

Method with Batching for Stochastic Finite-Sum Variational Inequalities in Non-Euclidean Setting*

Alexander Pichugin¹, Maksim Pechin¹, Aleksandr Beznosikov^{1,2,3}, Vasilii Novitskii¹, and Alexander Gasnikov^{3,1,2}

¹ Moscow Institute of Physics and Technology, Moscow, Russia,

² Ivannikov Institute for System Programming of RAS, Moscow, Russia

³ Innopolis University, Innopolis, Russia

Abstract. Variational inequalities are a universal optimization paradigm that incorporate classical minimization and saddle point problems. Nowadays more and more tasks require to consider stochastic formulations of optimization problems. In this paper, we present an analysis of a method that gives optimal convergence estimates for monotone stochastic finite-sum variational inequalities. In contrast to the previous works, our method supports batching, does not lose the oracle complexity optimality and uses an arbitrary Bregman distance to take into account geometry of the problem. Paper provides experimental confirmation to algorithm's effectiveness.

Keywords: stochastic optimization, variational inequalities, finite-sum problems, batching, Bregman distance

1 Introduction

In this paper, we consider the following variational inequality (VI) problem:

$$\text{Find } x^* \in \mathcal{X} \text{ such that } \langle F(x^*), x - x^* \rangle + g(x) - g(x^*) \geq 0, \text{ for all } x \in \mathcal{X}, \quad (1)$$

where F is some operator, g is a proper convex lower semicontinuous function with domain $\text{dom } g$. Such problem is the classical formulation of a variational inequality [46]. Moreover, we use the composite scheme for which g is responsible. These kinds of variational inequalities have a wide usage.

Application of variational inequalities. Variational inequalities found significant implementations across various disciplines. In particular, they are used for classical problems such as equilibrium theory, games and economics

*The work of A. Pichugin and M. Pechin was supported by a grant for research centers in the field of artificial intelligence, provided by the Analytical Center for the Government of the Russian Federation in accordance with the subsidy agreement (agreement identifier 000000D730321P5Q0002) and the agreement with the Moscow Institute of Physics and Technology dated November 1, 2021 No. 70-2021-00138.

[37,20,16,43]. Later, the machine learning community discovered a wide range of implications such as reinforcement learning [38,21], GANs (notably, Wasserstein GANs [6]), adversarial training [30], supervised learning, unsupervised learning, discriminative clustering [22], matrix factorization [7] and robust optimization [8]. Some of these problems specifically require to utilize a stochastic setup of variational inequalities [30,44,9]. Also stochastic variational inequalities are an instrumental in addressing several fundamental challenges, including balance problems in finances, management, and engineering contexts [45,47,53].

Our contribution and related works. The simplest method for solving variational inequalities is the classical Gradient method. Unfortunately, it converges only for the cases of strongly monotone F , in the non-strongly monotone case convergence is lacking [14]. A big improvement in algorithms for VIs happened when the Extragradient method was proposed [25]. It used the idea of the extrapolation for the Gradient method. Remarkably, it solved the convergence issue of the Gradient method. Another big step was done in [42], where the Optimistic method was proposed. It is worth noticing a significant difference in these two methods, as Extragradient required us to call the oracle twice, while the algorithm from [42] is a single-call scheme. Later, A. Nemirovsky proposed exploiting the Bregman divergence for variational inequalities [33]. This approach allowed to take into account generalized geometry that could be non-Euclidean. As the result, Mirror Prox algorithm was created. Next, there was a number of stochastic modifications of Extragradient and Mirror Prox.

Juditsky et al. [24] for the first time considered the stochastic version of Mirror Prox. This paper examines the general stochastic case with bounded variance. Another approach would be to consider a frequently used in practice finite-sum stochastic setup of the problem, where we can achieve better convergence rate due to its specifics. In particular, one of the earliest solution only for the saddle point problem in the Euclidean setup was considered by B. Paliannappan and F. Bach [39]. They based their method on the SVRG (Stochastic Variance reduction) algorithm [23] and added the Catalyst envelope acceleration. As mentioned, the result was made for strongly convex-concave saddle point problems, which corresponds to strongly monotone variational inequations. Other algorithms utilizing the variance reduction technique (VR) for solving finite-sum variational inequalities with strongly monotone operators with better convergence rate are L-SVRGDA [10] and SVRE [13]. The idea of SVRE was continued in [51] with the introduction of VR-AGDA. Here results were obtained in Polyak-Lojasiewicz conditions. Moreover, in [48], strongly convex-concave saddle point problems were considered. In this work, methods were designed with the Catalyst envelope acceleration. Achieved rate is comparable to [39]. In addition, it is worth mentioning the work [12]. This paper deals with the Bregman setup and adapts the variance reduction technique. Unfortunately, this solution is limited by additional assumptions on the operator F and by considering the matrix games setup. Either papers that fruitfully implement VR techniques suitable for finite-sum setup are [15,5,4]; among others in the non-Euclidean case [28]. Malitsky

and Alacaoglu in [1] applied variance reduction to Mirror Prox using double loop SVRG technique [5].

Works above have one common problem – proposed algorithms are sensitive to the size of mini-batches, apart from paper [24]. However, convergence rate in [24] does not reach the lower bound [19], because the general case is considered. There are papers solve this problem for the strongly monotone operators [26] and for the monotone operators [40] in the Euclidean setup .

We compare the mentioned algorithms and their convergence results in Table 1.

In the present paper, our approach is a natural extension of the results from [40]. On the contrary, here we consider an arbitrary non-Euclidean space. Correspondingly, our method is based on the following ideas: optimistic scheme [42] with negative momentum [27] to deal with variance reduction and avoid double update like in Extragradient; on double-loop VR scheme to taking into account the Bregman divergence [1]. The oracle complexity of the new method is independent of batching, as long as the size of the batch does not exceed the square root of the full sample size.

2 Main part

In this paper, we assume that $\mathcal{X} \subseteq \mathbb{R}^n$ is a normed vector space with a dual space \mathcal{X}^* and primal-dual norm pair $\|\cdot\|$ and $\|\cdot\|_*$. As $\|\cdot\|$, we are using the norm $\|\cdot\|_p$ for $p \in [1, 2]$.

We expect a stochastic nature of a problem: $F(x) = \mathbb{E}_{\xi \sim \mathcal{D}}[F_\xi(x)]$. As noted above, due to the fact that distribution \mathcal{D} is non-trivial or even unknown, which appears in a number of applications, we can represent this as finite-sum approximation:

$$F(x) = \frac{1}{M} \sum_{m=1}^M F_m(x), \quad (2)$$

using the Monte-Carlo approach [29]. In machine learning problems, the term “empirical risk” is often encountered. Note that calls of the full operator are expensive in practice. Thus, in order to avoid frequent computing operator F , one can use calls of single F_m or mini-batches of them.

Desire to work in arbitrary geometry requires us to introduce an alternative way to measure distance on the space rather than default Euclidean distance.

Definition 1. *Let $h : \mathcal{X} \rightarrow \mathbb{R} \cup \{+\infty\}$ be a convex lower semicontinuous function such that $\text{dom } g \subseteq \text{dom } h$ and h is differentiable on $\text{dom } h$, h is 1-strongly convex on $\text{dom } g$. We define the Bregman distance $V : \text{dom } g \times \text{dom } \partial h \rightarrow \mathbb{R}_+$ associated with h by*

$$V(u, v) = h(u) - h(v) - \langle \nabla h(v), u - v \rangle.$$

Here we used the following definition.

Table 1: Summary complexities for finding an ε -solution for monotone stochastic (finite-sum) variational inequality (1)+(2) with Lipschitz operators. Convergence is measured by gap function. *Notation:* L and L_2 are Lipschitz constants for F and F_m in terms of $\|\cdot\|_*$ and $\|\cdot\|_2$ norms, respectively (see Assumptions 3a, 3b), M = size of dataset, b = batch size per iteration.

Reference	Complexity	Non-Euclidean
Nemirovski et al. [33] ⁽¹⁾	$\mathcal{O}\left(M\frac{L}{\varepsilon}\right)$	✓
Juditsky et al. [24] ⁽⁴⁾	$\mathcal{O}\left(\frac{L}{\varepsilon} + \frac{1}{\varepsilon^2}\right)$	✓
Palaniappan & Bach [39] ^(2,3,4)	$\tilde{\mathcal{O}}\left(b\frac{L_2^2}{\varepsilon^2}\right)$	
Palaniappan & Bach [39] ^(2,3,4)	$\tilde{\mathcal{O}}\left(\sqrt{bM}\frac{L_2}{\varepsilon}\right)$	
Chavdarova et al. [13] ⁽³⁾	$\tilde{\mathcal{O}}\left(\frac{bL_2^2}{\varepsilon^2}\right)$	
Carmon et al. [12] ^(2,4)	$\tilde{\mathcal{O}}\left(\sqrt{bM}\frac{L}{\varepsilon}\right)$	✓
Yang et al. [51] ^(2,3,4)	$\tilde{\mathcal{O}}\left(b^{\frac{1}{3}}M^{\frac{2}{3}}\frac{L_2^3}{\varepsilon^3}\right)$	
Alacaoglu & Malitsky [1]	$\mathcal{O}\left(\sqrt{bM}\frac{L}{\varepsilon}\right)$	✓
Tomini et al. [48] ^(2,3,4)	$\tilde{\mathcal{O}}\left(\sqrt{bM}\frac{L_2}{\varepsilon}\right)$	
Beznosikov et al. [10] ^(3,4)	$\tilde{\mathcal{O}}\left(b\frac{L_2^2}{\varepsilon^2}\right)$	
Kovalev et al. [26] ^(3,4)	$\tilde{\mathcal{O}}\left(\sqrt{M}\frac{L_2}{\varepsilon}\right)$	
Pichugin et al. [40] ⁽⁵⁾	$\tilde{\mathcal{O}}\left(\sqrt{M}\frac{L_2}{\varepsilon}\right)$	
This paper ⁽⁵⁾	$\mathcal{O}\left(\sqrt{M}\frac{L_2}{\varepsilon}\right)$	✓
This paper ⁽⁵⁾	$\tilde{\mathcal{O}}\left(\sqrt{M}\frac{L}{\varepsilon}\right)$	✓
Han et al.(lower bounds) [19]	$\Omega\left(\sqrt{M}\frac{L_2}{\varepsilon}\right)$	

(1) deterministic methods, similar results were also obtained in [49,32],

(2) for saddle point problems only,

(3) only for μ -strongly monotone operators. To compare algorithms, we apply regularization trick,

(4) only for bounded domain,

(5) for $b \leq \sqrt{M}$.

Definition 2. *Function f is 1-strongly convex if it can be lower bounded by a quadratic function of the form:*

$$f(y) \geq f(x) + \langle y - x, \nabla f(x) \rangle + \frac{1}{2} \|y - x\|^2$$

for all x and y in $\text{dom } f$.

We also provide a couple of examples of h and their corresponding V :

- For $h(x) = \frac{1}{2} \|x\|_2^2$ we have $V(x, y) = \frac{1}{2} \|x - y\|_2^2$.
- The entropy function

$$h(x) = \sum_{i=1}^n x_i \log x_i \tag{3}$$

in a probabilistic simplex

$$\Delta^n = \{x \in \mathbb{R}^n \mid x_i \geq 0 \sum_{i=1}^n x_i = 1\} \tag{4}$$

generates KL-divergence

$$V(x, y) = \sum_{i=1}^n x_i \log \frac{x_i}{y_i}.$$

Let us highlight the remarkable property of the Bregman distance which we use for the analysis of convergence of our algorithm. Since h is 1-strongly convex with respect to norm $\|\cdot\|$, we have for all $u, v \in \mathcal{X}$

$$V(u, v) \geq \frac{1}{2} \|u - v\|^2. \tag{5}$$

To further analyze the problem, we introduce the following assumptions:

Assumption 1 *The solution (maybe not unique) for the problem (1)+(2) exists.*

Assumption 2 *The operator F is monotone, i.e. for all $u, v \in \mathcal{X}$ we have*

$$\langle F(u) - F(v); u - v \rangle \geq 0.$$

Assumption 3a *The operator F is L_2 -Lipschitz and for all $m \in [1, M]$ F_m is $L_{2,m}$ -Lipschitz, i.e. for all $u, v \in \mathcal{X}$ we have*

$$\begin{aligned} \|F(u) - F(v)\|_2 &\leq L_2 \|u - v\|_2, \\ \|F_m(u) - F_m(v)\|_2 &\leq L_{2,m} \|u - v\|_2. \end{aligned}$$

In addition, we define \bar{L}_2 such that $\bar{L}_2^2 = \frac{1}{M} \sum_{m=1}^M L_{2,m}^2$. With this notation, for all $u, v \in \mathcal{X}$ we have

$$\frac{1}{M} \sum_{m=1}^M \|F_m(u) - F_m(v)\|_2^2 \leq \bar{L}_2^2 \|u - v\|_2^2.$$

Assumption 3b *The operator F is L -Lipschitz and for all $m \in [1, M]$ F_m is L -Lipschitz, i.e. for all $u, v \in \mathcal{X}$ we have*

$$\begin{aligned} \|F(u) - F(v)\|_* &\leq L\|u - v\|, \\ \|F_m(u) - F_m(v)\|_* &\leq L\|u - v\|. \end{aligned}$$

Note that $\|\cdot\| \leq \|\cdot\|_2$ for all $p \in [1, 2]$. Thus, we can state that $L \leq \bar{L}_2$. In addition, we would like to point out that these meanings often differ dramatically.

2.1 Algorithm

Now let us state the following algorithm:

Algorithm 1 Optimistic Method with Momentum and Batching

- 1: **Parameters:** stepsize $\eta > 0$, momentum $\gamma > 0$, probability $p \in (0; 1)$, batch size $b \in \{1, \dots, M\}$, number of iterations K
 - 2: **Initialization:** choose $x_0^{-1} = x_0^0 = w_0^0 = x_{-1}^k = w_{-1}^k \in \mathcal{X}$ for all $k \in [-1, K - 1]$
 - 3: **for** $s = 0, 1, \dots, S - 1$ **do**
 - 4: **for** $k = 0, 1, \dots, K - 1$ **do**
 - 5: Sample j_1^k, \dots, j_b^k independently from $\{1, \dots, M\}$ uniformly at random
 - 6: $B^k = \{j_1^k, \dots, j_b^k\}$
 - 7: $\Delta^k = \frac{1}{b} \sum_{j \in B^k} (F_j(x_s^k) - F_j(w_s) + (F_j(x_s^k) - F_j(x_s^{k-1}))) + F(w_s)$
 - 8: $x_s^{k+1} = \operatorname{argmin}_{x \in \mathcal{X}} \left(g(x) + \frac{1}{\eta}(1 - \gamma)V(x, x_s^k) + \frac{1}{\eta} \cdot \gamma V(x, \bar{w}_s) + \langle \Delta^k, x \rangle \right)$
 - 9: **end for**
 - 10: $w_{s+1} = \frac{1}{K} \sum_{k=1}^K x_s^k$
 - 11: $\nabla h(\bar{w}_{s+1}) = \frac{1}{K} \sum_{k=1}^K \nabla h(x_s^k)$
 - 12: $x_{s+1}^0 = x_s^K$
 - 13: $x_{s+1}^{-1} = x_s^{K-1}$
 - 14: **end for**
-

Note that in line 7 of Algorithm 1 we use the variance reduction technique in $(F_j(x_s^k) - F_j(w_s) + F(w_s))$ part. The key difference from Stochastic Gradient Descent (SGD) is that instead of using $g_k = F_j(x^k)$, variance reduced methods use the approximation $g_k = F(w^k) + F_j(x^k) - F_j(w^k)$. This helps to decrease

"variance" $\mathbb{E}\|g_k - F(x_k)\|^2$ comparing to SGD in the case of "good" choice of ω^k [23,3,1]. We also add batching in lines 5 and 7.

In addition, in line 7 of Algorithm 1 we use $(F_j(x_s^k) - F_j(x_s^{k-1}))$ term to implement so-called optimistic scheme, slightly different from the original option [42]. Our update is a modification of Forward-Reflected-Backward approach [31], where we set $\alpha \equiv 1$ in $x^{k+1} = \eta F_j(x^k) + \eta\alpha [F_j(x^k) - F_j(x^{k-1})]$.

While for the minimization problems it is usual to apply positive (heavy-ball) momentum [41], the opposite approach turns out to be suitable for variational inequalities. This effect was noticed earlier [18,52,2] and appeared now in the theory of stochastic methods for VIs. Hence, in line 8 we also apply the negative momentum in $(-\frac{1}{\eta}\gamma V(x, x_s^k) + \frac{1}{\eta} \cdot \gamma V(x, \bar{w}_s))$ part. For illustration, in the Euclidean case the negative momentum would have looked like $\gamma(w^k - x^k)$.

The snapshot point updates in line 10 similar to [1] and SVRG [23]. Due to the Bregman setup is that we have the additional point \bar{w}_{s+1} that averages in the dual space.

Now we are ready to proof convergence of Algorithm 1.

2.2 Analysis

In order to calculate a convergence rate, we use the gap function [33,24] as a standard convergence criterion for such problems:

$$\text{Gap}(z) := \sup_{u \in \mathcal{C}} [\langle F(u), z - u \rangle + g(x) - g(u)]. \quad (6)$$

Here \mathcal{C} is a compact subset of \mathcal{X} used to handle the case if $\text{dom } g$ is unbounded [see Lemma 1 from [36]].

To begin working with update, we propose the following lemmas.

Lemma 1. [See Lemma 3.2 from [1]]

Let g be a proper convex lower semicontinuous function. Denote

$$x^\dagger = \underset{x \in \mathcal{C}}{\text{argmin}} (g(x) + \langle u, x \rangle + \gamma V(x, x_1) + (1 - \gamma)V(x, x_2)).$$

Then, for all x it delivers

$$\begin{aligned} g(x) - g(x^\dagger) + \langle u, x - x^\dagger \rangle \\ \geq V(x, x^\dagger) + \gamma(V(x^\dagger, x_1) - V(x, x_1)) + (1 - \gamma)(V(x^\dagger, x_2) - V(x, x_2)). \end{aligned}$$

Proof. The proof of this lemma directly follows from the first order optimality condition [11].

Hereafter, in the next lemma, we estimate the variance of Δ^k .

Lemma 2a. [See Lemma 2 from [40]] For any step s from 0 to $S - 1$ and for any k from 0 to $K - 1$ of Algorithm 1 under the Assumption 3a the following inequality holds:

$$\mathbb{E} \left[\|\Delta^k - \mathbb{E}_k [\Delta^k]\|_2^2 \right] \leq \frac{2\bar{L}_2^2}{b} \mathbb{E} \left[\|x_s^k - w_s\|^2 + \|x_s^k - x_s^{k-1}\|^2 \right],$$

where $\mathbb{E}_k [\Delta^k]$ is equal to

$$\mathbb{E}_k [\Delta^k] = 2F(x_s^k) - F(x_s^{k-1}). \quad (7)$$

Proof. We start from line 7 of Algorithm 1,

$$\begin{aligned} \mathbb{E}_k \left[\|\Delta^k - \mathbb{E}_k [\Delta^k]\|_2^2 \right] &= \mathbb{E}_k \left[\left\| \frac{1}{b} \sum_{j \in B^k} (F_j(x_s^k) - F_j(w_s) + (F_j(x_s^k) - F_j(x_s^{k-1}))) \right. \right. \\ &\quad \left. \left. + F(w_s) - (2F(x_s^k) - F(x_s^{k-1})) \right\|_2^2 \right]. \end{aligned}$$

With Cauchy–Schwarz inequality (43), we have

$$\begin{aligned} &\mathbb{E}_k \left[\|\Delta^k - \mathbb{E}_k [\Delta^k]\|_2^2 \right] \\ &\leq 2\mathbb{E}_k \left[\left\| \frac{1}{b} \sum_{j \in B^k} (F_j(x_s^k) - F_j(w_s)) - (F(x_s^k) - F(w_s)) \right\|_2^2 \right] \\ &\quad + 2\mathbb{E}_k \left[\left\| \frac{1}{b} \sum_{j \in B^k} (F_j(x_s^k) - F_j(x_s^{k-1})) - (F(x_s^k) - F(x_s^{k-1})) \right\|_2^2 \right]. \end{aligned}$$

Using that we choose j_1^k, \dots, j_s^k in B^k independently and uniformly, one can note

$$\begin{aligned} &\mathbb{E}_k \left[\left\langle (F_j(x_s^k) - F_j(w_s)) - (F(x_s^k) - F(w_s)), (F_j(x_s^k) - F_j(w_s)) \right. \right. \\ &\quad \left. \left. - (F(x_s^k) - F(w_s)) \right\rangle \right] \\ &= \mathbb{E}_k \left[\left\langle \mathbb{E}_{j_i^k} \left[(F_{j_i^k}(x_s^k) - F_{j_i^k}(w_s)) - (F(x_s^k) - F(w_s)) \right], \right. \right. \\ &\quad \left. \left. \mathbb{E}_{j_i^k} \left[(F_{j_i^k}(x_s^k) - F_{j_i^k}(w_s)) - (F(x_s^k) - F(w_s)) \right] \right\rangle \right] \\ &= 0. \end{aligned}$$

Hence, we get

$$\begin{aligned} &\mathbb{E}_k \left[\|\Delta^k - \mathbb{E}_k [\Delta^k]\|_2^2 \right] \\ &\leq 2\mathbb{E}_k \left[\sum_{j \in B^k} \frac{1}{b^2} \|(F_j(x_s^k) - F_j(w_s)) - (F(x_s^k) - F(w_s))\|_2^2 \right] \\ &\quad + 2\mathbb{E}_k \left[\sum_{j \in B^k} \frac{1}{b^2} \|(F_j(x_s^k) - F_j(x_s^{k-1})) - (F(x_s^k) - F(x_s^{k-1}))\|_2^2 \right] \end{aligned}$$

$$\begin{aligned}
 &= \frac{2}{b^2} \mathbb{E}_k \left[\sum_{j \in B^k} \|(F_j(x_s^k) - F_j(w_s)) - (F(x_s^k) - F(w_s))\|_2^2 \right] \\
 &\quad + \frac{2}{b^2} \mathbb{E}_k \left[\sum_{j \in B^k} \|(F_j(x_s^k) - F_j(x_s^{k-1})) - (F(x_s^k) - F(x_s^{k-1}))\|_2^2 \right] \\
 &\leq \frac{2}{b^2} \mathbb{E}_k \left[\sum_{j \in B^k} \|F_j(x_s^k) - F_j(w_s)\|_2^2 \right] + \frac{2}{b^2} \mathbb{E}_k \left[\sum_{j \in B^k} \|F_j(x_s^k) - F_j(x_s^{k-1})\|_2^2 \right].
 \end{aligned}$$

In the last step, we used the fact that $\mathbb{E}\|X - \mathbb{E}X\|^2 = \mathbb{E}\|X\|^2 - \|\mathbb{E}X\|^2$. Next, we again take into account that j_1^k, \dots, j_s^k in B^k are chosen uniformly,

$$\begin{aligned}
 &\mathbb{E}_k \left[\|\Delta^k - \mathbb{E}_k[\Delta^k]\|_2^2 \right] \\
 &\leq \frac{2}{b} \mathbb{E}_k \left[\mathbb{E}_{j \sim \text{u.a.r. } \{1, \dots, M\}} \left[\|F_j(x_s^k) - F_j(w_s)\|_2^2 + \|F_j(x_s^k) - F_j(x_s^{k-1})\|_2^2 \right] \right] \\
 &= \frac{2}{Mb} \sum_{j=1}^M \left(\|F_j(x_s^k) - F_j(w_s)\|_2^2 + \|F_j(x_s^k) - F_j(x_s^{k-1})\|_2^2 \right).
 \end{aligned}$$

We take the full expectation of both parts:

$$\mathbb{E} \left[\|\Delta^k - \mathbb{E}_k[\Delta^k]\|_2^2 \right] \leq \frac{2}{Mb} \mathbb{E} \left[\sum_{j=1}^M \left(\|F_j(x_s^k) - F_j(w_s)\|_2^2 + \|F_j(x_s^k) - F_j(x_s^{k-1})\|_2^2 \right) \right].$$

Using L_2 -Lipschitzness of F (Assumption 3a) and the fact that $\|\cdot\|_2 \leq \|\cdot\|$, we can rewrite it as

$$\begin{aligned}
 \mathbb{E} \left[\|\Delta^k - \mathbb{E}_k[\Delta^k]\|_2^2 \right] &\leq \frac{2\bar{L}_2^2}{b} \mathbb{E} \left[\|x_s^k - w_s\|_2^2 + \|x_s^k - x_s^{k-1}\|_2^2 \right] \\
 &\leq \frac{2\bar{L}_2^2}{b} \mathbb{E} \left[\|x_s^k - w_s\|^2 + \|x_s^k - x_s^{k-1}\|^2 \right].
 \end{aligned}$$

This finishes the proof.

Lemma 2b. *For any step s from 0 to $S-1$ and for any k from 0 to $K-1$ of Algorithm 1 under the Assumption 3b the following inequality holds:*

$$\mathbb{E} \left[\|\Delta^k - \mathbb{E}_k[\Delta^k]\|_*^2 \right] \leq \frac{2(1 + C \ln n)\bar{L}^2}{b} \mathbb{E} \left[\|x_s^k - w_s\|^2 + \|x_s^k - x_s^{k-1}\|^2 \right],$$

where $\mathbb{E}_k[\Delta^k]$ is equal to

$$\mathbb{E}_k[\Delta^k] = 2F(x_s^k) - F(x_s^{k-1}).$$

Proof. Similar to the proof of Lemma 2a we start from line 7 of Algorithm 1,

$$\begin{aligned} & \mathbb{E} \left[\left\| \Delta^k - \mathbb{E} [\Delta^k] \right\|_*^2 \right] \\ &= \mathbb{E}_k \left[\left\| \frac{1}{b} \sum_{j \in B^k} (F_j(x_s^k) - F_j(w_s) + (F_j(x_s^k) - F_j(x_s^{k-1}))) + F(w_s) \right. \right. \\ & \quad \left. \left. - (2F(x_s^k) - F(x_s^{k-1})) \right\|_*^2 \right]. \end{aligned}$$

For shortness we introduce the stochastic variables

$$\xi^j = (F_j(x_s^k) - F_j(w_s) + (F_j(x_s^k) - F_j(x_s^{k-1}))) + F(w_s) - (2F(x_s^k) - F(x_s^{k-1})),$$

and declare the following properties of them:

1. $\mathbb{E} \xi^j = 0 \quad \forall j \in B^k$,
2. ξ^j is independent,
3. $\|\xi^j\| < \sigma = \sqrt{2\bar{L} (\|x_s^k - w_s\|^2 + \|x_s^k - x_s^{k-1}\|^2)}$. (Here we used the Assumption 3b).

We also bring in the sequence $\{u^j\}_{j \in B^k}$ defined as

$$u^{j+1} = \operatorname{argmin}_{y \in B_1(0)} \left(V(u, u^j) + \left\langle \frac{\Omega}{\sigma\sqrt{b}} \xi^j, y \right\rangle \right) \quad (8)$$

with a starting point $u^1 = 0$. In addition, we introduce $\max_{u \in B_1(0)} V(u, u^1) =$

$$\max_{u \in B_1(0)} V(u, 0) = \frac{\Omega^2}{2}.$$

With these denotations and the definition of $\|\cdot\|_*$, it is correct that

$$\left\| \sum_{j \in B^k} \frac{\Omega}{\sigma\sqrt{b}} \xi^j \right\|_* = \max_{u \in B_1(0)} \left\langle \sum_{j \in B^k} \frac{\Omega}{\sigma\sqrt{b}} \xi^j, u \right\rangle.$$

To estimate the right side, we can write the optimality condition [11] for (8):

$$\left\langle \nabla h(u^{j+1}) - \nabla h(u^j) - \frac{\Omega}{\sigma\sqrt{b}} \xi^{j+1}, y - u^{j+1} \right\rangle \geq 0.$$

We apply the three point identity (45) to obtain

$$V(u, u^j) - V(y, u^{j+1}) - V(u^{j+1}, u^j) - \left\langle \frac{\Omega}{\sigma\sqrt{b}} \xi^{j+1}, y - u^{j+1} \right\rangle \geq 0. \quad (9)$$

Using Young's inequality (44) and the inequality (5), we achieve

$$\left\langle \frac{\Omega}{\sigma\sqrt{b}} \xi^{j+1}, y - u^{j+1} \right\rangle = \left\langle \frac{\Omega}{\sigma\sqrt{b}} \xi^{j+1}, y - u^j \right\rangle + \left\langle \frac{\Omega}{\sigma\sqrt{b}} \xi^{j+1}, u^j - u^{j+1} \right\rangle$$

$$\begin{aligned}
 &\geq \left\langle \frac{\Omega}{\sigma\sqrt{b}} \xi^{j+1}, y - u^j \right\rangle + \frac{1}{2} \left\| \frac{\Omega}{\sigma\sqrt{b}} \xi^{j+1} \right\|_*^2 - \frac{1}{2} \|u^{j+1} - u^j\|^2 \\
 &\geq \left\langle \frac{\Omega}{\sigma\sqrt{b}} \xi^{j+1}, y - u^j \right\rangle + \frac{1}{2} \left\| \frac{\Omega}{\sigma\sqrt{b}} \xi^{j+1} \right\|_*^2 - V(u^{j+1}, u^j).
 \end{aligned} \tag{10}$$

Combining (9) with (10), we get

$$\left\langle \frac{\Omega}{\sigma\sqrt{b}} \xi^{j+1}, y \right\rangle \leq V(y, u^j) - V(y, u^{j+1}) + \left\langle \frac{\Omega}{\sigma\sqrt{b}} \xi^{j+1}, u^j \right\rangle + \frac{1}{2} \left\| \frac{\Omega}{\sigma\sqrt{b}} \xi^{j+1} \right\|_*^2.$$

Now we take sum for all j to estimate

$$\begin{aligned}
 \left\| \sum_{j \in B^k} \frac{\Omega}{\sigma\sqrt{b}} \xi^j \right\|_* &\leq \max_{u \in B_1(0)} V(u, u^1) + \frac{1}{2} \sum_{j \in B^k} \left\| \frac{\Omega}{\sigma\sqrt{b}} \xi^j \right\|_*^2 + \sum_{j \in B^k} \left\langle \frac{\Omega}{\sigma\sqrt{b}} \xi^j, u^j \right\rangle \\
 &\leq \frac{\Omega^2}{2} + \frac{1}{2} \sum_{j \in B^k} \left(\frac{\Omega}{\sqrt{b}} \right)^2 + \sum_{j \in B^k} \left\langle \frac{\Omega}{\sigma\sqrt{b}} \xi^j, u^j \right\rangle.
 \end{aligned}$$

Here we used the third property of ξ^j . Then,

$$\left\| \sum_{j \in B^k} \xi^j \right\|_* \leq \sigma\Omega\sqrt{b} + \sum_{j \in B^k} \langle \xi^j, u^j \rangle.$$

After taking square of both sides, one can get

$$\left\| \sum_{j=1}^b \xi^j \right\|_*^2 \leq \sigma^2 \Omega^2 b + \sum_{j \in B^k} \langle \xi^j, u^j \rangle^2 + 2\sigma\Omega\sqrt{b} \sum_{j \in B^k} \langle \xi^j, u^j \rangle.$$

Now we take the expectation of both sides

$$\mathbb{E} \left\| \sum_{j \in B^k} \xi^j \right\|_*^2 \leq \mathbb{E} [\sigma^2 \Omega^2 b] + \mathbb{E} \sum_{j \in B^k} \langle \xi^j, u^j \rangle^2 + 2\Omega\sqrt{b} \sum_{j \in B^k} \mathbb{E} \sigma \langle \xi^j, u^j \rangle. \tag{11}$$

We use the fact that for all j from B^K ξ^j and u^j are independent to get

$$\sum_{j \in B^k} \mathbb{E} \sigma \langle \xi^j, u^j \rangle = 0. \tag{12}$$

Combining (11) and (12), one can achieve

$$\mathbb{E} \left\| \sum_{j \in B^k} \xi^j \right\|_*^2 \leq \mathbb{E} [\sigma^2 \Omega^2 b] + \mathbb{E} \sum_{j \in B^k} \langle \xi^j, u^j \rangle^2.$$

Now what is left is to apply the Young's inequality (44) and the first property of ξ^j to get

$$\begin{aligned} \mathbb{E} \left\| \sum_{j \in B^k} \xi^j \right\|_*^2 &\leq \mathbb{E} [\sigma^2 \Omega^2 b] + \mathbb{E} \sum_{j \in B^k} [\|\xi^j\|_*^2 \cdot \|u^j\|^2] \\ &\leq \mathbb{E} [\sigma^2 \Omega^2 b] + \mathbb{E} [b\sigma^2] \\ &\leq \mathbb{E} [\sigma^2 b(\Omega^2 + 1)]. \end{aligned}$$

In order to perform the next step, we need to evaluate Ω . This kind of result may be found in Table 2 of [17]. For convinience, it is represented below.

$Q = B_p^n(1)$	$1 \leq p \leq a$	$a \leq p \leq 2$	$2 \leq p \leq \infty$
$\ \cdot\ $	$\ \cdot\ _1$	$\ \cdot\ _p$	$\ \cdot\ _2$
$h(x)$	$\frac{1}{2(a-1)} \ x\ _a^2$	$\frac{1}{2(p-1)} \ x\ _p^2$	$\frac{1}{2} \ x\ _2^2$
Ω^2	$\mathcal{O}(\ln n)$	$\mathcal{O}((p-1)^{-1})$	$\mathcal{O}(n^{\frac{1}{2}-\frac{1}{p}})$

Table 2: Examples of prox-functions for ball-shaped sets Q in various norms

In particular, we can choose $h(x) = \frac{1}{2(a-1)} \|x\|_a^2$ for $a = \frac{2 \ln n}{2 \ln n - 1}$. At this case for $p \in [1, 2]$ it is also true that $\Omega^2 = \mathcal{O}(\ln n)$, where n is dimensionality of \mathcal{X} . After plugging constant \sqrt{C} in the definition of $\mathcal{O}(\cdot)$ bound, substituting $\sigma = \sqrt{2\bar{L} (\|x_s^k - w_s\|^2 + \|x_s^k - x_s^{k-1}\|^2)}$, we finally get

$$\mathbb{E} \left\| \frac{1}{b} \sum_{j \in B^k} \xi^j \right\|_*^2 \leq \frac{2\bar{L}(1 + C \ln n)}{b} \mathbb{E} [\|x_s^k - w_s\|^2 + \|x_s^k - x_s^{k-1}\|^2].$$

That allows us to finish the proof.

Lemma 3a. [see Lemma 2.4 from [1]] Let $\mathcal{F} = (\mathcal{F}_k)_{k \geq 0}$ be a filtration and (u^k) a stochastic process adapted to \mathcal{F} with $\mathbb{E}[u^{k+1} | \mathcal{F}_k] = 0$. Then for any $K \in \mathbb{N}$, $x^0 \in X$ and any compact set $\mathcal{C} \subseteq X$

$$\mathbb{E} \left[\max_{x \in \mathcal{C}} \sum_{k=0}^{K-1} \langle u^{k+1}, x \rangle \right] \leq \max_{x \in \mathcal{C}} \frac{1}{2} \|x^0 - x\|_2^2 + \frac{1}{2} \sum_{k=0}^{K-1} \mathbb{E} \|u^{k+1}\|_2^2.$$

Proof. See proof from Lemma 3b in the case of substitution $V(x, y) = \frac{1}{2} \|x - y\|^2$.

Lemma 3b. [see Lemma 3.5 from [1]] Let $\mathcal{F} = (\mathcal{F}_k)_{k \geq 0}$ be a filtration and (u^k) a stochastic process adapted to \mathcal{F} with $\mathbb{E}[u^{k+1} | \mathcal{F}_k] = 0$. Then for any $K \in \mathbb{N}$, $x^0 \in X$ and any compact set $\mathcal{C} \subseteq X$

$$\mathbb{E} \left[\max_{x \in \mathcal{C}} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \langle u^{k+1}, x \rangle \right] \leq \max_{x \in \mathcal{C}} V(x, x_0) + \frac{1}{2} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \mathbb{E} \|u^{k+1}\|_*^2$$

Proof. Let us define

$$z^{k+1} = \operatorname{argmin}_{x \in \operatorname{dom} g} \{ \langle -u^{k+1}, x \rangle + V(x, z^k) \}.$$

We use the the first order optimality condition [11] for this notation to get

$$\langle \nabla h(z^{k+1}) - \nabla h(z^k) - u^{k+1}, x - z^{k+1} \rangle \geq 0.$$

We apply the three point identity (45) to obtain

$$V(x, z^k) - V(x, z^{k+1}) - V(z^{k+1}, z^k) - \langle u^{k+1}, x - z^{k+1} \rangle \geq 0. \quad (13)$$

Applying the Young's inequality (44), we achieve

$$\begin{aligned} \langle u^{k+1}, x - z^{k+1} \rangle &= \langle u^{k+1}, x - z^k \rangle + \langle u^{k+1}, z^k - z^{k+1} \rangle \\ &\geq \langle u^{k+1}, x - z^k \rangle + \frac{1}{2} \|u^{k+1}\|_*^2 - \frac{1}{2} \|z^{k+1} - z^k\|^2 \\ &\geq \langle u^{k+1}, x - z^k \rangle + \frac{1}{2} \|u^{k+1}\|_*^2 - V(z^{k+1}, z^k). \end{aligned} \quad (14)$$

Combining (13) with (14), we get

$$\langle u^{k+1}, x \rangle \leq V(x, z^k) - V(x, z^{k+1}) + \langle u^{k+1}, z^k \rangle + \frac{1}{2} \|u^{k+1}\|_*^2.$$

We sum this inequality for all k from 0 to $K-1$, take maximum, expectation and use the fact that $\mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=1}^{K-1} \langle u^{k+1}, z^k \rangle \right] = 0$. That result finishes the proof.

Finally, we introduce a simple technical lemma.

Lemma 4. (see (25) from [1]) In the conditions of Algorithm 1 it is applied that

$$V(x, \bar{\omega}_s) - V(x_s^{k+1}, \bar{\omega}_s) = \frac{1}{K} \sum_{j=0}^{K-1} \left(V(x, x_{s-1}^j) - V(x_s^{k+1}, x_{s-1}^j) \right).$$

Proof. Proof follows from the note that for any x, y the expression $V(x, z) - V(y, z)$ is linear in terms of $\nabla h(z)$.

Now we present the theorem to state the convergence of Algorithm 1.

Theorem 1. Consider the problem (1)+(2) under Assumptions 1, 2 and 3a. Let $\{x_B^k\}$ be the sequence generated by Algorithm 1 with tuning of $\eta, \theta, \alpha, \beta, \gamma$ as follows:

$$0 < \gamma = p \leq \frac{1}{16}, \quad (15)$$

$$\eta = \min \left\{ \frac{\sqrt{\gamma b}}{8L_2}, \frac{1}{8L_2} \right\}. \quad (16)$$

Then for $x_S = \frac{1}{KS} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} x_B^k$ it holds that

$$\mathbb{E}[\text{Gap}(x_S)] \leq \frac{(2 + K\gamma)}{\eta KS} \max_{x \in \mathcal{C}} \{V(x, x_0)\}.$$

Proof. We start from using Lemma 1 after plugging parameters $u = \eta \Delta^k$, $x^\dagger = x_s^{k+1}$, $x_1 = \bar{\omega}_s$, $x_2 = x_s^k$, and get

$$\begin{aligned} V(x, x_s^{k+1}) &\leq \eta \langle \Delta^k, x - x_s^{k+1} \rangle - \gamma V(x_s^{k+1}, \bar{\omega}_s) + \gamma V(x, \bar{\omega}_s) \\ &\quad - (1 - \gamma) (V(x_s^{k+1}, x_s^k) - V(x, x_s^k)) + \eta g(x) - \eta g(x_s^{k+1}). \end{aligned}$$

Then, after applying (7) from Lemma 2a and simple algebra, we get

$$\begin{aligned} V(x, x_s^{k+1}) &\leq -\eta \langle \mathbb{E}_k [\Delta^k], x_s^{k+1} - x \rangle + \eta \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^{k+1} - x_s^k \rangle \\ &\quad + \eta \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^k - x \rangle - \gamma V(x_s^{k+1}, \bar{\omega}_s) + \gamma V(x, \bar{\omega}_s) \\ &\quad - (1 - \gamma) (V(x_s^{k+1}, x_s^k) - V(x, x_s^k)) + \eta g(x) - \eta g(x_s^{k+1}) \\ &\leq -\eta \langle F(x_s^k) + F(x_s^k) - F(x_s^{k-1}), x_s^{k+1} - x \rangle \\ &\quad + \eta \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^{k+1} - x_s^k \rangle + \eta \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^k - x \rangle \\ &\quad - \gamma V(x_s^{k+1}, \bar{\omega}_s) + \gamma V(x, \bar{\omega}_s) \\ &\quad - (1 - \gamma) (V(x_s^{k+1}, x_s^k) - V(x, x_s^k)) + \eta g(x) - \eta g(x_s^{k+1}) \\ &\leq -\eta \langle F(x_s^k) - F(x_s^{k+1}) + F(x_s^k) - F(x_s^{k-1}), x_s^{k+1} - x \rangle \\ &\quad - \eta \langle F(x_s^{k+1}), x_s^{k+1} - x \rangle + \eta \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^{k+1} - x_s^k \rangle \\ &\quad + \eta \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^k - x \rangle - \gamma V(x_s^{k+1}, \bar{\omega}_s) + \gamma V(x, \bar{\omega}_s) \\ &\quad - (1 - \gamma) (V(x_s^{k+1}, x_s^k) - V(x, x_s^k)) + \eta g(x) - \eta g(x_s^{k+1}). \end{aligned}$$

By simple rearrangements, we obtain

$$\begin{aligned} &\eta (g(x_s^{k+1}) - g(x)) + \eta \langle F(x_s^{k+1}), x_s^{k+1} - x \rangle \\ &\leq -\eta \langle F(x_s^k) - F(x_s^{k+1}) + F(x_s^k) - F(x_s^{k-1}), x_s^{k+1} - x \rangle \\ &\quad + \eta \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^{k+1} - x_s^k \rangle \\ &\quad + \eta \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^k - x \rangle - \gamma V(x_s^{k+1}, \bar{\omega}_s) + \gamma V(x, \bar{\omega}_s) \end{aligned}$$

$$\begin{aligned}
 & - (1 - \gamma)(V(x_s^{k+1}, x_s^k) - V(x, x_s^k)) - V(x, x_s^{k+1}) \\
 \leq & -\eta \langle F(x_s^k) - F(x_s^{k+1}), x_s^{k+1} - x \rangle - \eta \langle F(x_s^k) - F(x_s^{k-1}), x_s^{k+1} - x \rangle \\
 & + \eta \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^{k+1} - x_s^k \rangle \\
 & + \eta \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^k - x \rangle - \gamma V(x_s^{k+1}, \bar{\omega}_s) + \gamma V(x, \bar{\omega}_s) \\
 & - (1 - \gamma)(V(x_s^{k+1}, x_s^k) - V(x, x_s^k)) - V(x, x_s^{k+1}) \\
 \leq & -\eta \langle F(x_s^k) - F(x_s^{k+1}), x_s^{k+1} - x \rangle - \eta \langle F(x_s^k) - F(x_s^{k-1}), x_s^{k+1} - x_s^k \rangle \\
 & - \eta \langle F(x_s^k) - F(x_s^{k-1}), x_s^k - x \rangle + \eta \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^{k+1} - x_s^k \rangle \\
 & + \eta \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^k - x \rangle - \gamma V(x_s^{k+1}, \bar{\omega}_s) + \gamma V(x, \bar{\omega}_s) \\
 & - (1 - \gamma)(V(x_s^{k+1}, x_s^k) - V(x, x_s^k)) - V(x, x_s^{k+1}).
 \end{aligned}$$

Next, we apply Lemma 4 and get

$$\begin{aligned}
 & \eta(g(x_s^{k+1}) - g(x)) + \eta \langle F(x_s^{k+1}), x_s^{k+1} - x \rangle \\
 & \leq -\eta \langle F(x_s^k) - F(x_s^{k+1}), x_s^{k+1} - x \rangle - \eta \langle F(x_s^k) - F(x_s^{k-1}), x_s^{k+1} - x_s^k \rangle \\
 & \quad - \eta \langle F(x_s^k) - F(x_s^{k-1}), x_s^k - x \rangle + \eta \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^{k+1} - x_s^k \rangle \\
 & \quad + \eta \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^k - x \rangle + \frac{\gamma}{K} \sum_{j=0}^{K-1} \left(V(x, x_{s-1}^j) - V(x_s^{k+1}, x_{s-1}^j) \right) \\
 & \quad - (1 - \gamma)(V(x_s^{k+1}, x_s^k) - V(x, x_s^k)) - V(x, x_s^{k+1}) \\
 = & -\eta \langle F(x_s^k) - F(x_s^{k+1}), x_s^{k+1} - x \rangle - \eta \langle F(x_s^k) - F(x_s^{k-1}), x_s^{k+1} - x_s^k \rangle \\
 & - \eta \langle F(x_s^k) - F(x_s^{k-1}), x_s^k - x \rangle + \eta \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^{k+1} - x_s^k \rangle \\
 & + \eta \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^k - x \rangle + \frac{\gamma}{K} \sum_{j=0}^{K-1} V(x, x_{s-1}^j) - \frac{\gamma}{K} \sum_{j=0}^{K-1} V(x_s^{k+1}, x_{s-1}^j) \\
 & - (1 - \gamma)V(x_s^{k+1}, x_s^k) + (1 - \gamma)V(x, x_s^k) - V(x, x_s^{k+1}). \tag{17}
 \end{aligned}$$

Applying the fact (5) to $\|x_s^{k+1} - x_{s-1}^j\|^2$, Jensen's inequality and line 10 of Algorithm 1, we state

$$\begin{aligned}
 -\frac{\gamma}{K} \sum_{j=0}^{K-1} V(x_s^{k+1}, x_{s-1}^j) & \leq -\frac{\gamma}{2K} \sum_{j=0}^{K-1} \|x_s^{k+1} - x_{s-1}^j\|^2 \\
 & \leq -\frac{\gamma}{2} \left\| \frac{1}{K} \sum_{j=0}^{K-1} (x_s^{k+1} - x_{s-1}^j) \right\|^2 \\
 & \leq -\frac{\gamma}{2} \|x_s^{k+1} - \omega_s\|^2. \tag{18}
 \end{aligned}$$

Substituting the fact (18) to (17), we estimate

$$\eta(g(x_s^{k+1}) - g(x)) + \eta \langle F(x_s^{k+1}), x_s^{k+1} - x \rangle$$

$$\begin{aligned}
&\leq -\eta\langle F(x^k) - F(x_s^{k+1}), x_s^{k+1} - x \rangle - \eta\langle F(x_s^k) - F(x_s^{k-1}), x_s^{k+1} - x_s^k \rangle \\
&\quad - \eta\langle F(x_s^k) - F(x_s^{k-1}), x_s^k - x \rangle + \eta\langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^{k+1} - x_s^k \rangle \\
&\quad + \eta\langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^k - x \rangle + \frac{\gamma}{K} \sum_{j=0}^{K-1} V(x, x_{s-1}^j) - \frac{\gamma}{2} \|x_s^{k+1} - \omega_s\|^2 \\
&\quad - (1 - \gamma)V(x_s^{k+1}, x_s^k) + (1 - \gamma)V(x, x_s^k) - V(x, x_s^{k+1}). \tag{19}
\end{aligned}$$

For the next step we apply the Young's inequality (44), Assumption 3a, the choice of step (16) and the fact that $\|\cdot\|_2 \leq \|\cdot\|$. This allows us to achieve

$$\begin{aligned}
-\eta\langle F(x_s^k) - F(x_s^{k-1}), x_s^{k+1} - x_s^k \rangle &\leq 2\eta^2 \|F(x_s^k) - F(x_s^{k-1})\|_2^2 + \frac{1}{8} \|x_s^{k+1} - x_s^k\|_2^2 \\
&\leq \bar{L}_2^2 \eta^2 \|x_s^k - x_s^{k-1}\|_2^2 + \frac{1}{8} \|x_s^{k+1} - x_s^k\|_2^2 \\
&\leq \frac{1}{32} \|x_s^k - x_s^{k-1}\|_2^2 + \frac{1}{8} \|x_s^{k+1} - x_s^k\|_2^2 \\
&\leq \frac{1}{32} \|x_s^k - x_s^{k-1}\|^2 + \frac{1}{8} \|x_s^{k+1} - x_s^k\|^2. \tag{20}
\end{aligned}$$

We substitute (20) to (19) to get

$$\begin{aligned}
&\eta(g(x_s^{k+1}) - g(x)) + \langle F(x_s^{k+1}), x_s^{k+1} - x \rangle \\
&\quad \leq -\eta\langle F(x^k) - F(x_s^{k+1}), x_s^{k+1} - x \rangle - \eta\langle F(x_s^k) - F(x_s^{k-1}), x_s^k - x \rangle \\
&\quad + \frac{1}{32} \|x_s^k - x_s^{k-1}\|^2 + \frac{1}{4} V(x_s^{k+1}, x_s^k) + \eta\langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^{k+1} - x_s^k \rangle \\
&\quad + \eta\langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^k - x \rangle + \frac{\gamma}{K} \sum_{j=0}^{K-1} V(x, x_{s-1}^j) - \frac{\gamma}{2} \|x_s^{k+1} - \omega_s\|^2 \\
&\quad - (1 - \gamma)V(x_s^{k+1}, x_s^k) + (1 - \gamma)V(x, x_s^k) - V(x, x_s^{k+1}).
\end{aligned}$$

Next, we sum for all k from 0 to $K - 1$ and obtain

$$\begin{aligned}
&\eta \sum_{k=0}^{K-1} \left[(g(x_s^{k+1}) - g(x)) + \eta\langle F(x_s^{k+1}), x_s^{k+1} - x \rangle \right] \\
&\quad \leq -\eta \sum_{k=0}^{K-1} \langle F(x_s^k) - F(x_s^{k+1}), x_s^{k+1} - x \rangle \\
&\quad - \eta \sum_{k=0}^{K-1} \langle F(x_s^k) - F(x_s^{k-1}), x_s^k - x \rangle \\
&\quad + \frac{1}{32} \sum_{k=0}^{K-1} \|x_s^k - x_s^{k-1}\|^2 + \frac{1}{4} \sum_{k=0}^{K-1} V(x_s^{k+1}, x_s^k) \\
&\quad + \eta \sum_{k=0}^{K-1} \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^{k+1} - x_s^k \rangle
\end{aligned}$$

$$\begin{aligned}
& + \eta \sum_{k=0}^{K-1} \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^k - x \rangle + \frac{\gamma}{K} \sum_{k=0}^{K-1} \sum_{j=0}^{K-1} V(x, x_{s-1}^j) \\
& - \frac{\gamma}{2} \sum_{k=0}^{K-1} \|x_s^{k+1} - \omega_s\|^2 - (1-\gamma) \sum_{k=0}^{K-1} V(x_s^{k+1}, x_s^k) \\
& + (1-\gamma) \sum_{k=0}^{K-1} V(x, x_s^k) - \sum_{k=0}^{K-1} V(x, x_s^{k+1}) \\
\leq & \eta \langle F(x_s^{-1}) - F(x_s^0), x_s^0 - x \rangle - \eta \langle F(x_s^{K-1}) - F(x_s^K), x_s^K - x \rangle \\
& + \frac{1}{32} \sum_{k=0}^{K-1} \|x_s^k - x_s^{k-1}\|^2 + \frac{1}{4} \sum_{k=0}^{K-1} V(x_s^{k+1}, x_s^k) \\
& + \eta \sum_{k=0}^{K-1} \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^{k+1} - x_s^k \rangle \\
& + \eta \sum_{k=0}^{K-1} \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^k - x \rangle \\
& + \gamma \sum_{k=0}^{K-1} V(x, x_{s-1}^k) - \frac{\gamma}{2} \sum_{k=0}^{K-1} \|x_s^{k+1} - \omega_s\|^2 \\
& - (1-\gamma) \sum_{k=0}^{K-1} V(x_s^{k+1}, x_s^k) + (1-\gamma) \sum_{k=0}^{K-1} V(x, x_s^k) \\
& - \sum_{k=0}^{K-1} V(x, x_s^{k+1}).
\end{aligned}$$

We take sum for all s from 0 to $S-1$ and apply lines 12 and 13 of Algorithm 1 to achieve

$$\begin{aligned}
& \eta \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[(g(x_s^{k+1}) - g(x)) + \eta \langle F(x_s^{k+1}), x_s^{k+1} - x \rangle \right] \\
& \leq \eta \sum_{s=0}^{S-1} \langle F(x_s^{-1}) - F(x_s^0), x_s^0 - x \rangle \\
& \quad - \eta \sum_{s=0}^{S-1} \langle F(x_s^{K-1}) - F(x_s^K), x_s^K - x \rangle \\
& \quad + \frac{1}{32} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^k - x_s^{k-1}\|^2 + \frac{1}{4} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^{k+1}, x_s^k) \\
& \quad + \eta \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^{k+1} - x_s^k \rangle
\end{aligned}$$

$$\begin{aligned}
& + \eta \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^k - x \rangle \\
& + \gamma \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x, x_{s-1}^k) - \frac{\gamma}{2} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^{k+1} - \omega_s\|^2 \\
& - (1 - \gamma) \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^{k+1}, x_s^k) + (1 - \gamma) \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x, x_s^k) \\
& - \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x, x_s^{k+1}) \\
& = \eta \sum_{s=0}^{S-1} \langle F(x_s^{-1}) - F(x_s^0), x_s^0 - x \rangle \\
& - \eta \sum_{s=0}^{S-1} \langle F(x_{s+1}^{-1}) - F(x_{s+1}^0), x_{s+1}^0 - x \rangle \\
& + \frac{1}{32} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_B^k - x_s^{k-1}\|^2 + \frac{1}{4} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^{k+1}, x_B^k) \\
& + \eta \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^{k+1} - x_s^k \rangle \\
& + \eta \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^k - x \rangle \\
& + \gamma \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x, x_{s-1}^k) - \frac{\gamma}{2} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^{k+1} - \omega_s\|^2 \\
& - (1 - \gamma) \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^{k+1}, x_s^k) + (1 - \gamma) \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x, x_s^k) \\
& - \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x, x_s^{k+1}). \tag{21}
\end{aligned}$$

Simple algebra and line 2 from Algorithm 1 allows us to get

$$\begin{aligned}
& \eta \sum_{s=0}^{S-1} \langle F(x_s^{-1}) - F(x_s^0), x_s^0 - x \rangle - \eta \sum_{s=0}^{S-1} \langle F(x_{s+1}^{-1}) - F(x_{s+1}^0), x_{s+1}^0 - x \rangle \\
& = \eta \langle F(x_0^S) - F(x_S^{-1}), x_S^0 - x \rangle - \eta \langle F(x_0^0) - F(x_0^{-1}), x - x_0^0 \rangle \\
& = \eta \langle F(x_0^S) - F(x_S^{-1}), x_S^0 - x \rangle \tag{22}
\end{aligned}$$

Applying the Young's inequality (44), Assumption 3a, the choice of step (16) and the fact that $\|\cdot\|_2 \leq \|\cdot\|$, we get

$$\begin{aligned}
 \eta \langle F(x_S^0) - F(x_S^{-1}), x_S^0 - x \rangle &\leq \frac{\eta^2}{2} \|F(x_S^0) - F(x_S^{-1})\|_2^2 + \frac{1}{2} \|x_S^0 - x\|_2^2 \\
 &\leq \frac{\eta^2 L_2^2}{2} \|x_S^0 - x_S^{-1}\|_2^2 + \frac{1}{2} \|x_S^0 - x\|_2^2 \\
 &\leq \frac{1}{128} \|x_S^0 - x_S^{-1}\|_2^2 + \frac{1}{2} \|x_S^0 - x\|_2^2 \\
 &\leq \frac{1}{128} \|x_S^0 - x_S^{-1}\|^2 + \frac{1}{2} \|x_S^0 - x\|^2. \tag{23}
 \end{aligned}$$

Inequalities (23) and (22) allow us to turn (21) into

$$\begin{aligned}
 &\eta \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[(g(x_s^{k+1}) - g(x)) + \eta \langle F(x_s^{k+1}), x_s^{k+1} - x \rangle \right] \\
 &\leq \frac{1}{128} \|x_S^0 - x_S^{-1}\|^2 + \frac{1}{2} \|x_S^0 - x\|^2 + \frac{1}{4} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^{k+1}, x_s^k) \\
 &\quad + \frac{1}{32} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_B^k - x_s^{k-1}\|^2 \\
 &\quad + \eta \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^{k+1} - x_s^k \rangle \\
 &\quad + \eta \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^k - x \rangle + \gamma \sum_{s=1}^{S-1} \sum_{k=0}^{K-1} V(x, x_{s-1}^k) \\
 &\quad - \frac{\gamma}{2} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^{k+1} - \omega_s\|^2 - (1-\gamma) \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^{k+1}, x_s^k) \\
 &\quad + (1-\gamma) \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x, x_s^k) - \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x, x_s^{k+1}) \\
 &\leq \frac{1}{128} \|x_{S-1}^K - x_{S-1}^{K-1}\|^2 + V(x, x_S^0) + \frac{1}{4} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^{k+1}, x_s^k) \\
 &\quad + \frac{1}{32} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_B^k - x_s^{k-1}\|^2 \\
 &\quad + \eta \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^{k+1} - x_s^k \rangle \\
 &\quad + \eta \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^k - x \rangle + \gamma \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x, x_{s-1}^k)
 \end{aligned}$$

$$\begin{aligned}
& -\frac{\gamma}{2} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^{k+1} - \omega_s\|^2 - (1-\gamma) \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^{k+1}, x_s^k) \\
& + (1-\gamma) \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x, x_s^k) - \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x, x_s^{k+1}). \quad (24)
\end{aligned}$$

In this step, we again used lines 12 and 13 from Algorithm 1 and the fact (5). This result and line 12 of Algorithm 1 gives us

$$\begin{aligned}
& \gamma \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x, x_{s-1}^k) + (1-\gamma) \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x, x_s^k) - \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x, x_s^{k+1}) + V(x, x_S^0) \\
& \leq \gamma \sum_{k=0}^{K-1} V(x, x_{-1}^k) + \gamma \left(\sum_{s=0}^{S-2} \sum_{k=0}^{K-1} V(x, x_s^k) - \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x, x_s^k) \right) \\
& \quad + \left(\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x, x_s^k) - \sum_{s=0}^{S-1} \sum_{k=1}^K V(x, x_s^k) \right) + V(x, x_S^0) \\
& \leq \gamma \sum_{k=0}^{K-1} V(x, x_{-1}^k) - \gamma \sum_{k=0}^{K-1} V(x, x_{S-1}^k) + \sum_{s=0}^{S-1} V(x, x_s^0) \\
& \quad - \sum_{s=0}^{S-1} V(x, x_s^K) + V(x, x_S^0) \\
& = \gamma K V(x, x_0^0) - \gamma \sum_{k=0}^{K-1} V(x, x_{S-1}^k) \\
& \quad + \sum_{s=0}^{S-1} V(x, x_s^0) - \sum_{s=1}^S V(x, x_s^0) + V(x, x_S^0) \\
& = \gamma K V(x, x_0^0) - \gamma \sum_{k=0}^{K-1} V(x, x_{S-1}^k) + V(x, x_0^0) - V(x, x_S^0) + V(x, x_S^0) \\
& \leq (1 + \gamma K) V(x, x_0^0).
\end{aligned}$$

We apply this result to (24) and get

$$\begin{aligned}
& \eta \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[(g(x_s^{k+1}) - g(x)) + \eta \langle F(x_s^{k+1}), x_s^{k+1} - x \rangle \right] \\
& \leq \frac{1}{128} \|x_{S-1}^K - x_{S-1}^{K-1}\|^2 + \frac{1}{32} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^k - x_s^{k-1}\|^2 \\
& \quad + \frac{1}{4} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^{k+1}, x_s^k) \\
& \quad + \eta \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^{k+1} - x_s^k \rangle \right]
\end{aligned}$$

$$\begin{aligned}
& + \eta \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^k - x \rangle \right] \\
& - \frac{\gamma}{2} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^{k+1} - \omega_s\|^2 - (1 - \gamma) \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^{k+1}, x_s^k) \\
& + (1 + \gamma K)V(x, x_0^0).
\end{aligned}$$

After that we take maximum and expectation of both sides, we used fact that the maximum of sum is less then the sum of maximums. In addition, here we use property of conditional expectation to obtain $\mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^k \rangle \right] \right] = 0$. Thus,

$$\begin{aligned}
& \mathbb{E} \left[\max_{x \in \mathcal{C}} \left\{ \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\eta(g(x_s^{k+1}) - g(x)) + \eta \langle F(x_s^{k+1}), x_s^{k+1} - x \rangle \right] \right\} \right] \\
& \leq \frac{1}{128} \mathbb{E} \left[\|x_{S-1}^K - x_{S-1}^{K-1}\|^2 \right] + \frac{1}{32} \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^k - x_s^{k-1}\|^2 \right] \\
& \quad + \left(\gamma - \frac{3}{4} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^{k+1}, x_s^k) \right] \\
& \quad + \eta \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^{k+1} - x_s^k \rangle \right] \right] \\
& \quad + \eta \mathbb{E} \left[\max_{x \in \mathcal{C}} \left\{ \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^k - x \rangle \right] \right\} \right] \\
& \quad - \frac{\gamma}{2} \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^{k+1} - \omega_s\|^2 \right] + \max_{x \in \mathcal{C}} \Phi_0(x) \\
& \leq \frac{1}{128} \mathbb{E} \left[\|x_{S-1}^K - x_{S-1}^{K-1}\|^2 \right] + \frac{1}{32} \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^k - x_s^{k-1}\|^2 \right] \\
& \quad + \left(\gamma - \frac{3}{4} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^{k+1}, x_s^k) \right] \\
& \quad + \eta \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^{k+1} - x_s^k \rangle \right] \right] \\
& \quad + \eta \mathbb{E} \left[\max_{x \in \mathcal{C}} \left\{ \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\langle \mathbb{E}_k [\Delta^k] - \Delta^k, x \rangle \right] \right\} \right]
\end{aligned}$$

$$-\frac{\gamma}{2}\mathbb{E}\left[\sum_{s=0}^{S-1}\sum_{k=0}^{K-1}\|x_s^{k+1}-\omega_s\|^2\right]+(1+\gamma K)\max_{x\in\mathcal{C}}V(x,x_0^0). \quad (25)$$

Using the Young's inequality (44), we can estimate

$$\mathbb{E}\left[\eta\langle\mathbb{E}_k[\Delta^k]-\Delta^k,x_s^{k+1}-x_s^k\rangle\right]\leq 2\eta^2\mathbb{E}\left[\|\mathbb{E}_k[\Delta^k]-\Delta^k\|_2^2\right]+\frac{1}{8}\mathbb{E}\left[\|x_s^{k+1}-x_s^k\|_2^2\right].$$

Then

$$\begin{aligned} \eta\mathbb{E}\left[\sum_{s=0}^{S-1}\sum_{k=0}^{K-1}\left[\langle\mathbb{E}_k[\Delta^k]-\Delta^k,x_s^{k+1}-x_s^k\rangle\right]\right] \\ \leq 2\eta^2\sum_{s=0}^{S-1}\sum_{k=0}^{K-1}\mathbb{E}\|\mathbb{E}_k[\Delta^k]-\Delta^k\|_2^2+\frac{1}{8}\sum_{s=0}^{S-1}\sum_{k=0}^{K-1}\left[\mathbb{E}\left[\|x_s^{k+1}-x_s^k\|_2^2\right]\right]. \end{aligned} \quad (26)$$

$\mathbb{E}_k[\Delta^k]-\Delta^k$ is a stochastic process with $\mathbb{E}[\mathbb{E}_k[\Delta^k]-\Delta^k]=0$. Therefore, according to Lemma 3a

$$\begin{aligned} \eta\mathbb{E}\left[\max_{x\in\mathcal{C}}\left\{\sum_{s=0}^{S-1}\sum_{k=0}^{K-1}\langle\mathbb{E}_k[\Delta^k]-\Delta^k,x\rangle\right\}\right] \\ \leq\frac{1}{2}\max_{x\in\mathcal{C}}\|x-x_0^0\|_2^2+\frac{\eta^2}{2}\sum_{s=0}^{S-1}\sum_{k=0}^{K-1}\mathbb{E}\|\mathbb{E}_k[\Delta^k]-\Delta^k\|_2^2. \end{aligned} \quad (27)$$

Applying (26) and (27) to (25), we get

$$\begin{aligned} \mathbb{E}\left[\max_{x\in\mathcal{C}}\left\{\sum_{s=0}^{S-1}\sum_{k=0}^{K-1}\left[\eta(g(x_s^{k+1})-g(x))+\eta\langle F(x_s^{k+1}),x_s^{k+1}-x\rangle\right]\right\}\right] \\ \leq\frac{1}{128}\mathbb{E}\left[\|x_{S-1}^K-x_{S-1}^{K-1}\|^2\right]+\frac{1}{32}\mathbb{E}\left[\sum_{s=0}^{S-1}\sum_{k=0}^{K-1}\|x_s^k-x_s^{k-1}\|^2\right] \\ +\left(\gamma-\frac{3}{4}\right)\mathbb{E}\left[\sum_{s=0}^{S-1}\sum_{k=0}^{K-1}V(x_s^{k+1},x_s^k)\right] \\ +2\eta^2\sum_{s=0}^{S-1}\sum_{k=0}^{K-1}\mathbb{E}\|\mathbb{E}_k[\Delta^k]-\Delta^k\|_2^2 \\ +\frac{1}{8}\sum_{s=0}^{S-1}\sum_{k=0}^{K-1}\left[\mathbb{E}\left[\|x_s^{k+1}-x_s^k\|_2^2\right]\right] \\ +\frac{1}{2}\max_{x\in\mathcal{C}}\|x-x_0^0\|_2^2+\frac{\eta^2}{2}\sum_{s=0}^{S-1}\sum_{k=0}^{K-1}\mathbb{E}\|\mathbb{E}_k[\Delta^k]-\Delta^k\|_2^2 \end{aligned}$$

$$\begin{aligned}
 & -\frac{\gamma}{2}\mathbb{E}\left[\sum_{s=0}^{S-1}\sum_{k=0}^{K-1}\|x_s^{k+1}-\omega_s\|^2\right]+(1+\gamma K)\max_{x\in\mathcal{C}}V(x,x_0^0) \\
 & \leq\frac{1}{128}\mathbb{E}\left[\|x_{S-1}^K-x_{S-1}^{K-1}\|^2\right]+\frac{1}{32}\mathbb{E}\left[\sum_{s=0}^{S-1}\sum_{k=0}^{K-1}\|x_s^k-x_s^{k-1}\|^2\right] \\
 & \quad +\left(\gamma-\frac{3}{4}\right)\mathbb{E}\left[\sum_{s=0}^{S-1}\sum_{k=0}^{K-1}V(x_s^{k+1},x_B^k)\right] \\
 & \quad +\frac{5\eta^2}{2}\sum_{s=0}^{S-1}\sum_{k=0}^{K-1}\mathbb{E}\|\mathbb{E}_k[\Delta^k]-\Delta^k\|_2^2 \\
 & \quad +\frac{1}{8}\sum_{s=0}^{S-1}\sum_{k=0}^{K-1}\left[\mathbb{E}\left[\|x_s^{k+1}-x_B^k\|_2^2\right]\right]+\frac{1}{2}\max_{x\in\mathcal{C}}\|x-x_0^0\|_2^2 \\
 & \quad -\frac{\gamma}{2}\mathbb{E}\left[\sum_{s=0}^{S-1}\sum_{k=0}^{K-1}\|x_s^{k+1}-\omega_s\|^2\right]+(1+\gamma K)\max_{x\in\mathcal{C}}V(x,x_0^0).
 \end{aligned} \tag{28}$$

The Cauchy–Schwarz(43) inequality allows us to estimate

$$\begin{aligned}
 & -\frac{\gamma}{2}\mathbb{E}\left[\sum_{s=0}^{S-1}\sum_{k=0}^{K-1}\|x_s^{k+1}-\omega_s\|^2\right] \\
 & \leq-\frac{\gamma}{4}\mathbb{E}\left[\sum_{s=0}^{S-1}\sum_{k=0}^{K-1}\|x_s^k-\omega_s\|^2\right]+\frac{\gamma}{2}\mathbb{E}\left[\sum_{s=0}^{S-1}\sum_{k=0}^{K-1}\|x_s^{k+1}-x_s^k\|^2\right]
 \end{aligned} \tag{29}$$

Applying Lemma 2a, the fact (5), property $\|\cdot\|_2\leq\|\cdot\|$ and (29) to (28), we achieve

$$\begin{aligned}
 & \mathbb{E}\left[\max_{x\in\mathcal{C}}\left\{\sum_{s=0}^{S-1}\sum_{k=0}^{K-1}\left[\eta(g(x_s^{k+1})-g(x))+\eta\langle F(x_s^{k+1}),x_s^{k+1}-x\rangle\right]\right\}\right] \\
 & \leq\frac{1}{128}\mathbb{E}\left[\|x_{S-1}^K-x_{S-1}^{K-1}\|^2\right]+\frac{1}{32}\mathbb{E}\left[\sum_{s=0}^{S-1}\sum_{k=0}^{K-1}\|x_s^k-x_s^{k-1}\|^2\right] \\
 & \quad +\left(\gamma-\frac{3}{4}\right)\mathbb{E}\left[\sum_{s=0}^{S-1}\sum_{k=0}^{K-1}V(x_s^{k+1},x_s^k)\right] \\
 & \quad +\frac{5\eta^2}{2}\sum_{s=0}^{S-1}\sum_{k=0}^{K-1}\left[\frac{2\bar{L}_2^2}{b}\mathbb{E}\left[\|x_s^k-w_s\|^2+\|x_s^k-x_s^{k-1}\|^2\right]\right] \\
 & \quad +\frac{1}{8}\sum_{s=0}^{S-1}\sum_{k=0}^{K-1}\left[\mathbb{E}\left[\|x_s^{k+1}-x_s^k\|_2^2\right]\right]+\frac{1}{2}\max_{x\in\mathcal{C}}\|x-x_0^0\|_2^2 \\
 & \quad -\frac{\gamma}{2}\mathbb{E}\left[\sum_{s=0}^{S-1}\sum_{k=0}^{K-1}\|x_s^{k+1}-\omega_s\|^2\right]+(1+\gamma K)\max_{x\in\mathcal{C}}V(x,x_0^0)
 \end{aligned}$$

$$\begin{aligned}
&\leq \frac{1}{64} \mathbb{E} [V(x_{S-1}^K, x_{S-1}^{K-1})] + \frac{1}{16} \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^k, x_s^{k-1}) \right] \\
&\quad + \left(\gamma - \frac{3}{4} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^{k+1}, x_s^k) \right] + \frac{1}{4} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\mathbb{E} [V(x_s^{k+1}, x_s^k)] \right] \\
&\quad + \frac{5\eta^2}{2} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\frac{\bar{L}_2^2}{b} \mathbb{E} [2\|x_s^k - w_s\|^2 + V(x_s^k, x_s^{k-1})] \right] + \max_{x \in \mathcal{C}} V(x, x_0^0) \\
&\quad - \frac{\gamma}{4} \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^k - \omega_s\|^2 \right] + \frac{\gamma}{2} \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^{k+1} - x_s^k\|^2 \right] \\
&\quad + (1 + \gamma K) \max_{x \in \mathcal{C}} V(x, x_0^0) \\
&\leq \frac{1}{64} \mathbb{E} [V(x_{S-1}^K, x_{S-1}^{K-1})] + \left(\frac{1}{16} + \frac{5\bar{L}_2^2 \eta^2}{2b} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^k, x_s^{k-1})^2 \right] \\
&\quad + \left(\gamma - \frac{1}{2} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^{k+1}, x_s^k) \right] \\
&\quad + \frac{5\bar{L}_2^2 \eta^2}{b} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\mathbb{E} \|x_s^k - w_s\|^2 \right] + \max_{x \in \mathcal{C}} V(x, x_0^0) \\
&\quad - \frac{\gamma}{4} \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^k - \omega_s\|^2 \right] + \gamma \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^{k+1}, x_s^k) \right] \\
&\quad + (1 + \gamma K) \max_{x \in \mathcal{C}} V(x, x_0^0). \tag{30}
\end{aligned}$$

Next, we apply the fact (5), inequality (15), lines 12 and 13 of Algorithm 1, the choice of the step (16) and the initialization (line 2) to get

$$\begin{aligned}
&\frac{1}{64} \mathbb{E} [V(x_{S-1}^K, x_{S-1}^{K-1})] + \left(\frac{1}{16} + \frac{5\bar{L}_2^2 \eta^2}{2b} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^k, x_s^{k-1}) \right] \\
&\quad + \left(\gamma - \frac{1}{2} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^{k+1}, x_s^k) \right] + \gamma \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^{k+1}, x_s^k) \right] \\
&\leq \frac{1}{64} \mathbb{E} [V(x_{S-1}^K, x_{S-1}^{K-1})] + \left(\frac{1}{16} + \frac{5\bar{L}_2^2 \eta^2}{2b} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=1}^{K-1} V(x_s^k, x_s^{k-1}) \right] \\
&\quad + \left(\frac{1}{16} + \frac{5\bar{L}_2^2 \eta^2}{2b} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} V(x_s^0, x_s^{-1}) \right] \\
&\quad + \left(2\gamma - \frac{1}{2} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=1}^{K-1} V(x_s^k, x_s^{k-1}) \right] + \left(\gamma - \frac{1}{2} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} V(x_s^k, x_s^{K-1}) \right]
\end{aligned}$$

$$\begin{aligned}
&\leq \left(\frac{1}{16} + \frac{5\bar{L}_2^2\eta^2}{2b} + 2\gamma - \frac{1}{2} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=1}^{K-1} V(x_s^k, x_s^{k-1}) \right] \\
&\quad + \frac{1}{64} \mathbb{E} [V(x_{S-1}^K, x_{S-1}^{K-1})] + \left(\frac{1}{16} + \frac{5\bar{L}_2^2\eta^2}{2b} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} V(x_s^0, x_s^{-1}) \right] \\
&\quad + \left(2\gamma - \frac{1}{2} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} V(x_s^k, x_s^{K-1}) \right] \\
&\leq \left(\frac{1}{16} + \frac{5\bar{L}_2^2\eta^2}{2b} + 2\gamma - \frac{1}{2} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=1}^{K-1} V(x_s^k, x_s^{k-1}) \right] \\
&\quad + \frac{1}{64} \mathbb{E} [V(x_S^0, x_S^{-1})] + \left(\frac{1}{16} + \frac{5\bar{L}_2^2\eta^2}{2b} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} V(x_s^0, x_s^{-1}) \right] \\
&\quad + \left(2\gamma - \frac{1}{2} \right) \mathbb{E} \left[\sum_{s=1}^S V(x_s^0, x_s^{-1}) \right] \\
&\leq \left(\frac{1}{16} + \frac{5\bar{L}_2^2\eta^2}{2b} + 2\gamma - \frac{1}{2} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=1}^{K-1} V(x_s^k, x_s^{k-1}) \right] \\
&\quad + \left(\frac{1}{16} + \frac{5\bar{L}_2^2\eta^2}{2b} + 2\gamma - \frac{1}{2} \right) \mathbb{E} \left[\sum_{s=1}^S V(x_s^0, x_s^{-1}) \right] \\
&\quad + \left(\frac{1}{16} + \frac{5\bar{L}_2^2\eta^2}{2b} \right) \mathbb{E} \left[V(x_0^0, x_0^{-1}) \right] \\
&\leq \left(\frac{1}{16} + \frac{5}{2} \cdot \frac{1}{64} \cdot \frac{1}{16} + \frac{2}{16} - \frac{1}{2} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=1}^{K-1} V(x_s^k, x_s^{k-1}) \right] \\
&\quad + \left(\frac{1}{16} + \frac{5}{2} \cdot \frac{1}{64} \cdot \frac{1}{16} + \frac{2}{16} - \frac{1}{2} \right) \mathbb{E} \left[\sum_{s=1}^S V(x_s^0, x_s^{-1}) \right] \\
&\quad + \left(\frac{1}{16} + \frac{5\bar{L}_2^2\eta^2}{2b} \right) \mathbb{E} \left[V(x_0^0, x_0^{-1}) \right] \leq 0, \tag{31}
\end{aligned}$$

and

$$\begin{aligned}
&\frac{5\bar{L}_2^2\eta^2}{b} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\mathbb{E} \|x_s^k - w_s\|^2 \right] - \frac{\gamma}{4} \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^k - \omega_s\|^2 \right] \\
&\leq \left(\frac{5\bar{L}_2^2\eta^2}{b} - \frac{\gamma}{4} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^{k+1} - \omega_s\|^2 \right] \\
&\leq \left(\frac{5\gamma}{64} - \frac{\gamma}{4} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^{k+1} - \omega_s\|^2 \right] \leq 0. \tag{32}
\end{aligned}$$

Applying (31) and (32) to (30), we achieve

$$\begin{aligned}
& \mathbb{E} \left[\max_{x \in \mathcal{C}} \left\{ \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\eta(g(x_s^{k+1}) - g(x)) + \eta \langle F(x_s^{k+1}), x_s^{k+1} - x \rangle \right] \right\} \right] \\
& \leq \frac{5\bar{L}_2^2 \eta^2}{b} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\mathbb{E} \|x_s^k - w_s\|^2 \right] + \max_{x \in \mathcal{C}} V(x, x_0^0) \\
& \quad - \frac{\gamma}{4} \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^{k+1} - \omega_s\|^2 \right] + (1 + \gamma K) \max_{x \in \mathcal{C}} V(x, x_0^0) \\
& \leq \max_{x \in \mathcal{C}} \{V(x, x_0^0)\} + (1 + \gamma K) \max_{x \in \mathcal{C}} V(x, x_0^0) \\
& = (2 + \gamma K) \max_{x \in \mathcal{C}} \{V(x, x_0)\}. \tag{33}
\end{aligned}$$

Next we use (6) and the fact that

$$\begin{aligned}
& \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\eta(g(x_s^{k+1}) - g(x)) + \eta \langle F(x_s^{k+1}), x_s^{k+1} - x \rangle \right] \\
& \geq \eta \left(g \left(\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} x_s^{k+1} \right) - g(x) \right) \\
& \quad + \eta \left\langle F \left(\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} x_s^{k+1} \right), \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} x_s^{k+1} - x \right\rangle \\
& \geq \eta K S (g(x_S) - g(x) + \langle F(x_S), x_S - x \rangle)
\end{aligned}$$

to get

$$\begin{aligned}
& \eta K S \mathbb{E} [\text{Gap}(x_S)] \\
& \leq \mathbb{E} \left[\max_{x \in \mathcal{C}} \left\{ \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\eta(g(x_s^{k+1}) - g(x)) + \eta \langle F(x_s^{k+1}), x_s^{k+1} - x \rangle \right] \right\} \right]. \tag{34}
\end{aligned}$$

Using (34) with (33), we achieve the result

$$\eta K S \mathbb{E} [\text{Gap}(x_S)] \leq (2 + K\gamma) \max_{x \in \mathcal{C}} V(x, x_0^0).$$

Then,

$$\mathbb{E} [\text{Gap}(x_S)] \leq \frac{(2 + K\gamma)}{\eta K S} \max_{x \in \mathcal{C}} V(x, x_0^0). \tag{35}$$

This finishes the proof.

Now we are ready to obtain complexity of Algorithm 1. We choice $\gamma = \frac{1}{K}$ to achieve $O(1)$ in the numerator of estimation (35). In addition, it is natural to take $K = \frac{M}{3b}$ due to triple operator call in line 7.

Corollary 1. Let $K = \frac{M}{3b}$, $\gamma = \frac{1}{K} = \frac{3b}{M}$ and $\eta = \min \left\{ \frac{\sqrt{\gamma b}}{2L_2}, \frac{1}{8L_2} \right\} = \min \left\{ \frac{\sqrt{b}}{2\bar{L}_2} \cdot \sqrt{\frac{3b}{M}}, \frac{1}{8L_2} \right\}$ and $b \leq \frac{\sqrt{M}\bar{L}_2}{L_2}$. Then the total complexity of Algorithm 1 to reach ϵ -accuracy is $\mathcal{O}(M + \frac{L\sqrt{M}}{\epsilon})$.

Proof. From Theorem 1 it follows

$$\begin{aligned} \mathbb{E}[\text{Gap}(x_S)] &\leq \frac{(2 + K\gamma)}{\eta KS} \max_{x \in \mathcal{C}} \{V(x, x_0)\} \\ &\leq \left(\frac{2 + K \cdot \frac{1}{K}}{\frac{b}{2L_2} \cdot \sqrt{\frac{3}{M}} \cdot \frac{M}{3b} \cdot S} + \frac{2 + K \cdot \frac{1}{K}}{\frac{1}{8L_2} \cdot \frac{M}{3b} \cdot S} \right) \max_{x \in \mathcal{C}} \{V(x, x_0)\} \\ &\leq \left(\frac{\bar{L}_2 \cdot 6\sqrt{3}}{\sqrt{MS}} + \frac{L_2 \cdot 72b}{NS} \right) \max_{x \in \mathcal{C}} \{V(x, x_0)\} \\ &= \mathcal{O} \left(\frac{\bar{L}_2}{\sqrt{MS}} + \frac{L_2 \cdot b}{MS} \right). \end{aligned}$$

With $b \leq \frac{\sqrt{M}\bar{L}_2}{L_2}$ we have $\mathbb{E}[\text{Gap}(x_S)] = \mathcal{O} \left(\frac{\bar{L}_2}{\sqrt{MS}} \right)$. One outer iteration requires $\mathcal{O}(M)$ evaluations of F_m . Hence, the final complexity of $\mathcal{O}(M + \frac{\bar{L}_2\sqrt{M}}{\epsilon})$.

Now let us consider the case when Lipschitzness of the operator is given in the form of Assumption 3b.

Theorem 2. Consider the problem (1)+(2) under Assumptions 1, 2 and 3b. Let $\{x_s^k\}$ be the sequence generated by Algorithm 1 with tuning of $\eta, \theta, \alpha, \beta, \gamma$ as follows:

$$0 < \gamma = p \leq \frac{1}{16},$$

$$\eta = \min \left\{ \frac{\sqrt{\gamma b}}{8L\sqrt{1 + C \ln n}}, \frac{1}{8L\sqrt{1 + C \ln n}} \right\}. \quad (36)$$

Then for $x_S = \frac{1}{KS} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} x_s^k$ it holds that

$$\mathbb{E}[\text{Gap}(x_S)] \leq \frac{(2 + K\gamma)}{\eta KS} \max_{x \in \mathcal{C}} \{V(x, x_0)\}.$$

Proof. The beginning of the proof is similar to the proof of Theorem 1. With given assumptions we replace (20) and (23) with

$$-\eta \langle F(x_s^k) - F(x_s^{k-1}), x_s^{k+1} - x_s^k \rangle \leq 2\eta^2 \|F(x_s^k) - F(x_s^{k-1})\|_*^2 + \frac{1}{8} \|x_s^{k+1} - x_s^k\|^2$$

$$\begin{aligned}
&\leq 2L^2\eta^2\|x_s^k - x_s^{k-1}\|^2 + \frac{1}{8}\|x_s^{k+1} - x_s^k\|^2 \\
&\leq \frac{1}{32}\|x_s^k - x_s^{k-1}\|^2 + \frac{1}{8}\|x_s^{k+1} - x_s^k\|^2
\end{aligned}$$

and

$$\begin{aligned}
\eta\langle F(x_S^0) - F(x_S^{-1}), x_S^0 - x \rangle &\leq \frac{\eta^2}{2}\|F(x_S^0) - F(x_S^{-1})\|_*^2 + \frac{1}{2}\|x_S^0 - x\|^2 \\
&\leq \frac{\eta^2 L^2}{2}\|x_S^0 - x_S^{-1}\|^2 + \frac{1}{2}\|x_S^0 - x\|^2 \\
&\leq \frac{1}{128}\|x_S^0 - x_S^{-1}\|^2 + \frac{1}{2}\|x_S^0 - x\|^2.
\end{aligned}$$

These inequalities are correct due to Young's inequality (44), Assumption 3b and the choice of step (36).

In a similar way to (25) from Theorem 1 we get the result

$$\begin{aligned}
&\mathbb{E} \left[\max_{x \in \mathcal{C}} \left\{ \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\eta(g(x_s^{k+1}) - g(x)) + \eta\langle F(x_s^{k+1}), x_s^{k+1} - x \rangle \right] \right\} \right] \\
&\leq \frac{1}{128}\mathbb{E} \left[\|x_{S-1}^K - x_{S-1}^{K-1}\|^2 \right] + \frac{1}{32}\mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^k - x_s^{k-1}\|^2 \right] \\
&\quad + \left(\gamma - \frac{3}{4} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^{k+1}, x_s^k) \right] \\
&\quad + \eta \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^{k+1} - x_s^k \rangle \right] \right] \\
&\quad + \eta \mathbb{E} \left[\max_{x \in \mathcal{C}} \left\{ \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\langle \mathbb{E}_k [\Delta^k] - \Delta^k, x \rangle \right] \right\} \right] \\
&\quad - \frac{\gamma}{2} \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^{k+1} - \omega_s\|^2 \right] + (1 + \gamma K) \max_{x \in \mathcal{C}} V(x, x_0^0).
\end{aligned} \tag{37}$$

Similar to (26), we get

$$\begin{aligned}
&\eta \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\langle \mathbb{E}_k [\Delta^k] - \Delta^k, x_s^{k+1} - x_s^k \rangle \right] \right] \\
&\leq 2\eta^2 \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \mathbb{E} \|\mathbb{E}_k [\Delta^k] - \Delta^k\|_*^2 + \frac{1}{8} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\mathbb{E} [\|x_s^{k+1} - x_s^k\|^2] \right].
\end{aligned} \tag{38}$$

$\mathbb{E}_k [\Delta^k] - \Delta^k$ is a stochastic process with $\mathbb{E}[\mathbb{E}_k [\Delta^k] - \Delta^k] = 0$. Therefore, according to Lemma 3b

$$\begin{aligned} & \eta \mathbb{E} \left[\max_{x \in \mathcal{C}} \left\{ \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x \rangle \right\} \right] \\ & \leq V(x, x_0^0) + \frac{\eta^2}{2} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \mathbb{E} \|\mathbb{E}_k [\Delta^k] - \Delta^k\|_*^2. \end{aligned} \quad (39)$$

Applying (38) and (39) to (37), we get

$$\begin{aligned} & \mathbb{E} \left[\max_{x \in \mathcal{C}} \left\{ \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\eta (g(x_s^{k+1}) - g(x)) + \eta \langle F(x_s^{k+1}), x_s^{k+1} - x \rangle \right] \right\} \right] \\ & \leq \frac{1}{128} \mathbb{E} [\|x_{S-1}^K - x_{S-1}^{K-1}\|^2] + \frac{1}{32} \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^k - x_s^{k-1}\|^2 \right] \\ & \quad + \left(\gamma - \frac{3}{4} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^{k+1}, x_s^k) \right] \\ & \quad + 2\eta^2 \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \mathbb{E} \|\mathbb{E}_k [\Delta^k] - \Delta^k\|_*^2 + \frac{1}{8} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\mathbb{E} [\|x_s^{k+1} - x_s^k\|^2] \right] \\ & \quad + V(x, x_0^0) + \frac{\eta^2}{2} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \mathbb{E} \|\mathbb{E}_k [\Delta^k] - \Delta^k\|_*^2 \\ & \quad - \frac{\gamma}{2} \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^{k+1} - \omega_s\|^2 \right] + (1 + \gamma K) \max_{x \in \mathcal{C}} V(x, x_0^0) \\ & \leq \frac{1}{128} \mathbb{E} [\|x_{S-1}^K - x_{S-1}^{K-1}\|^2] + \frac{1}{32} \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^k - x_s^{k-1}\|^2 \right] \\ & \quad + \left(\gamma - \frac{3}{4} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^{k+1}, x_s^k) \right] \\ & \quad + \frac{5\eta^2}{2} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \mathbb{E} \|\mathbb{E}_k [\Delta^k] - \Delta^k\|_*^2 \\ & \quad + \frac{1}{8} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\mathbb{E} [\|x_s^{k+1} - x_s^k\|^2] \right] + \frac{1}{2} \max_{x \in \mathcal{C}} \|x - x_0^0\|^2 \\ & \quad - \frac{\gamma}{2} \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^{k+1} - \omega_s\|^2 \right] + (1 + \gamma K) \max_{x \in \mathcal{C}} V(x, x_0^0). \end{aligned} \quad (40)$$

Applying Lemma 2b, the fact (5), property $\|\cdot\|_2 \leq \|\cdot\|$ and (29) to (40), we achieve

$$\begin{aligned}
& \mathbb{E} \left[\max_{x \in \mathcal{C}} \left\{ \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\eta(g(x_s^{k+1}) - g(x)) + \eta \langle F(x_s^{k+1}), x_s^{k+1} - x \rangle \right] \right\} \right] \\
& \leq \frac{1}{128} \mathbb{E} \left[\|x_{S-1}^K - x_{S-1}^{K-1}\|^2 \right] + \frac{1}{32} \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^k - x_s^{k-1}\|^2 \right] \\
& \quad + \left(\gamma - \frac{3}{4} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^{k+1}, x_s^k) \right] \\
& \quad + \frac{5\eta^2}{2} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\frac{2(1 + C \ln n)L^2}{b} \mathbb{E} [\|x_s^k - w_s\|^2 + \|x_s^k - x_s^{k-1}\|^2] \right] \\
& \quad + \frac{1}{8} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\mathbb{E} [\|x_s^{k+1} - x_s^k\|_2^2] \right] + \frac{1}{2} \max_{x \in \mathcal{C}} \|x - x_0\|_2^2 \\
& \quad - \frac{\gamma}{2} \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^{k+1} - \omega_s\|^2 \right] + (1 + \gamma K) \max_{x \in \mathcal{C}} V(x, x_0^0) \\
& \leq \frac{1}{64} \mathbb{E} [V(x_{S-1}^K, x_{S-1}^{K-1})] + \frac{1}{16} \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^k, x_s^{k-1}) \right] \\
& \quad + \left(\gamma - \frac{3}{4} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^{k+1}, x_s^k) \right] + \frac{1}{4} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\mathbb{E} [V(x_s^{k+1}, x_s^k)] \right] \\
& \quad + \frac{5\eta^2}{2} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\frac{L^2}{b} \mathbb{E} \left[2(1 + C \ln n) \|x_s^k - w_s\|^2 \right. \right. \\
& \quad \left. \left. + \frac{2(1 + C \ln n)}{2} V(x_s^k, x_s^{k-1}) \right] \right] + \max_{x \in \mathcal{C}} V(x, x_0^0) \\
& \quad - \frac{\gamma}{4} \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^k - \omega_s\|^2 \right] + \frac{\gamma}{2} \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^{k+1} - x_s^k\|^2 \right] \\
& \quad + (1 + \gamma K) \max_{x \in \mathcal{C}} V(x, x_0^0) \\
& \leq \frac{1}{64} \mathbb{E} [V(x_{S-1}^K, x_{S-1}^{K-1})] \\
& \quad + \left(\frac{1}{16} + \frac{10(1 + C \ln n)L^2\eta^2}{4b} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^k, x_s^{k-1})^2 \right] \\
& \quad + \left(\gamma - \frac{1}{2} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^{k+1}, x_s^k) \right]
\end{aligned}$$

$$\begin{aligned}
& + \frac{10(1 + C \ln n)L^2\eta^2}{2b} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\mathbb{E} \|x_s^k - w_s\|^2 \right] + \max_{x \in \mathcal{C}} V(x, x_0^0) \\
& - \frac{\gamma}{4} \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^k - \omega_s\|^2 \right] + \gamma \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^{k+1}, x_s^k) \right] \\
& + (1 + \gamma K) \max_{x \in \mathcal{C}} V(x, x_0^0).
\end{aligned}$$

The following steps are similar up to (30). We use the fact (5), inequality (15), lines 12 and 13 of Algorithm 1, the choice of the step (16) and similar calculations to (31) and (32) to get

$$\begin{aligned}
& \frac{1}{64} \mathbb{E} \left[V(x_{S-1}^K, x_{S-1}^{K-1}) \right] + \left(\frac{1}{16} + \frac{10(1 + C \ln n)L^2\eta^2}{4b} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^k, x_s^{k-1}) \right] \\
& + \left(\gamma - \frac{1}{2} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^{k+1}, x_s^k) \right] + \gamma \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} V(x_s^{k+1}, x_s^k) \right] \\
& \leq \left(\frac{1}{16} + \frac{10(1 + C \ln n)L^2\eta^2}{4b} + 2\gamma - \frac{1}{2} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=1}^{K-1} V(x_s^k, x_s^{k-1}) \right] \\
& + \left(\frac{1}{16} + \frac{10(1 + C \ln n)L^2\eta^2}{4b} + 2\gamma - \frac{1}{2} \right) \mathbb{E} \left[\sum_{s=1}^S V(x_s^0, x_s^{-1}) \right] \\
& + \left(\frac{1}{16} + \frac{10(1 + C \ln n)L^2\eta^2}{4b} \right) \mathbb{E} \left[V(x_0^0, x_0^{-1}) \right] \\
& \leq \left(\frac{1}{16} + \frac{10}{4} \cdot \frac{1}{64} \cdot \frac{1}{16} + \frac{2}{16} - \frac{1}{2} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=1}^{K-1} V(x_s^k, x_s^{k-1}) \right] \\
& + \left(\frac{1}{16} + \frac{10}{4} \cdot \frac{1}{64} \cdot \frac{1}{16} + \frac{2}{16} - \frac{1}{2} \right) \mathbb{E} \left[\sum_{s=1}^S V(x_s^0, x_s^{-1}) \right] \\
& + \left(\frac{1}{16} + \frac{5(1 + C \ln n)L^2\eta^2}{2b} \right) \mathbb{E} \left[V(x_0^0, x_0^{-1}) \right] \leq 0,
\end{aligned}$$

and

$$\begin{aligned}
& \frac{10(1 + C \ln n)L^2\eta^2}{2b} \sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \left[\mathbb{E} \|x_s^k - w_s\|^2 \right] - \frac{\gamma}{4} \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^k - \omega_s\|^2 \right] \\
& \leq \left(\frac{10(1 + C \ln n)L^2\eta^2}{2b} - \frac{\gamma}{4} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^{k+1} - \omega_s\|^2 \right] \\
& \leq \left(\frac{5\gamma}{64} - \frac{\gamma}{4} \right) \mathbb{E} \left[\sum_{s=0}^{S-1} \sum_{k=0}^{K-1} \|x_s^{k+1} - \omega_s\|^2 \right] \leq 0.
\end{aligned}$$

Steps similar to the proof of Theorem 1 finish the proof.

Now we proceed similarly to Corollary 1.

Corollary 2. Let $K = \frac{M}{3b}$, $\gamma = \frac{1}{K} = \frac{3b}{M}$ and $\eta = \min \left\{ \frac{\sqrt{\gamma b}}{2\sqrt{1+C \ln n}L}, \frac{1}{8L\sqrt{1+C \ln n}} \right\} = \min \left\{ \frac{b}{2\sqrt{1+C \ln n}L} \cdot \sqrt{\frac{3}{M}}, \frac{1}{8L\sqrt{1+C \ln n}} \right\}$ and $b \leq \sqrt{M}$. Then the total complexity of the Algorithm 1 to reach ϵ -accuracy is $\mathcal{O}(M + \frac{L\sqrt{M}}{\epsilon})$.

Proof. From Theorem 2 it follows

$$\begin{aligned} \mathbb{E}[\text{Gap}(x_S)] &\leq \frac{(2 + K\gamma)}{\eta K S} \max_{x \in \mathcal{C}} \{V(x, x_0)\} \\ &\leq \left(\frac{2 + K \cdot \frac{1}{K}}{\frac{b}{2L} \cdot \sqrt{\frac{3}{M}} \cdot \frac{M}{3b} \cdot S} + \frac{2 + K \cdot \frac{1}{K}}{\frac{1}{8L} \cdot \frac{M}{3b} \cdot S} \right) \sqrt{1 + C \ln n} \max_{x \in \mathcal{C}} \{V(x, x_0)\} \\ &\leq \left(\frac{L \cdot 6\sqrt{3}}{\sqrt{MS}} + \frac{L \cdot 72b}{MS} \right) \sqrt{1 + C \ln n} \cdot \max_{x \in \mathcal{C}} \{V(x, x_0)\} \\ &= \tilde{\mathcal{O}} \left(\frac{L}{\sqrt{MS}} + \frac{L \cdot b}{MS} \right). \end{aligned}$$

With $b \leq \sqrt{M}$ we have $\mathbb{E}[\text{Gap}(x_S)] = \tilde{\mathcal{O}} \left(\frac{L}{\sqrt{MS}} \right)$. On epoch requires $\mathcal{O}(N)$ evaluations F_m . Hence, the final complexity of $\tilde{\mathcal{O}}(M + \frac{L\sqrt{M}}{\epsilon})$.

As mentioned before, $L \leq \bar{L}_2$. Thus, we can state $\frac{L}{\sqrt{MS}} \leq \frac{\bar{L}_2}{\sqrt{MS}}$.

3 Experiments

In this section, empirical performance of Algorithm 1 is shown. Similar to the previous works [1,40], we consider matrix games

$$\min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}} \langle Ax, y \rangle \quad (41)$$

with matrix $A \in \mathbb{R}^{n \times n}$ and use simplex constraints (4) in the entropic setup:

$$\mathcal{X} = \Delta^n, \quad \mathcal{Y} = \Delta^n.$$

For our experiments we choose the policeman and burglar matrix [34] and the first test matrix from [35].

Note that the problem (41) does not have the finite-sum form as (1)+(2), but it can be rewritten as

$$A = \sum_{i=1}^n A_i \text{ or } A = \sum_{j=1}^n A_{\cdot j}. \quad (42)$$

where $A_{i\cdot}$ is the i -th row of A and $A_{\cdot j}$ is the j -th column of A – see details in Section 5.1.2 from [1]. In this formulation the problem already has the form of a finite sum.

In order to calculate the update (in particular, we need to compute ∇h in the closed form) we introduce $u = (u^x, u^y)$ and $v = (v^x, v^y)$. With this notation we can adapt (3) in a similar to the previous works [2,40,12] way:

$$\begin{aligned} \nabla h(x_s^{k+1}) &= (1 - \gamma)\nabla h(x_s^k) + \gamma\nabla h(\bar{u}^s) \\ &\quad - \eta \left(\frac{1}{b} \sum_{i \in B^k} \left[\frac{1}{r_i} A_{i\cdot} (2y_{s,i}^k - v_{s,i} - y_{s,i}^{k-1}) \right] + A^\top v_s \right) \\ &= (1 - \gamma)\nabla h(x_s^k) + \gamma\nabla h(\bar{u}^s) - \eta A^\top v_s \\ &\quad - \frac{\eta}{b} \sum_{i \in B^k} \left[\frac{1}{r_i} A_{i\cdot} \|2y_s^k - v_s - y_s^{k-1}\| \cdot \text{sign}(2y_{s,i}^k - v_{s,i} - y_{s,i}^{k-1}) \right]. \end{aligned}$$

Likewise, for the second component one can write

$$\begin{aligned} \nabla h(y_s^{k+1}) &= (1 - \gamma)\nabla h(y_s^k) + \gamma\nabla h(\bar{v}^s) + \eta A u_s \\ &\quad + \frac{\eta}{b} \sum_{j \in B^k} \left[\frac{1}{c_j} A_{\cdot j} \|2x_s^k - u_s - x_s^{k-1}\| \cdot \text{sign}(2x_{s,j}^k - u_{s,j} - x_{s,j}^{k-1}) \right]. \end{aligned}$$

Now we are ready to implement the Algorithm 1 and compare our method with other methods. In particular, we choose Algorithm 2 from [1], Algorithm 1+2 from [12] — state-of-the-art competitors (see Table 1). The parameters of the algorithms are taken from the provided papers.

We run all methods with various batch sizes and use theoretical parameters (see Theorems 1 and 2 for our method and Section 6 from [1] for competitors). Duality gap (6) is used as the convergence measure, in the matrix games setup it can be simply computed as $[\max_i(A^\top x)_i - \min_j(Ay)_j]$ for simplex constraints. The comparison criterion is the number of oracle calls (one call is computationally equal to calculations of $A_{i\cdot}y_i$ and $A_{\cdot j}x_j$). The results are reflected in Figures 1 and 2. It is clearly shown that Algorithm 1 outperforms competitors and its convergence rate is insensitive to the size of batch (unlike the other methods). Additionally, we also notice that on given test matrices it is possible to improve convergence rates by adjusting parameters. Using grid search we achieve the parameters that guaranteed the best convergence. With these parameters we run algorithms and plot Figures 3 and 4. The results of this experiment are similar: our algorithm provides better results than the competitors and the convergence rate does not depend on the batch size.

4 References

1. Alacaoglu, A., Malitsky, Y.: Stochastic variance reduction for variational inequality methods. arXiv preprint arXiv:2102.08352 (2021)

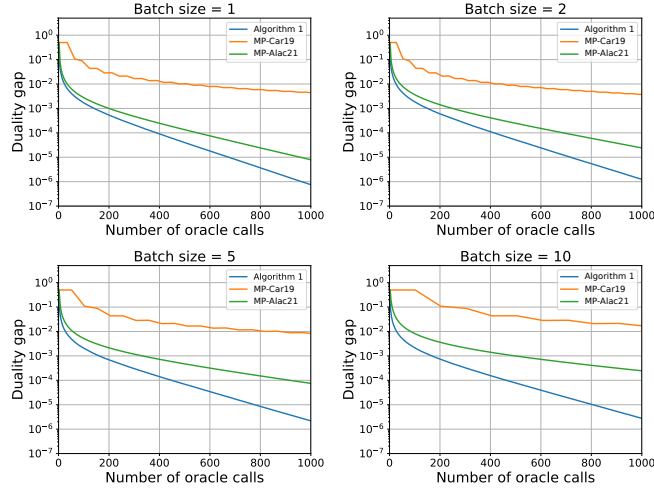


Fig. 1: Comparison of computational complexities for Algorithm 1, MP-Car19 (Algorithm 1+2 from [12]), and MP-Alac21 (Algorithm 2 [1]) with different batch sizes on test matrix from [34].

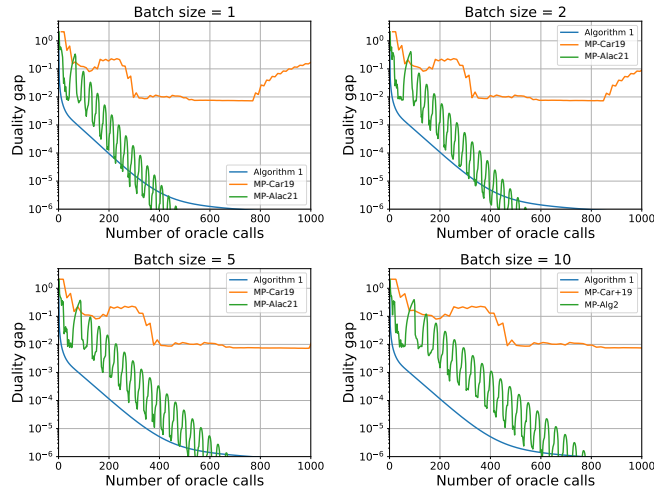


Fig. 2: Comparison of computational complexities for Algorithm 1, MP-Car19 (Algorithm 1+2 from [12]), and MP-Alac21 (Algorithm 2 [1]) with different batch sizes on policeman and burglar matrix from [34].

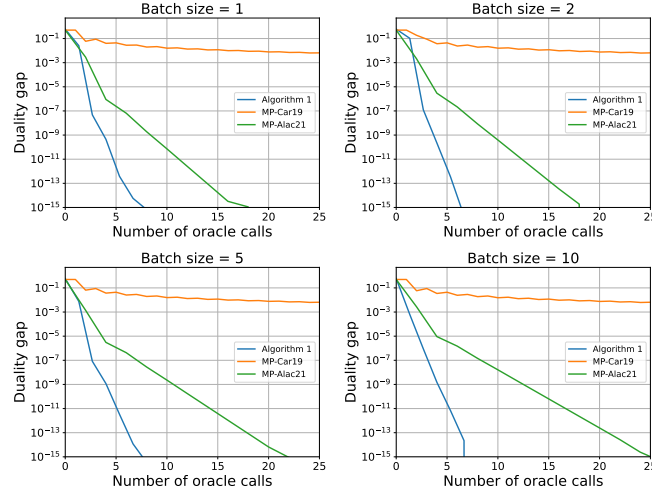


Fig. 3: Comparison of computational complexities for Algorithm 1 (with adjusted parameters), MP-Car19 (Algorithm 1+2 from [12]), and MP-Alac21 (Algorithm 2 [1]) with different batch sizes on first test matrix from [34].

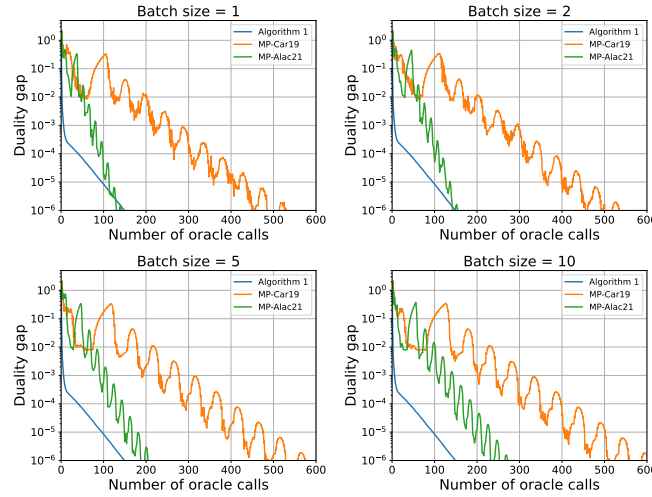


Fig. 4: Comparison of computational complexities for Algorithm 1 (with adjusted parameters), MP-Car19 (Algorithm 1+2 from [12]), and MP-Alac21 (Algorithm 2 [1]) with different batch sizes on policeman and burglar matrix from [34].

2. Alacaoglu, A., Malitsky, Y., Cevher, V.: Forward-reflected-backward method with variance reduction. *Computational optimization and applications* **80**(2), 321–346 (2021)
3. Allen-Zhu, Z.: Katyusha: The first direct acceleration of stochastic gradient methods. In: *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing*. pp. 1200–1205 (2017)
4. Allen-Zhu, Z.: Katyusha x: Practical momentum method for stochastic sum-of-nonconvex optimization. *arXiv preprint arXiv:1802.03866* (2018)
5. Allen-Zhu, Z., Yuan, Y.: Improved svrg for non-strongly-convex or sum-of-nonconvex objectives. In: *International conference on machine learning*. pp. 1080–1089. PMLR (2016)
6. Arjovsky, M., Chintala, S., Bottou, L.: Wasserstein generative adversarial networks. In: *International conference on machine learning*. pp. 214–223. PMLR (2017)
7. Bach, F., Mairal, J., Ponce, J.: Convex sparse matrix factorizations. *arXiv preprint arXiv:0812.1869* (2008)
8. Ben-Tal, A., El Ghaoui, L., Nemirovski, A.: *Robust optimization*, vol. 28. Princeton university press (2009)
9. Beznosikov, A., Dvinskikh, D., Semenov, A., Gasnikov, A.: Bregman proximal method for efficient communications under similarity. *arXiv preprint arXiv:2311.06953* (2023)
10. Beznosikov, A., Gorbunov, E., Berard, H., Loizou, N.: Stochastic gradient descent-ascent: Unified theory and new efficient methods. In: *International Conference on Artificial Intelligence and Statistics*. pp. 172–235. PMLR (2023)
11. Boyd, S.P., Vandenberghe, L.: *Convex optimization*. Cambridge university press (2004)
12. Carmon, Y., Jin, Y., Sidford, A., Tian, K.: Variance reduction for matrix games. *arXiv preprint arXiv:1907.02056* (2019)
13. Chavdarova, T., Gidel, G., Fleuret, F., Lacoste-Julien, S.: Reducing noise in gan training with variance reduced extragradient. vol. 32, pp. 393–403 (2019)
14. Chavdarova, T., Hsieh, Y.P., Jordan, M.I.: Continuous-time analysis for variational inequalities: An overview and desiderata. *arXiv preprint arXiv:2207.07105* (2022)
15. Defazio, A., Bach, F., Lacoste-Julien, S.: Saga: A fast incremental gradient method with support for non-strongly convex composite objectives. In: Ghahramani, Z., Welling, M., Cortes, C., Lawrence, N., Weinberger, K. (eds.) *Advances in Neural Information Processing Systems*. vol. 27, pp. 1646–1654. Curran Associates, Inc. (2014)
16. Facchinei, F., Pang, J.S.: *Finite-Dimensional Variational Inequalities and Complementarity Problems*. Springer Series in Operations Research, Springer (2003)
17. Gasnikov, A.: Universal gradient descent. *arXiv preprint arXiv:1711.00394* (2017)
18. Gidel, G., Hemmat, R.A., Pezeshki, M., Le Priol, R., Huang, G., Lacoste-Julien, S., Mitliagkas, I.: Negative momentum for improved game dynamics. In: *The 22nd International Conference on Artificial Intelligence and Statistics*. pp. 1802–1811. PMLR (2019)
19. Han, Y., Xie, G., Zhang, Z.: Lower complexity bounds of finite-sum optimization problems: The results and construction. *arXiv preprint arXiv:2103.08280* (2021)
20. Harker, P.T., Pang, J.S.: *Finite-dimensional variational inequality and nonlinear complementarity problems: a survey of theory, algorithms and applications*. Mathematical programming (1990)
21. Jin, Y., Sidford, A.: Efficiently solving mdps with stochastic mirror descent. In: *International Conference on Machine Learning*. pp. 4890–4900. PMLR (2020)

22. Joachims, T.: A support vector method for multivariate performance measures. In: Proceedings of the 22nd international conference on Machine learning. pp. 377–384 (2005)
23. Johnson, R., Zhang, T.: Accelerating stochastic gradient descent using predictive variance reduction. vol. 26, pp. 315–323 (2013)
24. Juditsky, A., Nemirovski, A., Tauvel, C.: Solving variational inequalities with stochastic mirror-prox algorithm. *Stochastic Systems* **1**(1), 17–58 (2011)
25. Korpelevich, G.M.: The extragradient method for finding saddle points and other problems. *Matecon* **12**, 35–49 (1977; Russian original: *Economika Mat Metody*, 12(4):747–756, 1976)
26. Kovalev, D., Beznosikov, A., Borodich, E., Gasnikov, A., Scutari, G.: Optimal gradient sliding and its application to distributed optimization under similarity. arXiv preprint arXiv:2205.15136 (2022)
27. Kovalev, D., Beznosikov, A., Sadiev, A., Pershiianov, M., Richtárik, P., Gasnikov, A.: Optimal algorithms for decentralized stochastic variational inequalities. arXiv preprint arXiv:2202.02771 (2022)
28. Lan, G., Li, Z., Zhou, Y.: A unified variance-reduced accelerated gradient method for convex optimization. *Advances in Neural Information Processing Systems* **32** (2019)
29. Leluc, R.: Monte Carlo Methods and Stochastic Approximation: Theory and Applications to Machine Learning. Ph.D. thesis, Institut Polytechnique de Paris (2023)
30. Madry, A., Makelov, A., Schmidt, L., Tsipras, D., Vladu, A.: Towards deep learning models resistant to adversarial attacks. arXiv preprint arXiv:1706.06083 (2017)
31. Malitsky, Y., Tam, M.K.: A forward-backward splitting method for monotone inclusions without cocoercivity. *SIAM Journal on Optimization* **30**(2), 1451–1472 (2020)
32. Mokhtari, A., Ozdaglar, A., Pattathil, S.: A unified analysis of extra-gradient and optimistic gradient methods for saddle point problems: Proximal point approach. In: International Conference on Artificial Intelligence and Statistics. pp. 1497–1507. PMLR (2020)
33. Nemirovski, A.: Prox-method with rate of convergence $o(1/t)$ for variational inequalities with lipschitz continuous monotone operators and smooth convex-concave saddle point problems. *SIAM Journal on Optimization* **15**(1), 229–251 (2004)
34. Nemirovski, A.: Mini-course on convex programming algorithms (2013)
35. Nemirovski, A., Juditsky, A., Lan, G., Shapiro, A.: Robust stochastic approximation approach to stochastic programming. *SIAM Journal on optimization* **19**(4), 1574–1609 (2009)
36. Nesterov, Y.: Dual extrapolation and its applications to solving variational inequalities and related problems. *Mathematical Programming* **109**(2), 319–344 (2007)
37. Neumann, J.V., Morgenstern, O.: Theory of games and economic behavior. Princeton University Press (1944)
38. Omidshafiei, S., Pazis, J., Amato, C., How, J.P., Vian, J.: Deep decentralized multi-task multi-agent reinforcement learning under partial observability. In: International Conference on Machine Learning. pp. 2681–2690. PMLR (2017)
39. Palaniappan, B., Bach, F.: Stochastic variance reduction methods for saddle-point problems. In: Advances in Neural Information Processing Systems. pp. 1416–1424 (2016)
40. Pichugin, A., Pechin, M., Beznosikov, A., Gasnikov, A.: Optimal analysis of method with batching for monotone stochastic finite-sum variational inequalities. arXiv preprint arXiv:2401.07858 (2024)

41. Polyak, B.: Introduction to optimization. Optimization Software (1987)
42. Popov, L.D.: A modification of the arrow-hurwicz method for search of saddle points. Mathematical notes of the Academy of Sciences of the USSR **28**, 845–848 (1980)
43. Rudeva, A., Shananin, A.: Variational inequalities for economic equilibrium in the model with the deficit of the working capital. Moscow University Computational Mathematics and Cybernetics **31**(4), 170–178 (2007)
44. Sadiev, A., Danilova, M., Gorbunov, E., Horváth, S., Gidel, G., Dvurechensky, P., Gasnikov, A., Richtárik, P.: High-probability bounds for stochastic optimization and variational inequalities: the case of unbounded variance. In: International Conference on Machine Learning. pp. 29563–29648. PMLR (2023)
45. Shapiro, A., Dentcheva, D., Ruszczyński, A.: Lectures on stochastic programming: modeling and theory. SIAM (2021)
46. Stampacchia, G.: Formes bilinéaires coercitives sur les ensembles convexes. Académie des Sciences de Paris **258**, 4413–4416 (1964)
47. Sun, H.L., Chen, X.J.: Two-stage stochastic variational inequalities: theory, algorithms and applications. Journal of the Operations Research Society of China **9**(1), 1–32 (2021)
48. Tominin, V., Tominin, Y., Borodich, E., Kovalev, D., Gasnikov, A., Dvurechensky, P.: On accelerated methods for saddle-point problems with composite structure. arXiv preprint arXiv:2103.09344 (2021)
49. Tseng, P.: A modified forward-backward splitting method for maximal monotone mappings. SIAM Journal on Control and Optimization **38**(2), 431–446 (2000)
50. Vandenberghe, L.: Generalized distances and mirror descent. Lecture notes (2022)
51. Yang, J., Kiyavash, N., He, N.: Global convergence and variance-reduced optimization for a class of nonconvex-nonconcave minimax problems. arXiv preprint arXiv:2002.09621 (2020)
52. Yoon, T., Ryu, E.K.: Accelerated algorithms for smooth convex-concave minimax problems with $\mathcal{O}(1/k^2)$ rate on squared gradient norm. In: International Conference on Machine Learning. pp. 12098–12109. PMLR (2021)
53. You, Z., Zhang, H.: A prediction–correction admm for multistage stochastic variational inequalities. Journal of Optimization Theory and Applications **199**(2), 693–731 (2023)

A Basic facts

Lemma 4. For all $u, v \in \mathbb{R}^n$ the Cauchy–Schwarz inequality holds:

$$|\langle u, v \rangle|^2 \leq \langle u, u \rangle \cdot \langle v, v \rangle. \quad (43)$$

Lemma 4. For all $u, v \in \mathbb{R}^n$, for any $p > 0$ and $q > 0$ such that $\frac{1}{p} + \frac{1}{q} = 1$ the Young’s inequality holds:

$$\langle u, v \rangle \leq \frac{u^p}{p} + \frac{v^q}{q}. \quad (44)$$

Lemma 4 ([50]). For all u, v and $w \in \mathcal{X}$ the three point identity holds:

$$V(u, w) = V(u, v) + V(v, w) + \langle \nabla h(v) - \nabla h(w), u - v \rangle. \quad (45)$$