

Stochastic proximal splitting algorithm for composite minimization

Andrei Patrascu · Paul Irofti

Received: date / Accepted: date

Abstract Supported by the recent contributions in multiple branches, the first-order splitting algorithms became central for structured nonsmooth optimization. In the large-scale or noisy contexts, when only stochastic information on the objective function is available, the extension of proximal gradient schemes to stochastic oracles is heavily based on the tractability of the proximal operator corresponding to nonsmooth component, which has been deeply analyzed in the literature. However, there remained some questions about the difficulty of the composite models where the nonsmooth term is not proximally tractable anymore. Therefore, in this paper we tackle composite optimization problems, where the access only to stochastic information on both smooth and nonsmooth components is assumed, using a stochastic proximal first-order scheme with stochastic proximal updates. We provide sublinear $\mathcal{O}(\frac{1}{k})$ convergence rates (in expectation of squared distance to the optimal set) under the strong convexity assumption on the objective function. Also, linear convergence is achieved for convex feasibility problems. The empirical behavior is illustrated by numerical tests on parametric sparse representation models.

1 Introduction

In this paper we consider the following convex composite optimization problem:

$$\min_{x \in \mathbb{R}^n} f(x) + h(x),$$

A. Patrascu and P. Irofti are with the Research Center for Logic, Optimization and Security (LOS), Department of Computer Science, Faculty of Mathematics and Computer Science, University of Bucharest. (e-mails: andrei.patrascu@fmi.unibuc.ro, paul@irofti.net). The research of A. Patrascu was supported by a grant of the Romanian Ministry of Education and Research, CNCS - UEFISCDI, project number PN-III-P1-1.1-PD-2019-1123, within PNCDI III. Also, the research work of P. Irofti was supported by a grant of the Romanian Ministry of Education and Research, CNCS - UEFISCDI, project number PN-III-P1-1.1-PD-2019-0825, within PNCDI III.

where f is the smooth component and h is a proper, convex, lower-semicontinuous function. In literature, many applications from statistics [22] or signal processing often motivate noisy contexts allowing access only to stochastic first order information on smooth function f , having regularizer h as a typical proximally-tractable convex function. By proximally-tractable we mean that the proximal map of a given function is computable in closed form or, at most, in linear time. Therefore, in these situations the following stochastic model

$$\min_{x \in \mathbb{R}^n} \mathbb{E}[f(x; \xi)] + h(x),$$

where ξ is a random variable, express better the previous real assumptions. However, the recent dimensionality inflation of machine learning [8,26,33] and signal processing [26,27] models gave birth to optimization problems with complicated regularizers or complicated many constraints. As practical examples we recall: parametric sparse representation [27], group lasso [8,26,33], CUR-like factorization [31], graph trend filtering [24,29], dictionary learning [27,32]. Motivated by all these models, in this paper instead of assuming proximal-tractable regularizer h , we consider that h is expressed as an expectation of stochastic proximally-tractable components $h(\cdot; \xi)$ (i.e. $h(x) = \mathbb{E}[h(x; \xi)]$). Thus we focus on the following model:

$$\min_{x \in \mathbb{R}^n} F(x) := \mathbb{E}_{\xi \in \Omega}[f(x; \xi)] + \mathbb{E}_{\xi \in \Omega}[h(x; \xi)], \quad (1)$$

where ξ is a random variable associated with probability space (\mathbb{P}, Ω) . Functions $f(\cdot; \xi) : \mathbb{R}^n \rightarrow \mathbb{R}$ are smooth with Lipschitz gradients and $h(\cdot; \xi) : \mathbb{R}^n \mapsto (-\infty, +\infty]$ are proper convex and lower-semicontinuous, $\mathbb{E}[\cdot]$ is the expectation over respective random variable. In general, many existing primal schemes encounter computational difficulties when a large (possibly infinite) number of constraints are present, since they are based on full projections onto complicated feasible set.

Contributions. (i) We analyze a stochastic first-order splitting scheme relying on stochastic gradients and stochastic proximal updates, which naturally generalize the widely known Stochastic Gradient Descent (SGD) and Stochastic Proximal Point (SPP) algorithms toward composite models with untractable regularizations; (ii) we provide $\mathcal{O}(\frac{1}{k})$ iteration complexity estimates which were previously unknown for this type of schemes. In particular, for convex feasibility problems, the analysis yields naturally linear convergence rates.

We briefly recall further the milestone results from stochastic optimization literature with focus on the complexity of stochastic first-order methods.

1.1 Previous work

Great attention has been given in the last decade to the behaviour of stochastic first order schemes, with special focus in stochastic gradient descent (SGD), on a variety of models under different convexity properties, see [10–13,16,22,25]. Since the analysis of SGD naturally require a typical smoothness, appropriate extensions are necessary to attack nonsmooth models. These extensions are embodied by the stochastic

proximal point (SPP) algorithm, which has been recently analyzed using various differentiability assumptions, see [1,5,9,18,23,28,30] and has shown surprising analytical and empirical performances. In [28] is considered the typical stochastic learning model involving the expectation of random particular components $f(\cdot; \xi)$ defined by the composition of a smooth function and a linear operator, i.e.: $f(x; \xi) = \ell(a_\xi^T x)$, where $a_\xi \in \mathbb{R}^n$. Their complexity analysis requires smoothness and strong convexity to obtain in the quadratic mean and an $\mathcal{O}\left(\frac{1}{k^\gamma}\right)$ convergence rate, using vanishing step-size. The generalization of these convergence guarantees is undertaken in [18], where no linear composition structure is required and an (in)finite number of constraints are included in the stochastic model, i.e. $h(\cdot; \xi) = \mathbb{I}_{X_\xi}(\cdot)$. However, the analysis of [18] requires strong convexity and Lipschitz gradient continuity for each functional component $f(\cdot; \xi)$. Note that our analysis surpasses this restriction and provides a natural generalization of [18] to nonsmooth composite models. Further, in [5] a general asymptotic convergence analysis of slightly modified SPP scheme has been provided, under mild convexity assumptions on a finitely constrained stochastic problem. In [30] similar SPP algorithms are developed, which are also tailored for complicated feasible sets. The authors focus on mild convex models deriving optimal $\mathcal{O}(1/\sqrt{k})$ rates. Recently, in [1], the authors analyze SPP schemes for stochastic models with "shared minimizers" obtaining linear convergence results, for variable stepsize SPP. Also, without shared minimizers assumption, they obtain for SPP $\mathcal{O}\left(\frac{1}{\sqrt{k}}\right)$ in convex Lipschitz continuous case and, furthermore, $\mathcal{O}\left(\frac{1}{k}\right)$ in strongly convex case.

Splitting first-order schemes received significantly attention due to their natural insight and simplicity in contexts when a sum of two components are minimized (see [4,14]). However, when the information access is limited only to stochastic samples of the two components, extending the existing guarantees is not straightforward. Notice that overall, on one hand, SPP avoids any splitting in composite models by treating the constrained smooth problems as black box expectations (see [1,18,28]). On the other hand, until recently the composite nonsmooth models assumed proximal-tractable h in order to extend proof arguments of the results related to stochastic smooth optimization [22]. Therefore, only recently the full stochastic composite models with stochastic regularizers have been properly tackled [24], where almost sure asymptotic convergence is established for a stochastic splitting scheme, where each iteration represents a typical proximal gradient update with respect to stochastic samples of f and h . In our paper we analyze the nonasymptotic behaviour of this scheme and point the relations with other algorithms from the literature.

The stochastic splitting schemes are also related to the model-based methods developed in [6]. Here, the authors developed a unified algorithmic framework, which generates stochastic algorithms, for different models arising in learning applications. They assume their composite objective function to be the sum of a (weakly convex) stochastic component with bounded gradients and simple (proximally tractable) convex regularization. Although their framework is algorithmically more general, our analysis avoid these boundedness assumptions, allows objectives with a component having Lipschitz continuous gradient and do not require proximal tractability on regularizations.

Notations. We use notation $[m] = \{1, \dots, m\}$. For $x, y \in \mathbb{R}^n$ denote the scalar product $\langle x, y \rangle = x^T y$ and Euclidean norm by $\|x\| = \sqrt{x^T x}$. The projection operator onto set X is denoted by π_X and the distance from x to the set X is denoted $\text{dist}_X(x) = \min_{z \in X} \|x - z\|$. The indicator function of a set X is denoted:

$$\mathbb{I}_X(x) = \begin{cases} 0, & \text{if } x \in X \\ \infty, & \text{otherwise} \end{cases}. \text{ We use notations } \partial h(x; \xi) \text{ for the subdifferential set}$$

and $g_h(x; \xi)$ for a subgradient of $h(\cdot; \xi)$ at x . In differentiable case, $\nabla f(\cdot; \xi)$ is the gradient of component ξ . Finally, we use $\mathbb{E}[\cdot]$ for (conditional) expectation operator.

1.2 Preliminaries

We denote the set of optimal solutions with X^* and x^* any optimal point for (1).

Assumption 1 *The objective function of our main problem (1) satisfies:*

(i) *The function $f(\cdot; \xi)$ has L_f -Lipschitz gradient, i.e. there exists $L_f > 0$ such that:*

$$\|\nabla f(x; \xi) - \nabla f(y; \xi)\| \leq L_f \|x - y\|, \quad \forall x, y \in \mathbb{R}^n, \xi \in \Omega.$$

and f is σ_f -strongly convex, i.e. there exists $\sigma_f > 0$ satisfying:

$$f(x) \geq f(y) + \langle \nabla f(y), x - y \rangle + \frac{\sigma_f}{2} \|x - y\|^2 \quad \forall x, y \in \mathbb{R}^n. \quad (2)$$

(ii) *There exists subgradient mappings $g_F(\cdot; \xi) : \text{dom } F(\cdot; \xi) \times \Omega \mapsto \mathbb{R}^n$ such that $g_F(x; \xi) \in \partial F(x; \xi)$, $\forall \xi \in \Omega$ and $\mathbb{E}[g_F(x; \xi)] \in \partial F(x)$.*

(iii) *$F(\cdot; \xi)$ has bounded gradients on the optimal set: there exists $\mathcal{S}_F^* \geq 0$ such that $\mathbb{E}[\|g_F(x^*; \xi)\|^2] \leq \mathcal{S}_F^* < \infty$ for all $x^* \in X^*$.*

(iv) *For any $g_F(x^*) \in \partial F(x^*)$ there exists bounded subgradients $g_F(x^*; \xi) \in \partial F(x^*; \xi)$, such that $\mathbb{E}[g_F(x^*; \xi)] = g_F(x^*)$ and $\mathbb{E}[\|g_F(x^*; \xi)\|^2] < \mathcal{S}_F^*$. Moreover, for simplicity, we assume throughout the paper $g_F(x^*) = \mathbb{E}[g_F(x^*; \xi)] = 0$.*

The first part of the above assumption is natural in the stochastic optimization problems. The Assumption (ii) guarantee the existence of a subgradient mapping. The third part Assumption (iii) is standard in the literature related to proximal stochastic algorithms. Notice that most stochastic subgradient algorithms require bounded gradients on the entire domain [10,13,16], which is more restrictive than condition (iii). The assumption (iv) needs a more consistent discussion, detailed in [17]. However, for completeness we include it in the remark below.

Remark 1 Denote the set $\mathbb{E}[\partial F(\cdot; \xi)] = \{\mathbb{E}[g_F(\cdot; \xi)] \mid g_F(\cdot; \xi) \in \partial F(\cdot; \xi)\}$. In general, for convex functions it can be easily shown $\mathbb{E}[\partial F(x; \xi)] \subseteq \partial F(x)$ for all $x \in \text{dom}(F)$ (see [21]). However, (iv) is guaranteed by the stronger equality:

$$\mathbb{E}[\partial h(x; \xi)] = \partial h(x). \quad (3)$$

Discrete case. Let us consider finite discrete domains $\Omega = \{1, \dots, m\}$. Then [20, Theorem 23.8] guarantees that the finite sum objective function from (1) satisfy (3) if $\bigcap_{\xi \in \Omega} \text{ri}(\text{dom}(h(\cdot; \xi))) \neq \emptyset$. The $\text{ri}(\text{dom}(\cdot))$ can be further relaxed to $\text{dom}(\cdot)$ for

polyhedral components. In particular, let X_1, \dots, X_m be finitely many closed convex satisfying qualification condition: $\bigcap_{i=1}^m \text{ri}(X_i) \neq \emptyset$, then also (3) holds, i.e. $\mathcal{N}_X(x) = \sum_{i=1}^m \mathcal{N}_{X_i}(x)$ (see (by [20, Corrolary 23.8.1])). Again, $\text{ri}(X_i)$ can be relaxed to the set itself for polyhedral sets. As pointed by [3], the (bounded) linear regularity property of $\{X_i\}_{i=1}^m$ implies the intersection qualification condition.

Under support of these arguments, observe that (iv) can be easily checked for most finite-sum examples arisen in convex learning problems.

Continuous case. In the nondiscrete case, sufficient conditions for (3) are discussed in [21]. Based on the arguments from [21], an assumption equivalent to (iv) is considered in [24] under the name of 2-integrable representation of x^* (definition in [24, Section B]). On short, if $h(\cdot; \xi)$ is normal convex integrand with full domain then x^* admits an 2-integrable representation $\{g_h(x^*; \xi)\}_{\xi \in \Omega}$, and implicitly (iv) holds.

We mention that deriving a more complicated result similar to Lemma 3 we could avoid assumption (iv). However, since (iv) facilitates the simplicity and naturality of our results and while our target applications are not excluded, we assume throughout the paper that (iv) holds.

Let closed convex sets $\{X_\xi\}_{\xi \in \Omega}$ and $X = \bigcap_{\xi \in \Omega} X_\xi$, then a favorable "conditioning" property for most projection methods is the linear regularity property: there exists $\kappa > 0$ such that

$$\mathbb{E}[\text{dist}_{X_\xi}^2(x)] \geq \kappa \text{dist}_X^2(x) \quad \forall x \in \mathbb{R}^n. \quad (4)$$

Given some smoothing parameter $\mu > 0$, define Moreau envelope of $h(x; \xi)$ and the prox operator as follows:

$$\begin{aligned} h_\mu(x; \xi) &:= \min_{z \in \mathbb{R}^n} h(z; \xi) + \frac{1}{2\mu} \|z - x\|^2 \\ \text{prox}_{h, \mu}(x; \xi) &:= \arg \min_{z \in \mathbb{R}^n} h(z; \xi) + \frac{1}{2\mu} \|z - x\|^2. \end{aligned}$$

The approximate $h_\mu(\cdot; \xi)$ inherits the same convexity properties of $h(\cdot; \xi)$ and additionally has Lipschitz continuous gradient with constant $\frac{1}{\mu}$, see [19]. In particular, when $h(x; \xi) = \mathbb{I}_{X_\xi}(x)$ the prox operator becomes the projection operator $\text{prox}_{h, \mu}(x; \xi) = \pi_{X_\xi}(x)$.

2 Stochastic Splitting Proximal Gradient Algorithm

In the following section we present the Stochastic Splitting Proximal Gradient (SSPG) algorithm and analyze its nonasymptotic convergence towards the optimal set of the original problem (1). The asymptotic convergence of vanishing stepsize SSPG have been analyzed in [24].

Let $x^0 \in \mathbb{R}^n$ be a starting point and $\{\mu_k\}_{k \geq 0}$ be a nonincreasing positive sequence of stepsizes.

Stochastic Splitting Proximal Gradient algorithm (SSPG): For $k \geq 0$ compute

1. Choose randomly $\xi_k \in \Omega$ w.r.t. probability distribution \mathbb{P}
2. Update:

$$\begin{aligned} y^k &= x^k - \mu_k \nabla f(x^k; \xi_k) \\ x^{k+1} &= \text{prox}_{h, \mu_k}(y^k; \xi_k) \end{aligned}$$

3. If the stopping criterion holds, then **STOP**, otherwise $k = k + 1$.

The SSPG iteration $x^{k+1} = \text{prox}_{h, \mu_k}(x^k - \mu_k \nabla f(x^k; \xi_k); \xi_k)$ is mainly a Stochastic Proximal Gradient iteration based on stochastic proximal maps [24]. Thus, the first step of algorithm SSPG consists of a varying-stepsize (vanilla) stochastic gradient update, while the second step rely on a stochastic proximal update, or equivalently a gradient step in the direction of the randomly sampled gradient of expected Moreau envelope $h_\mu(\cdot)$. Further results will state that a diminishing stepsize is an appropriate choice to obtain convergence in expectation. By varying our central model, this general SSPG scheme recovers several well-known stochastic first order algorithms.

(i) In the smooth case ($h = 0$), SSPG reduces to vanilla SGD [10]:

$$x^{k+1} = x^k - \mu_k \nabla f(x^k; \xi_k).$$

(ii) By considering proximal-tractable regularizers (i.e. $h(\cdot; \xi) = h(\cdot)$) or simple convex sets (i.e. $h(\cdot; \xi) = \mathbb{I}_X(\cdot)$, with $\pi_X(\cdot)$ computable in closed form), then we recover Proximal (or Projected, respectively) SGD [22]:

$$x^{k+1} = \text{prox}_{h, \mu_k}(x^k - \mu_k \nabla f(x^k; \xi_k)) \quad \text{or} \quad x^{k+1} = \pi_X(x^k - \mu_k \nabla f(x^k; \xi_k)).$$

(iii) For nonsmooth objective functions, when $f = 0$, SSPG is equivalent with SPP iteration [1,18,28]:

$$x^{k+1} = \text{prox}_{h, \mu_k}(x^k; \xi_k).$$

(iv) For CFPs, i.e. $f = 0$, $h(\cdot; \xi) = \mathbb{I}_{X_\xi}(\cdot)$, the SSPG iteration is simplified since $y^k = x^k$ and $x^{k+1} = \text{prox}_{h, \mu_k}(x^k; \xi_k)$. Thus it recovers Randomized Alternating Projections scheme [2]:

$$x^{k+1} = \pi_{X_{\xi_k}}(x^k).$$

Therefore, the below results will implicitly represent unifying convergence rates for these algorithms, under stated assumptions.

3 Iteration complexity in expectation

In this section we assume that the function f is strongly convex and derive convergence rates for this particular case. We recall a few elementary inequalities that will be used in our proofs. The proof of the following lemma can be found in [15].

Lemma 1 ([15]) *Under Assumption 1(i), the following inequality holds:*

$$f(x; \xi) \leq f(y; \xi) + \langle \nabla f(y; \xi), x - y \rangle + \frac{L_f}{2} \|x - y\|^2 \quad \forall x, y \in \mathbb{R}^n.$$

A second elementary bound the we will use is: let $v, x, y \in \mathbb{R}^n$ and $\nu > 0$, then

$$\langle v, x - y \rangle + \frac{1}{2\nu} \|x - y\|^2 \geq -\frac{\nu}{2} \|v\|^2. \quad (5)$$

We denote the history of index choices by $\Xi_k = \{\xi_0, \dots, \xi_{k-1}\}$.

Lemma 2 *Let Assumption 1 hold and $\mu_k \leq \frac{1}{2L_f}$. Then the sequence $\{x^k\}_{k \geq 0}$ generated by SSPG satisfies:*

$$\begin{aligned} \mathbb{E}[\|x^{k+1} - x^*\|^2] &\leq (1 - \sigma_f \mu_k) \mathbb{E}[\|x^k - x^*\|^2] \\ &\quad + 2\mu_k \mathbb{E} \left[F(x^*) - F(x^{k+1}; \xi_k) - \frac{1}{4\mu_k} \|x^{k+1} - x^k\|^2 \right]. \end{aligned}$$

Proof First notice that from the optimality conditions of the subproblem $\min_z h(z; \xi) + \frac{1}{2\mu} \|z - y\|^2$, the following relation holds:

$$g_h(x^{k+1}; \xi_k) + \frac{1}{\mu_k} (x^{k+1} - y^k) = 0. \quad (6)$$

Further we obtain the main recurrence:

$$\begin{aligned} \|x^{k+1} - x^*\|^2 &= \|x^k - x^*\|^2 + 2\langle x^{k+1} - x^k, x^k - x^* \rangle + \|x^{k+1} - x^k\|^2 \\ &= \|x^k - x^*\|^2 + 2\langle x^{k+1} - x^k, x^{k+1} - x^* \rangle - \|x^{k+1} - x^k\|^2 \\ &= \|x^k - x^*\|^2 + 2\langle \mu_k \nabla f(x^k; \xi_k) + \mu_k g_h(x^{k+1}; \xi_k), x^* - x^{k+1} \rangle - \|x^{k+1} - x^k\|^2 \\ &\leq \|x^k - x^*\|^2 + 2\mu_k \langle \nabla f(x^k; \xi_k), x^* - x^{k+1} \rangle - \|x^{k+1} - x^k\|^2 \\ &\quad + 2\mu_k [h(x^*; \xi_k) - h(x^{k+1}; \xi_k)] \\ &= \|x^k - x^*\|^2 - 2\mu_k \left(\langle \nabla f(x^k; \xi_k), x^{k+1} - x^k \rangle + \frac{1}{4\mu_k} \|x^{k+1} - x^k\|^2 \right. \\ &\quad \left. + h(x^{k+1}; \xi_k) \right) + 2\mu_k \langle \nabla f(x^k; \xi_k), x^* - x^k \rangle - \frac{1}{2} \|x^{k+1} - x^k\|^2 + 2\mu_k h(x^*; \xi_k). \end{aligned} \quad (7)$$

By using in (7) the stepsize bound $\mu_k \leq \frac{1}{2L_f}$ and Lemma 1, we obtain:

$$\begin{aligned} \|x^{k+1} - x^*\|^2 &\stackrel{\mu_k \leq \frac{1}{2L_f}}{\leq} \|x^k - x^*\|^2 \\ &\quad - 2\mu_k \left(\langle \nabla f(x^k; \xi_k), x^{k+1} - x^k \rangle + \frac{L_f}{2} \|x^{k+1} - x^k\|^2 + h(x^{k+1}; \xi_k) \right) \\ &\quad + 2\mu_k \langle \nabla f(x^k; \xi_k), x^* - x^k \rangle - \frac{1}{2} \|x^{k+1} - x^k\|^2 + 2\mu_k h(x^*; \xi_k) \\ &\stackrel{\text{Lemma 1}}{\leq} \|x^k - x^*\|^2 - 2\mu_k \left(F(x^{k+1}; \xi_k) - f(x^k; \xi_k) \right) \\ &\quad + 2\mu_k \langle \nabla f(x^k; \xi_k), x^* - x^k \rangle - \frac{1}{2} \|x^{k+1} - x^k\|^2 + 2\mu_k h(x^*; \xi_k). \end{aligned}$$

Further we take expectation w.r.t. ξ_k in both sides and use the strong convexity property (2) (with $x = x^k$ and $y = x^*$) to finally derive the following:

$$\begin{aligned}
\mathbb{E}[\|x^{k+1} - x^*\|^2 \mid \Xi_k] &\leq \|x^k - x^*\|^2 + 2\mu_k \mathbb{E}[(f(x^k; \xi_k) - F(x^{k+1}; \xi_k)) \mid \Xi_k] \\
&\quad + 2\mu_k \langle \nabla f(x^k), x^* - x^k \rangle - \frac{1}{2} \mathbb{E}[\|x^{k+1} - x^k\|^2 \mid \Xi_k] + 2\mu_k h(x^*) \\
&\leq \|x^k - x^*\|^2 + 2\mu_k \mathbb{E}[f(x^k; \xi_k) - F(x^{k+1}; \xi_k) \mid \Xi_k] \\
&\quad + 2\mu_k \left(F(x^*) - f(x^k) - \frac{\sigma_f}{2} \|x^k - x^*\|^2 \right) - \frac{1}{2} \mathbb{E}[\|x^{k+1} - x^k\|^2 \mid \Xi_k] \\
&= (1 - \sigma_f \mu_k) \|x^k - x^*\|^2 + 2\mu_k \mathbb{E} \left[F(x^*) - F(x^{k+1}; \xi_k) - \frac{1}{4\mu_k} \|x^{k+1} - x^k\|^2 \mid \Xi_k \right].
\end{aligned}$$

Finally, by taking full expectation $\mathbb{E}[\cdot]$, over the entire index history, we obtain the above result.

Further we present some lower bounds on the second term from the right hand side, which will allow the synthesis of general complexity estimates over stochastic and deterministic contexts.

Lemma 3 *Given $\mu > 0$, let Assumption 1 hold. Then F satisfies the following relations: given $k \geq 0$*

- (i) $\mathbb{E} \left[F(x^{k+1}; \xi_k) - F^* + \frac{1}{4\mu_k} \|x^{k+1} - x^k\|^2 \right] \geq -\mu_k \mathbb{E} [\|g_F(x^*; \xi)\|^2].$
- (ii) *Let $\{X_\xi\}_{\xi \in \Omega}$ be convex sets satisfying linear regularity with constant κ , such that $X := \bigcap_{\xi \in \Omega} X_\xi \neq \emptyset$. Also let $h(\cdot; \xi) = \mathbb{I}_{X_\xi}(\cdot)$, then*

$$\begin{aligned}
&\mathbb{E} \left[F(x^{k+1}; \xi_k) - F^* + \frac{1}{4\mu_k} \|x^{k+1} - x^k\|^2 \right] \\
&\geq -2\mu_k \mathbb{E} [\|\nabla f(x^*; \xi)\|^2] - \frac{4\mu_k}{\kappa} \|\nabla f(x^*)\|^2 + \frac{\kappa}{16\mu_k} \mathbb{E} [\text{dist}_X^2(x^k)]
\end{aligned}$$

Proof In order to prove (i) let $x^* \in X^*$ and $g_F(x^*; \xi) \in \partial F(x^*; \xi)$, by convexity of $f(\cdot; \xi)$ we have:

$$\begin{aligned}
&F(x^{k+1}; \xi_k) - F^* + \frac{1}{4\mu_k} \|x^{k+1} - x^k\|^2 \\
&\geq \langle g_F(x^*; \xi_k), x^{k+1} - x^* \rangle + \frac{1}{4\mu_k} \|x^{k+1} - x^k\|^2 \\
&\geq \langle g_F(x^*; \xi_k), x^k - x^* \rangle + \langle g_F(x^*; \xi_k), x^{k+1} - x^k \rangle + \frac{1}{4\mu_k} \|x^{k+1} - x^k\|^2. \quad (8)
\end{aligned}$$

On one hand, based on Assumption 1 (iv), we observe that the first term of the right hand side has null expectation:

$$\begin{aligned}
\mathbb{E} [\langle g_F(x^*; \xi_k), x^k - x^* \rangle] &= \mathbb{E} [\langle \mathbb{E} [g_F(x^*; \xi_k) \mid \Xi_k], x^k - x^* \rangle] \\
&\stackrel{\text{Assump. 1(iv)}}{=} \mathbb{E} [\langle g_F(x^*), x^k - x^* \rangle] = 0. \quad (9)
\end{aligned}$$

On the other hand, for any $x^* \in X^*$, the second term can be lower bounded by (5):

$$\langle g_F(x^*; \xi_k), x^{k+1} - x^k \rangle + \frac{1}{4\mu_k} \|x^{k+1} - x^k\|^2 \geq -\mu_k \|g_F(x^*; \xi_k)\|^2. \quad (10)$$

We take full expectation (over the entire index history) in (8) and use relations (9)-(10) to get the first part (i).

To derive the second part (ii), we proceed slightly different:

$$\begin{aligned} & F(x^{k+1}; \xi_k) - F^* + \frac{1}{4\mu_k} \|x^{k+1} - x^k\|^2 \\ &= f(x^{k+1}; \xi_k) - f^* + \frac{1}{8\mu_k} \|x^{k+1} - x^k\|^2 + h(x^{k+1}; \xi_k) + \frac{1}{8\mu_k} \|x^{k+1} - x^k\|^2 \\ &\geq f(x^{k+1}; \xi_k) - f^* + \frac{1}{8\mu_k} \|x^{k+1} - x^k\|^2 + \min_z h(z; \xi_k) + \frac{1}{8\mu_k} \|z - x^k\|^2 \\ &\geq \langle \nabla f(x^*; \xi_k), x^{k+1} - x^k \rangle + \frac{1}{8\mu_k} \|x^{k+1} - x^k\|^2 + h_{4\mu_k}(x^k; \xi_k) \\ &= \langle \nabla f(x^*; \xi_k), x^k - x^* \rangle + \\ &\quad \langle \nabla f(x^*; \xi_k), x^{k+1} - x^k \rangle + \frac{1}{8\mu_k} \|x^{k+1} - x^k\|^2 + h_{4\mu_k}(x^k; \xi_k). \end{aligned} \quad (11)$$

We proceed to bound each term of right hand side, similarly as in (i). By using the optimality condition (i.e. $\langle \nabla f(x^*), z - x^* \rangle \geq 0$ for all $z \in X$), the expectation of the first term can be lower bounded as follows:

$$\begin{aligned} & \mathbb{E} [\langle \nabla f(x^*; \xi_k), x^k - x^* \rangle] \\ &= \mathbb{E} [\langle \mathbb{E} [\nabla f(x^*; \xi_k) \mid \Xi_k], x^k - x^* \rangle] \\ &= \mathbb{E} [\langle \nabla f(x^*), x^k - x^* \rangle] \\ &= \mathbb{E} [\langle \nabla f(x^*), \pi_X(x^k) - x^* \rangle + \langle \nabla f(x^*), x^k - \pi_X(x^k) \rangle] \\ &\stackrel{O.C.}{\geq} \mathbb{E} [\langle \nabla f(x^*), x^k - \pi_X(x^k) \rangle] \\ &\stackrel{C.S.}{\geq} \mathbb{E} [-\|f(x^*)\| \text{dist}_X(x^k)] \geq -\|f(x^*)\| \sqrt{\mathbb{E} [\text{dist}_X^2(x^k)]}, \end{aligned} \quad (12)$$

where in the second inequality we used Cauchy-Schwartz inequality and in the last one we used $\mathbb{E}[Y] \leq \sqrt{\mathbb{E}[Y^2]}$. After taking full expectation in (11) and using (5), (12) and the linear regularity property $h_\mu(x) := \frac{1}{2\mu} \mathbb{E} [\text{dist}_{X_\xi}^2(x)] \geq \frac{\kappa}{2\mu} \text{dist}_X^2(x)$

yields the following:

$$\begin{aligned}
& \mathbb{E} \left[F(x^{k+1}; \xi_k) - F^* + \frac{1}{4\mu_k} \|x^{k+1} - x^k\|^2 \right] \\
& \stackrel{(5)+(12)}{\geq} -2\mu_k \mathbb{E} [\|\nabla f(x^*; \xi)\|^2] - \|f(x^*)\| \sqrt{\mathbb{E} [\text{dist}_X^2(x^k)]} + \mathbb{E}[h_{4\mu_k}(x^k)] \\
& \stackrel{\text{lin.reg.}}{\geq} -2\mu_k \mathbb{E} [\|\nabla f(x^*; \xi)\|^2] - \|f(x^*)\| \sqrt{\mathbb{E} [\text{dist}_X^2(x^k)]} + \frac{\kappa}{8\mu_k} \mathbb{E}[\text{dist}_X^2(x^k)] \\
& \stackrel{(5)}{\geq} -2\mu_k \mathbb{E} [\|\nabla f(x^*; \xi)\|^2] - \frac{4\mu_k}{\kappa} \|\nabla f(x^*)\|^2 + \frac{\kappa}{16\mu_k} \mathbb{E}[\text{dist}_X^2(x^k)].
\end{aligned}$$

In the last inequality we also used (5).

Next we present the main recurrences which will finally generate our nonasymptotic convergence rates.

Theorem 2 *Let Assumptions 1 hold, then the SSPG sequence $\{x^k\}_{k \geq 0}$ satisfies:*

(i) *Let $\mu_k \leq \frac{1}{2L_f}$ for all $k \geq 0$, then:*

$$\mathbb{E}[\|x^{k+1} - x^*\|^2] \leq (1 - \sigma_f \mu_k) \mathbb{E}[\|x^k - x^*\|^2] + \mu_k^2 \Sigma,$$

where $\Sigma = 2\mathbb{E} [\|g_F(x^*; \xi)\|^2]$.

(ii) *In particular, let $h(\cdot; \xi) = \mathbb{I}_{X_\xi}(\cdot)$ such that $\{X_\xi\}_{\xi \in \Omega}$ are linearly regular sets with constant κ . Then the following recurrence holds:*

$$\mathbb{E}[\|x^{k+1} - x^*\|^2] \leq (1 - \sigma_f \mu_k) \mathbb{E}[\|x^k - x^*\|^2] + \mu_k^2 \Sigma - \frac{\kappa}{8} \mathbb{E}[\text{dist}_X^2(x^k)],$$

where $\Sigma = 4\mathbb{E} [\|\nabla f(x^*; \xi)\|^2] + \frac{8}{\kappa} \|\nabla f(x^*)\|^2$.

Proof The proof result straightforwardly from Lemmas 2 and 3.

Remark 2 Consider deterministic setting $F(\cdot; \xi) = F(\cdot)$ and $\mu_k = \frac{1}{2L_f}$, then SSPG becomes the proximal gradient algorithm and Theorem 2(i) holds with $g_F(x^*; \xi) = g_F(x^*) = 0$, implying that $\Sigma = 0$. Thus the well-known iteration complexity estimate $\mathcal{O}\left(\frac{L_f}{\sigma_f} \log(1/\epsilon)\right)$ [4,14] of proximal gradient algorithm is recovered up to a constant from Theorem 2.

The above recurrences generates immediately the following convergence rates.

Corollary 3 *Under Assumption 1, the following convergence rates hold:*

(i) *Let $\mu_k = \frac{1}{k^\gamma}$, $\gamma \in (0, 1)$ then: $\mathbb{E}[\|x^k - x^*\|^2] \leq \mathcal{O}\left(\frac{1}{k^\gamma}\right)$*

(ii) *Let $\mu_k = \frac{1}{k}$, then :* $\mathbb{E}[\|x^k - x^*\|^2] \leq \begin{cases} \mathcal{O}\left(\frac{1}{k}\right) & \text{if } \mu_0 \sigma_f > e - 1 \\ \mathcal{O}\left(\frac{\ln k}{k}\right) & \text{if } \mu_0 \sigma_f > e - 1 \\ \mathcal{O}\left(\frac{1}{k}\right)^{2\ln(1+\mu_0 \sigma_f)} & \text{if } \mu_0 \sigma_f < e - 1. \end{cases}$

(iii) For constant stepsize $\mu_k = \mu > 0$, the recurrence from Theorem 2(i) implies:

$$\mathbb{E}[\|x^k - x^*\|^2] \leq (1 - \mu\sigma_f)^k \|x^0 - x^*\|^2 + \frac{\mu}{\sigma_f} \Sigma, \quad (13)$$

where $\Sigma = 2\mathbb{E}[\|g_F(x^*; \xi)\|^2]$.

(iv) Moreover, consider convex feasibility problem where $f = 0$, $h(\cdot; \xi) = \mathbb{I}_{X_\xi}(\cdot)$ with κ -linearly regular sets $\{X_\xi\}_{\xi \in \Omega}$ and constant stepsize $\mu_k = \mu$. Then the SSPG sequence $\{x^k\}_{k \geq 0}$ converges linearly as follows:

$$\mathbb{E}[\text{dist}_X^2(x^k)] \leq \left(1 - \frac{\kappa}{8}\right)^k \text{dist}_X^2(x^0).$$

Proof The proof for the first two results (i) and (ii) follows similar lines with [18, Corollary 15]. However, for completeness we present it in appendix.

For (iii), notice that Theorem 2(i) straightforwardly implies:

$$\begin{aligned} \mathbb{E}[\|x^k - x^*\|^2] &\leq (1 - \sigma_f \mu) \mathbb{E}[\|x^{k-1} - x^*\|^2] + \mu^2 \Sigma \\ &\leq (1 - \mu\sigma_f)^k \|x^0 - x^*\|^2 + \mu^2 \Sigma \sum_{i=0}^{k-1} (1 - \mu\sigma_f)^i \\ &\leq (1 - \mu\sigma_f)^k \|x^0 - x^*\|^2 + \frac{\mu}{\sigma_f} \Sigma [1 - (1 - \mu\sigma_f)^k] \\ &\leq (1 - \mu\sigma_f)^k \|x^0 - x^*\|^2 + \frac{\mu}{\sigma_f} \Sigma. \end{aligned}$$

Now, let CFP settings $f = 0$, $h(\cdot; \xi) = \mathbb{I}_{X_\xi}(\cdot)$ hold. Theorem 2(ii) states that:

$$\mathbb{E}[\|x^{k+1} - x^*\|^2] \leq \mathbb{E}[\|x^k - x^*\|^2] - \frac{\kappa}{8} \mathbb{E}[\text{dist}_X^2(x^k)] \quad \forall x^* \in X^*.$$

Since in this case, $X = \bigcap_{\xi} X_\xi = X^*$, then by choosing $x^* = \pi_X(x^k)$ and using in the LHS that $\|x^{k+1} - \pi_X(x^k)\| \geq \text{dist}_X(x^{k+1})$, then we obtain:

$$\mathbb{E}[\text{dist}_X^2(x^{k+1})] \leq \mathbb{E}[\|x^{k+1} - \pi_X(x^k)\|^2] \leq \left(1 - \frac{\kappa}{8}\right) \mathbb{E}[\text{dist}_X^2(x^k)],$$

which yields the linear convergence rate of SSPG.

Remark 3 Although sublinear $\mathcal{O}(1/k)$ convergence rates for strongly convex objectives are typically obtained in literature for many first-order stochastic schemes [10,18], these rates are novel due to their generalization potential.

Regarding second rate (ii), although it expresses a geometric decrease of the initial residual term, this rate states that, after $\mathcal{O}\left(\frac{1}{\mu\sigma_f} \log\left(\frac{1}{\epsilon}\right)\right)$ iterations, the sequence $\{x^k\}_{k \geq 0}$ will remain (in expectation) in a bounded neighborhood of the optimal point $\{x : \|x - x^*\|^2 \leq \frac{\mu}{\sigma_f} \Sigma\}$. This fact suggests that only sufficiently small constant stepsizes guarantee the convergence of SSPG sequence.

In the convex feasibility setting, SSPG reduces to Randomized Alternating Projections algorithm for which the obtained linear rate obeys the rates from the literature up to a constant. However, we believe that using some refinements of proof arguments there might be obtained optimal rates w.r.t. the constants.

4 Application to Parametric Sparse Representation

Sparse representation [7] starts with the given signal $y \in \mathbb{R}^m$ and aims to find the sparse signal $x \in \mathbb{R}^n$ by projecting y to a much smaller subspace through the over-complete dictionary $T \in \mathbb{R}^{m \times n}$. Among the many Dictionary Learning techniques, we focus on the multi-parametric cospase model proposed in [27], since it could efficiently illustrate the empirical capabilities of SSPG. The multi-parametric cospase representation problem is given by:

$$\begin{aligned} \min_x \quad & \|Tx - y\|_2^2 \\ \text{s.t.} \quad & \|\Delta x\|_1 \leq \delta, \end{aligned} \quad (14)$$

where T and x correspond to the dictionary and, respectively, the resulting sparse representation, with sparsity being imposed on a scaled subspace Δx with $\Delta \in \mathbb{R}^{p \times n}$. In pursuit of (1), we move to the exact penalty problem $\min_x \frac{1}{2m} \|Tx - y\|_2^2 + \lambda \|\Delta x\|_1$. In order to limit the solution norm we further regularize the unconstrained objective using an ℓ_2 term as follows:

$$\min_x \quad \frac{1}{2m} \|Tx - y\|_2^2 + \lambda \|\Delta x\|_1 + \frac{\alpha}{2} \|x\|_2^2. \quad (15)$$

The decomposition which puts the above formulation into model (1) consists of:

$$f(x; \xi) = \frac{1}{2} (T_\xi x - y_\xi)_2^2 + \frac{\alpha}{2} \|x\|_2^2 \quad (16)$$

where T_ξ represents line ξ of matrix T , and

$$h(x; \xi) = m\lambda |\Delta_\xi x|. \quad (17)$$

To compute the SSPG iteration for the sparse representation problem, we note that the proximal operator $\text{prox}_{h,\mu}(x; \xi) = \arg \min_{z \in \mathbb{R}^n} m\lambda |\Delta_\xi z| + \frac{1}{2\mu} \|z - x\|^2$ is given by

$$\text{prox}_{h,\mu}(x; \xi) = \begin{cases} x - \frac{\Delta_\xi x}{\|\Delta_\xi\|^2} \Delta_\xi^T & \text{if } \frac{|\Delta_\xi x|}{m\lambda \|\Delta_\xi\|^2} \leq \mu \\ x - \mu m \lambda \text{sgn}(\Delta_\xi x) \Delta_\xi^T & \text{if } \frac{|\Delta_\xi x|}{m\lambda \|\Delta_\xi\|^2} > \mu \end{cases}$$

Also, the gradient of the smooth function $f(\cdot; \xi)$ is given by

$$\nabla f(x; \xi) = (T_\xi^T T_\xi + \alpha I_n) x - T_\xi^T y_\xi$$

and with that we are ready to formulate the resulting particular variant of SSPG.

SSPG - Sparse Representation (SSPG-SR): For $k \geq 0$ compute

1. Choose randomly $\xi_k \in \Omega$ w.r.t. probability distribution \mathbb{P}
2. Update:

$$\begin{aligned} y^k &= [I_n - \mu_k (T_{\xi_k}^T T_{\xi_k} + \alpha I_n)] x^k + \mu_k T_{\xi_k}^T y_{\xi_k}, \quad \beta_k = \frac{\Delta_{\xi_k} y^k}{\|\Delta_{\xi_k}\|^2} \\ x^{k+1} &= \begin{cases} y^k - \beta_k \Delta_{\xi_k}^T & \text{if } |\beta_k| \leq m\lambda\mu \\ y^k - \mu m \lambda \text{sgn}(\beta_k) \Delta_{\xi_k}^T & \text{if } |\beta_k| > m\lambda\mu \end{cases} \end{aligned}$$

3. If the stopping criterion holds, then **STOP**, otherwise $k = k + 1$.

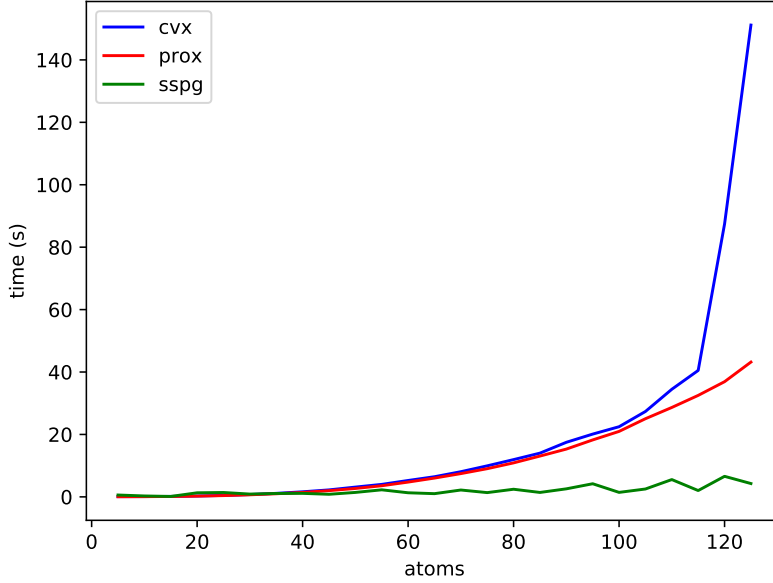


Fig. 1 Increasing problem dimension. Features are 4 times the number of atoms. ($\alpha = 0.5$, $\lambda = 5$, $\varepsilon = 10^{-6}$)

We use randomly generated data with a standard normal distribution.¹ In the following we present numerical simulations that show-case the application of SSPG to the multi-parametric sparse representation problem and compare it to CVX and Proximal Gradient method (denoted `prox`) [14]. Notice that at each iteration Proximal Gradient method requires, given current point x_{PG}^k , the computation of:

- gradient $\nabla f(x_{PG}^k) = \frac{1}{m} T^T (T x_{PG}^k - y)$
- proximal update $z(x_{PG}^k) = \arg \min_z \lambda \|\Delta z\|_1 + \frac{L}{2} \|z - (x_{PG}^k - \frac{1}{L} \nabla f(x_{PG}^k))\|^2$

Our first experiment investigates the impact of the problem size on execution time. We vary the atoms in the dictionary starting from $n = 5$ up to $n = 125$ in increments of 5 and maintain a ratio of 4 features per atom such that $m = 4n$. First CVX is used to determine x^* within a $\varepsilon = 10^{-6}$ margin. Then Proximal Gradient method and SSPG-SR are executed until x^* is reached with the same approximation error. Figure 1 depicts the results. The experiment is stopped after $n = 125$ when CVX execution times become too large. We can clearly see a similar yet gentler increasing trend for Proximal Gradient method (`prox`) and indeed other experiments have shown it to reach an impas shortly after. SSPG-SR on the other hand maintains a steady pace throughout the tests.

Next, we focus on the iterations of Proximal Gradient method and SSPG-SR. To this end we design a similar experiment consisting of 4 tests where we vary the

¹ Data generating code available at <https://github.com/pirofti/SSPG>

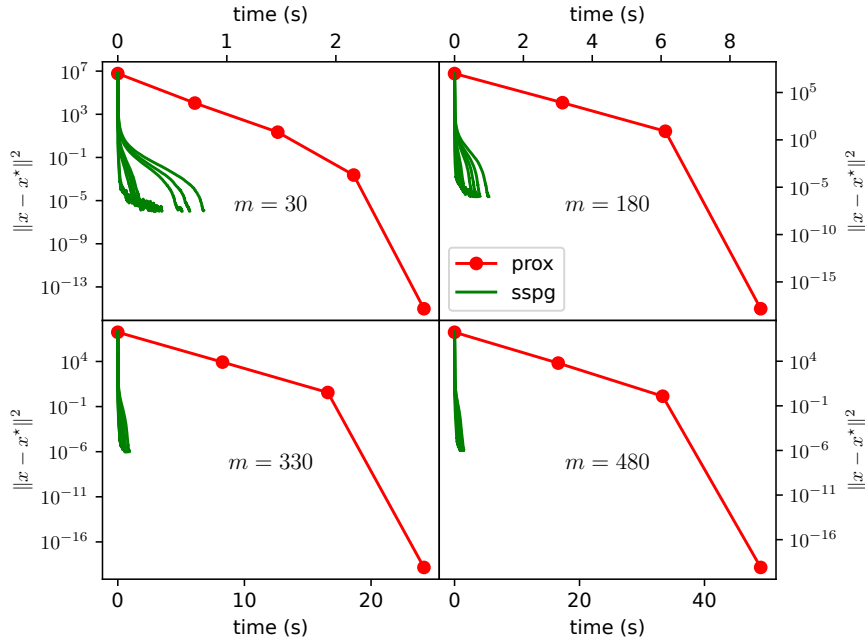


Fig. 2 Proximal Gradient method versus 10 rounds of SSPG iterations. Features are 6 times the number of atoms. ($\alpha = 0.2$, $\lambda = 5 \cdot 10^{-4}$, $\varepsilon = 10^{-6}$)

problem size by increasing the number of features with a ratio of 6 features per atom ($m = 6n$). As before, we use CVX to determine x^* with $\varepsilon = 10^{-6}$. With identical parameters and initialization we execute 10 rounds of SSPG-SR and compare them to the Proximal Gradient method iterations. We time the execution of each iteration and record its progress $\|x^k - x^*\|$. The result is shown in Figure 2. In all 4 panels it is clearly seen that the 10 SSPG-SR rounds perform much faster than Proximal Gradient method at least in a first phase. Further increases of the problem dimension lead to a stall in Proximal Gradient method. While SSPG-SR reaches a solution in less than a second in all the tests, Proximal Gradient method goes from 5 seconds in the first, to 9, 24 and 50 in the last. We note that even though the stopping criterion was the same for both methods ($\varepsilon = 10^{-6}$), Proximal Gradient method gets us much closer to x^* because, in our parameter setting, the SR problem is well-conditioned.

5 Conclusion

In this paper we presented preliminary guarantees for stochastic gradient schemes with stochastic proximal updates, which unify some well-known schemes in the literature. For future work, would be interesting to analyze convergence rate on (non)convex functions satisfying more relaxed convexity conditions (i.e. quadratic growth) and the empirical behaviour of SSPG scheme under various stepsize choices.

References

1. Hilal Asi and John C Duchi. Stochastic (approximate) proximal point methods: Convergence, optimality, and adaptivity. *SIAM Journal on Optimization*, 29(3):2257–2290, 2019.
2. Heinz Bauschke, Frank Deutsch, Hein Hundal, and Sung-Ho Park. Accelerating the convergence of the method of alternating projections. *Transactions of the American Mathematical Society*, 355(9):3433–3461, 2003.
3. Heinz H Bauschke, Jonathan M Borwein, and Wu Li. Strong conical hull intersection property, bounded linear regularity, jameson’s property (g), and error bounds in convex optimization. *Mathematical Programming*, 86(1):135–160, 1999.
4. Amir Beck and Marc Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM journal on imaging sciences*, 2(1):183–202, 2009.
5. Pascal Bianchi. Ergodic convergence of a stochastic proximal point algorithm. *SIAM Journal on Optimization*, 26(4):2235–2260, 2016.
6. Damek Davis and Dmitriy Drusvyatskiy. Stochastic model-based minimization of weakly convex functions. *SIAM Journal on Optimization*, 29(1):207–239, 2019.
7. Michael Elad. *Sparse and redundant representations: from theory to applications in signal and image processing*. Springer Science & Business Media, 2010.
8. David Hallac, Jure Leskovec, and Stephen Boyd. Network lasso: Clustering and optimization in large graphs. In *Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*, pages 387–396, 2015.
9. Jayash Koshal, Angelia Nedic, and Uday V Shanbhag. Regularized iterative stochastic approximation methods for stochastic variational inequality problems. *IEEE Transactions on Automatic Control*, 58(3):594–609, 2012.
10. Eric Moulines and Francis R Bach. Non-asymptotic analysis of stochastic approximation algorithms for machine learning. In *Advances in Neural Information Processing Systems*, pages 451–459, 2011.
11. Angelia Nedić. Random projection algorithms for convex set intersection problems. In *49th IEEE Conference on Decision and Control (CDC)*, pages 7655–7660. IEEE, 2010.
12. Angelia Nedić. Random algorithms for convex minimization problems. *Mathematical programming*, 129(2):225–253, 2011.
13. 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.
14. Yu Nesterov. Gradient methods for minimizing composite functions. *Mathematical Programming*, 140(1):125–161, 2013.
15. Yurii Nesterov. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer Science & Business Media, 2013.
16. Lam M Nguyen, Phuong Ha Nguyen, Marten van Dijk, Peter Richtárik, Katya Scheinberg, and Martin Takáč. Sgd and hogwild! convergence without the

- bounded gradients assumption. *arXiv preprint arXiv:1802.03801*, 2018.
17. Andrei Pătrașcu. New nonasymptotic convergence rates of stochastic proximal point algorithm for stochastic convex optimization. *Optimization*, pages 1–29, 2020.
 18. Andrei Patrascu and Ion Necoara. Nonasymptotic convergence of stochastic proximal point methods for constrained convex optimization. *The Journal of Machine Learning Research*, 18(1):7204–7245, 2017.
 19. R Tyrrell Rockafellar and Roger J-B Wets. *Variational analysis*, volume 317. Springer Science & Business Media, 2009.
 20. R.T. Rockafellar. *Convex Analysis*. Princeton University Press, Princeton, New Jersey, 1988.
 21. R.T. Rockafellar and R.J.-B. Wets. On the interchange of subdifferentiation and conditional expectation for convex functionals. *Stochastics*, 7(1):173–182, 1982.
 22. Lorenzo Rosasco, Silvia Villa, and Bang Công Vũ. Convergence of stochastic proximal gradient algorithm. *Applied Mathematics & Optimization*, pages 1–27, 2019.
 23. Ernest K Ryu and Stephen Boyd. Stochastic proximal iteration: a non-asymptotic improvement upon stochastic gradient descent. *Author website, early draft*, 2016.
 24. Adil Salim, Pascal Bianchi, and Walid Hachem. Snake: a stochastic proximal gradient algorithm for regularized problems over large graphs. *IEEE Transactions on Automatic Control*, 64(5):1832–1847, 2019.
 25. Shai Shalev-Shwartz, Yoram Singer, Nathan Srebro, and Andrew Cotter. Pegasos: Primal estimated sub-gradient solver for svm. *Mathematical programming*, 127(1):3–30, 2011.
 26. Wei Shi, Qing Ling, Gang Wu, and Wotao Yin. A proximal gradient algorithm for decentralized composite optimization. *IEEE Transactions on Signal Processing*, 63(22):6013–6023, 2015.
 27. Florin Stoican and Paul Irofti. Aiding dictionary learning through multi-parametric sparse representation. *Algorithms*, 12(7):131, 2019.
 28. Panos Toulis, Dustin Tran, and Edo Airoldi. Towards stability and optimality in stochastic gradient descent. In *Artificial Intelligence and Statistics*, pages 1290–1298, 2016.
 29. Rohan Varma, Harlin Lee, Jelena Kovacevic, and Yuejie Chi. Vector-valued graph trend filtering with non-convex penalties. *IEEE Transactions on Signal and Information Processing over Networks*, 2019.
 30. Mengdi Wang and Dimitri P Bertsekas. Stochastic first-order methods with random constraint projection. *SIAM Journal on Optimization*, 26(1):681–717, 2016.
 31. Xiao Wang, Shuxiong Wang, and Hongchao Zhang. Inexact proximal stochastic gradient method for convex composite optimization. *Computational Optimization and Applications*, 68(3):579–618, 2017.
 32. Yael Yankelevsky and Michael Elad. Dual graph regularized dictionary learning. *IEEE Transactions on Signal and Information Processing over Networks*, 2(4):611–624, 2016.
 33. Wenliang Zhong and James Kwok. Accelerated stochastic gradient method for composite regularization. In *Artificial Intelligence and Statistics*, pages 1086–1094, 2014.

6 Appendix

Proof (of Corollary 3) For simplicity denote $\theta_k = (1 - \mu_k \sigma_f)$, then Theorem 2 implies that:

$$\mathbb{E} [\|x^{k+1} - x^*\|^2] \leq \left(\prod_{i=0}^k \theta_i \right) \|x^0 - x^*\|^2 + \Sigma \sum_{i=0}^k \left(\prod_{j=i+1}^k \theta_j \right) \mu_i^2.$$

By using the Bernoulli inequality $1 - tx \leq \frac{1}{1+tx} \leq (1+x)^{-t}$ for $t \in [0, 1], x \geq 0$, then we have:

$$\prod_{i=l}^u \theta_i = \prod_{i=l}^u \left(1 - \frac{\mu_0}{i^\gamma} \sigma_f\right) \leq \prod_{i=l}^u (1 + \mu_0 \sigma_f)^{-1/i^\gamma} = (1 + \mu_0 \sigma_f)^{-\sum_{i=l}^u \frac{1}{i^\gamma}}. \quad (18)$$

On the other hand, if we use the lower bound

$$\sum_{i=l}^u \frac{1}{i^\gamma} \geq \int_l^{u+1} \frac{1}{\tau^\gamma} d\tau = \varphi_{1-\gamma}(u+1) - \varphi_{1-\gamma}(l). \quad (19)$$

then we can finally derive:

$$\begin{aligned} \sum_{i=0}^k \left(\prod_{j=i+1}^k \theta_j \right) \mu_i^2 &= \sum_{i=0}^m \left(\prod_{j=i+1}^k \theta_j \right) \mu_i^2 + \sum_{i=m+1}^k \left(\prod_{j=i+1}^k \theta_j \right) \mu_i^2 \\ &\stackrel{(18)+(19)}{\leq} \sum_{i=0}^m (1 + \mu_0 \sigma_f)^{\varphi_{1-\gamma}(i+1) - \varphi_{1-\gamma}(k)} \mu_i^2 + \mu_{m+1} \sum_{i=m+1}^k \left[\prod_{j=i+1}^k (1 - \mu_j \sigma_f) \right] \mu_i \\ &\leq (1 + \mu_0 \sigma_f)^{\varphi_{1-\gamma}(m) - \varphi_{1-\gamma}(k)} \sum_{i=0}^m \mu_i^2 \\ &\quad + \frac{\mu_{m+1}}{\sigma_f} \sum_{i=m+1}^k \left[\prod_{j=i+1}^k (1 - \mu_j \sigma_f) \right] (1 - (1 - \sigma_f \mu_i)) \\ &= (1 + \mu_0 \sigma_f)^{\varphi_{1-\gamma}(m) - \varphi_{1-\gamma}(k)} \mu_0^2 \sum_{i=0}^m \frac{1}{i^{2\gamma}} \\ &\quad + \frac{\mu_{m+1}}{\sigma_f} \sum_{i=m+1}^k \left[\prod_{j=i+1}^k (1 - \mu_j \sigma_f) - \prod_{j=i}^k (1 - \mu_j \sigma_f) \right] \\ &\leq (1 + \mu_0 \sigma_f)^{\varphi_{1-\gamma}(m) - \varphi_{1-\gamma}(k)} \frac{m^{1-2\gamma} - 1}{1 - 2\gamma} + \frac{\mu_{m+1}}{\sigma_f} \left[1 - \prod_{j=m+1}^k (1 - \mu_j \sigma_f) \right] \\ &\leq (1 + \mu_0 \sigma_f)^{\varphi_{1-\gamma}(m) - \varphi_{1-\gamma}(k)} \varphi_{1-2\gamma}(m) + \frac{\mu_{m+1}}{\sigma_f}. \end{aligned}$$

By denoting the second constant $\tilde{\theta}_0 = \frac{1}{1+\mu_0\sigma_f}$, then the last relation implies the following bound:

$$\mathbb{E} [\|x^{k+1} - x^*\|^2] \leq \tilde{\theta}_0^{\varphi_{1-\gamma}(k)} \|x^0 - x^*\|^2 + \tilde{\theta}_0^{\varphi_{1-\gamma}(k) - \varphi_{1-\gamma}(m)} \varphi_{1-2\gamma}(m) \Sigma + \frac{\mu_{m+1}}{\sigma_f} \Sigma.$$

Denote $r_k^2 = \mathbb{E}[\|x^k - x^*\|^2]$. To derive an explicit convergence rate order we analyze upper bounds on function ϕ .

(i) First assume that $\gamma \in (0, \frac{1}{2})$. This implies that $1 - 2\gamma > 0$ and that:

$$\varphi_{1-2\gamma} \left(\left\lfloor \frac{k}{2} \right\rfloor \right) \leq \varphi_{1-2\gamma} \left(\frac{k}{2} \right) = \frac{\left(\frac{k}{2}\right)^{1-2\gamma} - 1}{1 - 2\gamma} \leq \frac{\left(\frac{k}{2}\right)^{1-2\gamma}}{1 - 2\gamma}. \quad (20)$$

On the other hand, by using the inequality $e^{-x} \leq \frac{1}{1+x}$ for all $x \geq 0$, we obtain:

$$\begin{aligned} \tilde{\theta}_0^{\varphi_{1-\gamma}(k) - \varphi_{1-\gamma}(\frac{k-2}{2})} \varphi_{1-2\gamma} \left(\frac{k}{2} \right) &= e^{(\varphi_{1-\gamma}(k) - \varphi_{1-\gamma}(\frac{k-2}{2})) \ln \tilde{\theta}_0} \varphi_{1-2\gamma} \left(\frac{k}{2} \right) \\ &\leq \frac{\varphi_{1-2\gamma} \left(\frac{k}{2} \right)}{1 + [\varphi_{1-\gamma}(k) - \varphi_{1-\gamma}(\frac{k}{2} - 1)] \ln \frac{1}{\tilde{\theta}_0}} \stackrel{(20)}{\leq} \frac{\frac{k^{1-2\gamma}}{2^{1-2\gamma}(1-2\gamma)}}{\frac{1}{1-\gamma} [k^{1-\gamma} - (\frac{k}{2} - 1)^{1-\gamma}] \ln \frac{1}{\tilde{\theta}_0}} \\ &= \frac{\frac{k^{1-2\gamma}}{2^{1-2\gamma}(1-2\gamma)}}{\frac{k^{1-\gamma}}{1-\gamma} [1 - (\frac{1}{6})^{1-\gamma}] \ln \frac{1}{\tilde{\theta}_0}} = \frac{1-\gamma}{1-2\gamma} \frac{2^\gamma k^{-\gamma}}{2^{1-2\gamma} [1 - (\frac{1}{6})^{1-\gamma}] \ln \frac{1}{\tilde{\theta}_0}} = \mathcal{O} \left(\frac{1}{k^\gamma} \right). \end{aligned}$$

Therefore, in this case, the overall rate will be given by:

$$r_{k+1}^2 \leq \tilde{\theta}_0^{\mathcal{O}(k^{1-\gamma})} r_0^2 + \mathcal{O} \left(\frac{1}{k^\gamma} \right) \approx \mathcal{O} \left(\frac{1}{k^\gamma} \right).$$

If $\gamma = \frac{1}{2}$, then the definition of $\varphi_{1-2\gamma}(\frac{k}{2})$ provides that:

$$r_{k+1}^2 \leq \tilde{\theta}_0^{\mathcal{O}(\sqrt{k})} r_0^2 + \tilde{\theta}_0^{\mathcal{O}(\sqrt{k})} \mathcal{O}(\ln k) + \mathcal{O} \left(\frac{1}{\sqrt{k}} \right) \approx \mathcal{O} \left(\frac{1}{\sqrt{k}} \right).$$

When $\gamma \in (\frac{1}{2}, 1)$, it is obvious that $\varphi_{1-2\gamma}(\frac{k}{2}) \leq \frac{1}{2\gamma-1}$ and therefore the order of the convergence rate changes into:

$$r_{k+1}^2 \leq \tilde{\theta}_0^{\mathcal{O}(k^{1-\gamma})} [r_0^2 + \mathcal{O}(1)] + \mathcal{O} \left(\frac{1}{k^\gamma} \right) \approx \mathcal{O} \left(\frac{1}{k^\gamma} \right).$$

(ii) Lastly, if $\gamma = 1$, by using $\tilde{\theta}_0^{\ln k+1} \leq \left(\frac{1}{k}\right)^{\ln \frac{1}{\tilde{\theta}_0}}$ we obtain the second part of our result.