

# On the Convergence Rate of Stochastic Mirror Descent for Nonsmooth Nonconvex Optimization

Siqi Zhang\*

Niao He\*

## Abstract

In this paper, we investigate the non-asymptotic stationary convergence behavior of Stochastic Mirror Descent (SMD) for nonconvex optimization. We focus on a general class of nonconvex nonsmooth stochastic optimization problems, in which the objective can be decomposed into a relatively weakly convex function (possibly non-Lipschitz) and a simple non-smooth convex regularizer. We prove that SMD, without the use of mini-batch, is guaranteed to converge to a stationary point in a convergence rate of  $\mathcal{O}(1/\sqrt{t})$ . The efficiency estimate matches with existing results for stochastic subgradient method, but is evaluated under a stronger stationarity measure. Our convergence analysis applies to both the original SMD and its proximal version, as well as the deterministic variants, for solving relatively weakly convex problems.

## 1 Introduction

In this paper, we consider the composite nonsmooth nonconvex stochastic optimization problems with the following general form

$$\min_{x \in X} T(x) := f(x) + r(x) = \mathbb{E}_{\xi} [F(x; \xi)] + r(x) \quad (1.1)$$

where  $X \subseteq \mathcal{X}$  is a nonempty closed convex subset of a finite-dimensional Euclidean space  $\mathcal{X}$  equipped with a norm  $\|\cdot\|$ ,  $f(x) : X \rightarrow \mathbb{R}$  is a nonsmooth nonconvex function,  $r(x) : X \rightarrow \mathbb{R}$  is a simple nonsmooth convex regularizer.

Throughout, we assume that  $f(x)$  is  $\rho$ -relatively weakly convex, i.e., the function  $f(x) + \rho\omega(x)$  is convex for some  $\rho > 0$  and some function  $\omega(x) : X \rightarrow \mathbb{R}$  that is continuously differentiable and 1-strongly convex with respect to the norm  $\|\cdot\|$  defined on  $\mathcal{X}$ . In the case when  $\omega(x) = \frac{1}{2}\|x\|^2$  and  $\|\cdot\|$  is an inner product induced norm,  $f(x)$  is also called  $\rho$ -weakly convex. Weak convexity is a special yet very common case of nonconvex functions, which contains all convex functions and Lipschitz smooth functions. The composite form of the optimization problem covers a wide spectrum of regularized problems in machine learning, including the nonlinear least square, sparse logistic regression (Liu et al., 2009; Shen and Gu, 2018), sparse recovery (Chen and Gu, 2014), and robust phase retrieval (Davis et al., 2017).

When the function  $f(x)$  is convex, the proximal variant of Stochastic Mirror Descent (SMD) is one of the most widely used algorithms for solving the above composite problem; see, e.g., Duchi et al. (2010), He (2015), and Beck (2017). SMD performs the recurrence at each iteration:

$$x_{t+1} = \operatorname{argmin}_{x \in X} \{ \langle F'(x_t, \xi_t), x \rangle + r(x) + \frac{1}{\alpha_t} D_{\psi}(x, x_t) \} \quad (1.2)$$

---

\*Department of Industrial and Enterprise Systems Engineering (ISE), University of Illinois at Urbana-Champaign (UIUC), Urbana, IL 61801, USA. Emails: [siqiz4@illinois.edu](mailto:siqiz4@illinois.edu), [niahe@illinois.edu](mailto:niahe@illinois.edu). This work is supported by NSF CCF-1755829.

where  $\alpha_t > 0$  is the stepsize,  $F'(x_t, \xi_t)$  is an unbiased estimator of the subgradient of  $f(x)$  at  $x_t$ , the term  $D_\psi(x, x_t) := \psi(x) - \psi(x_t) - \langle \nabla \psi(x_t), x - x_t \rangle$  stands for some Bregman divergence generated by a 1-strongly convex and continuously differentiable function  $\psi(x)$  defined on  $X$ . Note that when the Bregman divergence is set to be the simple Euclidean distance, *i.e.*,  $D_\psi(x, x') = \frac{1}{2}\|x - x'\|_2^2$  with  $\psi(x) = \frac{1}{2}\|x\|_2^2$ , SMD reduces to proximal stochastic subgradient method (SGD). When there is no regularizer, *i.e.*,  $r(x) = 0$ , this reduces to the original SMD (Nemirovski and Yudin, 1983). The non-asymptotic convergences of SMD algorithm and its variants have been extensively studied in the convex regime; see e.g., Nemirovski et al. (2009) for analysis of the original SMD, and Duchi et al. (2010) for the proximal SMD. It is well-known that SMD achieves an optimal  $\mathcal{O}(1/\sqrt{t})$  convergence rate for solving general composite nonsmooth convex problems with unimprovable constant factors. However, the non-asymptotic convergence behavior of SMD is far from fully understood when moving to the nonconvex regime.

## 1.1 Related Works

There have been several recent works discussing the convergences of SGD or SMD for nonconvex problems. We mainly focus on the purely stochastic case, where  $f(x)$  can only be accessed through stochastic oracles. The special case where  $f(x)$  consists of a finite sum of components is beyond the scope of this work.

The seminal work by Ghadimi and Lan (2013) provides the first non-asymptotic convergence analysis of SGD for unconstrained smooth nonconvex objectives, *i.e.*, problem (1.1) with  $X = \mathbb{R}^n, r(x) = 0$  and smooth  $f(x)$ . They show that a modified SGD, called *Randomized Stochastic Gradient (RSG)*, requires  $\mathcal{O}(1/\epsilon^2)$  number of iterations to generate an  $\epsilon$ -stationary point such that  $\mathbb{E}[\|\nabla f(x)\|_2^2] \leq \epsilon$ . Later, Ghadimi et al. (2016) addressed the general constrained composite problem (1.1) with smooth  $f(x)$  and proposed a modified mini-batch SMD method, called *Randomized Stochastic Projected Gradient (RSPG)*, that requires using a mini-batch of size  $\mathcal{O}(1/\epsilon)$  samples to estimate the gradient at each iteration. They showed that the algorithm achieved the same  $\mathcal{O}(1/\epsilon^2)$  overall sample complexity to achieve an  $\epsilon$ -stationary point, in terms of the generalized projected gradient, *i.e.*,  $\mathbb{E}[\|g_X(x)\|_2^2] \leq \epsilon$ <sup>1</sup>. It is worth mentioning that although the RSPG algorithm utilizes the mirror descent framework, the analysis in Ghadimi et al. (2016) only applies to the Euclidean setting and smooth objectives.

To overcome the mini-batch requirement for solving constrained nonconvex problems, Davis and Grimmer (2017) proposed the *Proximally Guided Stochastic subGradient (PGSG)* method, by combining proximal point algorithm and SGD in a nested framework - iteratively solving subproblems arising from proximal point algorithm through SGD routines. The algorithm solves problem (1.1) with  $r(x) = 0$  and  $\rho$ -weakly convex  $f(x)$ , and attains the  $\tilde{\mathcal{O}}(1/\epsilon^2)$ <sup>2</sup> sample complexity to get an  $\epsilon$ -stationary point, measured by the squared distance from zero to the Fréchet subdifferential set, *i.e.*,  $\mathbb{E}[\text{dist}^2(0, \partial_F(f + \delta_X)(x))] \leq \epsilon$ . More recently, Davis and Drusvyatskiy (2018) considered the same weakly convex setting and showed that even the basic SGD and its proximal variant converge and exhibit an  $\mathcal{O}(1/\epsilon^2)$  sample complexity.

However, none of these works have considered or addressed the convergence behavior of SMD in the non-Euclidean setting. We point out that a recent work by Zhou et al. (2017) investigated the asymptotic convergence of SMD, but is only limited to a very special class of nonconvex problems that ensures global convergence. This paper aims to close this fundamental theoretical gap and establish the non-asymptotic stationary convergence analysis of Stochastic Mirror Descent for nonconvex problems. A detailed comparison of this work and previous ones is summarized in Table 1.

<sup>1</sup> Here the generalized project gradient is defined as  $g_X(x_t) := \frac{1}{\alpha_t}(x_t - x_{t+1})$ , where  $x_{t+1}$  is defined in (1.2).

<sup>2</sup>  $\tilde{\mathcal{O}}(\cdot)$  means the complexity neglects its logarithmic terms in its expression.

Table 1: Summary of algorithms

Method	<i>RSG</i> (Ghadimi and Lan, 2013)	<i>RSPG</i> (Ghadimi et al., 2016)	<i>PGSG</i> (Davis and Grimmer, 2017)	<i>PSG</i> (Davis and Drusvyatskiy, 2018)	<i>SMD</i> (this paper)
Convexity	NC	NC + C	WC	WC + C	RWC + C
Smoothness	Lip-smooth	Lip-smooth	Lip-continuous	Lip-continuous	(relative) Lip-continuous
Constraint	$\mathbb{R}^n$	closed convex	closed convex	$\mathbb{R}^n$	closed convex
Stationary	$\mathbb{E}[\ \nabla f(x)\ _2^2]$	$\mathbb{E}[\ g_X(x)\ _2^2]$	$\mathbb{E}[\ \mathcal{G}_{1/(2\rho)}(x)\ ^2]$	$\mathbb{E}[\ \mathcal{G}_{1/(2\rho)}(x)\ ^2]$	$\mathbb{E}[\Delta_{1/(2\rho)}(x)]$
Complexity	$\mathcal{O}(1/\epsilon^2)$	$\mathcal{O}(1/\epsilon^2)$	$\tilde{\mathcal{O}}(1/\epsilon^2)$	$\mathcal{O}(1/\epsilon^2)$	$\mathcal{O}(1/\epsilon^2)$
Batch size	1	$\mathcal{O}(1/\epsilon)$	1 (per inner iter)	1	1
Setting	Euclidean	Euclidean	Euclidean	Euclidean	Non-Euclidean

C = Convex, NC = Nonconvex, WC = Weakly Convex, RWC = Relatively Weakly Convex

Lip-smooth = Lipschitz Smooth, Lip-continuous = Lipschitz Continuous

$g_X(x_t) = \alpha_t^{-1} \|x_t - x_{t+1}\|$ ;  $\mathcal{G}_{1/(2\rho)}(x)$  and  $\Delta_{1/(2\rho)}(x)$  are defined in (1.3) and (1.4)

## 1.2 Contribution

In this paper, we establish the non-asymptotic stationary convergence rate analysis of SMD for constrained stochastic composite optimization problems in the general form of (1.1), where the objective is  $\rho$ -relatively weakly convex. We consider SMD with distance generating function setting to be exactly  $\psi(x) = \omega(x)$ . Our results apply to the basic SMD and its proximal version as well as the deterministic variants. More specifically, the main contributions can be summarized as follows.

Firstly, inspired by Davis and Drusvyatskiy (2018), we construct a new measure of stationary convergence called *Bregman gradient mapping* based on the *Bregman proximal operator*:

$$\mathcal{G}_\lambda(x) := \frac{1}{\lambda}(x - \text{prox}_{\lambda T}(x)) \quad (1.3)$$

where the Bregman proximal operator  $\text{prox}_{\lambda T}(x) := \text{argmin}_{y \in X} \{T(y) + \frac{1}{\lambda} D_\omega(y, x)\}$ . When  $T(x)$  is  $\rho$ -relatively weakly convex, the stationary measure is well-defined as long as  $\lambda < \rho^{-1}$ . Note that this is very distinct from the notion of generalized projection gradient used in Ghadimi et al. (2016). We also define another measure induced by Bregman divergence, called *Bregman stationarity*,

$$\Delta_\lambda(x) := \frac{1}{\lambda^2} \cdot (D_\omega(x, \text{prox}_{\lambda T}(x)) + D_\omega(\text{prox}_{\lambda T}(x), x)). \quad (1.4)$$

This quantity provides a stronger convergence criterion since  $\|\mathcal{G}_\lambda(x)\|^2 \leq \Delta_\lambda(x)$ . When the distance generating function  $\omega(x) = \frac{1}{2}\|x\|_2^2$ , both stationary measures reduce to the one used in Davis and Grimmer (2017) and Davis and Drusvyatskiy (2018) based on the gradient of the Moreau envelope of the objective. Later, we provide detailed analysis of these stationary measures and its relations to the gradient of Bregman Moreau envelope and traditional stationary measure, *i.e.*,  $\text{dist}(0, \partial(T + \delta_X)(x))$  in this problem.

As a main result, we show that SMD converges to a  $\epsilon$ -stationary point such that  $\mathbb{E}[\Delta_{1/(2\rho)}(x)] \leq \epsilon$  within  $\mathcal{O}(1/\epsilon^2)$  iterations. The rate matches with that of stochastic subgradient method recently estab-

lished in [Davis and Drusvyatskiy \(2018\)](#) and implies that using mini-batch is not necessary for SMD to converge for relatively weakly convex problems. This appears to be the first non-asymptotic convergence result for SMD in the nonconvex, nonsmooth regime, to the best of our knowledge. We provide a unified and simplified convergence analysis that apply to both plain SMD and its proximal variant. In contrast, [Davis and Drusvyatskiy \(2018\)](#) requires different analysis for the projected and proximal versions of stochastic subgradient method.

Lastly, we extend these results to a much weaker condition by assuming only relative continuity of the objective function. We show that similar convergence result can be obtained under this relaxed assumption.

### 1.3 Paper Organization

The paper is organized as follows. In Section 2, we introduce the concepts of relative weak convexity and Bregman stationarity measures. We also review some important properties of Bregman Moreau envelope and Bregman proximal operator. In Section 3, we present the SMD algorithm and its stationary convergence guarantee. Finally, in Section 4, we further extend the results to relative Lipschitz continuous problems.

## 2 Relatively Weak Convexity and Bregman Stationarity

In this section, we first introduce the concept of relatively weakly convex functions and discuss some important properties and calculus of this family of nonconvex functions.

### 2.1 Relatively Weakly Convex Functions

Let  $X$  be a closed convex set and  $\mathcal{X}$  be its embedding Euclidean space associated with some norm  $\|\cdot\|$ . Let  $\omega(x) : X \rightarrow \mathbb{R}$  be a reference function that is continuously differentiable and 1-strongly convex on  $X$  with respect to the given norm  $\|\cdot\|$ , i.e.,  $\omega(x) - \omega(y) - \langle \nabla \omega(y), x - y \rangle \geq \frac{1}{2} \|x - y\|^2$  for any  $x, y \in X$ . This induces the *Bregman divergence*, denoted by  $D_\omega(x, y)$ :

$$D_\omega(x, y) = \omega(x) - \omega(y) - \langle \nabla \omega(y), x - y \rangle. \quad (2.1)$$

It follows immediately that  $D_\omega(x, y) \geq \frac{1}{2} \|x - y\|^2$ .

**Definition 2.1** (Relatively Weak Convexity) *A function  $f(x) : X \rightarrow \mathbb{R}$  is said to be  $\rho$ -relatively weakly convex on  $X$  with respect to the reference function  $\omega(x)$  if  $f(x) + \rho\omega(x)$  is convex on  $X$ . We denote  $f(x)$  as  $(\rho, \omega(\cdot))$ -RWC.*

The above definition generalizes the traditional notion of weak convexity introduced in [Vial \(1983\)](#) and extensively studied in existing works ([Drusvyatskiy, 2017](#); [Davis and Grimmer, 2017](#); [Davis and Drusvyatskiy, 2018](#)). In the case when  $\omega(x) = \frac{1}{2} \|x\|^2$  with some inner product induced norm  $\|\cdot\|$ , such as the Euclidean norm, the function  $f(x)$  is called  $\rho$ -weakly convex. Obviously, any  $\rho$ -weakly convex function is also  $(\rho, \omega(\cdot))$ -weakly convex, for any reference  $\omega(x)$  that is 1-strongly convex with respect to the norm  $\|\cdot\|$ . However, the class of relatively weakly convex functions can be much broader. For example, the function  $f(x) = -\sum_{i=1}^n x_i \log(x_i)$  is relatively weakly convex, but not weakly convex.

In what follows, we will provide some equivalent characterizations of relatively weakly convexity.

**Proposition 2.1** *Let  $X \subseteq U$ , where  $U$  is a convex open set. The following statements are equivalent:*

- (i)  $f(x)$  is  $(\rho, \omega(\cdot))$ -RWC on  $U$ .

(ii) For any fixed  $y \in U$ ,  $f_\rho(x; y) := f(x) + \rho D_\omega(x, y)$  is convex in  $x \in U$ .

(iii) For any fixed  $y \in U$ , there exists  $g \in \mathcal{X}$ , such that

$$f(x) \geq f(y) + \langle g, x - y \rangle - \rho D_\omega(x, y), \forall x \in U. \quad (2.2)$$

*Proof.* (i) $\Rightarrow$ (ii) is straightforward.  $(\rho, \omega(\cdot))$ -RWC implies that  $f(x) + \rho\omega(x)$  is convex. Hence,

$$f_\rho(x; y) = f(x) + \rho[\omega(x) - \omega(y) - \langle \nabla\omega(y), x - y \rangle] = [f(x) + \rho\omega(x)] - [\rho\omega(y) + \langle \nabla\omega(y), x - y \rangle]$$

is equal to the sum of a convex function and an affine function, thus convex. (ii) $\Rightarrow$ (iii) is also straightforward. Since  $f_\rho(x; y)$  is convex in  $x \in U$ , subgradients exist on  $U$ . Let  $g \in \partial f_\rho(y; y)$ . We have

$$f(x) + \rho D_\omega(x, y) \geq f(y) + \rho D_\omega(y, y) + \langle g, x - y \rangle, \forall x \in U$$

Rearranging the terms, we obtain the third statement. Lastly, we show (iii) $\Rightarrow$ (i). Invoking the definition of Bregman divergence, (2.2) implies that for any  $y \in U$ , there exists  $g \in \mathcal{X}$

$$\begin{aligned} f(x) &\geq f(y) + \langle g, x - y \rangle - \rho[\omega(x) - \omega(y) - \langle \nabla\omega(y), x - y \rangle], \forall x \in U \\ \Leftrightarrow [f(x) + \rho\omega(x)] &\geq [f(y) + \rho\omega(y)] + \langle g + \rho\nabla\omega(y), x - y \rangle, \forall x \in U \end{aligned}$$

This implies that  $f(x) + \rho\omega(x)$  is convex, *i.e.*,  $f(x)$  is  $(\rho, \omega(\cdot))$ -RWC. ■

In fact, the above results also provide a valid subdifferential set of relatively weakly convex functions and the construction of subgradients.

**Definition 2.2** (Subgradient and Subdifferential Set) *A vector  $g \in \mathcal{X}$  is a subgradient of  $f(x)$  at  $x \in X$  if  $g \in \partial f_\rho(x; x)$ . The subdifferential set of  $f(x)$  at  $x$ , denoted as  $\partial f(x)$ , contains all subgradients at  $x$ . Note that  $\partial f(x) = \partial(f + \rho\omega)(x) - \rho\nabla\omega(x)$ .*

For  $\rho$ -relatively weakly convex functions, the above subdifferential set is always well-defined and non-empty. When the function is weakly convex (thus locally Lipschitz), this set is also equivalent to the Fréchet subdifferential set and the Clarke differential set (Davis and Grimmer, 2017).

**Examples.** A major class of relatively weakly convex functions is the family of smooth functions with Lipschitz continuous gradients. Suppose  $f(x)$  is continuously differentiable and has  $\rho$ -Lipschitz continuous gradient, *i.e.*,  $\|\nabla f(x) - \nabla f(y)\|_* \leq \rho\|x - y\|$ , where  $\|\cdot\|_*$  is the dual norm of  $\|\cdot\|$ , then by fundamental theorem of calculus, this implies that  $|f(x) - f(y) - \langle \nabla f(y), x - y \rangle| \leq \frac{\rho}{2}\|x - y\|^2$ . Hence, it follows

$$f(x) \geq f(y) - \langle \nabla f(y), x - y \rangle - \frac{\rho}{2}\|x - y\|^2 \geq f(y) - \langle \nabla f(y), x - y \rangle - \rho D_\omega(x, y),$$

so  $f(x)$  is relatively weakly convex. In the case when both  $f(x)$  and  $\omega(x)$  is twice differentiable, relative weak convexity is equivalent to say  $\nabla^2 f(x) \succeq -\rho\nabla^2 \omega(x)$ . Hence, the family of relatively weakly convex functions also include functions that are not necessarily Lipschitz smooth, *e.g.*, the relatively smooth functions (Lu et al., 2018). Moreover, the following proposition gives some calculus and more examples of relatively weakly convex functions.

**Proposition 2.2** *Let  $X$  be a nonempty closed convex set.*

- (a) Suppose  $f_1 : X \rightarrow \mathbb{R}$  is  $(\rho_1, \omega_1(\cdot))$ -RWC and  $f_2 : X \rightarrow \mathbb{R}$  is  $(\rho_2, \omega_2(\cdot))$ -RWC on  $X$ , then  $f_1 + f_2$  is  $(\rho_1 + \rho_2, \bar{\omega}(\cdot))$ -RWC on  $X$ , where  $\bar{\omega}(x) = (\rho_1 + \rho_2)^{-1}(\rho_1\omega_1(x) + \rho_2\omega_2(x))$  is differentiable and 1-strongly convex on  $X$ .
- (b) Suppose  $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$  is  $(\rho_i, \omega(x))$ -RWC for  $i \in I$ , and  $\rho := \sup_{i \in I} \rho_i < \infty$ , then the supreme function  $f(x) := \sup_{i \in I} f_i(x)$  is also  $(\rho, \omega(x))$ -RWC.
- (c) Suppose  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is closed convex and  $L_f$ -Lipschitz continuous such that  $|f(u) - f(v)| \leq L_f \|u - v\|, \forall u, v \in \mathbb{R}^d$ , and suppose  $g : \mathbb{R}^n \rightarrow \mathbb{R}^d$  is  $L_g$ -relatively smooth with respect to  $\omega(x)$  such that for any  $x, y \in \mathbb{R}^n$

$$\|g(x) - g(y) - \langle \nabla g(y), x - y \rangle\| \leq L_g \cdot D_\omega(x, y).$$

Then the composition  $f \circ g : \mathbb{R}^n \rightarrow \mathbb{R}$  is  $(L_f L_g, \omega(\cdot))$ -RWC.

*Proof.* Parts (a) and (b) of Proposition 2.2 are straightforward. Part (c) of Proposition 2.2 is because that for any  $x, y, w \in \partial f(g(y))$ , it holds that

$$\begin{aligned} f \circ g(x) &\geq f \circ g(y) + \langle w, g(x) - g(y) \rangle \\ &= f \circ g(y) + \langle w, \nabla g(y)^T(x - y) \rangle + \langle w, g(x) - g(y) - g(y)^T(x - y) \rangle \\ &\geq f \circ g(y) + \langle w, \nabla g(y)^T(x - y) \rangle - L_g \|w\|_* \cdot D_\omega(x, y) \\ &\geq f \circ g(y) + \langle \nabla g(y)w, x - y \rangle - L_f L_g D_\omega(x, y) \end{aligned}$$

The first inequality is due to the convexity of  $f$ ; the third inequality is due to Hölders inequality; and the last inequality is due to the Lipschitz continuity of  $f$ .  $\blacksquare$

## 2.2 Bregman Moreau Envelope and Bregman Proximal Operator

We now revisit the basic properties of Bregman divergence and introduce the stationary measures based on Bregman Moreau envelope. We first list a few important properties of Bregman divergence that will be heavily used in the rest of the paper.

**Lemma 2.1** (Properties of Bregman Divergence, Section 9.2.1, Beck (2017))

- (a) The Bregman divergence satisfies the three-point identity:

$$D_\omega(x, y) + D_\omega(y, z) = D_\omega(x, z) + \langle \nabla \omega(z) - \nabla \omega(y), x - y \rangle, \forall x, y, z \in X \quad (2.3)$$

- (b) Suppose  $\phi(x)$  is convex and  $z^+ = \operatorname{argmin}_{x \in X} \{\phi(x) + \frac{1}{\alpha} D_\omega(x, z)\}$  for some  $\alpha > 0$ , then we have

$$\phi(x) + \frac{1}{\alpha} D_\omega(x, z) \geq \phi(z^+) + \frac{1}{\alpha} D_\omega(z^+, z) + \frac{1}{\alpha} D_\omega(x, z^+), \forall x \in X. \quad (2.4)$$

Below we provide the definitions of Bregman Moreau envelope and Bregman proximal operator, which are natural extensions of Moreau envelope and proximal operator by replacing Euclidean distance with Bregman divergence (Bauschke et al., 2006). Because of the asymmetry of Bregman divergence, we should

be careful when extending Moreau envelope directly to the Bregman case. We consider the (left) Bregman envelope and the (left) Bregman proximal operator here<sup>1</sup>.

**Definition 2.3** (Bregman Moreau envelope and proximal operator) *Given positive number  $\lambda > 0$  and a function  $T(x)$ , for a vector  $z \in X$ , we define its Bregman Moreau envelope as*

$$T_\lambda(z) := \min_{x \in X} \left\{ T(x) + \frac{1}{\lambda} D_\omega(x, z) \right\} \quad (2.5)$$

and the corresponding Bregman proximal operator

$$\text{prox}_{\lambda T}(z) := \operatorname{argmin}_{x \in X} \left\{ T(x) + \frac{1}{\lambda} D_\omega(x, z) \right\} \quad (2.6)$$

It is obvious that when the function  $T(\cdot)$  is convex, then the proximal operator is always well-defined and unique for any positive number  $\lambda > 0$ . In fact, this holds true for any  $(\rho, \omega(\cdot))$ -RWC functions as long as  $0 < \lambda < \rho^{-1}$ . More specifically, we have

**Lemma 2.2** (Uniqueness of Bregman proximal operator) *Suppose a function  $T(x)$  is  $(\rho, \omega(\cdot))$ -RWC on  $X$  and  $0 < \lambda < \rho^{-1}$ . Then for any input  $z \in X$ , the function  $T(x) + \frac{1}{\lambda} D_\omega(x, z)$  is  $(\lambda^{-1} - \rho)$ -strongly convex. Moreover, the Bregman proximal operator  $\text{prox}_{\lambda T}(z)$  is unique.*

The result follows directly from the definition of relative weak convexity. Same as the Euclidean case, one can also show that the Bregman Moreau envelope is differentiable.

**Lemma 2.3** (Gradient of Bregman Moreau envelope) *Suppose  $T(x)$  is a proper closed function and  $(\rho, \omega(\cdot))$ -RWC on  $X$ , and  $0 < \lambda < \rho^{-1}$ . Suppose the above DGF  $\omega(x)$  is also twice continuously differentiable. Then the Bregman Moreau envelope  $T_\lambda(z)$  is differentiable, and its gradient is given by*

$$\nabla T_\lambda(z) = \frac{1}{\lambda} \nabla^2 \omega(z) (z - \text{prox}_{\lambda T}(z)). \quad (2.7)$$

The result follows immediately from Propositions 3.10 and 3.12 in [Bauschke et al. \(2006\)](#) by using the convexity of  $\lambda T(x) + \omega(x)$ . For sake of simplicity, we do not repeat the details here.

## 2.3 Stationarity Measures

Since the major goal of solving a general nonsmooth nonconvex problem is to find a stationary point, we are mainly interested in analyzing the stationary convergence of the SMD algorithm. For the general constrained composite problem in the form of (1.1), a stationary point  $x^*$  can often be described as such that  $0 \in \partial(T + \delta_X)(x^*)$ , or equivalently,  $\text{dist}_{\|\cdot\|}(0, \partial(T + \delta_X)(x^*)) = 0$ . Here we use  $\text{dist}_{\|\cdot\|}(x, Q) := \inf_{y \in Q} \|y - x\|$  to characterize the distance between a point  $x \in X$  and a set  $Q$  under a specific norm  $\|\cdot\|$  and we use  $\delta_X(\cdot)$  to denote the indicator function of the set  $X$ .

Inspired by [Davis and Drusvyatskiy \(2018\)](#), a natural option to measure the stationarity of a candidate solution  $x \in X$  is by evaluating the difference of  $x$  and its proximity:

$$\mathcal{G}_\lambda(x) := \frac{1}{\lambda} (x - \text{prox}_{\lambda T}(x)), \text{ where } \lambda \in (0, \rho^{-1}). \quad (2.8)$$

---

<sup>1</sup> In fact, there are also ‘‘right’’ versions of the Bregman envelope and proximal operator, with some different properties ([Bauschke et al., 2006](#)), but here we will focus on the left version.



In the case when the DGF is  $\omega(x) = \frac{1}{2}\|x\|_2^2$ , it follows immediately from Lemma 2.3 that  $\mathcal{G}_\lambda(x) = \nabla T_\lambda(x)$ , *i.e.*, the gradient of the Moreau envelope of  $T(x)$  (Davis and Drusvyatskiy, 2018). Under such a case, invoking the definition of proximal operator, and denoting  $\hat{x} := \text{prox}_{\lambda T}(x)$ , we have  $\mathcal{G}_\lambda(x) \in \partial(T + \delta_X)(\hat{x})$ , so one can show that

$$\text{dist}_{\|\cdot\|_2}^2(0, \partial(T + \delta_X)(\hat{x})) \leq \|\mathcal{G}_\lambda(x)\|_2^2 = \|\nabla T_\lambda(x)\|_2^2. \quad (2.9)$$

Hence, the magnitude of  $\mathcal{G}_\lambda(x)$  provides an upper bound for the distance from the origin to the subdifferential set  $\partial f(x)$ , and can be used to measure the progress of iterations. In our case, for general choices of distance generating functions  $\omega(x)$ , this also makes sense. From Lemma 2.3, suppose  $\omega(x)$  is twice continuously differentiable, we have

$$\mathcal{G}_\lambda(x) = (\nabla^2 \omega(x))^{-1} \nabla T_\lambda(x),$$

which can be viewed as a rescaled gradient of the Bregman Moreau envelope. Then with the assumption that  $\omega(x)$  is 1-strongly convex with respect to  $\|\cdot\|$ -norm, we have  $\|\mathcal{G}_\lambda(x)\| \leq \|\nabla T_\lambda(x)\|$ . Moreover, from the definition of Bregman proximal operator  $\hat{x} = \text{prox}_{\lambda T}(x)$ , we have

$$0 \in \partial(T + \delta_X)(\hat{x}) + \frac{1}{\lambda}(\nabla \omega(\hat{x}) - \nabla \omega(x)) \approx \partial(T + \delta_X)(\hat{x}) + \nabla^2 \omega(x) \mathcal{G}_\lambda(x) \quad (2.10)$$

where the approximation is based on the first-order Taylor expansion of  $\nabla \omega(\cdot)$ . Hence when  $\|\mathcal{G}_\lambda(x)\|$  is small, it indicates that the origin is near the set  $\partial(T + \delta_X)(x)$ , *i.e.*,  $\hat{x}$  is close to a stationary point.

To better capture the geometry of the non-Euclidean setup, we propose to measure the stationarity of a candidate solution through evaluating the Bregman divergence the solution and its proximity:

$$\Delta_\lambda(x) := \frac{1}{\lambda^2} \cdot (D_\omega(x, \text{prox}_{\lambda T}(x)) + D_\omega(\text{prox}_{\lambda T}(x), x)). \quad (2.11)$$

We call this the *Bregman stationarity* measure. It follows immediately from the 1-strongly convexity of  $\omega(\cdot)$  that  $\|\mathcal{G}_\lambda(x)\|^2 \leq \Delta_\lambda(x)$ . Hence, the measure  $\Delta_\lambda(x)$  yields a stronger convergence criterion than using the squared norm of the Bregman gradient mapping. Further, suppose the distance generating function  $\omega(x)$  has  $M$ -Lipschitz continuous gradient, then we have

$$\text{dist}_{\|\cdot\|}^2(0, \partial(T + \delta_X)(\hat{x})) \leq \frac{1}{\lambda^2} \|\nabla \omega(x) - \nabla \omega(\hat{x})\|^2 \leq \frac{M}{\lambda^2} \langle \nabla \omega(x) - \nabla \omega(\hat{x}), x - \hat{x} \rangle = M \cdot \Delta_\lambda(x). \quad (2.12)$$

Hence, the measure defined by  $\Delta_\lambda(x)$  provides a valid characterization of the stationarity of a candidate solution in terms of the norm  $\|\cdot\|$ . In particular, for the  $\ell_1$ -setup with  $\|\cdot\| = \|\cdot\|_1$ , the squared distance defined by  $\ell_1$ -norm in (2.12) could be of order  $\mathcal{O}(n)$  larger than that defined by  $\ell_2$ -norm in (2.9).

### 3 Stationary Convergence of Stochastic Mirror Descent (SMD)

In this section, we formally describe the problem setting and assumptions, and revisit the stochastic mirror descent algorithm. We will then discuss its convergence behavior in terms of the previously defined stationarity measure.



### 3.1 Problem Setting and Assumptions

We consider the general composite stochastic optimization problem:

$$\min_{x \in X} T(x) := f(x) + r(x) = \mathbb{E}_\xi [F(x; \xi)] + r(x) \quad (3.1)$$

under the following assumptions:

**Assumption 3.1** *We assume that*

- (i) *The set  $X \subseteq \mathcal{X}$  is a closed convex subset of a finite-dimensional Euclidean space  $\mathcal{X}$ .*
- (ii) *The function  $f(x)$  is  $(\rho, \omega(\cdot))$ -RWC on  $X$ , for some function  $\omega(\cdot)$  that is continuously differentiable and 1-strongly convex (1-SC) on  $X$  with respect to the norm  $\|\cdot\|$  defined on the Euclidean space  $\mathcal{X}$ .*
- (iii) *There exists a stochastic oracle that outputs a random vector  $G(x, \xi)$  given input  $x \in X$ , such that*

$$\mathbb{E}_\xi [G(x, \xi)] \in \partial f(x) \quad (3.2)$$

where  $\partial f(x)$  is the subdifferential set of  $f(x)$  at  $x$ . Moreover, we assume there exists a constant  $L > 0$ , such that  $\forall x \in X$

$$\mathbb{E}[\|G(x, \xi)\|_*^2] \leq L^2; \quad (3.3)$$

This is sometimes called  $L$ -stochastically continuity of  $f(x)$  (Lu, 2017).

- (iv) *The function  $r(x) : X \rightarrow \mathbb{R}$  is proper, closed, convex, nonnegative and perhaps nonsmooth.*
- (v) *The optimal objective value, denoted as  $T_{\min}$ , exists and  $T_{\min} > -\infty$ .*

Note that here we assume the term  $r(\cdot)$  to be nonnegative, which is a common assumption for proximal algorithms in the literature; see e.g., Duchi et al. (2010) and Beck (2017). This assumption is also satisfied by a wide range of regularizations used in practical applications.

### 3.2 Stochastic Mirror Descent (SMD)

We now formally present the SMD algorithm as outlined in Algorithm 1 below. Here we are going to use  $\omega(x)$  as the distance generating function for the Bregman divergence used in the SMD algorithm. For the sake of generality, we will adopt the proximal variant of SMD, which has been extensively studied for convex problems; see, e.g., Duchi et al. (2010), He (2015), and Beck (2017). The only modifications we make is that when generating an output solution after  $N$  iterations, we will randomly pick one from the sequence  $\{x_0, x_1, \dots, x_{N-1}\}$  according to a fixed distribution based on the stepsizes.

We emphasize that the SMD algorithm significantly differs from the RSPG algorithm proposed in Ghadimi et al. (2016) in several aspects: first, we don't need to use mini-batch samples to construct the subgradient estimator; second, the stepsize has to be decaying or in the order of  $\mathcal{O}(1/\sqrt{N})$  rather than a large constant; third, the probability mass function for selecting a random output is much simpler.

### 3.3 Convergence Results

Below we present the stationary convergence result of SMD.

---

**Algorithm 1** Stochastic Mirror Descent (SMD)

---

Input  $x_0, N, \{\alpha_t\}_{t=0}^{N-1}$

**for**  $t = 0$  to  $N - 1$  **do**

    Obtain  $G_t := G(x_t, \xi_t)$  from the stochastic oracle

    Update  $x_{t+1} = \operatorname{argmin}_{x \in X} \left\{ \langle G_t, x \rangle + r(x) + \frac{1}{\alpha_t} D_\omega(x, x_t) \right\}$

**end for**

Output  $x_R$  from  $\{x_0, \dots, x_{N-1}\}$  with probability as  $P(R = i) = \frac{\alpha_i}{\sum_{t=0}^{N-1} \alpha_t}$ , ( $i = 0, 1, \dots, N - 1$ )

---

**Theorem 3.1** (Stationary Convergence of SMD) *Let  $x_R$  be the output of the SMD algorithm after  $N$  iterations with non-increasing stepsize  $\alpha_t > 0, t = 0, \dots, N - 1$ . Then we have for any  $\hat{\rho}$  such that  $\hat{\rho} > \rho$ ,*

$$\mathbb{E}[\Delta_{1/\hat{\rho}}(x_R)] \leq \frac{\hat{\rho}}{\hat{\rho} - \rho} \cdot \frac{T_{1/\hat{\rho}}(x_0) - T_{\min} + \hat{\rho}\alpha_0 r(x_0) + \frac{\hat{\rho}L^2}{2} \sum_{t=0}^{N-1} \alpha_t^2}{\sum_{t=0}^{N-1} \alpha_t}, \quad (3.4)$$

where  $\Delta_{1/\hat{\rho}}(x_R)$  is as defined in (2.11) and the expectation is taken with respect to  $R$  and  $(\xi_0, \dots, \xi_{N-1})$ .

The above theorem provides the first characterization of the non-asymptotic convergence behavior of the SMD algorithm in expectation. As discussed in previous section, the stationary measure  $\Delta_{1/\hat{\rho}}(x)$  with  $\hat{\rho} > \rho$  provides a meaningful way to evaluate the stationarity of a candidate solution and also captures the underlying geometry of the non-Euclidean setup. It is worth mentioning that this result generalizes the recent convergence results (Davis and Drusvyatskiy, 2018) for stochastic projected subgradient method and stochastic proximal subgradient method in a unified sense. In Davis and Drusvyatskiy (2018), the authors develop two different results and analysis for the projected and proximal versions of stochastic subgradient method. For the proximal version, their convergence result requires  $\hat{\rho} \in (\rho, 2\rho]$  and the stepsize  $\alpha_t \leq 1/\hat{\rho}$  for algebraic purposes. However, such requirements are not needed in our analysis.

In particular, if we select the stepsize to be a constant and set  $\hat{\rho} = 2\rho$ , our result yields

**Corollary 3.1** *For a fixed number of iterations  $N$ , by setting the stepsize to be a constant  $\alpha_t \equiv \frac{c}{\sqrt{N}}, t = 0, 1, \dots, N - 1$ , for some positive scalar  $c > 0$ , the solution  $x_R$  generated by the SMD algorithm satisfies*

$$\mathbb{E}[\Delta_{1/(2\rho)}(x_R)] \leq 2 \cdot \left( \frac{T_{1/(2\rho)}(x_0) - T_{\min} + \rho c^2 L^2}{c\sqrt{N}} + \frac{r(x_0)}{N} \right). \quad (3.5)$$

We can further optimize the choice of stepsize and obtain

**Corollary 3.2** *Suppose that  $T_{\min}$  is known and assume that we can initialize SMD with  $x_0$  such that  $r(x_0) = 0$ . Then by setting the stepsize to be  $\alpha_t \equiv \frac{c}{\sqrt{N}}, t = 0, 1, \dots, N - 1$ , such that*

$$c = \sqrt{\frac{T_{1/(2\rho)}(x_0) - T_{\min}}{\rho L^2}}, \quad (3.6)$$

we further have

$$\mathbb{E}[\|\mathcal{G}_{1/(2\rho)}(x_R)\|^2] \leq \mathbb{E}[\Delta_{1/(2\rho)}(x_R)] \leq \frac{4L\sqrt{\rho(T_{1/(2\rho)}(x_0) - T_{\min})}}{\sqrt{N}} \quad (3.7)$$

The above corollaries imply that the SMD algorithm converges to a stationary point in the rate of  $\mathcal{O}(1/\sqrt{N})$ . In other words, to obtain an  $\epsilon$ -stationary solution such that  $\mathbb{E}[\Delta_{1/(2\rho)}(x_R)] \leq \epsilon$ , the iteration

complexity and sample complexity for SMD is at most  $O\left(\frac{\rho L^2(T_{1/(2\rho)}(x_0) - T_{\min})}{\epsilon^2}\right)$ . The order of sample complexity, *i.e.*,  $\mathcal{O}(1/\epsilon^2)$  matches with that of existing algorithms, such as the RSPG algorithm (Ghadimi et al., 2016), the PGSG algorithm (Davis and Grimmer, 2017), and the proximal stochastic subgradient algorithm (Davis and Drusvyatskiy, 2018) for solving nonsmooth nonconvex optimization.

### 3.4 Convergence Analysis

In this section, we provide the detailed proof for Theorem 3.1.

*Proof.* For sake of simplicity, in what follows, we will denote  $\hat{x} := \text{prox}_{T/\hat{\rho}}(x)$  for any  $x \in X$ . First, by the definition of Bregman envelope, we have  $T_{1/\hat{\rho}}(x_{t+1}) = T(\hat{x}_{t+1}) + \hat{\rho}D_\omega(\hat{x}_{t+1}, x_{t+1})$ . The optimality of  $\hat{x}_{t+1}$  implies

$$T_{1/\hat{\rho}}(x_{t+1}) \leq T(\hat{x}_t) + \hat{\rho}D_\omega(\hat{x}_t, x_{t+1}). \quad (3.8)$$

Recall the definition of  $x_{t+1}$  and apply the three-point property introduced in Lemma 2.1(b) and equation (2.4) by setting  $z = x_t$ ,  $z^+ = x_{t+1}$ ,  $x = \hat{x}_t$  and  $\alpha = \alpha_t$ ,  $\phi(x) = \langle G_t, x \rangle + r(x)$ . We have

$$\alpha_t[\langle G_t, \hat{x}_t - x_{t+1} \rangle + r(\hat{x}_t) - r(x_{t+1})] \geq D_\omega(\hat{x}_t, x_{t+1}) + D_\omega(x_{t+1}, x_t) - D_\omega(\hat{x}_t, x_t) \quad (3.9)$$

Combing equations (3.8) and (3.9), we have

$$\begin{aligned} & \mathbb{E}\left[T_{1/\hat{\rho}}(x_{t+1})\right] \\ & \leq \mathbb{E}\left[T(\hat{x}_t) + \hat{\rho}\alpha_t\langle G_t, \hat{x}_t - x_{t+1} \rangle + \hat{\rho}\alpha_t(r(\hat{x}_t) - r(x_{t+1})) + \hat{\rho}D_\omega(\hat{x}_t, x_t) - \hat{\rho}D_\omega(x_{t+1}, x_t)\right] \\ & = \mathbb{E}\left[T_{1/\hat{\rho}}(x_t) + \hat{\rho}\alpha_t\langle G_t, \hat{x}_t - x_{t+1} \rangle + \hat{\rho}\alpha_t(r(\hat{x}_t) - r(x_{t+1})) - \hat{\rho}D_\omega(x_{t+1}, x_t)\right] \\ & = \mathbb{E}\left[T_{1/\hat{\rho}}(x_t)\right] + \hat{\rho}\alpha_t\mathbb{E}\left[\langle G_t, \hat{x}_t - x_t \rangle + (r(\hat{x}_t) - r(x_t))\right] + \hat{\rho}\mathbb{E}\left[\alpha_t(r(x_t) - r(x_{t+1}))\right] \\ & \quad + \hat{\rho}\mathbb{E}\left[\alpha_t\langle G_t, x_t - x_{t+1} \rangle - D_\omega(x_{t+1}, x_t)\right] \end{aligned} \quad (3.10)$$

where the first equality comes from the definition of  $T_{1/\hat{\rho}}(x_t)$ .

Next, invoking the  $(\rho, \omega(\cdot))$ -relatively weakly convexity of the function  $f(x)$ , we have

$$\mathbb{E}[\langle G_t, \hat{x}_t - x_t \rangle] \leq f(\hat{x}_t) - f(x_t) + \rho D_\omega(\hat{x}_t, x_t) \quad (3.11)$$

where the expectation is taking over  $\xi_t | \xi_0, \dots, \xi_{t-1}$ . Combine with  $r(x)$ , this implies that the second term in equation (3.10) can be bounded by

$$\mathbb{E}\left[\langle G_t, \hat{x}_t - x_t \rangle + r(\hat{x}_t) - r(x_t)\right] \leq T(\hat{x}_t) - T(x_t) + \rho D_\omega(\hat{x}_t, x_t) \quad (3.12)$$

Moreover, it is easy to see that the last term in equation (3.10) can also be bounded as follows,

$$\begin{aligned} \hat{\rho}\mathbb{E}\left[\alpha_t\langle G_t, x_t - x_{t+1} \rangle - D_\omega(x_{t+1}, x_t)\right] & \leq \hat{\rho}\mathbb{E}\left[\alpha_t\langle G_t, x_t - x_{t+1} \rangle - \frac{1}{2}\|x_{t+1} - x_t\|^2\right] \\ & \leq \frac{1}{2}\hat{\rho}\alpha_t^2 \cdot \mathbb{E}[\|G_t\|_*^2] \leq \frac{1}{2}\hat{\rho}\alpha_t^2 L^2 \end{aligned} \quad (3.13)$$

Here the first inequality is due to the fact that  $D_\omega(x, y) \geq \frac{1}{2}\|x - y\|^2$  and the second inequality is due to

Young's inequality. Hence, combining (3.10) with (3.12) and (3.13), we end up with

$$\begin{aligned} & \mathbb{E}\left[T_{1/\hat{\rho}}(x_{t+1})\right] \\ & \leq \mathbb{E}\left[T_{1/\hat{\rho}}(x_t) + \hat{\rho}\alpha_t(T(\hat{x}_t) - T(x_t) + \rho D_\omega(\hat{x}_t, x_t)) + \hat{\rho}\alpha_t(r(x_t) - r(x_{t+1})) + \frac{\hat{\rho}\alpha_t^2 L^2}{2}\right] \end{aligned} \quad (3.14)$$

Therefore, by telescoping the sum from  $t = 0$  to  $N - 1$ , and moving terms around, we further arrive at

$$\sum_{t=0}^{N-1} \mathbb{E}\left[\alpha_t(T(x_t) - T(\hat{x}_t) - \rho D_\omega(\hat{x}_t, x_t))\right] \quad (3.15)$$

$$\leq \frac{1}{\hat{\rho}}(T_{1/\hat{\rho}}(x_0) - T_{1/\hat{\rho}}(x_N)) + \sum_{t=0}^{N-1} \alpha_t(r(x_t) - r(x_{t+1})) + \frac{L^2}{2} \sum_{t=0}^{N-1} \alpha_t^2 \quad (3.16)$$

$$\leq \frac{1}{\hat{\rho}}(T_{1/\hat{\rho}}(x_0) - T_{\min}) + \alpha_0 r(x_0) + \frac{L^2}{2} \sum_{t=0}^{N-1} \alpha_t^2 \quad (3.17)$$

The last inequality is because that the stepsize  $\alpha_t$  is non-increasing and the function  $r(x)$  is nonnegative, which leads to

$$\sum_{t=0}^{N-1} \alpha_t(r(x_t) - r(x_{t+1})) = \alpha_0 r(x_0) - \alpha_{N-1} r(x_N) + \sum_{t=0}^{N-2} (\alpha_{t+1} - \alpha_t) r(x_{t+1}) \leq \alpha_0 r(x_0).$$

Finally, let us divide both sides of equation (3.17) by  $\sum_{t=0}^{N-1} \alpha_t$ , we finally obtain

$$\frac{\sum_{t=0}^{N-1} \mathbb{E}\left[\alpha_t(T(x_t) - T(\hat{x}_t) - \rho D_\omega(\hat{x}_t, x_t))\right]}{\sum_{t=0}^{N-1} \alpha_t} \leq \frac{T_{1/\hat{\rho}}(x_0) - T_{\min} + \hat{\rho}\alpha_0 r(x_0) + \rho D_\omega(\hat{x}_t, x_t)}{\hat{\rho} \sum_{t=0}^{N-1} \alpha_t}. \quad (3.18)$$

Invoking the definition of  $x_R$  in the algorithm, this implies that

$$\text{LHS} = \mathbb{E}\left[T(x_R) - T(\hat{x}_R) - \rho D_\omega(\hat{x}_R, x_R)\right]. \quad (3.19)$$

Recall that  $\hat{x}_R$  is the minimizer of the problem,  $\operatorname{argmin}_{x \in X} \{T(x) + \hat{\rho} D_\omega(x, x_R)\}$ , and the objective is  $(\hat{\rho} - \rho)$ -relatively strongly convex. It follows from Lemma 2.1 that

$$T(x_R) - [T(\hat{x}_R) + \hat{\rho} D_\omega(\hat{x}_R, x_R)] \geq (\hat{\rho} - \rho) D_\omega(x_R, \hat{x}_R) \quad (3.20)$$

Hence, we can further derive that

$$\begin{aligned} \text{LHS} &= \mathbb{E}\left[T(x_R) - T(\hat{x}_R) - \rho D_\omega(\hat{x}_R, x_R)\right] \\ &= \mathbb{E}\left[T(x_R) - [T(\hat{x}_R) + \hat{\rho} D_\omega(\hat{x}_R, x_R)] + (\hat{\rho} - \rho) D_\omega(\hat{x}_R, x_R)\right] \\ &\geq \mathbb{E}\left[(\hat{\rho} - \rho) D(x_R, \hat{x}_R) + (\hat{\rho} - \rho) D_\omega(\hat{x}_R, x_R)\right] \\ &= \frac{(\hat{\rho} - \rho)^2}{\hat{\rho}^2} \mathbb{E}\left[\Delta_{1/\hat{\rho}}(x_R)\right]. \end{aligned} \quad (3.21)$$

Here second inequality is from (3.20) and the last inequality is simply using the definition of  $\Delta_{1/\hat{\rho}}(x_R)$ . Combining with equation (3.18), we arrive at the desired result as stated in the theorem. ■

## 4 Extension to Relatively Continuous Nonconvex Problems

In this section, we further extend the previous stationary convergence results of SMD under relaxed assumptions of the Lipschitz continuity of the function  $f(x)$ . A standard condition for applying the SMD algorithm to stochastic nonsmooth problems is to assume that the stochastic gradient has bounded moments, i.e.,  $\max_{x \in X} \mathbb{E}[\|G(x, \xi)\|_*^2] \leq L^2$ . Recent works (e.g., Lu (2017)) show that such an assumption is not always satisfied in practice, particularly for those objectives without Lipschitz continuity. Here we generalize the convergence results to a broader class of nonsmooth nonconvex functions that are possibly non-Lipschitz continuous.

Let  $\omega(x) : X \rightarrow \mathbb{R}$  be a reference function that is differentiable and 1-strongly convex on  $X$  with respect to the given norm  $\|\cdot\|$  and  $D_\omega(x, y)$  be the Bregman divergence induced by  $\omega(x)$ .

### Definition 4.1 (Stochastically (Fréchet) Relatively Continuous Functions)

A function  $f(x)$  is called  $L$ -Stochastically relatively continuous with respect to  $\omega(x)$  on a set  $X$ , denoted as  $(L, \omega(\cdot))$ -SRC, for some positive constant  $L > 0$ , if for any  $x \in X$  and any unbiased estimator  $G(x, \xi)$  of the subgradient of  $f(\cdot)$  at  $x$ , satisfy  $\mathbb{E}[G(x, \xi)] \in \partial f(x)$ , and

$$\mathbb{E}[\|G(x, \xi)\|_*^2] \leq \frac{L^2 D_\omega(y, x)}{\frac{1}{2}\|y - x\|^2}, \forall y \neq x. \quad (4.1)$$

**Lemma 4.1** (Binomial Property of SRC Functions, Lu (2017)) *Let  $f(x)$  be a  $(L, \omega(\cdot))$ -SRC function and  $x \in X$ . Define the random vector*

$$M(x, \xi) := \|G(x, \xi)\|_* \cdot \max_{y \in X, y \neq x} \frac{\|y - x\|}{\sqrt{2D_\omega(y, x)}}. \quad (4.2)$$

Then it holds that

- (a)  $\mathbb{E}[M^2(x, \xi)] \leq L^2$ , and
- (b) For any  $\alpha > 0$ ,  $\langle \alpha G(x, \xi), x - y \rangle - D_\omega(y, x) \leq \frac{1}{2}\alpha^2 M^2(x, \xi)$ .

**Theorem 4.1** *Suppose that  $f(x)$  is  $(\rho, \omega(\cdot))$ -RWC and  $(L, \omega(\cdot))$ -SRC as defined. Let  $x_R$  be the output of SMD algorithm for solving the problem (1.1) within a fixed iteration number  $N > 0$ , and constant stepsize  $\alpha_t = c/\sqrt{N}$ , where  $c > 0$ . We have*

$$\mathbb{E}[\Delta_{1/(2\rho)}(x_R)] \leq 2 \cdot \left( \frac{T_{1/(2\rho)}(x_0) - T_{\min} + \rho c^2 L^2}{c\sqrt{N}} + \frac{r(x_0)}{N} \right). \quad (4.3)$$

*Proof.* The theorem can be proved with a slight modification of the proof as detailed in the previous section. When the function  $T(x)$  is  $(L, \omega(\cdot))$ -SRC, we can still prove the equation (3.13) by directly apply Lemma 4.1. ■

**Remark 4.1** *Note that the reference function  $\omega(\cdot)$  used to define either the relatively weak convexity or the stochastically relative continuity is the same as the distance generating function used in the SMD algorithm as well as in the Bregman Moreau envelope.*

**Remark 4.2** (Deterministic Setting) *The results developed in Sections 3 and 4 also apply to deterministic nonconvex problems and the deterministic Mirror Descent algorithm. Particularly, suppose we have access to a subgradient oracle that returns  $g(x) \in \partial f(x)$  for any input  $x$ , and  $\|g(x)\|_*^2 \leq \frac{L^2 D_\omega(y,x)}{\frac{1}{2}\|y-x\|^2}, \forall y \neq x$ . Suppose the Mirror Descent algorithm performs the updates:*

$$x_{t+1} = \operatorname{argmin}_{x \in X} \left\{ \langle g(x_t), x \rangle + r(x) + \frac{1}{\alpha_t} D_\omega(x, x_t) \right\}$$

for  $t = 0, 1, \dots, N - 1$ , and outputs a solution  $x_R$  such that  $x_R = \operatorname{argmin}_{t=0, \dots, N-1} \{\Delta_{1/(2\rho)}(x_t)\}$ . Then with constant stepsize  $\alpha_t = c/\sqrt{N}$ , where  $c > 0$ , we have

$$\Delta_{1/(2\rho)}(x_R) \leq 2 \cdot \left( \frac{T_{1/(2\rho)}(x_0) - T_{\min} + \rho c^2 L^2}{c\sqrt{N}} + \frac{r(x_0)}{N} \right). \quad (4.4)$$

In other words, the number of subgradient evaluations needed to obtain an  $\epsilon$ -stationary solution such that  $\Delta_{1/(2\rho)}(x) \leq \epsilon$ , is at most  $\mathcal{O}(1/\epsilon^2)$ . This result seems to be also the first non-asymptotic convergence result for the deterministic Mirror Descent algorithm and its proximal variant.

## 5 Conclusion

In this paper, we establish the first non-asymptotic convergence analysis of Stochastic Mirror Descent (SMD) for solving a general class of nonconvex nonsmooth optimization problems, under relaxed conditions of weak convexity and continuity. Our analysis applies to many variants in the family of SMD algorithms, and indicates that using mini-batch is not necessary for stationary convergence of SMD. We also show that using non-Euclidean setup could yield stronger stationarity guarantees. For future work, we will investigate the convergence behaviors of other algorithms in the SMD family under different settings, both in theory and in real applications.

## References

- Heinz H. Bauschke, Patrick L. Combettes, and Dominikus Noll. Joint minimization with alternating bregman proximity operators. *Pacific Journal of Optimization*, 2(3):401–424, 2006.
- Amir Beck. *First-Order Methods in Optimization*. Society for Industrial and Applied Mathematics, Philadelphia, PA, 2017.
- Laming Chen and Yuantao Gu. The convergence guarantees of a non-convex approach for sparse recovery. *IEEE Transactions on Signal Processing*, 62(15):3754–3767, 2014.
- Damek Davis and Dmitriy Drusvyatskiy. Stochastic subgradient method converges at the rate  $o(k^{-1/4})$  on weakly convex functions. *arXiv preprint arXiv:1802.02988*, 2018.
- Damek Davis and Benjamin Grimmer. Proximally guided stochastic subgradient method for nonsmooth, nonconvex problems. *arXiv preprint arXiv:1707.03505*, 2017.
- Damek Davis, Dmitriy Drusvyatskiy, and Courtney Paquette. The nonsmooth landscape of phase retrieval. *arXiv preprint arXiv:1711.03247*, 2017.
- Dmitriy Drusvyatskiy. The proximal point method revisited. *arXiv preprint arXiv:1712.06038*, 2017.

- John Duchi, Shai Shalev-Shwartz, Yoram Singer, and Ambuj Tewari. Composite objective mirror descent. In *Proc. of the 23th Annual Conference on Learning Theory*, pages 14–26, 2010.
- Saeed Ghadimi and Guanghui Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.
- Saeed Ghadimi, Guanghui Lan, and Hongchao Zhang. Mini-batch stochastic approximation methods for nonconvex stochastic composite optimization. *Mathematical Programming*, 155(1-2):267–305, 2016.
- Niao He. *Saddle Point Techniques in Convex Composite and Error-in-Measurement Optimization*. Ph.D Thesis. Georgia Institute of Technology, 2015.
- Jun Liu, Jianhui Chen, and Jieping Ye. Large-scale sparse logistic regression. In *Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 547–556, 2009.
- Haihao Lu. “Relative-continuity” for non-lipschitz non-smooth convex optimization using stochastic (or deterministic) mirror descent. *arXiv preprint arXiv:1710.04718*, 2017.
- Haihao Lu, Robert M Freund, and Yurii Nesterov. Relatively smooth convex optimization by first-order methods, and applications. *SIAM Journal on Optimization*, 28(1):333–354, 2018.
- Arkadii. S. Nemirovski and David. B. Yudin. *Problem Complexity and Method Efficiency in Optimization*. Wiley, New York, 1983.
- Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on Optimization*, 19(4):1574–1609, 2009.
- Xinyue Shen and Yuantao Gu. Nonconvex sparse logistic regression with weakly convex regularization. *IEEE Transactions on Signal Processing*, 66(12):3199–3211, 2018.
- Jean-Philippe Vial. Strong and weak convexity of sets and functions. *Mathematics of Operations Research*, 8(2): 231–259, 1983.
- Zhengyuan Zhou, Panayotis Mertikopoulos, Nicholas Bambos, Stephen Boyd, and Peter W Glynn. Stochastic mirror descent in variationally coherent optimization problems. In *Advances in Neural Information Processing Systems*, pages 7043–7052, 2017.