pm \bar{x}) \geq \kappa(1 - 2p_{\text{fail}})\|x - \bar{x}\|\|x + \bar{x}\|.$$

*Proof.*   By expanding the difference, we find that

$$\hat{f}_P(x) - \hat{f}_P(\pm \bar{x}) = (1 - p_{\text{fail}})(f_P(x) - f_P(\bar{x})) + p_{\text{fail}}\mathbb{E}_{a, \xi}\left[|(a^T x)^2 - (a^T \bar{x})^2 - \xi| - |\xi|\right]$$

$$\geq (1 - p_{\text{fail}})f_P(x) - p_{\text{fail}}\mathbb{E}_a\left[|(a^T x)^2 - (a^T \bar{x})^2|\right]$$

$$\geq (1 - 2p_{\text{fail}})f_P(x).$$

Only the sharpness inequality is left to prove, but this is simply a consequence of the sharpness of $f_P$, which was proved in Eldar & Mendelson (2014, Corollary 3.7).                                              □

Therefore, by Corollary 5.13 we arrive at the complete characterization of the stationary points of $\hat{f}_P$.

THEOREM B.13 The set of stationary points of $\hat{f}_P$ are precisely

$$\{\pm x\} \cup \{0\} \cup \{x \mid \langle x, \bar{x} \rangle = 0, \text{ and } \exists \zeta \in \partial \hat{\varphi}(\lambda(X)), \max_i \{\zeta_i\} = 0\}.$$

The exact location of those stationary points orthogonal to $\bar{x}$ depends on the structure of the convex function $\hat{\varphi}$, which in turn depends the distribution of the noise $\delta \cdot \xi$. We will not attempt to characterize such $\hat{\varphi}$.

By the sharpness of $\hat{f}_P$, a quantitative version of Theorem B.13 immediately follows from Theorem B.3. When coupled together with a concentration inequality like that in Theorem 5.1, such a theorem would imply concentration of the subdifferential graphs of $\hat{f}_S$ and $\hat{f}_P$. We omit these straightforward details.