

---

# Stochastic Optimization for DC Functions and Non-smooth Non-convex Regularizers with Non-asymptotic Convergence

---

Yi Xu<sup>1</sup> Qi Qi<sup>1</sup> Qihang Lin<sup>2</sup> Rong Jin<sup>3</sup> Tianbao Yang<sup>1</sup>

## Abstract

Difference of convex (DC) functions cover a broad family of non-convex and possibly non-smooth and non-differentiable functions, and have wide applications in machine learning and statistics. Although deterministic algorithms for DC functions have been extensively studied, stochastic optimization that is more suitable for learning with big data remains under-explored. In this paper, we propose new stochastic optimization algorithms and study their first-order convergence theories for solving a broad family of DC functions. We improve the existing algorithms and theories of stochastic optimization for DC functions from both practical and theoretical perspectives. Moreover, we extend the proposed stochastic algorithms for DC functions to solve problems with a general non-convex non-differentiable regularizer, which does not necessarily have a DC decomposition but enjoys an efficient proximal mapping. To the best of our knowledge, this is the first work that gives the first non-asymptotic convergence for solving non-convex optimization whose objective has a general non-convex non-differentiable regularizer.

## 1. Introduction

In this paper, we consider a family of non-convex non-smooth optimization problems that can be written in the following form:

$$\min_{\mathbf{x} \in \mathbb{R}^d} g(\mathbf{x}) + r(\mathbf{x}) - h(\mathbf{x}), \quad (1)$$

---

<sup>1</sup>Department of Computer Science, University of Iowa, Iowa City, IA 52242, USA <sup>2</sup>Department of Management Sciences, University of Iowa, Iowa City, IA 52242, USA <sup>3</sup>Machine Intelligence Technology, Alibaba Group, Bellevue, WA 98004, USA. Correspondence to: Yi Xu <yi-xu@uiowa.edu>, Tianbao Yang <tianbao-yang@uiowa.edu>.

where  $g(\cdot)$  and  $h(\cdot)$  are real-valued lower-semicontinuous convex functions,  $r(\cdot)$  is a proper lower-semicontinuous function. We include the component  $r$  in order to capture non-differentiable functions that usually play the role of regularization, e.g., the indicator function of a convex set  $\mathcal{X}$  where  $r(\mathbf{x}) = \delta_{\mathcal{X}}(\mathbf{x})$  is zero if  $\mathbf{x} \in \mathcal{X}$  and infinity otherwise, and a non-differential regularizer such as the convex  $\ell_1$  norm  $\|\mathbf{x}\|_1$  or the non-convex  $\ell_0$  norm and  $\ell_p$  norm  $\|\mathbf{x}\|_p^p$  with  $p \in (0, 1)$ . We do not necessarily impose smoothness condition on  $g(\mathbf{x})$  or  $h(\mathbf{x})$  and the convexity condition on  $r(\mathbf{x})$ .

A special class of the problem (1) is the one with  $r(\mathbf{x})$  being a convex function - also known as difference of convex (DC) functions. We would like to mention that even the family of DC functions is broader enough to cover many interesting non-convex problems that are well-studied, including an additive composition of a smooth non-convex function and a non-smooth convex function, weakly convex functions, etc. We postpone this discussion to Section 2 after we formally introduce the definitions of smooth functions and weakly convex functions.

In the literature, deterministic algorithms for DC problems have been studied extensively since its introduction by Pham Dinh Tao in 1985 and are continuously receiving attention from the community (Khamaru & Wainwright, 2018; Wen et al., 2018). Please refer to (Thi & Dinh, 2018) for a survey on this subject. Although stochastic optimization (SO) algorithms for the special cases of DC functions mentioned above (smooth non-convex functions, weakly convex functions) have been studied recently (Davis & Grimmer, 2017; Davis & Drusvyatskiy, 2018b;a; Drusvyatskiy & Paquette, 2018; Chen et al., 2018b; Lan & Yang, 2018; Allen-Zhu, 2017; Chen & Yang, 2018; Allen-Zhu & Hazan, 2016; Reddi et al., 2016b;a; Zhang & He, 2018), a comprehensive study of SO algorithms with a broader applicability to the DC functions and the problem (1) with a non-smooth non-convex regularizer  $r(\mathbf{x})$  still remain rare.

The papers by (Mairal, 2013), (Nitanda & Suzuki, 2017) and (Thi et al., 2017) are the most related works dedicated to the stochastic optimization of special DC functions. Mairal (2013) studied a special case of problem (1) with  $h$  is smooth and proposed a stochastic majorization-minimization algo-

rithm enjoying an asymptotic convergence result for finding a stationary point. Thi et al. (2017) considered a special class of DC problems and they reformulated the problem into (1) such that  $h$  is a sum of  $n$  convex functions, and  $g$  is a quadratic function and  $r$  is the first component of the DC decomposition of the regularizer. Then, they proposed a stochastic variant of the classical DCA (Difference-of-Convex Algorithm) and established an asymptotic convergence result for finding a critical point. To our knowledge, the paper by (Nitanda & Suzuki, 2017) is the probably the first result that gives non-asymptotic convergence for finding an approximate critical point of a special class of DC problems, in which both  $g$  and  $h$  can be stochastic functions and  $r = 0$ . Their algorithm consists of multiple stages of solving a convex objective that is constructed by linearizing  $h(\mathbf{x})$  and adding a quadratic regularization. However, their algorithm and convergence theory have the following drawbacks. First, at each stage, they need to compute an unbiased stochastic gradient denoted by  $\mathbf{v}(\mathbf{x})$  of  $\nabla h(\mathbf{x})$  such that  $\mathbb{E}[\|\mathbf{v}(\mathbf{x}) - \nabla h(\mathbf{x})\|^2] \leq \epsilon^2$ , where  $\epsilon$  is the accuracy level imposed on the returned solution in terms of the gradient's norm. In reality, one has to resort to mini-batching technique by using a large number of samples to ensure this condition, which is impractical and not user-friendly. An user has to worry about what is the size of the mini-batch in order to find a sufficiently accurate solution while keeping the computational costs minimal. Second, for each constructed convex subproblem, their theory requires running a stochastic algorithm that solves each subproblem to the accuracy level of  $\epsilon$ , which could waste a lot of computations at earlier stages. Third, their convergence analysis requires that  $r(\mathbf{x}) = 0$  and  $g(\mathbf{x})$  is a smooth function with a Lipschitz continuous gradient. In addition, they obtained fast convergence result of the problem under Polyak-Łojasiewicz condition, which is not considered in this paper.

**Our Contributions - I.** In Section 3, we propose new stochastic optimization algorithms and establish their convergence results for solving the DC class of the problem (1) that improves the algorithm and theory in (Nitanda & Suzuki, 2017) from several perspectives. It is our intention to address the aforementioned drawbacks of their algorithm and theory. In particular, (i) our algorithm only requires unbiased stochastic (sub)-gradients of  $g(\mathbf{x})$  and  $h(\mathbf{x})$  without a requirement on the small variance of the used stochastic (sub)-gradients; (ii) we do not need to solve each constructed subproblem to the accuracy level of  $\epsilon$ . Instead, we allow the accuracy for solving each constructed subproblem to grow slowly without sacrificing the overall convergence rate; (iii) we improve the convergence theory significantly. First, our convergence analysis does not require  $g(\mathbf{x})$  to be smooth with a Lipschitz continuous gradient. Instead, we only require either  $g(\mathbf{x}) + r(\mathbf{x})$  or  $h(\mathbf{x})$  to be differentiable with a Hölder continuous gradient, under the former condition

$h(\mathbf{x})$  can be a non-smooth non-differentiable function and under the later condition  $r(\mathbf{x})$  and  $g(\mathbf{x})$  can be non-smooth non-differentiable functions. Second, the convergence rate is automatically adaptive to the Hölder continuity of the involved function without requiring the knowledge of the Hölder continuity to run the algorithm. Third, when adaptive stochastic gradient method is employed to solve each subproblem, we establish an adaptive convergence similar to existing theory of AdaGrad for convex problems (Duchi et al., 2011; Chen et al., 2018a) and weakly convex problems (Chen et al., 2018b), which is missing in (Nitanda & Suzuki, 2017).

**Our Contributions - II.** Moreover, in Section 4 we extend our algorithm and theory to the more general class of non-convex non-smooth problem (1), in which  $r(\mathbf{x})$  is a general non-convex non-differentiable regularizer that enjoys an efficient proximal mapping. Although such kind of non-smooth non-convex regularization has been considered in literature (Attouch et al., 2013; Bolte et al., 2014; Bot et al., 2016; Li & Lin, 2015; Yu et al., 2015; Yang, 2018; Liu et al., 2018; An & Nam, 2017; Zhong & Kwok, 2014), existing results are restricted to deterministic optimization and asymptotic or local convergence analysis. In addition, most of them consider a special case of our problem with  $g - h$  being a smooth non-convex function. To the best of our knowledge, this is the first work of stochastic optimization with a non-asymptotic first-order convergence result for tackling the non-convex objective (1) with a non-convex non-differentiable regularization and a smooth function  $g$  and a possibly non-smooth function  $h$  with a Hölder continuous gradient. Our algorithm and theory are based on using the Moreau envelope of  $r(\mathbf{x})$  that can be written as a DC function, which then reduces to the problem that is studied in Section 3. By using the algorithms and their convergence results established in Section 3 and carefully controlling the approximation parameter, we establish the first non-asymptotic convergence of stochastic optimization for solving the original non-convex problem with a non-convex non-differentiable regularizer. This non-asymptotic convergence result can be also easily extended to the deterministic optimization, which itself is novel and could be interesting to a broader community. A summary of our results is presented in Table 1.

## 2. Preliminaries

In this section, we present some preliminaries. Let  $\|\cdot\|_p$  denote the standard  $p$ -norm with  $p \geq 0$ . For a non-convex function  $f(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}$ , let  $\hat{\partial}f(\mathbf{x})$  denote the Fréchet subgradient and  $\partial f(\mathbf{x})$  denote the limiting subgradient, i.e.,

$$\hat{\partial}f(\bar{\mathbf{x}}) = \left\{ \mathbf{v} \in \mathbb{R}^d : \liminf_{\mathbf{x} \rightarrow \bar{\mathbf{x}}} \frac{f(\mathbf{x}) - f(\bar{\mathbf{x}}) - \mathbf{v}^\top (\mathbf{x} - \bar{\mathbf{x}})}{\|\mathbf{x} - \bar{\mathbf{x}}\|} \geq 0 \right\},$$

$$\partial f(\bar{\mathbf{x}}) = \{ \mathbf{v} \in \mathbb{R}^d : \exists \mathbf{x}_k \xrightarrow{f} \bar{\mathbf{x}}, v_k \in \hat{\partial}f(\mathbf{x}_k), \mathbf{v}_k \rightarrow \mathbf{v} \},$$

Table 1. Summary of our results for finding a (nearly)  $\epsilon$ -critical point of the problem (1), where  $g$  and  $h$  are assumed to be convex. HC refers to Hölder continuous gradient condition; SM refers to the smooth condition; CX means convex; NC means non-convex and NS means non-smooth; LP denotes Lipchitz continuous function; LB means lower bounded over  $\mathbb{R}^d$ ; FV means finite-valued over  $\mathbb{R}^d$ ; FVC means finite-valued over a compact set.  $\nu \in (0, 1]$  denotes the power constant of the involved function's Hölder continuity.  $n$  denotes the total number of components in a finite-sum problem. SPG denotes stochastic proximal gradient algorithm. SVRG denotes stochastic variance reduced gradient algorithm. AdaGrad denotes adaptive stochastic gradient method. Complexity for SPG and AdaGrad means iteration complexity, and for SVRG and AG means gradient complexity. **Better** results were obtained in the arXiv version after ICML.

$g$	$h$	$r$	Algorithms for subproblems	Complexity
-	HC	CX	SPG, AdaGrad	$O(1/\epsilon^{4/\nu})$
SM	HC	CX	SVRG	$O(n/\epsilon^{2/\nu})$
HC	-	CX, HC	SPG, AdaGrad	$O(1/\epsilon^{4/\nu})$
SM	-	CX, HC	SVRG	$O(n/\epsilon^{2/\nu})$
SM	HC	NC, NS, LP	SPG, AdaGrad	$O(1/\epsilon^{4(1+1/\nu)})$
SM	HC	NC, NS, FV, LB	SPG, AdaGrad	$O(1/\epsilon^{4(1+2/\nu)})$
SM	HC	NC, NS, LP	SVRG	$O(n/\epsilon^{2(1+1/\nu)})$
SM	HC	NC, NS, FV, LB	SVRG	$O(n/\epsilon^{2(1+2/\nu)})$
SM	HC	NC, NS, FVC	SVRG	$O(n/\epsilon^{2(1+2/\nu)})$

where the notation  $\mathbf{x} \xrightarrow{f} \bar{\mathbf{x}}$  means that  $\mathbf{x} \rightarrow \bar{\mathbf{x}}$  and  $f(\mathbf{x}) \rightarrow f(\bar{\mathbf{x}})$ . It is known that  $\hat{\partial}f(\mathbf{x}) \subseteq \partial f(\mathbf{x})$ . If  $f(\cdot)$  is differential at  $\mathbf{x}$ , then  $\hat{\partial}f(\mathbf{x}) = \{\nabla f(\mathbf{x})\}$ . Moreover, if  $f(\mathbf{x})$  is continuously differentiable on a neighborhood of  $\mathbf{x}$ , then  $\partial f(\mathbf{x}) = \{\nabla f(\mathbf{x})\}$ . When  $f$  is convex, the Fréchet and the limiting subgradient reduce to the subgradient in the sense of convex analysis:  $\partial f(\mathbf{x}) = \{\mathbf{v} \in \mathbb{R}^d : f(\mathbf{x}) \geq f(\mathbf{y}) + \mathbf{v}^\top(\mathbf{x} - \mathbf{y}), \forall \mathbf{y} \in \mathbb{R}^d\}$ . For simplicity, we use  $\|\cdot\|$  to denote the Euclidean norm (aka. 2-norm) of a vector. Let  $\text{dist}(\mathcal{S}_1, \mathcal{S}_2)$  denote the distance between two sets and  $[K] = \{1, \dots, K\}$ .

A function  $f(\mathbf{x})$  is smooth with a  $L$ -Lipchitz continuous gradient if it is differentiable and the following inequality holds  $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|, \forall \mathbf{x}, \mathbf{y}$ . A differentiable function  $f(\mathbf{x})$  has  $(L, \nu)$ -Hölder continuous gradient if there exists  $\nu \in (0, 1]$  such that  $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|^\nu, \forall \mathbf{x}, \mathbf{y}$ . Next, let us characterize the critical points of the considered problem (1) that are standard in the literature (Hiriart-Urruty, 1985; Horst & Thoai, 1999; Thi & Dinh, 2018; An & Nam, 2017), and introduce the convergence measure of an algorithm. First, let us consider the DC problem:

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) := g(\mathbf{x}) - h(\mathbf{x}), \quad (2)$$

where  $g : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  is a proper lower semicontinuous convex function and  $h : \mathbb{R}^d \rightarrow \mathbb{R}$  is convex. Any point  $\bar{\mathbf{x}}$  such that  $\partial h(\bar{\mathbf{x}}) \cap \partial g(\bar{\mathbf{x}}) \neq \emptyset$  is called a critical point of (2), which is a necessary condition for  $\bar{\mathbf{x}}$  to be a local minimizer. For an iterative optimization algorithm, it is hard to find an exactly critical point in a finite-number of iterations. Therefore, we find an  $\epsilon$ -critical point  $\mathbf{x}$  that satisfies

$$\text{dist}(\partial h(\mathbf{x}), \partial g(\mathbf{x})) \leq \epsilon. \quad (3)$$

Similarly, we can extend the above definition of critical points to the general problem (1) with  $r(\mathbf{x})$  being a

proper and lower semi-continuous (possibly non-convex) function (An & Nam, 2017). In particular, any point  $\bar{\mathbf{x}}$  such that  $\partial h(\bar{\mathbf{x}}) \cap \hat{\partial}(g+r)(\bar{\mathbf{x}}) \neq \emptyset$  is called a critical point of (1). When  $g$  is differentiable,  $\hat{\partial}(g+r)(\bar{\mathbf{x}}) = \nabla g(\bar{\mathbf{x}}) + \hat{\partial}r(\bar{\mathbf{x}})$  (Rockafellar & Wets, 1998)[Exercise 8.8], and when both  $g$  and  $r$  are convex and their domains cannot be separated  $\hat{\partial}(g+r)(\bar{\mathbf{x}}) = \partial g(\bar{\mathbf{x}}) + \partial r(\bar{\mathbf{x}})$  (Rockafellar & Wets, 1998)[Corollary 10.9]. An  $\epsilon$ -critical point of (1) is a point  $\mathbf{x}$  that satisfies  $\text{dist}(\partial h(\mathbf{x}), \hat{\partial}(g+r)(\mathbf{x})) \leq \epsilon$ . It is notable that when  $g+r$  is non-differentiable, finding an  $\epsilon$ -critical point could become a challenging task for an iterative algorithm even under the condition that  $r$  is a convex function. Let us consider the example of  $g = |x|, h = r = 0$ . As long as  $x \neq 0$ , we have  $\text{dist}(0, \partial|x|) = 1$ . To address this challenge when  $g+r$  is non-differentiable, we introduce the notion of nearly  $\epsilon$ -critical points. In particular, a point  $\mathbf{x}$  is called a nearly  $\epsilon$ -critical point of the problem (1) if there exists  $\bar{\mathbf{x}}$  such that

$$\|\mathbf{x} - \bar{\mathbf{x}}\| \leq O(\epsilon), \quad \text{dist}(\partial h(\bar{\mathbf{x}}), \hat{\partial}(g+r)(\bar{\mathbf{x}})) \leq \epsilon. \quad (4)$$

A similar notion of nearly critical points for non-smooth and non-convex optimization problems have been utilized in several recent works (Davis & Grimmer, 2017; Davis & Drusvyatskiy, 2018b;a; Chen et al., 2018b).

### Examples and Applications of DC functions.

*Example 1: Weakly convex functions.* Weakly convex functions have been recently studied in numerous papers (Davis & Grimmer, 2017; Davis & Drusvyatskiy, 2018b;a; Chen et al., 2018b; Zhang & He, 2018). A function  $f(\mathbf{x})$  is called  $\rho$ -weakly convex if  $f(\mathbf{x}) + \frac{\rho}{2}\|\mathbf{x}\|^2$  is a convex function. More generally,  $f(\mathbf{x})$  is called  $\rho$ -relative convex with respect to a strongly convex function  $\omega(\mathbf{x})$  if  $f(\mathbf{x}) + \rho\omega(\mathbf{x})$  is convex (Zhang & He, 2018). It is obvious that a weakly convex function  $f(\mathbf{x})$  is a DC function. Examples of weakly convex functions can be found in deep neural networks with a

smooth active function and a smooth/non-smooth loss function (Chen et al., 2018b), robust learning (Xu et al., 2018), robust phase retrieval (Davis & Drusvyatskiy, 2018a).

*Example 2: Non-Convex Sparsity-Promoting Regularizers.* Many non-convex sparsity-promoting regularizers in statistics can be written as a DC function, including log-sum penalty (LSP) (Candès et al., 2008), minimax concave penalty (MCP) (Zhang, 2010a), smoothly clipped absolute deviation (SCAD) (Fan & Li, 2001), capped  $\ell_1$  penalty (Zhang, 2010b), transformed  $\ell_1$  norm (Zhang & Xin, 2018). For detailed DC composition of these regularizers, please refer to (Wen et al., 2018; Gong et al., 2013). We also present the details in the supplement.

*Example 3: Least-squares Regression with  $\ell_{1-2}$  Regularization.* Recently, a non-convex regularization in the form of  $\lambda(\|\mathbf{x}\|_1 - \|\mathbf{x}\|_2)$  was proposed for least-squares regression or compressive sensing (Yin et al., 2015), which is naturally a DC function.

*Example 4: Positive-Unlabeled (PU) Learning.* In PU learning for binary classification, only positive data  $\{(\mathbf{z}_i, +1), i = 1, \dots, n_+\}$  are observed where  $\mathbf{z}_i \in \mathbb{R}^m$  denotes the feature vector of  $i$ -th positive example, conventional empirical risk minimization becomes problematic. A remedy to address this challenge is to use unlabeled data for computing an unbiased estimation of  $\mathbb{E}_{\mathbf{z}, y}[\ell(\mathbf{x}; \mathbf{z}, y)]$ , where  $y \in \{1, -1\}$  denotes the label. In particular, the objective in the following problem is an unbiased risk (Kiryo et al., 2017):  $\min_{\mathbf{x} \in \mathbb{R}^d} \frac{\pi_p}{n_+} \sum_{i=1}^{n_+} (\ell(\mathbf{x}; \mathbf{z}_i, 1) - \ell(\mathbf{x}; \mathbf{z}_i, -1)) + \frac{\sum_{j=1}^{n_u} \ell(\mathbf{x}; \mathbf{z}_j^u, -1)}{n_u}$ , where  $\{\mathbf{z}_i^u, i = 1, \dots, n_u\}$  is a set of unlabeled data, and  $\pi_p = \Pr(y = 1)$  is the prior probability of the positive class. It is obvious that if  $\ell(\mathbf{x}; \cdot)$  is a convex loss function in terms of  $\mathbf{x}$ , the above objective function is a DC function. In practice, an estimation of  $\pi_p$  is used.

**Examples of Non-Convex Non-Smooth Regularizers.** Finally, we present some examples of non-convex non-smooth regularizers  $r(\mathbf{x})$  that cannot be written as a DC function or whose DC decomposition is unknown. Thus, the algorithms and theories presented in Section 3 are not directly applicable, but the algorithms discussed in Section 4 are applicable when the proximal mapping of  $r(\mathbf{x})$  is efficient to compute. Examples include  $\ell_0$  norm (i.e., the number of non-zero elements of a vector) and  $\ell_p$  norm regularization for  $p \in (0, 1)$  (i.e.,  $\sum_{i=1}^d |x_i|^p$ ), whose proximal mapping can be efficiently computed (Attouch et al., 2013; Bolte et al., 2014). Let us consider a non-convex optimization problem with domain constraint  $\mathbf{x} \in \mathcal{C}$ , where  $\mathcal{C}$  is a non-convex set. Directly handling a non-convex constrained problem could be difficult. An alternative solution is to convert the constraint into a penalization  $r(\mathbf{x}) = \frac{\lambda}{2} \|\mathbf{x} - \mathbb{P}_{\mathcal{C}}(\mathbf{x})\|^2$  with  $\lambda > 0$  in the objective, where  $\mathbb{P}_{\mathcal{C}}(\cdot)$  denotes the projection of a point to the set  $\mathcal{C}$ . Note that when  $\mathcal{C}$  is a non-convex set,

$r(\mathbf{x})$  is a non-convex non-smooth function in general, and its proximal mapping enjoys a closed-form solution (Li & Pong, 2016).

### 3. New Stochastic Algorithms of DC functions

In this section, we present new stochastic algorithms for solving the problem (1) when  $r(\mathbf{x})$  is a convex function and their convergence results. We assume both  $g(\mathbf{x})$  and  $h(\mathbf{x})$  have a large number of components such that computing a stochastic gradient is much more efficient than computing a deterministic gradient. Without loss of generality, we assume  $g(\mathbf{x}) = \mathbb{E}_{\xi}[g(\mathbf{x}; \xi)]$  and  $h(\mathbf{x}) = \mathbb{E}_{\varsigma}[h(\mathbf{x}; \varsigma)]$ , and consider the following problem:

$$\min_{\mathbf{x} \in \mathbb{R}^d} F(\mathbf{x}) := \mathbb{E}_{\xi}[g(\mathbf{x}; \xi)] + r(\mathbf{x}) - \mathbb{E}_{\varsigma}[h(\mathbf{x}; \varsigma)]. \quad (5)$$

where  $g$  and  $h$  are real-valued lower-semicontinuous convex functions and  $r$  is a proper lower-semicontinuous convex function. It is notable that a special case of this problem is the finite-sum form:

$$\min_{\mathbf{x} \in \mathbb{R}^d} F(\mathbf{x}) := \frac{1}{n_1} \sum_{i=1}^{n_1} g_i(\mathbf{x}) + r(\mathbf{x}) - \frac{1}{n_2} \sum_{j=1}^{n_2} h_j(\mathbf{x}), \quad (6)$$

which allows us to develop faster algorithms for smooth functions by using variance reduction techniques.

Since we do not necessarily impose any smoothness assumption on  $g(\mathbf{x})$  and  $h(\mathbf{x})$ , we will postpone the particular assumptions for these functions in the statements of later theorems. For all algorithms presented below, we assume that the **proximal mapping** of  $r(\mathbf{x})$  can be efficiently computed, i.e.,  $\text{prox}_{\eta r}(\mathbf{y}) = \arg \min_{\mathbf{x} \in \mathbb{R}^d} \frac{1}{2\eta} \|\mathbf{x} - \mathbf{y}\|^2 + r(\mathbf{x})$  can be easily computed for any  $\eta > 0$ . But it is not necessary for developing subgradient methods when  $r$  is convex.

The basic idea of the proposed algorithm is similar to the stochastic algorithm proposed in (Nitanda & Suzuki, 2017). The algorithm consists of multiple stages of solving convex problems. At the  $k$ -th stage ( $k \geq 1$ ), given a point  $\mathbf{x}_k$ , a convex majorant function  $F_{\mathbf{x}_k}^{\gamma}(\mathbf{x})$  is constructed as following such that  $F_{\mathbf{x}_k}^{\gamma}(\mathbf{x}) \geq F(\mathbf{x}), \forall \mathbf{x}$  and  $F_{\mathbf{x}_k}^{\gamma}(\mathbf{x}_k) = F(\mathbf{x}_k)$ :

$$F_{\mathbf{x}_k}^{\gamma}(\mathbf{x}) = g(\mathbf{x}) + r(\mathbf{x}) - (h(\mathbf{x}_k) + \partial h(\mathbf{x}_k)^{\top}(\mathbf{x} - \mathbf{x}_k)) + \frac{\gamma}{2} \|\mathbf{x} - \mathbf{x}_k\|^2, \quad (7)$$

where  $\gamma > 0$  is a constant parameter. Then a stochastic algorithm is employed to optimize  $F_{\mathbf{x}_k}^{\gamma}$ . The key difference from the previous work lies at how to solve each convex majorant function. An important change introduced to our design is to make the proposed algorithms more efficient and more practical. Roughly speaking, we only require solving each function  $F_{\mathbf{x}_k}^{\gamma}(\mathbf{x})$  up to an accuracy level of  $c/k$  for some constant  $c > 0$ , i.e., finding  $\mathbf{x}_{k+1}$  such that

$$\mathbb{E}[F_{\mathbf{x}_k}^{\gamma}(\mathbf{x}_{k+1}) - \min_{\mathbf{x} \in \mathbb{R}^d} F_{\mathbf{x}_k}^{\gamma}(\mathbf{x})] \leq c/k. \quad (8)$$

In contrast, the results presented in (Nitanda & Suzuki,

**Algorithm 1** Stagewise Stochastic DC (SSDC) Algorithm

- 1: **Initialize:**  $\mathbf{x}_1 \in \text{dom}(r)$
- 2: **for**  $k = 1, \dots, K$  **do**
- 3:   Let  $F_k(\mathbf{x}) = F_{\mathbf{x}_k}^\gamma$  as defined in (7)
- 4:    $\mathbf{x}_{k+1} = \mathcal{A}(F_{\mathbf{x}_k}^\gamma, \Theta_k) \diamond \Theta_k$  denotes algorithm dependent parameters
- 5: **end for**

2017) require solving each convex problem up to an accuracy level of  $\epsilon$ , which is the expected accuracy level on the final solution. This change not only makes our algorithms more efficient by saving unnecessary computations but also more practical without requiring  $\epsilon$  to run the algorithm. We present a meta algorithm in Algorithm 1, in which  $\mathcal{A}$  refers to an appropriate stochastic algorithm for solving each convex majorant function. The Step 4 means that  $\mathcal{A}$  is employed for finding  $\mathbf{x}_{k+1}$  such that (8) is satisfied (or a more fine-grained condition is satisfied for a particular algorithm as discussed later), where  $\Theta_k$  denotes the algorithm dependent parameters (e.g., the number of iterations). Our convergence analysis also has its merits compared with the previous work (Nitanda & Suzuki, 2017). We will divide our convergence analysis into three parts. First, in subsection 3.1 we introduce a general convergence measure without requiring any smoothness assumptions of involved functions and conduct a convergence analysis of the proposed algorithm. Second, we analyze different stochastic algorithms and their convergence results in subsection 3.2, including an adaptive convergence result for using AdaGrad. Finally, we discuss the implications of these convergence results for solving the original problem in terms of finding a (nearly)  $\epsilon$ -stationary point in subsection 3.3.

### 3.1. A General Convergence Result

For any  $\gamma > 0$ , define  $P_\gamma(\mathbf{z}) = \arg \min_{\mathbf{x} \in \mathbb{R}^d} F_\mathbf{z}^\gamma(\mathbf{x})$ ,  $G_\gamma(\mathbf{z}) = \gamma(\mathbf{z} - P_\gamma(\mathbf{z}))$ . It is notable that  $P_\gamma(\mathbf{z})$  is well defined since  $F_\mathbf{z}^\gamma$  is strongly convex. The following proposition shows that when  $\mathbf{z} = P_\gamma(\mathbf{z})$ , then  $\mathbf{z}$  is a critical point of the original problem.

**Proposition 1.** *If  $\mathbf{z} = P_\gamma(\mathbf{z})$ , then  $\mathbf{z}$  is a critical point of the problem  $\min_{\mathbf{x} \in \mathbb{R}^d} g(\mathbf{x}) + r(\mathbf{x}) - h(\mathbf{x})$ .*

The above proposition implies that  $\|G_\gamma(\mathbf{z})\| = \gamma\|P_\gamma(\mathbf{z}) - \mathbf{z}\|$  can serve as a measure of convergence of an algorithm for solving the considered minimization problem. In subsection 3.3, we will discuss how the convergence in terms of  $\gamma\|P_\gamma(\mathbf{z}) - \mathbf{z}\|$  implies the standard convergence in terms of the (sub)gradient norm of the original problem. We use the following basic assumption for our analysis.

**Assumption 1.** *For an initial solution  $\mathbf{x}_1 \in \text{dom}(r)$ , assume that  $F(\mathbf{x}_1) - \inf_{\mathbf{x} \in \mathbb{R}^d} F(\mathbf{x}) \leq \Delta$  for some  $\Delta > 0$ .*

The theorems below are the main results of this subsection.

**Theorem 1.** *Suppose Assumption 1 holds and there exists an stochastic algorithm  $\mathcal{A}$  that when applied to  $F_{\mathbf{x}_k}^\gamma(\mathbf{x})$*

*can find a solution  $\mathbf{x}_{k+1}$  satisfying (8), then we have  $\mathbb{E}[\|G_\gamma(\mathbf{x}_\tau)\|^2] \leq (2\gamma\Delta + 2\gamma c(1 + \log(K)))/K$ , where  $\tau \in [K]$  is uniformly sampled.*

**Remark:** It is clear that when  $K \rightarrow \infty$ ,  $\gamma\|\mathbf{x}_\tau - P_\gamma(\mathbf{x}_\tau)\| \rightarrow 0$  in expectation, implying the convergence to a critical point. Note that the  $\log(K)$  factor will lead to an iteration complexity of  $O(\log(1/\epsilon)/\epsilon^4)$  for using stochastic (sub)gradient method. Nevertheless, such a logarithmic factor can be removed by exploiting non-uniform sampling under a slightly stronger condition of the problem.

**Theorem 2.** *Suppose there exists a stochastic algorithm  $\mathcal{A}$  that when applied to  $F_{\mathbf{x}_k}^\gamma(\mathbf{x})$  can find a solution  $\mathbf{x}_{k+1}$  satisfying (8), and there exists  $\Delta > 0$  such that  $\mathbb{E}[F(\mathbf{x}_k) - \min_{\mathbf{x}} F(\mathbf{x})] \leq \Delta$  for all  $k \in [K]$ , then we have  $\mathbb{E}[\|G_\gamma(\mathbf{x}_\tau)\|^2] \leq \frac{2\gamma(\Delta+c)(\alpha+1)}{K}$ , where  $\tau \in [K]$  is sampled according to probabilities  $p(\tau = k) = k^\alpha / \sum_{k=1}^K k^\alpha$  with  $\alpha \geq 1$ .*

**Remark:** Compared to Theorem 1, the condition  $\mathbb{E}[F(\mathbf{x}_k) - \min_{\mathbf{x}} F(\mathbf{x})] \leq \Delta$  for all  $k \in [K]$  is slightly stronger than Assumption 1. However, it can be easily satisfied if  $\mathbf{x}_k \in \text{dom}(r)$  resides in a bounded set (e.g., when  $r(\mathbf{x})$  is the indicator function of a bounded set), or if  $\mathbb{E}[F(\mathbf{x}_k)]$  is non-increasing (e.g., when using variance-reduction methods for the case that  $g(\mathbf{x})$  is smooth). In the following presentation, we assume this condition holds without explicitly mentioned.

### 3.2. Convergence Results of Different Algorithms

In this section, we will present the convergence results of Algorithm 1 for employing different stochastic algorithms to minimize  $F_k(\mathbf{x})$  at each stage. In particular, we consider three representative algorithms, namely stochastic proximal subgradient (SPG) method (Duchi et al., 2010; Zhao & Zhang, 2015), adaptive stochastic gradient (AdaGrad) method (Duchi et al., 2011; Chen et al., 2018a), and proximal stochastic gradient method with variance reduction (SVRG) (Xiao & Zhang, 2014). We refer to Algorithm 1 by using SPG, AdaGrad, SVRG for solving each subproblem as SSDC-SPG, SSDC-AdaGrad, SSDC-SVRG, respectively.

**SPG.** We make the additional assumptions about the problem for developing SPG, which are typical in the literature (Zhao & Zhang, 2015; Duchi et al., 2010).

**Assumption 2.** *Assume one of the following conditions: (i)  $g(\mathbf{x})$  is  $L$ -smooth and  $\mathbb{E}[\|(\nabla g(\mathbf{x}; \xi) - \partial h(\mathbf{x}; \varsigma)) - \mathbb{E}[\nabla g(\mathbf{x}; \xi) - \partial h(\mathbf{x}; \varsigma)]\|^2] \leq G^2$ . (ii)  $\mathbb{E}[\|\partial g(\mathbf{x}; \xi)\|^2] \leq G^2$ ,  $\mathbb{E}[\|\partial h(\mathbf{x}; \varsigma)\|^2] \leq G^2$  for  $\mathbf{x} \in \text{dom}(r)$ , and either  $r = \delta_{\mathcal{X}}(\mathbf{x})$  for a closed convex set  $\mathcal{X}$  or  $\|\partial r(\mathbf{x})\| \leq G$  for  $\mathbf{x} \in \text{dom}(r)$ .*

Without loss of generality, we consider minimizing  $F_{\mathbf{x}_1}^\gamma$  by

SPG. The key update of SPG is the following:

$$\mathbf{x}_{t+1} = \arg \min_{\mathbf{x} \in \Omega} \left\{ \mathbf{x}^\top \mathcal{G}(\mathbf{x}_t; \xi_t, \varsigma_t) + r(\mathbf{x}) + \frac{\gamma}{2} \|\mathbf{x} - \mathbf{x}_1\|^2 + \frac{1}{2\eta_t} \|\mathbf{x} - \mathbf{x}_t\|^2 \right\}, t = 1, \dots, T \quad (9)$$

where  $\mathcal{G}(\mathbf{x}_t; \xi_t, \varsigma_t) = \partial g(\mathbf{x}_t; \xi_t) - \partial h(\mathbf{x}_1; \varsigma_t)$ . For smooth  $g$ , we set  $\Omega = \mathbb{R}^d$ , and for non-smooth  $g$  we set  $\Omega = \mathcal{B}_{\mathbf{x}_1} = \{\mathbf{x} \in \mathbb{R}^d, \|\mathbf{x} - \mathbf{x}_1\| \leq 3G/\gamma\}$ . Restricting the solution to the ball  $\mathcal{B}_{\mathbf{x}_1}$  is to accommodate the proximal mapping of  $r(\mathbf{x})$  when  $g(\mathbf{x})$  is non-smooth. When using the subgradient of  $r(\mathbf{x})$  instead of the proximal mapping of  $r(\mathbf{x})$  in the update or  $r(\mathbf{x})$  is the indicator function of a bounded convex set, the projection onto  $\mathcal{B}_{\mathbf{x}_1}$  can be removed. The complete steps of the SPG algorithm are presented in Algorithm 3 in the supplement with the two options to handle smooth and non-smooth  $g$  separately. The convergence of Algorithm 1 when using SPG to solve each subproblem is stated below.

**Theorem 3.** *Suppose Assumption 2 (i) holds and Algorithm 3 with (9) ( $\Omega = \mathbb{R}^d$ ) is employed for solving  $F_k$  with  $\eta_t = 1/(L(t+1))$ ,  $\gamma \geq 3L$  and  $T_k = 4k/c$  iterations where  $c \in (0, 1]$ , then Algorithm 1 guarantees  $\mathbb{E}[\|G_\gamma(\mathbf{x}_\tau)\|^2] \leq \frac{8\gamma\Delta(\alpha+1)}{K} + \frac{8cG^2\gamma(\alpha+1)}{LK}$ . Similarly, Suppose Assumption 2 (ii) holds and Algorithm 3 with (9) ( $\Omega = \mathcal{B}_{\mathbf{x}_{k-1}}$ ) is employed for solving  $F_k$  with  $\eta_t = 4/(\gamma t)$ , and  $T_k = k/c$  iterations with  $c \in (0, 1]$ , then  $\mathbb{E}[\|G_\gamma(\mathbf{x}_\tau)\|^2] \leq \frac{8\gamma\Delta(\alpha+1)}{K} + \frac{448c(\alpha+1)}{K}$ , where  $\tau$  is sampled similarly as in Theorem 2.*

**Remark:** Let us consider the iteration complexity of using SPG for finding a solution that satisfies  $\mathbb{E}[\|G_\gamma(\mathbf{x}_\tau)\|^2] \leq \epsilon^2$ . For the non-smooth case, by setting  $\gamma < 1$  and  $c = \gamma$ , we need a total number of stages  $K = O(\gamma/\epsilon^2)$  and total iteration complexity  $\sum_{k=1}^K T_k = \sum_{k=1}^K k/\gamma = O(1/\epsilon^4)$ . For the smooth case, by setting  $c = 1$  we have  $K = O(\max(L, 1)/\epsilon^2)$  and total iteration complexity  $\sum_{k=1}^K T_k = \sum_{k=1}^K 4k = O(\max(L, 1)/\epsilon^4)$ .

**AdaGrad.** AdaGrad (Duchi et al., 2011) is an important algorithm in the literature of stochastic optimization, which uses adaptive step size for each coordinate. It has potential benefit of speeding up the convergence when the cumulative growth of stochastic gradient is slow. Next, we show that AdaGrad can be leveraged to solve each convex majorant function and yield adaptive convergence for the original problem. Similar to (Duchi et al., 2011; Chen et al., 2018a), we make the following assumption.

**Assumption 3.** *For any  $\mathbf{x} \in \text{dom}(r)$ , there exists  $G > 0$  such that  $\max(\|\partial g(\mathbf{x}; \xi)\|_\infty, \|\partial h(\mathbf{x}; \varsigma)\|_\infty) \leq G$ , either  $r = \delta_{\mathcal{X}}(\mathbf{x})$  for a closed convex set  $\mathcal{X}$  or  $\|\partial r(\mathbf{x})\| \leq G_r$ .*

The convergence result of Algorithm 1 by using AdaGrad to solve each problem is described by following theorem.

**Theorem 4.** *Suppose Assumption 3 holds and Algorithm 2 is employed for solving  $F_k$  with  $\eta_k = c/\sqrt{k}$ ,  $T_k$  being the minimum number that is larger than  $M_k \max\{a(2G +$*

**Algorithm 2** ADAGRAD( $F_{\mathbf{x}_1}^\gamma, \mathbf{x}_1, \eta$ )

- 1: **Initialize:**  $t = 1, \mathbf{g}_{1:0} = \emptyset, H_0 \in \mathbb{R}^{d \times d}, \Omega = \{\mathbf{x} \in \text{dom}(r) : \|\mathbf{x} - \mathbf{x}_1\| \leq \frac{2\sqrt{d}G + G_r}{\gamma}\}$
- 2: **while**  $T$  doesn't satisfy the condition in Thm. 4 **do**
- 3:   Compute  $\mathbf{g}_t = \partial g(\mathbf{x}_t; \xi_t) - \partial h(\mathbf{x}_1; \varsigma_t)$
- 4:   Update  $g_{1:t} = [g_{1:t-1}, \mathbf{g}_t], s_{t,i} = \|g_{1:t,i}\|_2$
- 5:   Set  $H_t = H_0 + \text{diag}(s_t)$
- 6:   Let  $\mathbf{x}_{t+1} = \arg \min_{\mathbf{x} \in \Omega} \mathbf{x}^\top \left( \frac{1}{t} \sum_{\tau=1}^t \mathbf{g}_\tau \right) + r(\mathbf{x}) + \frac{\gamma}{2} \|\mathbf{x} - \mathbf{x}_1\|^2 + \frac{1}{t\eta} \frac{1}{2} (\mathbf{x} - \mathbf{x}_1)^\top H_t (\mathbf{x} - \mathbf{x}_1)$
- 7: **end while**
- 8: **Output:**  $\hat{\mathbf{x}}_T = \sum_{t=1}^T \mathbf{x}_t / T$

$\max_i \|g_{1:T_k,i}^k\|, \sum_{i=1}^d \|g_{1:T_k,i}^k\|/a, \frac{G_r}{\eta_k} \|\mathbf{x}_1^k - \mathbf{x}_{T_k+1}^k\| \}$  where  $M_k \eta_k \geq 4/(a\gamma)$ , then Algorithm 1 guarantees  $\mathbb{E}[\|G_\gamma(\mathbf{x}_\tau)\|^2] \leq \frac{8\gamma\Delta(\alpha+1)}{K} + \frac{4\gamma^2 c^2 a(\alpha+1)(\alpha+1)}{K}$ , where  $g_{1:t,i}^k$  denotes the cumulative stochastic gradient of the  $i$ -th coordinate at the  $k$ -th stage, and  $\tau$  is sampled similarly as in Theorem 2.

**Remark:**  $a$  is a parameter used to balance the two involved terms for minimizing the value of  $T_k$ . It is obvious that the total number of iterations  $\sum_{k=1}^K T_k$  is adaptive to the data. Next, let us present more discussion on the iteration complexity. Note that  $M_k = O(\sqrt{k})$ . By the boundness of stochastic gradient  $\|g_{1:T_k,i}^k\| \leq O(\sqrt{T_k})$ , therefore  $T_k$  in the order of  $O(k)$  will satisfy the condition in Theorem 4. Thus in the worst case, the iteration complexity for finding  $\mathbb{E}[\|G_\gamma(\mathbf{x}_\tau)\|^2] \leq \epsilon^2$  is in the order of  $\sum_{k=1}^K O(k) \leq O(1/\epsilon^4)$ . We can show the potential advantage of adaptiveness similar to that in (Chen et al., 2018b). In particular, let us consider  $r = \delta_{\mathcal{X}}$  and thus  $G_r = 0$  in the above result. When the cumulative growth of stochastic gradient is slow, e.g., assuming  $\|g_{1:T_k,i}^k\| \leq O(T_k^\beta)$  with  $\beta < 1/2$ . Then  $T_k = O(k^{1/(2(1-\beta))})$  will work, and then the total number of iterations  $\sum_{k=1}^K T_k \leq K^{1+1/(2(1-\alpha))} \leq O(1/\epsilon^{2+1/(1-\alpha)})$ , which is better than  $O(1/\epsilon^4)$ .

**SVRG.** Next, we discuss SVRG (a variance reduction method) for solving each subproblem when it has a finite-sum form (6) and  $g$  is a smooth function. It is notable that the smoothness of  $h$  is not necessary for developing the SVRG algorithm since at each stage we linearize  $h(\mathbf{x})$ . We can use the proximal SVRG proposed in (Xiao & Zhang, 2014) to minimize  $F_{\mathbf{x}_{k-1}}^\gamma$ , which is presented in Algorithm 4 in the supplement, whose convergence result is stated below.

**Theorem 5.** *Suppose Assumption 1 and  $g$  is smooth, and SVRG (Algorithm 4) is employed for solving  $F_k$  with  $\eta_k = 0.05/L$ ,  $T_k \geq \max(2, 200L/\gamma)$ ,  $S_k = \lceil \log_2(k) \rceil$ , then Algorithm 1 guarantees  $\mathbb{E}[\|G_\gamma(\mathbf{x}_\tau)\|^2] \leq 12\gamma\Delta(\alpha+1)/K$ , where  $\tau$  is sampled similarly as in Theorem 2.*

**Remark:** For finding a solution such that  $\mathbb{E}[\|G_\gamma(\mathbf{x}_\tau)\|^2] \leq \epsilon^2$ , the total number of stages  $K = O(\gamma/\epsilon^2)$  and the total gradient complexity is  $\tilde{O}((n\gamma + L)/\epsilon^2)$ .

### 3.3. Finding a (nearly) $\epsilon$ -critical point

We summarize below the convergence results of the proposed algorithms for finding a (nearly)  $\epsilon$ -critical point.

**Theorem 6.** *Assume Algorithm 1 returns a solution  $\mathbf{x}_\tau$  such that  $\mathbb{E}[\|G_\gamma(\mathbf{x}_\tau)\|^2] \leq O(1/K)$  under appropriate conditions. Then if  $g(\mathbf{x}) + r(\mathbf{x})$  is differentiable and has  $(L, \nu)$ -Hölder continuous gradient, we have  $\mathbb{E}[\text{dist}(\partial h(\mathbf{x}_\tau), \nabla(g(\mathbf{x}_\tau) + r(\mathbf{x}_\tau)))] \leq O\left(\frac{1}{K^{\nu/2}} + \frac{1}{\sqrt{K}}\right)$ . If  $h(\mathbf{x})$  is differentiable and has  $(L, \nu)$ -Hölder continuous gradient, we have  $\mathbb{E}[\|\mathbf{x}_\tau - \mathbf{z}_\tau\|] \leq O(1/\sqrt{K})$  and  $\mathbb{E}[\text{dist}(\nabla h(\mathbf{z}_\tau), \partial(g(\mathbf{z}_\tau) + r(\mathbf{z}_\tau)))] \leq O\left(\frac{1}{K^{\nu/2}} + \frac{1}{\sqrt{K}}\right)$ , where  $\mathbf{z}_\tau = P_\gamma(\mathbf{x}_\tau)$ .*

**Remark:** Note that the convergence of the proposed algorithms can be automatically adaptive to the Hölder continuous of the involved functions without requiring the value of  $\nu$  for running the algorithm. Both SSDC-SPG and SSDC-AdaGrad have an iteration complexity (in the worst-case) of  $O(1/\epsilon^{4/\nu})$  for finding a (nearly)  $\epsilon$ -critical point. When the problem has a finite-sum structure (6) and  $g(\mathbf{x})$  is smooth, SSDC-SVRG has a gradient complexity of  $O(n/\epsilon^{2/\nu})$  for finding a (nearly)  $\epsilon$ -critical point.

## 4. Non-Smooth Non-Convex Regularization

In this section, we consider a more challenging class of problem (1) where  $r(\mathbf{x})$  is a proper non-smooth and non-convex lower-semicontinuous function that is not necessarily a DC function (e.g.,  $\ell_0$  norm). Even if  $r(\mathbf{x})$  is a DC function such that both components in its DC decomposition  $r(\mathbf{x}) = r_1(\mathbf{x}) - r_2(\mathbf{x})$  are non-differentiable functions without Hölder continuous gradients (e.g.,  $\ell_{1-2}$  regularization, capped  $\ell_1$  norm), the theories presented in this section are useful to derive non-asymptotic convergence results in terms of finding an  $\epsilon$ -critical point. Please note that in this case the results presented in section 3.3 are not applicable. Similarly, we assume  $r(\mathbf{x})$  is simple such that its proximal mapping exists and can be efficiently computed.

The problem is challenging due to the presence of non-smooth non-convex function  $r$ . To tackle this function, we introduce the Moreau envelope of  $r$ :

$$r_\mu(\mathbf{x}) = \min_{\mathbf{y} \in \mathbb{R}^d} \frac{1}{2\mu} \|\mathbf{y} - \mathbf{x}\|^2 + r(\mathbf{y}),$$

where  $\mu > 0$ . A nice property of the Moreau envelope function is that it can be written as a DC function:

$$r_\mu(\mathbf{x}) = \frac{1}{2\mu} \|\mathbf{x}\|^2 - \underbrace{\max_{\mathbf{y} \in \mathbb{R}^d} \frac{1}{\mu} \mathbf{y}^\top \mathbf{x} - \frac{1}{2\mu} \|\mathbf{y}\|^2 - r(\mathbf{y})}_{R_\mu(\mathbf{x})},$$

where  $R_\mu(\mathbf{x})$  is a convex function because it is the max of convex functions of  $\mathbf{x}$  (Boyd & Vandenberghe, 2004). The following properties about the Moreau envelope will be useful for our analysis.

**Lemma 1.**  $\frac{1}{\mu} \text{prox}_{\mu r}(\mathbf{x}) \subseteq \partial R_\mu(\mathbf{x})$ , and  $\frac{1}{\mu}(\mathbf{x} - \mathbf{v}) \subseteq \hat{\partial} r(\mathbf{v}), \forall \mathbf{v} \in \text{prox}_{\mu r}(\mathbf{x})$ .

Given the Moreau envelope of  $r$ , the key idea is to solve the following DC problem:

$$\min_{\mathbf{x} \in \mathbb{R}^d} g(\mathbf{x}) - h(\mathbf{x}) + \frac{1}{2\mu} \|\mathbf{x}\|^2 - R_\mu(\mathbf{x}). \quad (10)$$

By carefully controlling the value of  $\mu$  and combining the results presented in previous section, we are able to derive non-asymptotic convergence results for the original problem. It is worth mentioning that using the Moreau envelope of  $r$  and its DC decomposition for handling non-smooth non-convex function is first proposed in (Liu et al., 2018). However, their algorithms are deterministic and convergence results are only asymptotic. To formally state our non-asymptotic convergence results, we make the following assumptions.

**Assumption 4.** *Assume  $g$  and  $h$  are smooth, and one of the following conditions holds: (i)  $r$  is Lipschitz continuous; (ii)  $r$  is lower bounded and finite-valued over  $\mathbb{R}^d$ ; (iii)  $g(\mathbf{x}) - h(\mathbf{x}) + r_\mu(\mathbf{x})$  is level bounded for a small  $\mu < 1$ , and  $r$  is finite-valued on a compact set, and lower bounded over  $\mathbb{R}^d$ .*

**Remark:** The above assumptions on  $r$  capture many interesting non-convex non-smooth regularizers. For example,  $\ell_{1-2}$  regularization and capped  $\ell_1$  norm satisfy Assumption 4 (i). The  $\ell_0$  norm satisfies Assumption 4 (ii). A coercive function  $r$  usually satisfies Assumption 4 (iii), e.g.,  $\ell_p$  norm  $r(\mathbf{x}) = \sum_{i=1}^d |x_i|^p$  for  $p \in (0, 1)$ . In Appendix K, we further extend our results to handle a differentiable  $h$  that has only a Hölder-continuous gradient.

When employing the presented algorithms in last section to solve the problem (10), we let  $r = \frac{1}{2\mu} \|\mathbf{x}\|^2$ ,  $g \leftarrow g$  and  $h \leftarrow h + R_\mu$ . It is also notable that the new component  $R_\mu(\mathbf{x})$  is deterministic, whose subgradient can be computed according to Lemma 1. Thus the condition in Assumption 2 (i) is sufficient for running SPG, and the smoothness condition of  $g$  is sufficient for running SVRG. Now we are ready to present our results for solving the problem (1) with a non-smooth and non-convex  $r$ .

**Theorem 7.** *Suppose SSDC is employed for solving (10) and returns  $\mathbf{x}_\tau$ . Let  $\mathbf{w}_\tau = \text{prox}_{\mu r}(\mathbf{x}_\tau)$  be the final output, we have the following results to ensure  $\mathbb{E}[\text{dist}(\nabla h(\mathbf{w}_\tau), \nabla g(\mathbf{w}_\tau) + \hat{\partial} r(\mathbf{w}_\tau))] \leq \epsilon$ . (a) If Assumption 4 (i) and Assumption 2 (i) hold, then we can set  $\mu = \epsilon$ , use SPG for solving the subproblems, and have a total gradient complexity of  $O(1/\epsilon^8)$ . (b) If Assumption 4 (ii) and Assumption 2 (i) hold, then we can set  $\mu = \epsilon^2$ , use SPG for solving the subproblems, and have a total gradient complexity of  $O(1/\epsilon^{12})$ . (c) If  $g$  and  $h$  have a finite-sum form and are smooth, then we can use SVRG for solving the subproblems. Under assumption 4 (i), we can set  $\mu = \epsilon$  and have a total gradient complexity of  $O(n/\epsilon^4)$ . Under assumption 4 (ii) or (iii), we can set  $\mu = \epsilon^2$  and have a total gradient complexity of  $O(n/\epsilon^6)$ .*

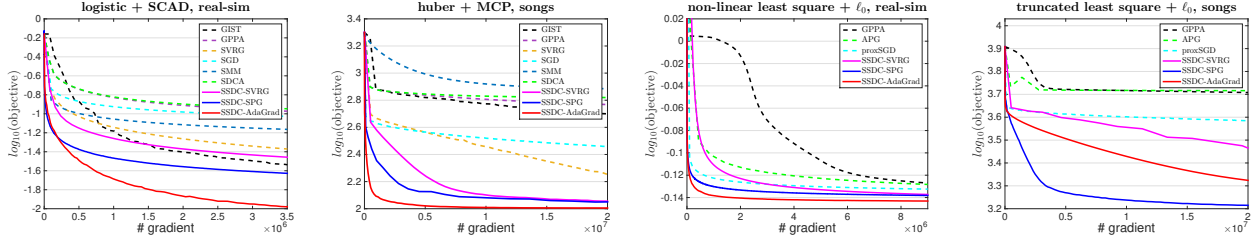


Figure 1. Learning with DC (left two) and non-DC regularizers (right two) on different datasets for classification and regression.

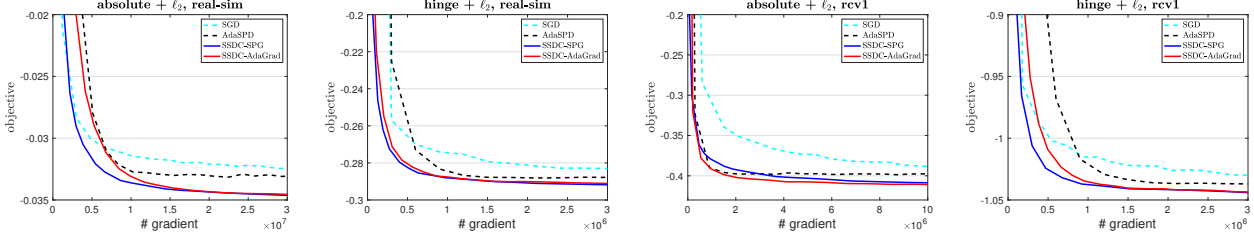


Figure 2. PU learning with different loss functions on different datasets.

## 5. Numerical Experiments

In this section, we perform some experiments for solving different tasks to demonstrate effectiveness of proposed algorithms by comparing with different baselines. We use very large-scale datasets from libsvm website in experiments, including real-sim ( $n = 72309$ ) and rcv1 ( $n=20242$ ) for classification, million songs ( $n = 463715$ ) for regression. For all algorithms, the initial stepsizes are tuned in the range of  $\{10^{-6}:1:4\}$ , and the same initial solution with all zero entries is used. The initial iteration number  $T_0$  of SSDC-SPG is tuned in  $\{10^{1:1:4}\}$ .

First, we compare SSDC algorithms with SDCA (Thi et al., 2017), SMM (Mairal, 2013), SGD (Davis & Drusvyatskiy, 2018b), SVRG (Reddi et al., 2016c), GIST (Gong et al., 2013) and GPPA (An & Nam, 2017) for learning with a DC regularizer: minimizing logistic loss with a SCAD regularizer for classification and huber loss with a MCP regularizer for regression. The parameter in Huber loss is set to be 1. The value of regularization parameter is set to be  $10^{-4}$ . We used the form of weight in SMM following (Mairal, 2013). Since these regularizers are weakly convex, SGD with step size  $\eta_0/\sqrt{t}$  is applicable (Davis & Drusvyatskiy, 2018b). We set the inner iteration number of SVRG as  $n$  following (Reddi et al., 2016c) and the same value is used as the inner iteration number  $T$  of SSDC-SVRG. We set the values of parameters in GIST with their suggested BB rule (Gong et al., 2013). Similar to (Thi et al., 2017), we tune the batch size of SDCA in a wide range and choose the one with the best performance. GIST and GPPA are deterministic algorithms that use all data points in each iteration. For fairness of comparison, we plot the objective in log scale versus the number of gradient computations in Figure 1 (left two).

Second, we consider minimizing  $\ell_0$  regularized non-linear least square loss function  $\frac{1}{n} \sum_{i=1}^n (y_i - \sigma(\mathbf{w}^\top \mathbf{x}_i))^2 + \lambda \|\mathbf{w}\|_0$  with a sigmod function  $\sigma(s) = \frac{1}{1+e^{-s}}$  for classifica-

tion and  $\ell_0$  regularized truncated least square loss function  $\frac{1}{2n} \sum_{i=1}^n \alpha \log(1 + (y_i - \mathbf{w}^\top \mathbf{x}_i)^2/\alpha) + \lambda \|\mathbf{w}\|_0$  (Xu et al., 2018) for regression. We compare the proposed algorithms with GPPA, APG (Li & Lin, 2015) and proximal version of SGD (proxSGD), where GPPA and APG are deterministic algorithms. We fix the truncation value as  $\alpha = \sqrt{10n}$ . The loss function in these two tasks are smooth and non-convex. The value of regularization parameter is fixed as  $10^{-6}$ . For APG, we implement both monotone and non-monotone versions following (Li & Lin, 2015), and then the better one is reported. Although the convergence guarantee of proxSGD remains unclear for the considered problems, we still include it for comparison. The results on two data sets are plotted in Figure 1 (right two).

The results of these two experiments indicate that the proposed stochastic algorithms outperform all deterministic baselines (GITS, GPPA, APG) on all tasks, which verify the necessity of using stochastic algorithms on large datasets. In addition, our algorithms especially SSDC-AdaGrad and SSDC-SPG also converge faster than stochastic algorithms SGD, SDCA, and non-convex SVRG verifying that our stochastic algorithms are more practical for the considered problems. We also see that in most cases SSDC-AdaGrad is more effective than SSDC-SPG and SSDC-SVRG.

Finally, we compare SSDC algorithms with two baselines AdaSPD (Nitanda & Suzuki, 2017) and SGD (Davis et al., 2018) for solving two  $\ell_2$  regularized positive-unlabeled (PU) learning problems (Du Plessis et al., 2015) with non-smooth losses, i.e., hinge loss and absolute loss. The  $\ell_2$  regularization parameter is set to be  $10^{-4}$ . For SGD, we use the standard stepsize  $\eta = \eta_0/\sqrt{t}$  (Ghadimi & Lan, 2013) with  $\eta_0$  tuned. The mini-batch size and the number of iterations of each stage of AdaSPD are simply set as  $10^4$ . The results on two classification datasets are plotted in Figure 2, which show that SSDC-SPG and SSDC-AdaGrad outperforms SGD and AdaSPD.



## Acknowledgements

The authors thank the anonymous reviewers for their helpful comments. Y. Xu and T. Yang are partially supported by National Science Foundation (IIS-1545995).

## References

- Allen-Zhu, Z. Natasha: Faster non-convex stochastic optimization via strongly non-convex parameter. In *International Conference on Machine Learning*, pp. 89–97, 2017.
- Allen-Zhu, Z. and Hazan, E. Variance reduction for faster non-convex optimization. In *International Conference on Machine Learning*, pp. 699–707, 2016.
- An, N. T. and Nam, N. M. Convergence analysis of a proximal point algorithm for minimizing differences of functions. *Optimization*, 66(1):129–147, 2017.
- Attouch, H., Bolte, J., and Svaiter, B. F. Convergence of descent methods for semi-algebraic and tame problems: proximal algorithms, forward–backward splitting, and regularized gauss–seidel methods. *Mathematical Programming*, 137(1):91–129, Feb 2013.
- Bolte, J., Sabach, S., and Teboulle, M. Proximal alternating linearized minimization for nonconvex and nonsmooth problems. *Mathematical Programming*, 146(1-2), August 2014.
- Bot, R. I., Csetnek, E. R., and László, S. C. An inertial forward–backward algorithm for the minimization of the sum of two nonconvex functions. *EURO Journal on Computational Optimization*, 4(1):3–25, Feb 2016.
- Boyd, S. and Vandenberghe, L. *Convex Optimization*. Cambridge University Press, 2004.
- Candès, E. J., Wakin, M. B., and Boyd, S. P. Enhancing sparsity by reweighted  $\ell_1$  minimization. *Journal of Fourier Analysis and Applications*, 14(5):877–905, Dec 2008.
- Chen, Z. and Yang, T. A variance reduction method for non-convex optimization with improved convergence under large condition number. *CoRR*, abs/1809.06754, 2018.
- Chen, Z., Xu, Y., Chen, E., and Yang, T. Sadagrad: Strongly adaptive stochastic gradient methods. In *International Conference on Machine Learning*, pp. 912–920, 2018a.
- Chen, Z., Yang, T., Yi, J., Zhou, B., and Chen, E. Universal stagewise learning for non-convex problems with convergence on averaged solutions. *CoRR*, abs/1808.06296, 2018b.
- Davis, D. and Drusvyatskiy, D. Stochastic model-based minimization of weakly convex functions. *CoRR*, abs/1803.06523, 2018a.
- Davis, D. and Drusvyatskiy, D. Stochastic subgradient method converges at the rate  $O(k^{-1/4})$  on weakly convex functions. *CoRR*, abs/1802.02988, 2018b.
- Davis, D. and Grimmer, B. Proximally guided stochastic subgradient method for nonsmooth, nonconvex problems. *arXiv preprint arXiv:1707.03505*, 2017.
- Davis, D., Drusvyatskiy, D., Kakade, S., and Lee, J. D. Stochastic subgradient method converges on tame functions. *CoRR*, abs/1804.07795, 2018.
- Drusvyatskiy, D. and Paquette, C. Efficiency of minimizing compositions of convex functions and smooth maps. *Mathematical Programming*, Jul 2018.
- Du Plessis, M., Niu, G., and Sugiyama, M. Convex formulation for learning from positive and unlabeled data. In *International Conference on Machine Learning*, pp. 1386–1394, 2015.
- Duchi, J., Hazan, E., and Singer, Y. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(Jul):2121–2159, 2011.
- Duchi, J. C., Shalev-Shwartz, S., Singer, Y., and Tewari, A. Composite objective mirror descent. In *Conference on Learning Theory*, pp. 14–26, 2010.
- Fan, J. and Li, R. Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association*, 96(456):1348–1360, 2001.
- Ghadimi, S. and Lan, G. Stochastic first- and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.
- Gong, P., Zhang, C., Lu, Z., Huang, J., and Ye, J. A general iterative shrinkage and thresholding algorithm for non-convex regularized optimization problems. In *International Conference on Machine Learning*, pp. 37–45, 2013.
- Hiriart-Urruty, J.-B. Generalized differentiability / duality and optimization for problems dealing with differences of convex functions. In Ponstein, J. (ed.), *Convexity and Duality in Optimization*, pp. 37–70, Berlin, Heidelberg, 1985. Springer Berlin Heidelberg.
- Horst, R. and Thoai, N. V. Dc programming: Overview. *Journal of Optimization Theory and Applications*, 103(1): 1–43, Oct 1999.
- Khamaru, K. and Wainwright, M. Convergence guarantees for a class of non-convex and non-smooth optimization problems. In *International Conference on Machine Learning*, pp. 2601–2610, 2018.

- Kiryo, R., Niu, G., du Plessis, M. C., and Sugiyama, M. Positive-unlabeled learning with non-negative risk estimator. In *Advances in Neural Information Processing Systems 30*, pp. 1675–1685, 2017.
- Lan, G. and Yang, Y. Accelerated stochastic algorithms for nonconvex finite-sum and multi-block optimization. *CoRR*, abs/1805.05411, 2018.
- Li, G. and Pong, T. K. Douglas-rachford splitting for non-convex optimization with application to nonconvex feasibility problems. *Mathematical Programming*, 159(1-2): 371–401, 2016.
- Li, H. and Lin, Z. Accelerated proximal gradient methods for nonconvex programming. In *Neural Information Processing Systems*, pp. 379–387, 2015.
- Liu, T., Pong, T. K., and Takeda, A. A successive difference-of-convex approximation method for a class of nonconvex nonsmooth optimization problems. *Mathematical Programming*, Sep 2018.
- Mairal, J. Stochastic majorization-minimization algorithms for large-scale optimization. In *Advances in Neural Information Processing Systems*, pp. 2283–2291, 2013.
- Nitanda, A. and Suzuki, T. Stochastic Difference of Convex Algorithm and its Application to Training Deep Boltzmann Machines. In *International Conference on Artificial Intelligence and Statistics*, pp. 470–478, 2017.
- Reddi, S. J., Hefny, A., Sra, S., Póczós, B., and Smola, A. Stochastic variance reduction for nonconvex optimization. In *International Conference on Machine Learning*, pp. 314–323, 2016a.
- Reddi, S. J., Sra, S., Póczos, B., and Smola, A. J. Fast incremental method for smooth nonconvex optimization. In *55th IEEE Conference on Decision and Control*, pp. 1971–1977, 2016b.
- Reddi, S. J., Sra, S., Póczos, B., and Smola, A. J. Proximal stochastic methods for nonsmooth nonconvex finite-sum optimization. In *Advances in Neural Information Processing Systems*, pp. 1145–1153, 2016c.
- Rockafellar, R. and Wets, R. J.-B. *Variational Analysis*. Springer Verlag, Heidelberg, Berlin, New York, 1998.
- Thi, H. A. L. and Dinh, T. P. Dc programming and dca: thirty years of developments. *Mathematical Programming*, 169(1):5–68, May 2018.
- Thi, H. A. L., Le, H. M., Phan, D. N., and Tran, B. Stochastic dca for the large-sum of non-convex functions problem and its application to group variable selection in classification. In *International Conference on Machine Learning*, pp. 3394–3403, 2017.
- Wen, B., Chen, X., and Pong, T. K. A proximal difference-of-convex algorithm with extrapolation. *Computational Optimization and Applications*, 69(2):297–324, Mar 2018.
- Xiao, L. and Zhang, T. A proximal stochastic gradient method with progressive variance reduction. *SIAM Journal on Optimization*, 24(4):2057–2075, 2014.
- Xu, Y., Zhu, S., Yang, S., Zhang, C., Jin, R., and Yang, T. Learning with non-convex truncated losses by SGD. *CoRR*, abs/1805.07880, 2018.
- Yang, L. Proximal gradient method with extrapolation and line search for a class of nonconvex and nonsmooth problems. *CoRR*, abs/1711.06831, 2018.
- Yin, P., Lou, Y., He, Q., and Xin, J. Minimization of 1-2 for compressed sensing. *SIAM J. Scientific Computing*, 37, 2015.
- Yu, Y., Zheng, X., Marchetti-Bowick, M., and Xing, E. P. Minimizing nonconvex non-separable functions. In *International Conference on Artificial Intelligence and Statistics*, 2015.
- Zhang, C.-H. Nearly unbiased variable selection under minimax concave penalty. *The Annals of Statistics*, 38: 894 – 942, 2010a.
- Zhang, S. and He, N. On the convergence rate of stochastic mirror descent for nonsmooth nonconvex optimization. *arXiv preprint arXiv:1806.04781*, 2018.
- Zhang, S. and Xin, J. Minimization of transformed  $l_1$  penalty: theory, difference of convex function algorithm, and robust application in compressed sensing. *Mathematical Programming*, 169(1):307–336, 2018.
- Zhang, T. Analysis of multi-stage convex relaxation for sparse regularization. *J. Mach. Learn. Res.*, 11:1081–1107, March 2010b. ISSN 1532-4435.
- Zhao, P. and Zhang, T. Stochastic optimization with importance sampling for regularized loss minimization. In *International Conference on Machine Learning*, pp. 1–9, 2015.
- Zhong, W. and Kwok, J. T. Gradient descent with proximal average for nonconvex and composite regularization. In *AAAI Conference on Artificial Intelligence*, pp. 2206–2212, 2014.