# Subgradient methods for sharp weakly convex functions

Damek Davis[*]    Dmitriy Drusvyatskiy[†]    Kellie J. MacPhee[‡]

Courtney Paquette[§]

**Abstract**

Subgradient methods converge linearly on a convex function that grows sharply away from its solution set. In this work, we show that the same is true for sharp functions that are only weakly convex, provided that the subgradient methods are initialized within a fixed tube around the solution set. A variety of statistical and signal processing tasks come equipped with good initialization, and provably lead to formulations that are both weakly convex and sharp. Therefore, in such settings, subgradient methods can serve as inexpensive local search procedures. We illustrate the proposed techniques on phase retrieval and covariance estimation problems.

## 1 Introduction

Typical methods for statistics and signal processing tasks follow the two-step strategy: (1) find a moderately accurate solution $\hat{x}$ at a low sample complexity cost (e.g., using spectral initialization), and (2) refine $\hat{x}$ by an iterative "local search algorithm" that converges rapidly under natural statistical assumptions. For smooth problem formulations, the term "local search" almost universally refers to gradient descent or a close variant thereof; see e.g. [1, 2, 4, 6, 20, 21, 23, 37]. For nonsmooth and nonconvex problems, the meaning of local search is much less clear. In this work, we ask the following question.

> Is there a generic gradient-based local search procedure for nonsmooth and nonconvex problems, which converges linearly under standard regularity conditions?

Not surprisingly, our approach is rooted in subgradient methods for convex optimization. To motivate the discussion, consider the constrained optimization problem

$$\min_{x \in \mathcal{X}} \ g(x), \tag{1.1}$$

[*]School of Operations Research and Information Engineering, Cornell University, Ithaca, NY 14850, USA; `people.orie.cornell.edu/dsd95/`.

[†]Department of Mathematics, U. Washington, Seattle, WA 98195; `www.math.washington.edu/~ddrusv`. Research of Drusvyatskiy was supported by the AFOSR YIP award FA9550-15-1-0237 and by the NSF DMS 1651851 and CCF 1740551 awards.

[‡]Department of Mathematics, U. Washington, Seattle, WA 98195; `sites.math.washington.edu/~kmacphee`

[§]Industrial and Systems Engineering Department, Lehigh University, Bethlehem, PA 18015; `sites.math.washington.edu/~yumiko88/`. Research of Paquette was supported by NSF CCF 1740796.

1

where $g$ is an $L$-Lipschitz convex function on $\mathbb{R}^d$ and $\mathcal{X}$ is a closed convex set. Given a current iterate $x_k$, subgradient methods proceed as follows:

$$\left\{ \begin{array}{l} \text{Choose any } \zeta_k \in \partial g(x_k) \\[2mm] \text{Set } x_{k+1} = \text{proj}_{\mathcal{X}}\left( x_k - \alpha_k \cdot \frac{\zeta_k}{\|\zeta_k\|} \right) \end{array} \right\}.$$

Here, the symbol $\text{proj}_{\mathcal{X}}(y)$ denotes the nearest point of $\mathcal{X}$ to $y$ and $\{\alpha_k\}$ is a specified stepsize sequence. The choice of the sequence $\{\alpha_k\}$ determines the behavior of the scheme, and is the main distinguishing feature among subgradient methods. In this work, we will only be interested in subgradient methods that are linearly convergent. As usual, linear rates of convergence of iterative methods require some regularity conditions to hold. Here, the appropriate regularity condition is *sharpness* [3, 31] (or equivalently a global error bound): there exists a real $\mu > 0$ satisfying

$$g(x) - \min_{x \in \mathcal{X}} g \geq \mu \cdot \text{dist}(x; \mathcal{X}^*) \qquad \text{for all } x \in \mathcal{X},$$

where $\mathcal{X}^*$ denotes the set of minimizers of (1.1). Assuming sharpness holds, subgradient methods, with a judicious choice of $\{\alpha_k\}$, produce iterates that converge to $\mathcal{X}^*$ at the linear rate $\sqrt{1 - (\mu/L)^2}$. Results of this type date back to 60's and 70's [18, 19, 29, 30, 33], while some more recent approaches have appeared in [22, 35, 38].

Various contemporary problems lead to formulations that are indeed sharp, but are only *weakly convex* and *locally Lipschitz*. Recall that a function $g$ is $\rho$-weakly convex [25] if the perturbed function $x \mapsto g(x) + \frac{\rho}{2}\| \cdot \|^2$ is convex for some $\rho > 0$. Note that weakly convex functions need not be smooth nor convex. A quick computation (Lemma 2.1) shows that if $g$ is $\mu$-sharp and $\rho$-weakly convex, then there is a tube around the solution set $\mathcal{X}^*$ that contains no extraneous stationary points:

$$\mathcal{T} := \left\{ x \in \mathcal{X} : \text{dist}(x; \mathcal{X}^*) \leq \frac{2\mu}{\rho} \right\}.$$

In this work, we show that the standard linearly convergent subgradient methods originally designed for convex problems, apply in this much greater generality, provided they are initialized within a slight contraction of the tube $\mathcal{T}$. The methods exhibit essentially the same linear rate of convergence as in the convex case, while the weak convexity constant $\rho$ only determines the validity of the initialization. We focus on three step-size rules: Polyak stepsize [18, 29], geometrically decaying step [19, 33], and constant stepsize [22, 35, 38]. As proof of concept, we illustrate the resulting algorithms on phase retrieval and covariance estimation problems.

Our current work sits within the broader scope of analyzing subgradient and proximal methods for weakly convex problems [9, 11, 13–16, 25, 26]; see also the recent survey [12]. In particular, the paper [9] proves a global sublinear rate of convergence, in terms of a natural stationarity measure, of a (stochastic) subgradient method on any weakly convex function. In contrast, here we are interested in subgradient methods that are locally linearly convergent under the additional sharpness assumption. The arguments we present are all quick modification of the proofs already available in the convex setting. Nonetheless, we believe that the drawn conclusions are interesting and powerful, opening the door to generic local search procedures for nonsmooth and nonconvex problems.

2

# 2  Notation

Throughout, we consider the Euclidean space $\mathbb{R}^d$, equipped with the inner-product $\langle \cdot, \cdot \rangle$ and the induced norm $\|x\| := \sqrt{\langle x, x \rangle}$. The *distance* and the *projection* of any point $y \in \mathbb{R}^d$ onto a set $\mathcal{X}$, are defined by

$$\text{dist}(y; \mathcal{X}) := \inf_{x \in \mathcal{X}} \|y - x\| \qquad \text{and} \qquad \text{proj}_{\mathcal{X}}(y) := \underset{x \in \mathcal{X}}{\text{argmin}} \ \|y - x\|,$$

respectively. Note that $\text{proj}_{\mathcal{X}}(y)$ is nonempty as long as $\mathcal{X}$ is a closed set. The *indicator function* of a set $\mathcal{X}$, denoted by $\delta_{\mathcal{X}}$, is defined to be zero on $\mathcal{X}$ and $+\infty$ off it.

## 2.1  Weakly convex functions

Our main focus is on those functions that are convex up to an additive quadratic perturbation. Namely, a function $g \colon \mathbb{R}^d \to \mathbb{R} \cup \{+\infty\}$ is called $\rho$-*weakly convex* (with $\rho \geq 0$) if the assignment $x \mapsto g(x) + \frac{\rho}{2}\|x\|^2$ is a convex function. The algorithms we consider will all use generalized derivative constructions. Variational analytic literature highlights a number of distinct subdifferentials (e.g. [24,27,32]); for weakly convex functions, all these constructions coincide. Consider a $\rho$-weakly convex function $g$. The *subdifferential* of $g$ at $x$, denoted $\partial g(x)$, is the set of all vectors $v \in \mathbb{R}^d$ satisfying

$$g(y) \geq g(x) + \langle v, y - x \rangle + o(\|y - x\|) \qquad \text{as } y \to x. \tag{2.1}$$

Though the condition (2.1) appears to lack uniformity with respect to the basepoint $x$, the subgradients of $g$ automatically satisfy the much stronger property [32, Theorem 12.17]:

$$g(y) \geq g(x) + \langle v, y - x \rangle - \frac{\rho}{2}\|y - x\|^2, \qquad \forall x, y \in \mathbb{R}^d, \ v \in \partial g(x). \tag{2.2}$$

Thus we may use the two conditions, (2.1) and (2.2), interchangeably for weakly convex functions. We note in passing that localizing condition (2.2) leads to so-called prox-regular functions, introduced in [28].

Weakly convex functions are widespread in applications and are typically easy to recognize. One common source is the composite problem class:

$$\min_x \ F(x) := h(c(x)), \tag{2.3}$$

where $h \colon \mathbb{R}^m \to \mathbb{R}$ is convex and $L$-Lipschitz, and $c \colon \mathbb{R}^d \to \mathbb{R}^m$ is a $C^1$-smooth map with $\beta$-Lipschitz gradient. An easy argument shows that $F$ is $L\beta$-weakly convex. This is a worst case estimate. In concrete circumstances, the composite function $F$ may have a much more favorable weak convexity constant $\rho$. The elements of the subdifferential $\partial F(x)$ are straightforward to compute through the chain rule [32, Theorem 10.6, Corollary 10.9]:

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

For a discussion of some recent uses of weakly convex functions in optimization, see the short survey [12]. Throughout the paper, we will use the following two running examples to illustrate our results.

3

**Example 2.1** (Phase retrieval). Phase retrieval is a common computational problem, with applications in diverse areas such as imaging, X-ray crystallography, and speech processing. For simplicity, we will focus on the version of the problem over the reals. The (real) phase retrieval problem seeks to determine a point $x$ satisfying the magnitude conditions,

$$|\langle a_i, x \rangle|^2 \approx b_i \quad \text{for } i = 1, \ldots, m,$$

where $a_i \in \mathbb{R}^d$ and $b_i \in \mathbb{R}$ are given. Note that we can only recover the optimal $x$ up to a universal sign change, since $|\langle a_i, x \rangle| = |\langle a_i, -x \rangle|$. In this work, we will focus on the following optimization formulation of the problem [10, 15, 17]:

$$\min_{x} \ \frac{1}{m} \sum_{i=1}^{m} |\langle a_i, x \rangle^2 - b_i|.$$

Clearly, this is an instance of (2.3). Indeed, under mild statistical assumptions on the way $a_i$ are generated, the formulation is $\rho$-weakly convex, for some numerical constant $\rho$ independent of $d$ and $m$ [15, Corollary 3.2]. Moreover, under an appropriate model of the noise in the measurements, the problem is sharp [15, Propostion 3]. It is worthwhile to mention that numerous other approaches to phase retrieval exist, based on different problem formulations; see for example [4, 5, 34, 36].

Experiment set-up: All of the experiments on phase retrieval will be generated according to the following procedure. In the *exact set-up*, we generate standard Gaussian measurements $a_i \sim N(0, I_{d \times d})$, for $i = 1, \ldots, m$, and generate the target signal $\bar{x} \sim N(0, I_{d \times d})$. We then set $b_i = \langle a_i, \bar{x} \rangle^2$ for each $i = 1, \ldots, m$. In the *corrupted set-up*, we generate $a_i$ and $\bar{x}$ as in the noiseless case. We then corrupt a proportion of the measurements with outliers. Namely, we set $b_i = (1 - z_i) \langle a_i, \bar{x} \rangle^2 + z_i |\zeta_i|$, where $z_i \sim \text{Bernoulli}(0.1)$ and $\zeta_i \sim \mathcal{N}(0, 100)$.

**Example 2.2** (Covariance matrix estimation). The problem of covariance estimation from quadratic measurements, introduced in [7], is a higher rank variant of phase retrieval. Let $a_1, \ldots, a_m \in \mathbb{R}^d$ be measurement vectors. The goal is to recover a low rank decomposition of a covariance matrix $\bar{X}\bar{X}^T$, with $\bar{X} \in \mathbb{R}^{d \times r}$, from quadratic measurements

$$b_i \approx a_i^T \bar{X}\bar{X}^T a_i = \text{Tr}(\bar{X}\bar{X}^T a_i a_i^T).$$

Note that we can only recover $\bar{X}$ up to multiplication by an orthogonal matrix. This problem arises in a variety of contexts, such as covariance sketching for data streams and spectrum estimation of stochastic processes. We refer the reader to [7] for details. In our examples, we will assume $m$ is even and will focus on the potential function

$$\min_{x} \ \frac{1}{m} \sum_{i=1}^{m} \left| \left\langle XX^T, a_{2i} a_{2i}^T - a_{2i-1} a_{2i-1}^T \right\rangle - (b_{2i} - b_{2i-1}) \right|. \tag{2.4}$$

Under exact measurements, i.e., $b_i = a_i^T \bar{X}\bar{X}^T a_i$ and under appropriate statistical assumptions on how $a_i$ are generated, the formulation (2.4) is $\rho$-weakly convex for a numerical constant $\rho$, independent of $d$ or $m$, and is sharp. Indeed, it is a simple consequence of two results, namely [7, Corollary 1] and [37, Lemma 5.4]. It is possible to show the objective is

also sharp when the measurements are corrupted by gross outliers. This guarantee is beyond the scope of our current work, and will appear in a different paper.

Experiment set-up: All of the experiments on covariance matrix estimation will be generated according to the following procedure. In the *exact set-up*, we generate standard Gaussian measurements $a_i \sim N(0, I_{d \times d})$ for $i = 1, \ldots, m$, and generate the target matrix $\bar{X} \in \mathbb{R}^{d \times r}$ as a standard Gaussian. We then set $b_i = \|\bar{X}^T a_i\|_F^2$ for each $i = 1, \ldots, m$. In the *corrupted set-up*, we generate $a_i$ and $\bar{X}$ as in the exact case. We then corrupt a proportion of the measurements with outliers. Namely, we set $b_i = (1 - z_i)\|\bar{X}^T a_i\|_F^2 + z_i|\zeta_i|$, where $\zeta_i \sim \mathcal{N}(0, 100)$ and $z_i \sim \text{Bernoulli}(0.1)$. All plots will show iteration counter $k$ versus the scaled Procrustes distance $\text{dist}(X_k, \mathcal{X}^*)/\|\bar{X}\| = \min_{\Omega^T \Omega = I} \|\Omega X_k - \bar{X}\|_F/\|\bar{X}\|_F$.

## 2.2   Setting of the paper

Throughout the manuscript, we make the following assumption.

**Assumption A.** Consider the optimization problem

$$\min_{x \in \mathcal{X}} \; g(x), \tag{2.5}$$

satisfying the following properties for some real $\mu, \rho > 0$.

1. **(Weak-convexity)** The function $g \colon \mathbb{R}^d \to \mathbb{R}$ is $\rho$-weakly convex, and the set $\mathcal{X} \subset \mathbb{R}^d$ is closed and convex. The set of minimizers $\mathcal{X}^* := \operatorname{argmin}_{x \in \mathcal{X}} g(x)$ is nonempty.

2. **(Sharpness)** The inequality

$$g(x) - \min_{\mathcal{X}} g \geq \mu \cdot \text{dist}(x; \mathcal{X}^*) \qquad \text{holds for all } x \in \mathcal{X}.$$

We will say that a point $\bar{x} \in \mathcal{X}$ is *stationary* for the target problem (2.5) if

$$g(x) - g(\bar{x}) \geq o(\|x - \bar{x}\|) \qquad \text{as } x \xrightarrow{\mathcal{X}} \bar{x}.$$

That is, $\bar{x}$ is *stationary* precisely when the zero vector is a subgradient of $g + \delta_{\mathcal{X}}$ at $\bar{x}$.

Shortly, we will discuss subgradient methods that converge linearly to $\mathcal{X}^*$ under appropriate initialization. As a first step, therefore, we must identify a neighborhood of $\mathcal{X}^*$ that is devoid of extraneous stationary points of (2.5). This is the content of the following lemma.

**Lemma 2.1** (Neighborhood with no stationary points)**.** *The problem* (2.5) *has no stationary points $x$ satisfying*

$$0 < \text{dist}(x; \mathcal{X}^*) < \frac{2\mu}{\rho}. \tag{2.6}$$

*Proof.* Fix a stationary point $x \in \mathcal{X} \setminus \mathcal{X}^*$ of (2.5). Choosing an arbitrary $\bar{x} \in \operatorname{proj}_{\mathcal{X}^*}(x)$, observe

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

Dividing through by $\text{dist}(x; \mathcal{X}^*)$, the result follows.  $\square$

In light of Lemma 2.1, for any $\gamma > 0$ define the following tube

$$\mathcal{T}_\gamma := \left\{ x \in \mathcal{X} : \operatorname{dist}(x; \mathcal{X}^*) < \gamma \cdot \frac{\mu}{\rho} \right\}$$

and the constant

$$L := \sup \left\{ \|\zeta\| : \zeta \in \partial g(x),\, x \in \mathcal{T}_1 \right\}. \tag{2.7}$$

Lemma 2.1 guarantees that the tubes $\mathcal{T}_\gamma$ contain no extraneous stationary points of the problem for any $\gamma \in (0, 2]$. Moreover, observe that $\mu$ and $L$ play reciprocal roles; consequently, the ratio $\tau := \mu/L$ should serve as a measure of conditioning. The following lemma verifies the inclusion $\tau \in [0, 1]$.

**Lemma 2.2** (Condition number). *The inclusion $\tau \in [0, 1]$ holds.*

*Proof.* Consider an arbitrary point $x \in \mathcal{T}_1 \setminus \mathcal{X}^*$ and choose a point $\bar{x} \in \operatorname{proj}_{\mathcal{X}^*}(x)$. By Lebourg's mean value theorem [8, Theorem 2.3.7], there exists a point $z$ in the open segment $(x, \bar{x})$ and a vector $\zeta \in \partial g(z)$ satisfying

$$g(x) - g(\bar{x}) = \langle \zeta, x - \bar{x} \rangle. \tag{2.8}$$

Trivially $z$ lies in $\mathcal{T}_1$, and therefore $\|\zeta\| \leq L$. Using this estimate and sharpness in (2.8) yields the guarantee

$$\mu \cdot \operatorname{dist}(x; \mathcal{X}^*) \leq g(x) - g(\bar{x}) \leq \|\zeta\| \cdot \|x - \bar{x}\| \leq L \cdot \operatorname{dist}(x; \mathcal{X}^*).$$

The result follows. □

To summarize, we will use the following symbols to describe the parameters of the problem class (2.5): $\rho$ is the weak convexity constant of $g$, $\mu$ is the sharpness constant of $g$, $L$ is the maximal subgradient norm at points in the tube $\mathcal{T}_1$, and $\tau$ is the condition measure $\tau = \mu/L \in [0, 1]$.

# 3 Polyak subgradient method

In this section, we consider the Polyak subgradient method for the problem (2.5). A preliminary version of this material applied to the phase retrieval problem appeared in [10]; we present the arguments here for the sake of completeness.

The Polyak subgradient method is summarized in Algorithm 1. This method requires knowing the optimal value $\min_{x \in \mathcal{X}} g(x)$. In a number of circumstances, this indeed is reasonable (e.g. exact penalty approach for solving nonlinear equations). The latter sections explore subgradient methods that do not require a known optimal value.

---
**Algorithm 1:** Polyak Subgradient Method

**Data:** $x_0 \in \mathbb{R}^d$
**Step** $k$: $(k \geq 0)$
    Choose $\zeta_k \in \partial g(x_k)$. **If** $\zeta_k = 0$, then exit algorithm.
    Set $x_{k+1} = \operatorname{proj}_{\mathcal{X}} \left( x_k - \dfrac{g(x_k) - \min_{\mathcal{X}} g}{\|\zeta_k\|^2} \zeta_k \right)$.

---

The following theorem shows that Algorithm 1, originally proposed for convex problems, enjoys the same linear convergence guarantees for functions that are only weakly convex, provided it is initialized within a certain tube of the optimal solution set.

**Theorem 3.1** (Linear rate). *Fix a real $\gamma \in (0, 1)$. Then Algorithm 1 initialized at any point $x_0 \in \mathcal{T}_\gamma$ produces iterates that converge Q-linearly to $\mathcal{X}^*$, that is*

$$\mathrm{dist}^2(x_{k+1}; \mathcal{X}^*) \leq \left(1 - (1 - \gamma)\tau^2\right) \mathrm{dist}^2(x_k; \mathcal{X}^*). \tag{3.1}$$

*Proof.* We proceed by induction. Suppose that the theorem holds up to iteration $k$. We will prove the inequality (3.1). To this end, choose $\bar{x} \in \mathrm{proj}_{\mathcal{X}^*}(x_k)$. 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 $\mathrm{dist}(x_k; \mathcal{X}^*) \leq \mathrm{dist}(x_0; \mathcal{X}^*)$ and therefore $x_k$ lies in $\mathcal{T}_\gamma$. Lemma 2.1 therefore guarantees $\zeta_k \neq 0$. We successively deduce, by non-expansiveness of $\mathrm{proj}_{\mathcal{X}}$, that

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

Combining the inclusion $x_k \in \mathcal{T}_\gamma$ with sharpness, we therefore deduce

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

The result follows. $\qquad\square$

As a numerical illustration, let us apply the Polyak subgradient method (Figure 1) to our two running examples, phase retrieval and covariance matrix estimation. Notice that a linear rate of convergence is observed in all experiments except for two, with the rate improving monotonically with an increasing number of measurements $m$. In the two exceptional experiments, the number of measurements $m$ is too small to guarantee that the initial point $x_0$ is within the basin of attraction, and the subgradient methods stagnates.

# 4 Subgradient method with constant step-size

Recall that the Polyak subgradient method (Algorithm 1) crucially relies on knowing the minimal value of the optimization problem (2.5). Henceforth, all the subgradient methods we
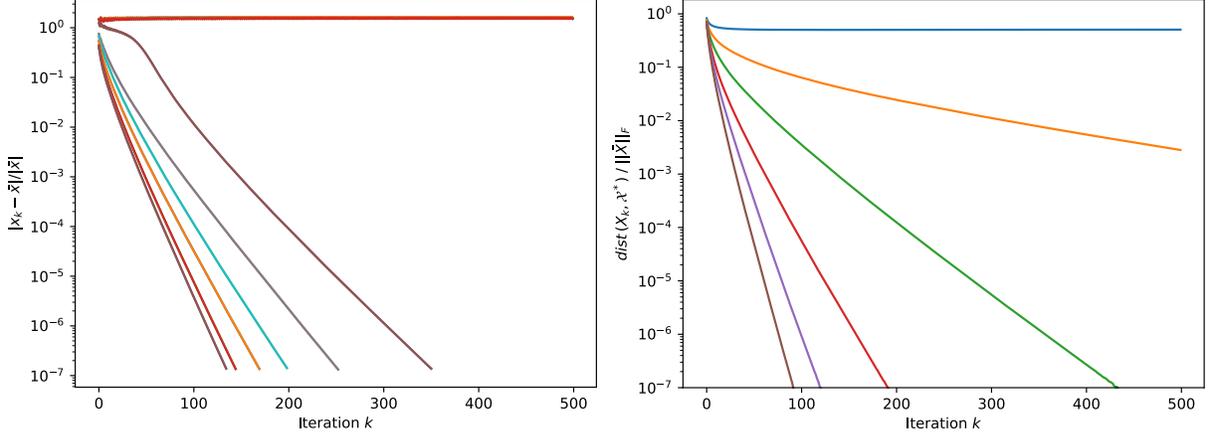
Figure 1: Polyak subgradient method. (Left) Phase retrieval with the exact set-up; $d = 5000$ and $m \in \{11000, 12225, 13500, 14750, 16000, 17250, 18500\}$. (Right) Covariance matrix estimation with the exact set-up; $d = 1000$, $r = 3$, and $m \in \{5000, 8000, 11000, 14000, 17000, 20000\}$. In both experiments, convergence rates uniformly improve with increasing $m$.

consider are agnostic to this value. That being said, they will require some estimates on the problem parameters $(\mu, \rho, L)$. We begin by analyzing a subgradient method with a constant step size (Algorithm 2). Constant-step schemes are often methods of choice in practice. We will show that when properly initialized, the subgradient method with constant stepsize generates iterates $x_k$ such that $\mathrm{dist}(x_k; \mathcal{X}^*)$ converges linearly up to a certain threshold.

---

**Algorithm 2:** Subgradient method with constant stepsize

> **Data:** Initial point $x_0 \in \mathbb{R}^d$ and stepsize $\alpha > 0$
> **Step** $k$: ($k \geq 0$)
>       Choose $\zeta_k \in \partial g(x_k)$. **If** $\zeta_k = 0$, then exit algorithm.
>       Set $x_{k+1} = \mathrm{proj}_{\mathcal{X}}\left(x_k - \alpha \cdot \frac{\zeta_k}{\|\zeta_k\|}\right)$.

---

The analysis we present fundamentally relies on the following estimate, often used in the analysis of subgradient methods. To simplify notation, for any point $x \in \mathbb{R}^d$, we set

$$E(x) := \mathrm{dist}^2(x; \mathcal{X}^*).$$

Whenever $x$ has an index $k$ as a subscript, we will set $E_k := E(x_k)$. The following lemma will feature in both the constant and geometrically decaying stepsize schemes.

**Lemma 4.1** (Basic recurrence). *Consider a point $x \in \mathcal{T}_1$ and a nonzero subgradient $\zeta \in \partial g(x)$, and define $x^+ := \mathrm{proj}_{\mathcal{X}}\left(x - \alpha \frac{\zeta}{\|\zeta\|}\right)$ for some $\alpha > 0$. Then the estimate holds:*

$$E(x^+) \leq \left(1 + \frac{\rho \alpha}{L}\right) E(x) - 2\alpha\tau\sqrt{E(x)} + \alpha^2. \tag{4.1}$$

8

*Proof.* Choose an arbitrary point $\bar{x} \in \text{proj}_{\mathcal{X}^*}(x)$. Observe

$$\|x^+ - \bar{x}\|^2 \leq \left\|(x - \bar{x}) - \alpha\tfrac{\zeta}{\|\zeta\|}\right\| = \|x - \bar{x}\|^2 + \tfrac{2\alpha}{\|\zeta\|} \cdot \langle \zeta, \bar{x} - x \rangle + \alpha^2$$

$$\leq \|x - \bar{x}\|^2 + \tfrac{2\alpha}{\|\zeta\|} \cdot \left(g(\bar{x}) - g(x) + \tfrac{\rho}{2}\|x - \bar{x}\|^2\right) + \alpha^2$$

$$\leq \left(1 + \tfrac{\alpha\rho}{\|\zeta\|}\right)\|x - \bar{x}\|^2 - \tfrac{2\alpha\mu}{\|\zeta\|} \cdot \|x - \bar{x}\| + \alpha^2.$$

Thus the inequality holds:

$$E(x^+) \leq \left(1 + \tfrac{\alpha\rho}{\|\zeta\|}\right) E(x) - \tfrac{2\alpha\mu}{\|\zeta\|} \cdot \sqrt{E(x)} + \alpha^2,$$

and consequently taking into account $\alpha/\|\zeta\| \geq \alpha/L$ we have

$$E(x^+) \leq \sup_{t \geq \alpha/L} \left\{(1 + \rho t)\, E(x) - 2\mu t \cdot \sqrt{E(x)} + \alpha^2\right\}.$$

Notice, the function inside the supremum is linear in $t$ with slope $s := \rho E(x) - 2\mu\sqrt{E(x)}$. The inclusion $x \in \mathcal{T}_1$ directly implies $s \leq 0$. Therefore the supremum on the right-hand-side is attained at $t = \tfrac{\alpha}{L}$, yielding the claimed estimate (4.1). $\qquad\square$

In light of Lemma 4.1, we can now prove that the quantities $E(x_k)$ converge linearly below a certain fixed threshold. The proof is a modification of that in [22, Section 4].

**Lemma 4.2** (Contraction inequality)**.** *Fix a constant $\alpha \in (0, \tfrac{\tau\mu}{\rho})$ and let $\{x_k\}_{k\geq 0}$ be the iterates generated by Algorithm 2. Define the quantity*

$$E^* := \left(\frac{\alpha L}{\mu + \sqrt{\mu^2 - \alpha\rho L}}\right)^2. \tag{4.2}$$

*Then whenever an iterate $x_k$ lies in $\mathcal{T}_1$, the estimate holds:*

$$E_{k+1} - E^* \leq q_k(E_k - E^*),$$

*where $q_k := 1 + \tfrac{\alpha}{L}\left(\rho - \tfrac{2\mu}{\sqrt{E_k} + \sqrt{E^*}}\right)$ satisfies $q_k < 1$.*

*Proof.* Looking back at the estimate (4.1), consider the following equation in the variable $e$:

$$e = \left(1 + \tfrac{\alpha\rho}{L}\right) e - 2\alpha\tau \cdot \sqrt{e} + \alpha^2. \tag{4.3}$$

An easy computation shows that the minimal positive solution to (4.3) is exactly $E^*$, defined in (4.2). Note that $E^*$ is well-defined by the inequality $\alpha \leq \tau \cdot \tfrac{\mu}{\rho}$.

Subtracting (4.3) from (4.1) yields the estimate

$$E_{k+1} - E^* \leq (1 + \tfrac{\alpha\rho}{L})(E_k - E^*) - 2\alpha\tau(\sqrt{E_k} - \sqrt{E^*})$$

$$= \left(1 + \tfrac{\alpha}{L}\left(\rho - \tfrac{2\mu}{\sqrt{E_k} + \sqrt{E^*}}\right)\right)(E_k - E^*).$$

Finally, notice

$$\rho - \frac{2\mu}{\sqrt{E_k} + \sqrt{E^*}} < \rho - \frac{2\mu}{2\mu/\rho} = 0.$$

This completes the proof of the lemma. $\qquad\square$

9

Iterating Lemma 4.2, we see that the quantities $E_k$ decrease to a value lower than $E^*$ at a linear rate. Figure 2 illustrates this behavior on our two running examples. It is also clear from the figure that the linear rate of convergence improves as $E_k$ tends to $E^*$. An explanation is immediate from the expression for $q_k$ in Lemma 4.2. Indeed, as $E_k$ decreases, so do the contraction factors $q_k$, and for $E_k \approx E^*$, we have $q_k \approx (1 - \tau^2) + \frac{\alpha\rho}{L}$. Thus as the step-size tends to zero, the limiting linear rate coincides with the ideal rate of $1 - \tau^2$.

Another interesting feature, apparent in Figure 2, is that even after $E_k$ becomes smaller than $E^*$, all the following values $E_k$ stay close to $E^*$. This is the content of the following theorem. The convex version of this theorem appears in [22, Theorem 2].
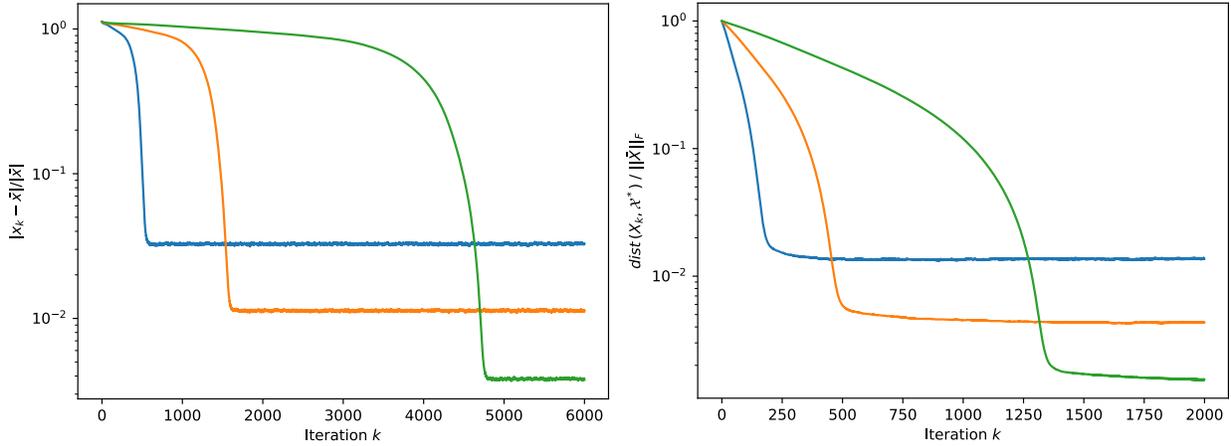


Figure 2: Constant step subgradient method. (Left) Phase retrieval with the corrupted set-up; $d = 1000$, $m = 3000$, and $\alpha \in \{1, 1/3, 1/9\}$. (Right) Covariance matrix estimation with the corrupted set-up; $d = 1000$, $r = 3$, $m = 10000$, and $\alpha \in \{1, 1/3, 1/9\}$. The lower curves curves correspond to smaller step-size in both experiments.

**Theorem 4.3** (Convergence of fixed stepsize subgradient method). *Fix a real $\gamma \in (0, 1)$ and a real $\alpha > 0$ satisfying*

$$0 < \alpha < \frac{\gamma\tau}{\sqrt{1 + 2\tau^2}} \cdot \frac{\mu}{\rho}. \tag{4.4}$$

*Let $x_k$ be the iterates generated by Algorithm 2 with stepsize $\alpha$ and initial point $x_0 \in \mathcal{T}_\gamma$. Define the constants*

$$E^* := \left( \frac{\alpha L}{\mu + \sqrt{\mu^2 - \alpha\rho L}} \right)^2 \qquad and \qquad D := \sqrt{\max\{E_0, 2\alpha^2 + E^*\}}.$$

*Then for each index $k$, the estimates hold:*

$$\sqrt{E_k} \le D \le \frac{\gamma\mu}{\rho} \qquad and \qquad E_k - E^* \le \max\left\{ q^k(E_0 - E^*), 2\alpha^2 \right\},$$

*where the coefficient $q := 1 + \frac{\alpha}{L}\left(\rho - \frac{\mu}{D}\right)$ satisfies $q \in (0, 1)$.*

*Proof.* We first verify the claims that are independent of the iteration counter. To this end, observe that (4.4) directly implies $\alpha \leq \tau \cdot \frac{\mu}{\rho}$, and therefore $E^*$ is well defined. Next, we show $D < \frac{\gamma\mu}{\rho}$. Indeed, noting $\sqrt{E^*} \leq \alpha\tau^{-1}$ and using (4.4), we deduce

$$D^2 = \max\{E_0, 2\alpha^2 + E^*\} \leq \max\{E_0, \alpha^2(2 + \tau^{-2})\} \leq \left(\frac{\gamma\mu}{\rho}\right)^2.$$

Next, we show the inclusion $q \in (0, 1)$. To this end, observe

$$-1 \leq -\frac{\tau\alpha}{D} < \tfrac{\alpha}{L}\left(\rho - \tfrac{\mu}{D}\right) \leq (1 - \gamma^{-1}) \cdot \frac{\rho\alpha}{L} < 0,$$

where the first inequality follows from the inequality, $\alpha \leq D$, and the third follows from the inequality, $D < \frac{\gamma\mu}{\rho}$. Thus we conclude $q \in (0, 1)$, as claimed.

We now proceed by induction. Fix an index $k$ and suppose as inductive hypothesis that for each index $i = 0, 1, \ldots, k$, the estimates hold:

$$\sqrt{E_i} \leq D \qquad \text{and} \qquad E_i - E^* \leq \max\left\{q^i(E_0 - E^*), 2\alpha^2\right\}.$$

Let us consider two cases. Suppose first $E_k \geq E^*$. Then by applying Lemma 4.2, we deduce

$$\begin{aligned}
E_{k+1} - E^* &\leq \left(1 + \tfrac{\alpha}{L}\left(\rho - \tfrac{2\mu}{\sqrt{E_k}+\sqrt{E^*}}\right)\right)(E_k - E^*) \\
&\leq \left(1 + \tfrac{\alpha}{L}\left(\rho - \tfrac{\mu}{D}\right)\right)(E_k - E^*) \\
&= q(E_k - E^*).
\end{aligned}$$

Suppose now that the second case, $E_k < E^*$, holds. Then Lemma 4.2 implies

$$\begin{aligned}
E_{k+1} &\leq E_k + \tfrac{\alpha\rho}{L}(E_k - E^*) - 2\alpha\tau(\sqrt{E_k} - \sqrt{E^*}) \\
&\leq \max_{E \in [0, E^*]}\{E + \tfrac{\alpha\rho}{L}(E - E^*) - 2\alpha\tau(\sqrt{E} - \sqrt{E^*})\} \\
&= \max\{E^*, \tfrac{\alpha}{L}(2\mu\sqrt{E^*} - \rho E^*)\}.
\end{aligned}$$

Subtracting $E^*$, we conclude

$$E_{k+1} - E^* \leq \max\{0, 2\tau\alpha\sqrt{E^*}\} \leq 2\alpha^2.$$

Thus in both cases, we have the estimate

$$E_{k+1} - E^* \leq \max\left\{q(E_k - E^*), 2\alpha^2\right\}.$$

In particular, we immediately deduce $\sqrt{E_{k+1}} \leq D$. Applying the inductive hypothesis, we conclude

$$E_{k+1} - E^* \leq \max\left\{q \cdot \max\{q^k(E_0 - E^*), 2\alpha^2\}, 2\alpha^2\right\} \leq \max\{q^{k+1}(E_0 - E^*), 2\alpha^2\}.$$

The theorem is proved. $\qquad\square$

# 5 Geometrically decaying step

In the last section, we showed linear convergence of the constant step size scheme up to a fixed tolerance $E^*$. To obtain a linearly convergent method to the true solution set, we will allow the step-size to decrease geometrically. The analogous strategy in the convex setting goes back to [19], and our argument follows the same strategy. The intuition for why one may expect linear convergence under such step sizes may be gleaned from the Polyak method under the optimal step size

$$\alpha_k = \frac{g(x_k) - \min g(x)}{\|\zeta_k\|}.$$

It is easy to verify that since $E_k$ tend to zero $Q$-linearly, the steps $\alpha_k$ tends to zero R-linearly. We implement such a geometrically decaying stepsize in Algorithm 3 and prove linear convergence of the method in Theorem 5.1.

---

**Algorithm 3:** Subgradient method with geometrically decreasing stepsize

**Data:** Real $\lambda > 0$ and $q \in (0, 1)$.
**Step** $k$: $(k \geq 0)$
    Choose $\zeta_k \in \partial g(x_k)$. **If** $\zeta_k = 0$, then exit algorithm.
    Set stepsize $\alpha_k = \lambda \cdot q^k$.
    Update iterate $x_{k+1} = \mathrm{proj}_{\mathcal{X}}\left(x_k - \alpha_k \frac{\zeta_k}{\|\zeta_k\|}\right)$.

---

**Theorem 5.1.** *Fix a real $\gamma \in (0, 1)$ and suppose $\tau \leq \sqrt{\frac{1}{2-\gamma}}$. Set*

$$\lambda := \frac{\gamma\mu^2}{\rho L} \quad and \quad q := \sqrt{1 - (1-\gamma)\tau^2}.$$

*Then the iterates $x_k$ generated by Algorithm 3, initialized at some point $x_0 \in \mathcal{T}_\gamma$, satisfy:*

$$\mathrm{dist}^2(x_k; \mathcal{X}^*) \leq \frac{\gamma^2\mu^2}{\rho^2}\left(1 - (1-\gamma)\tau^2\right)^k. \tag{5.1}$$

*Proof.* We will prove the result by induction. To this end, suppose the bound (5.1) holds for all $i = 0, \ldots, k$. Appealing to Lemma 4.1. and using the relation $\alpha_k = \lambda q^k$, we obtain

$$E_{k+1} \leq \left(1 + \frac{\rho\lambda q^k}{L}\right)E_k - 2\lambda\tau q^k\sqrt{E_k} + \lambda^2 q^{2k}. \tag{5.2}$$

Define the constant $M := \frac{\gamma\mu}{\rho}$. Recall the induction assumption guarantees $\sqrt{E_k} \leq Mq^k$. Let us therefore fix some value $R \in [0, M]$ satisfying $\sqrt{E_k} = Rq^k$. Inequality (5.2) then implies

$$E_{k+1} \leq \max_{R \in [0,M]}\left\{R^2 q^{2k} + \frac{\rho\lambda R^2}{L}q^{3k} - 2\lambda\tau Rq^{2k} + \lambda^2 q^{2k}\right\}.$$

Note that the expression inside the maximum is a convex quadratic in $R$ and therefore the maximum must occur either at $R = 0$ or $R = M$. We therefore deduce

$$E_{k+1} \leq q^{2k} \cdot \max\left\{\lambda^2, \ M^2 + \frac{\rho\lambda}{L}M^2 q^k - 2\lambda\tau M + \lambda^2\right\}. \tag{5.3}$$

To complete the induction, it is therefore sufficient to show

$$\lambda^2 \leq M^2 q^2 \qquad \text{and} \qquad M^2 + \frac{\rho\lambda}{L}M^2 q^k - 2\lambda\tau M + \lambda^2 \leq M^2 q^2. \tag{5.4}$$

First, we show that $M$ satisfies the first property. Note the equality $M = \frac{\lambda}{\tau}$. Hence, it suffices to show that $\tau \leq q$, Observe that the assumption $\tau \leq \sqrt{\frac{1}{2-\gamma}}$ directly implies

$$\tau^2 + (1-\gamma)\tau^2 \leq 1.$$

Rearranging yields $\tau^2 \leq 1 - (1-\gamma)\tau^2 = q^2$. Hence, the first condition in (5.4) holds.

Next we show that $M$ satisfies the second property in (5.4). Thus, rearranging the expression, we must establish

$$\left(1 + \frac{\rho\lambda}{L}q^k - q^2\right)M^2 - 2\lambda\tau M + \lambda^2 \leq 0. \tag{5.5}$$

We will show that the quadratic on the left-hand-side in $M$ has two real positive roots. To this end, a quick computation shows that the two roots are

$$\frac{\lambda\tau \pm \sqrt{\lambda^2\tau^2 - \lambda^2\left(1 + \frac{\rho\lambda}{L}q^k - q^2\right)}}{1 + \frac{\rho\lambda}{L}q^k - q^2} = \frac{\lambda}{\tau \mp \sqrt{\tau^2 - \left(1 + \frac{\rho\lambda}{L}q^k - q^2\right)}}.$$

To see that the discriminant is nonnegative, observe

$$\tau^2 - \left(1 + \frac{\rho\lambda}{L}q^k - q^2\right) \geq \tau^2 - \left(1 + \frac{\rho\lambda}{L} - q^2\right) = \tau^2 - (1 + \gamma\tau^2 - q^2) = 0.$$

Thus the convex quadratic in (5.5) has two real roots, and our choice $M = \lambda/\tau$ lies between them. Hence the condition (5.5) holds, and the inductive step is complete.

$\square$

We now illustrate the performance of Algorithm 3 on our two running examples in Figure 3. Empirically, we observed that $\lambda > 0$ and $q \in (0, 1)$ must be tuned for performance, which is what we did in the experiments. We observe linear convergence in all cases, and the convergence rate of the method improves monotonically as the chosen rate $q$ is decreased. While the Polyak scheme pictured in Figure 1 clearly outperforms all other methods, the geometrically decaying step scheme performs much better than the constant step scheme in Figure 2.

# References

[1] N. Boumal. Nonconvex phase synchronization. *SIAM J. Optim.*, 26(4):2355–2377, 2016.

[2] A. Brutzkus and A. Globerson. Globally optimal gradient descent for a convnet with gaussian inputs. *arXiv:1702.07966*, 2017.
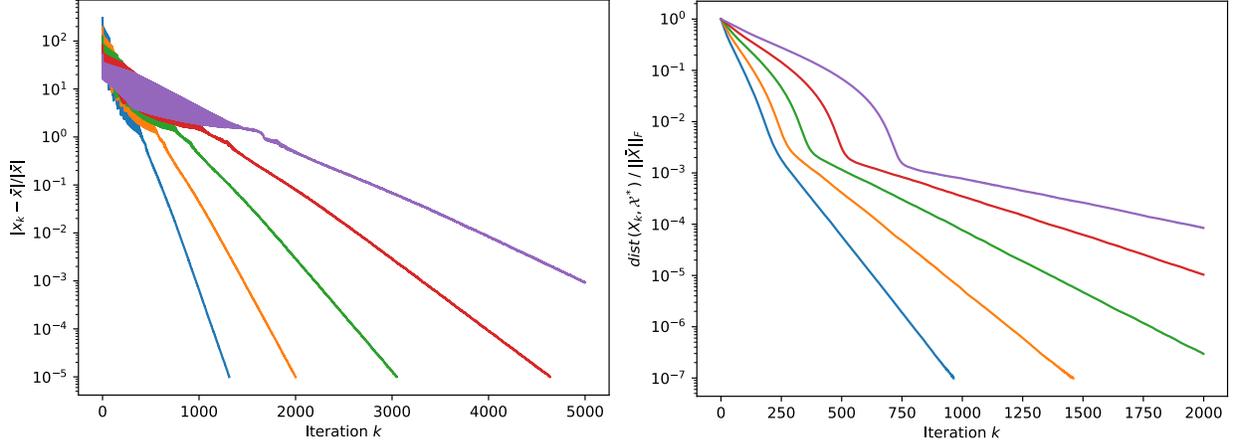
Figure 3: Geometrically decaying step size. (Left) Phase retrieval with the corrupted set-up; $d = 1000$, $m = 3000$, $q \in \{0.983, 0.989, 0.993, 0.996, 0.997\}$. (Right) Covariance matrix estimation with the corrupted set-up; $d = 1000$, $r = 3$, $m = 10000$, $q \in \{0.986, 0.991, 0.994, 0.996, 0.998\}$. The depicted rates uniformly improve with lower values of $q$, in both figures.

[3] J.V. Burke and M.C. Ferris. Weak sharp minima in mathematical programming. *SIAM J. Control Optim.*, 31(5):1340–1359, 1993.

[4] E.J. Candès, X. Li, and M. Soltanolkotabi. Phase retrieval via Wirtinger flow: theory and algorithms. *IEEE Trans. Inform. Theory*, 61(4):1985–2007, 2015.

[5] Y. Chen and E.J. Candès. Solving random quadratic systems of equations is nearly as easy as solving linear systems. *Comm. Pure Appl. Math.*, 70(5):822–883, 2017.

[6] Y. Chen and M.J Wainwright. Fast low-rank estimation by projected gradient descent: General statistical and algorithmic guarantees. *arXiv:1509.03025*, 2015.

[7] Yuxin Chen, Yuejie Chi, and Andrea J Goldsmith. Exact and stable covariance estimation from quadratic sampling via convex programming. *IEEE Transactions on Information Theory*, 61(7):4034–4059, 2015.

[8] F.H. Clarke. *Optimization and Nonsmooth Analysis*. Wiley Interscience, NY, 1983.

[9] D. Davis and D. Drusvyatskiy. Stochastic subgradient method converges at the rate $O(k^{-1/4})$ on weakly convex functions. *arXiv:1802.02988*, 2018.

[10] D. Davis, D. Drusvyatskiy, and C. Paquette. The nonsmooth landscape of phase retrieval. *arXiv:1711.03247*, 2017.

[11] D. Davis and B. Grimmer. Proximally guided stochastic method for nonsmooth, non-convex problems. *Preprint arXiv:1707.03505*, 2017.

[12] D. Drusvyatskiy. The proximal point method revisited. *To appear in SIAG/OPT Views and News, arXiv:1712.06038*, 2018.

[13] D. Drusvyatskiy and A.S. Lewis. Error bounds, quadratic growth, and linear convergence of proximal methods. *To appear in Math. Oper. Res., arXiv:1602.06661*, 2016.

14

[14] D. Drusvyatskiy and C. Paquette. Efficiency of minimizing compositions of convex functions and smooth maps. *Preprint arXiv:1605.00125*, 2016.

[15] J.C. Duchi and F. Ruan. Solving (most) of a set of quadratic equalities: Composite optimization for robust phase retrieval. *arXiv:1705.02356*, 2017.

[16] J.C. Duchi and F. Ruan. Stochastic methods for composite optimization problems. *Preprint arXiv:1703.08570*, 2017.

[17] Y.C. Eldar and S. Mendelson. Phase retrieval: stability and recovery guarantees. *Appl. Comput. Harmon. Anal.*, 36(3):473–494, 2014.

[18] I.I. Eremin. The relaxation method of solving systems of inequalities with convex functions on the left-hand side. *Dokl. Akad. Nauk SSSR*, 160:994–996, 1965.

[19] J.L. Goffin. On convergence rates of subgradient optimization methods. *Math. Program.*, 13(3):329–347, 1977.

[20] P. Jain, C. Jin, S.M. Kakade, and P. Netrapalli. Global convergence of non-convex gradient descent for computing matrix squareroot. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 479–488, 2017.

[21] P. Jane and P. Netrapalli. Fast exact matrix completion with finite samples. In Peter Grnwald, Elad Hazan, and Satyen Kale, editors, *Proceedings of The 28th Conference on Learning Theory*, volume 40 of *Proceedings of Machine Learning Research*, pages 1007–1034, Paris, France, 03–06 Jul 2015. PMLR.

[22] P.R. Johnstone and P. Moulin. Faster subgradient methods for functions with Hölderian growth. *arXiv:1704.00196*, 2017.

[23] R. Meka, P. Jain, and I.S. Dhillon. Guaranteed rank minimization via singular value projection. *arXiv:0909.5457*, 2009.

[24] B.S. Mordukhovich. *Variational Analysis and Generalized Differentiation I: Basic Theory*. Grundlehren der mathematischen Wissenschaften, Vol 330, Springer, Berlin, 2006.

[25] E. A. Nurminskii. The quasigradient method for the solving of the nonlinear programming problems. *Cybernetics*, 9(1):145–150, Jan 1973.

[26] E. A. Nurminskii. Minimization of nondifferentiable functions in the presence of noise. *Cybernetics*, 10(4):619–621, Jul 1974.

[27] J.-P. Penot. *Calculus without derivatives*, volume 266 of *Graduate Texts in Mathematics*. Springer, New York, 2013.

[28] R.A. Poliquin and R.T. Rockafellar. Prox-regular functions in variational analysis. *Trans. Amer. Math. Soc.*, 348:1805–1838, 1996.

[29] B.T. Poljak. Minimization of unsmooth functionals. *USSR Computational Mathematics and Mathematical Physics*, 9:14–29, 1969.

[30] B.T. Poljak. Subgradient methods: a survey of Soviet research. In *Nonsmooth optimization (Proc. IIASA Workshop, Laxenburg, 1977)*, volume 3 of *IIASA Proc. Ser.*, pages 5–29. Pergamon, Oxford-New York, 1978.

[31] B. Polyak. Sharp minima. *Institute of Control Sciences Lecture Notes, Moscow, USSR; Presented at the IIASA Workshop on Generalized Lagrangians and Their Applications, IIASA, Laxenburg, Austria*, (3):369–380, 1979.

[32] R.T. Rockafellar and R.J-B. Wets. *Variational Analysis.* Grundlehren der mathematischen Wissenschaften, Vol 317, Springer, Berlin, 1998.

[33] N.Z. Shor. The rate of convergence of the method of the generalized gradient descent with expansion of space. *Kibernetika (Kiev)*, (2):80–85, 1970.

[34] J. Sun, Q. Qu, and J. Wright. A geometric analysis of phase retrieval. *To appear in Found. Comp. Math., arXiv:1602.06664*, 2017.

[35] S. Supittayapornpong and M.J. Neely. Staggered time average algorithm for stochastic non-smooth optimization with $O(1/t)$ convergence. *arXiv:1607.02842*, 2016.

[36] Y.S Tan and R. Vershynin. Phase retreival via randomized kaczmarz: Theoretical guarantees. *arXiv:1605.08285*, 2017.

[37] S. Tu, R. Boczar, M. Simchowitz, M. Soltanolkotabi, and B. Recht. Low-rank solutions of linear matrix equations via procrustes flow. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48*, ICML'16, pages 964–973. JMLR.org, 2016.

[38] T. Yang and Q. Lin. RSG: Beating subgradient method without smoothness and strong convexity. *arXiv:1512.03107*, 2016.