

# Optimal Analysis of Method with Batching for Monotone Stochastic Finite-Sum Variational Inequalities

Alexander Pichugin<sup>1</sup>, Maksim Pechin<sup>1</sup>, Aleksandr Beznosikov<sup>1,2</sup>, and Alexander Gasnikov<sup>1,2</sup>

<sup>1</sup> Moscow Institute of Physics and Technology, Moscow, Russia,

<sup>2</sup> Institute for Information Transmission Problems, Moscow, Russia

**Abstract.** Variational inequalities are a universal optimization paradigm that is interesting in itself, but also incorporates classical minimization and saddle point problems. Modern realities encourage 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 and does not lose the oracle complexity optimality. The effectiveness of the algorithm, especially in the case of small but not single batches is confirmed experimentally.

**Keywords:** stochastic optimization · variational inequalities · finite-sum problems · batching

## 1 Introduction

In this paper, we consider the following variational inequality 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  $\text{dom}g$ . Variational inequalities are a popular optimization formulation. This is due to the fact that they incorporate many different problem statements that arise widely in various fields of applied science. To understand the importance of variational inequalities, let us consider two important examples.

*Example 1 (Minimization).* Consider the classic and widely known convex regularized minimization problem:

$$\min_{x \in \mathbb{R}^d} f(x) + g(x), \quad (2)$$

where  $f$  is typically a main target function, and  $g$  a regularizer. If we put  $F(x) := \nabla f(x)$ , then one can prove that  $x^* \in \text{dom}g$  is a solution for (1) if and only if  $x^* \in \text{dom}g$  is a solution for (2) [13, Section 1.4.1].

*Example 2 (Saddle point).* Consider the convex-concave saddle point problem:

$$\min_{x_1 \in \mathbb{R}^{d_1}} \max_{x_2 \in \mathbb{R}^{d_2}} f(x_1, x_2) + g_1(x_1) - g_2(x_2), \quad (3)$$

where  $g_1$  and  $g_2$  can also be interpreted as regularizers. If we define  $F(x) := F(x_1, x_2) := [\nabla_{x_1} f(x_1, x_2), -\nabla_{x_2} f(x_1, x_2)]$  and  $g(x) = g(x_1, x_2) = g_1(x_1) + g_2(x_2)$ , then it can be proved that  $x^* \in \text{dom } g$  is a solution for (1) if and only if  $z^* \in \text{dom } g$  is a solution for (3) [13, Section 1.4.1]. It gives that convex-concave saddle point problems (3) can be investigated via a reformulation in the form of the variational inequality (1).

While the minimization problems are usually considered separately from the paradigm of variational inequalities, the saddle point problems are very often studied in the generality of variational inequalities. But even excluding minimization problems, variational inequalities themselves [36] and their special case of saddle problems have plenty of real-world applications. First of all, it is important to note the classic examples of economics and game theory, which had their golden age in the middle of the last century [31,17,13]. It is also important to note also the classic but newer stories from robust optimization [6,30,27] and supervised/unsupervised learning [19,4,38,5,12,8]. And the cherry on top of the enduring popularity of variational inequalities are all-new applications from reinforcement learning [33,18], adversarial training [24], and generative models [15,10,14,25].

Meanwhile, every year modern applied problems become more and more complicated and very often have a stochastic nature:  $F(x) = \mathbb{E}_{\xi \sim \mathcal{D}}[F_\xi(x)]$ . Often such mathematical expectation cannot be represented in a closed form due to the fact that the distribution  $\mathcal{D}$  is non-trivial or even unknown, therefore one often resorts to the Monte Carlo approach, which approximates the expectation integral using a set of samples from the distribution  $\mathcal{D}$ :

$$F(x) = \frac{1}{M} \sum_{m=1}^M F_m(x). \quad (4)$$

In algorithms designed to solve stochastic problems [35,20,11,32], one often avoids calling the full operators/gradients, and resorts to the batching technique (e.g., we can choose one or more terms from (4)). Therefore, it is very important that the algorithm is robust to the use of any size batches, both in terms of theory and in terms of practice. But even for minimization problems, not all stochastic methods satisfy such properties [3, Section 5]. This brings us to the objective of this paper:

*Develop an algorithm for solving the monotone stochastic (finite-sum) variational inequality with Lipschitz summands (1)+(4). In terms of theory, the algorithm should be non-sensitive to the size of batches.*

**Our contribution and related works.** The theory of solving monotone variational inequalities with Lipschitz operators has an extensive history. A natural idea to get a method is to use a gradient descent with replacement of the

gradient  $\nabla f(x)$  by the operator  $F(x)$ . But such a method for monotone operators may diverge. A breakthrough was the creation of the ExtraGradient method [22], which was widely theoretically investigated [26,27] and subsequently received various modifications/analogues [34,37], including the stochastic case [21]. Regarding the methods for monotone stochastic finite-sum problems, one can highlight papers [7,1,2] where the authors adapt the variance reduction technique to variational inequalities. Moreover, the results from [1] are optimal [16]. But this is valid only for the case when each iteration we call 1 term from the sum (4). As soon as we use a non-single-batch, the results from works [7,1] become worse. See Table 1 for details.

We solve this issue and propose a method that is optimal for any batch size. Our method is based on the optimistic scheme with momentum from [23], but in this paper, the authors studied only strongly monotone variational inequalities, while we consider more general and more complicated from the point of view of theoretical analysis monotone ones.

Table 1: Summary complexities for finding an  $\varepsilon$ -solution for the monotone stochastic (finite-sum) variational inequality with Lipschitz terms (1)+(4). Convergence is measured by the gap function. *Notation:*  $\mu$  = constant of strong monotonicity of the operator  $F$ ,  $L$  = Lipschitz constants for all  $F_i$ ,  $n$  = the size of the local dataset,  $b$  = the batch size per iteration

Reference	Complexity	Additional assumptions
Nemirovski et al. [27] <sup>(1)</sup>	$\mathcal{O}\left(n\frac{L}{\varepsilon}\right)$	
Carmon et al. [7]	$\tilde{\mathcal{O}}\left(\sqrt{bn}\frac{L}{\varepsilon}\right)$	$x \rightarrow \langle F(x) + \nabla g(x), x - u \rangle$ is convex for any $u$ <b>or</b> bounded domain
Alacaoglu and Malitsky [1]	$\mathcal{O}\left(\sqrt{bn}\frac{L}{\varepsilon}\right)$	
Alacaoglu et al. [2]	$\mathcal{O}\left(n\frac{L}{\varepsilon}\right)$	
This paper	$\mathcal{O}\left(\sqrt{n}\frac{L}{\varepsilon}\right)$	
Han et al. [16]	$\Omega\left(\sqrt{n}\frac{L}{\varepsilon}\right)$	lower bounds

<sup>(1)</sup> deterministic methods, similar results were also obtained in [37,26]

## 2 Problem Setup and Assumptions

We consider the problem (1)+(4), where  $\mathcal{X}$  be a finite dimensional vector space with Euclidean inner product  $\langle \cdot, \cdot \rangle$  and induced norm  $\|\cdot\|$ . For a proper convex lower semicontinuous  $g$ , we denote domain as  $\text{dom}g = \{x : g(x) < +\infty\}$ . We

also assume that the function  $g$  is proximally friendly, i.e. the following operator  $\text{prox}_{\alpha g}(x) = \arg \min_{y \in \mathcal{X}} \{ \alpha g(y) + \frac{1}{2} \|y - x\|^2 \}$  with any  $\alpha > 0$  can be exactly calculated for free. For  $y = \text{prox}_{\alpha g}(x)$  it holds that

$$x - y \in \partial(\alpha g)(y). \quad (5)$$

By  $\mathbb{E}_k$  we denote a conditional expectation  $\mathbb{E}[\cdot | S^{k-1}, S^{k-2}, \dots, S^0]$ . Each operator  $F_i : \text{dom } g \rightarrow \mathbb{R}^d$  is single-valued. The key assumptions of our paper we give as follows

**Assumption 1**

- The solution (may be not unique) for the problem (1)+(4) exists.
- The operator  $F$  is monotone, i.e. for all  $u, v \in \mathbb{R}^d$  we have

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

- The operator  $F$  is  $L$ -Lipschitz, i.e. for all  $u, v \in \mathbb{R}^d$  we have

$$\|F(u) - F(v)\| \leq L \|u - v\|.$$

- Each operator  $F_m$  is  $L_m$ -Lipschitz. We define  $\bar{L}$  as  $\bar{L}^2 = \frac{1}{M} \sum_{m=1}^M L_m^2$ .

### 3 Main Part

We base our method on the non-distributed version of Algorithm 1 from [23]. This algorithm is based on the so-called optimistic modification [34] of the ExtraGradient method [22]. The essence of this modification is to use the value of the operator not only at the points of the current iteration ( $x^k$  and  $w^k$ ), but also from the past iteration ( $x^{k-1}$  and  $w^{k-1}$ ). We also use the variance reduction technique [20,1], for this we introduce an additional sequence  $w^k$ . The point  $w^k$  (the reference point) changes rarely (if  $p$  is small), and hence the full operator  $F(w^k)$  is almost not recalculated. We also apply the random negative momentum  $\gamma(w^k - x^k)$ , which is necessary to use the variance reduction approach in methods for variational inequalities [1,2].

We use the gap function as convergence criterion:

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

Here we do not take the maximum over the entire set  $\mathcal{X}$ , but over  $\mathcal{C}$  – a compact subset of  $\mathcal{X}$ . Thus, we can also consider unbounded sets  $\mathcal{X} \subseteq \mathbb{R}^d$ . This is permissible, since such a version of the criterion is valid if the solution  $x^*$  lies in  $\mathcal{C}$ ; for details see [30]. The following theorem gives the convergence of the proposed method.

---

**Algorithm 1** Optimistic Method with Momentum and Batching
 

---

- 1: **Parameters:** stepsize  $\eta > 0$ , momentum  $\gamma > 0$ , probability  $p \in (0, 1)$ , batchsize  $b \in \{1, \dots, M\}$ , number of iterations  $K$
  - 2: **Initialization:** choose  $x^0 = w^0 = x^{-1} = w^{-1} \in \mathbb{R}^d$
  - 3: **for**  $k = 0, 1, \dots$ , **do**
  - 4:   Sample  $j_1^k, \dots, j_b^k$  independently from  $\{1, \dots, M\}$  uniformly at random
  - 5:    $S^k = \{j_1^k, \dots, j_b^k\}$
  - 6:    $\Delta^k = \frac{1}{b} \sum_{j \in S^k} (F_j(x^k) - F_j(w^{k-1}) + (F_j(x^k) - F_j(x^{k-1}))) + F(w^{k-1})$
  - 7:    $x^{k+1} = \text{prox}_{\eta g}(x^k + \gamma(w^k - x^k) - \eta \Delta^k)$
  - 8:    $w^{k+1} = \begin{cases} x^{k+1}, & \text{with probability } p \\ w^k, & \text{with probability } 1 - p \end{cases}$
  - 9: **end for**
- 

**Theorem 1.** Consider the problem (1)+(4) under Assumptions 1. Let  $\{x^k\}$  be the sequence generated by Algorithm 1 with tuning of  $\eta, \theta, \alpha, \beta, \gamma$  as follows:

$$0 < p = \gamma \leq \frac{1}{16}, \quad \eta = \min \left\{ \frac{\sqrt{\gamma b}}{8\bar{L}}, \frac{1}{8L} \right\}$$

Then, given  $\varepsilon > 0$ , the number of iterations for  $\mathbb{E}[\text{Gap}(x^k)] \leq \varepsilon$  is

$$\mathcal{O} \left( \frac{1}{\sqrt{pb}} \frac{\bar{L}}{\varepsilon} + \frac{L}{\varepsilon} \right).$$

See the proof in Appendix A. Theorem gives iteration complexity of Algorithm 1, but it is primarily the oracle complexity (the number of calls of terms  $F_i$ ) that is of interest. Note that Theorem gives iterationionic complexity of Algorithm 1, but it is primarily the main point of interest is the oracle complexity (the number of calls to terms  $F_i$ ). One can note that at each iteration we call  $\mathcal{O}(b + pn)$  terms in average/expectation – each time we call the batch with size of  $b$  and with probability  $p$  we compute the full operator. From where we immediately get the optimal choice for  $p$ :

**Corollary 1.** Under the conditions of Theorem 1, if we choose  $p = \frac{b}{n}$ , then we get the following oracle complexities

$$\mathcal{O} \left( \sqrt{n} \frac{\bar{L}}{\varepsilon} + b \frac{L}{\varepsilon} \right).$$

With  $b \leq \frac{\bar{L}\sqrt{n}}{L}$ , we get  $\mathcal{O} \left( \sqrt{n} \frac{\bar{L}}{\varepsilon} \right)$ . And this results is optimal and unimproveable [16].

## 4 Experiments

In this section, we aim to test the performance of Algorithm 1 in practice.

We consider the bilinear problem:

$$\min_{x \in \Delta^d} \max_{y \in \Delta^d} x^\top Ay, \quad (7)$$

where  $\Delta^d$  is the unit simplex in  $\mathbb{R}^d$ . We use the same experimental setup as in [1], in particular we consider the policeman and burglar matrix from [28] and the first test matrix from [29]. Note that the problem (7) does not have the finite-sum form as (1)+(4), but we can rewrite  $A$  from (7) as follows  $A = \sum_{i=1}^d A_{i\cdot}$  or  $A = \sum_{i=1}^d A_{\cdot i}$ , where  $A_{i\cdot}$  is the  $i$ th row of  $A$  and  $A_{\cdot i}$  is the  $i$ th column of  $A$  – see details in Section 5.1.2 from [1].

For comparison, we take methods from Table 1. In particular, we choose Algorithms 1 and 2 from [1], Algorithm 1+2 from [7]. All algorithms are considered in the Euclidean setting with the projection to simplex from [9].

The parameters of all methods are tuned for the best convergence among the theoretical possible – see Section 6 from [1]. We run all methods with different batch sizes. We use duality gap (6) as the convergence measure, it can be simply computed as  $[\max_i (A^T x)_i - \min_j (Ay)_j]$  for simplex constraints. The comparison criterion is the number of operations (one operation is computationally equal to calculations of  $Ay$  and  $A^T x$ ).

The results are reflected in Figures 1 and 2. They show that Algorithm 1 outperforms all competitors.

Fig. 1: Comparison of computational complexities for Algorithm 1, EG-Ma122-1 (Algorithm 1 [1]), EG-Ma122-2 (Algorithm 2 [1]), EG-Car19 (Algorithm 1+2 from [7]) on (7) with policeman and burglar matrix from [28].

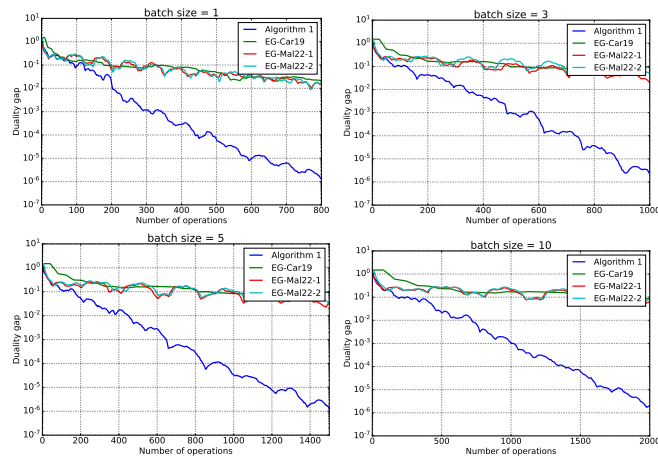
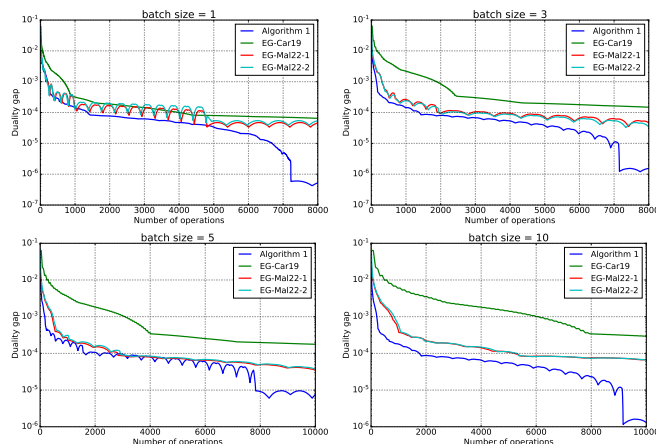


Fig. 2: Comparison of computational complexities for Algorithm 1, EG-Ma122-1 (Algorithm 1 [1]), EG-Ma122-2 (Algorithm 2 [1]), EG-Car19 (Algorithm 1+2 from [7]) on (7) with policeman and burglar matrix from [28].



## 5 Conclusion

We presented an algorithm for solving monotone stochastic (finite-sum) variational inequalities with Lipschitz terms. From the theoretical perspective, the algorithm is optimal for almost any batches size. As directions for future work, it would be interesting to get a version of Algorithm 1 in the non-Euclidean case with arbitrary Bregman divergence.

## Acknowledgements

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.

## References

1. Alacaoglu, A., Malitsky, Y.: Stochastic variance reduction for variational inequality methods. arXiv preprint arXiv:2102.08352 (2021)
2. Alacaoglu, A., Malitsky, Y., Cevher, V.: Forward-reflected-backward method with variance reduction. *Computational Optimization and Applications* **80** (11 2021). <https://doi.org/10.1007/s10589-021-00305-3>

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. Bach, F., Jenatton, R., Mairal, J., Obozinski, G.: Optimization with sparsity-inducing penalties. arXiv preprint arXiv:1108.0775 (2011)
5. Bach, F., Mairal, J., Ponce, J.: Convex sparse matrix factorizations. arXiv preprint arXiv:0812.1869 (2008)
6. Ben-Tal, A., Ghaoui, L.E., Nemirovski, A.: Robust Optimization. Princeton University Press (2009)
7. Carmon, Y., Jin, Y., Sidford, A., Tian, K.: Variance reduction for matrix games. arXiv preprint arXiv:1907.02056 (2019)
8. Chambolle, A., Pock, T.: A first-order primal-dual algorithm for convex problems with applications to imaging. *Journal of mathematical imaging and vision* **40**(1), 120–145 (2011)
9. Condat, L.: Fast projection onto the simplex and the  $l_1$  ball. *Mathematical Programming* **158**(1-2), 575–585 (2016)
10. Daskalakis, C., Ilyas, A., Syrgkanis, V., Zeng, H.: Training gans with optimism. arXiv preprint arXiv:1711.00141 (2017)
11. Defazio, A., Bach, F., Lacoste-Julien, S.: Saga: A fast incremental gradient method with support for non-strongly convex composite objectives. In: Advances in neural information processing systems. pp. 1646–1654 (2014)
12. Esser, E., Zhang, X., Chan, T.F.: A general framework for a class of first order primal-dual algorithms for convex optimization in imaging science. *SIAM Journal on Imaging Sciences* **3**(4), 1015–1046 (2010)
13. Facchinei, F., Pang, J.S.: Finite-Dimensional Variational Inequalities and Complementarity Problems. Springer Series in Operations Research, Springer (2003)
14. Gidel, G., Berard, H., Vignoud, G., Vincent, P., Lacoste-Julien, S.: A variational inequality perspective on generative adversarial networks. arXiv preprint arXiv:1802.10551 (2018)
15. Goodfellow, I.J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial networks (2014)
16. Han, Y., Xie, G., Zhang, Z.: Lower complexity bounds of finite-sum optimization problems: The results and construction. arXiv preprint arXiv:2103.08280 (2021)
17. Harker, P.T., Pang, J.S.: Finite-dimensional variational inequality and nonlinear complementarity problems: a survey of theory, algorithms and applications. *Mathematical programming* (1990)
18. Jin, Y., Sidford, A.: Efficiently solving MDPs with stochastic mirror descent. In: Proceedings of the 37th International Conference on Machine Learning (ICML). vol. 119, pp. 4890–4900. PMLR (2020)
19. Joachims, T.: A support vector method for multivariate performance measures. pp. 377–384 (01 2005). <https://doi.org/10.1145/1102351.1102399>
20. Johnson, R., Zhang, T.: Accelerating stochastic gradient descent using predictive variance reduction. vol. 26, pp. 315–323 (2013)
21. Juditsky, A., Nemirovski, A., Tauvel, C.: Solving variational inequalities with stochastic mirror-prox algorithm. *Stochastic Systems* **1**(1), 17–58 (2011)
22. 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)
23. 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)



24. Madry, A., Makelov, A., Schmidt, L., Tsipras, D., Vladu, A.: Towards deep learning models resistant to adversarial attacks. In: International Conference on Learning Representations (ICLR) (2018)
25. Mertikopoulos, P., Lecouat, B., Zenati, H., Foo, C.S., Chandrasekhar, V., Piliouras, G.: Optimistic mirror descent in saddle-point problems: Going the extra (gradient) mile. arXiv preprint arXiv:1807.02629 (2018)
26. 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)
27. 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)
28. Nemirovski, A.: Mini-course on convex programming algorithms (2013)
29. Nemirovski, A., Juditsky, A., Lan, G., Shapiro, A.: Robust stochastic approximation approach to stochastic programming. *SIAM Journal on optimization* **19**(4), 1574–1609 (2009)
30. Nesterov, Y.: Dual extrapolation and its applications to solving variational inequalities and related problems. *Mathematical Programming* **109**(2), 319–344 (2007)
31. Neumann, J.V., Morgenstern, O.: Theory of games and economic behavior. Princeton University Press (1944)
32. Nguyen, L.M., Liu, J., Scheinberg, K., Takáč, M.: SARAH: a novel method for machine learning problems using stochastic recursive gradient. In: International Conference on Machine Learning. pp. 2613–2621. PMLR (2017)
33. Omidshafiei, S., Pazis, J., Amato, C., How, J.P., Vian, J.: Deep decentralized multi-task multi-agent reinforcement learning under partial observability. In: Proceedings of the 34th International Conference on Machine Learning (ICML). vol. 70, pp. 2681–2690. PMLR (2017), <http://proceedings.mlr.press/v70/omidshafiei17a.html>
34. 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)
35. Robbins, H., Monro, S.: A Stochastic Approximation Method. *The Annals of Mathematical Statistics* **22**(3), 400 – 407 (1951)
36. Scutari, G., Palomar, D.P., Facchinei, F., Pang, J.S.: Convex optimization, game theory, and variational inequality theory. *IEEE Signal Processing Magazine* **27**(3), 35–49 (2010)
37. Tseng, P.: A modified forward-backward splitting method for maximal monotone mappings. *SIAM Journal on Control and Optimization* **38**(2), 431–446 (2000)
38. Xu, L., Neufeld, J., Larson, B., Schuurmans, D.: Maximum margin clustering. In: Advances in Neural Information Processing Systems. vol. 17, pp. 1537–1544 (2005)

## A Proof of Theorem 1

Before proving the theorem we introduce the following Lemmas.

**Lemma 1.** [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  $\mathbb{E}[u^{k+1}|F_k] = 0$ . Then for any  $K \in \mathbb{N}$ ,  $x^0 \in X$  and any compact set  $\mathcal{C} \subseteq X$

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

**Lemma 2.** Under Assumption 1 for iterates of Algorithm 1 the following inequality holds:

$$\mathbb{E} \left[ \|\Delta^k - \mathbb{E}_k [\Delta^k]\|^2 \right] \leq \frac{2\bar{L}^2}{b} \mathbb{E} \left[ \|x^k - w^{k-1}\|^2 + \|x^k - x^{k-1}\|^2 \right], \quad (8)$$

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

$$\mathbb{E}_k [\Delta^k] = 2F(x^k) - F(x^{k-1}). \quad (9)$$

*Proof.* We start from line 6 of Algorithm 1

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

With Cauchy–Schwarz inequality, we have

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

Using we choose  $j_1^k, \dots, j_b^k$  in  $S^k$  independently and uniformly, one can note that

$$\begin{aligned} & \mathbb{E}_k \left[ \left\langle (F_j(x^k) - F_j(w^{k-1})) - (F(x^k) - F(w^{k-1})), \right. \right. \\ & \quad \left. \left. (F_j(x^k) - F_j(w^{k-1})) - (F(x^k) - F(w^{k-1})) \right\rangle \right] \end{aligned}$$

$$\begin{aligned}
 &= \mathbb{E}_k \left[ \langle \mathbb{E}_{j_i^k} \left[ \left( F_{j_i^k} (x^k) - F_{j_i^k} (w^{k-1}) \right) - (F(x^k) - F(w^{k-1})) \right] \right. \\
 &\quad \left. \mathbb{E}_{j_t^k} \left[ \left( F_{j_t^k} (x^k) - F_{j_t^k} (w^{k-1}) \right) - (F(x^k) - F(w^{k-1})) \right] \rangle \right] \\
 &= 0.
 \end{aligned}$$

Hence, we get

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

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  $j_1^k, \dots, j_b^k$  in  $S^k$  are chosen uniformly

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

Since each operator  $F_j$  is  $L_j$ -Lipschitz, we can rewrite it as

$$\mathbb{E}_k \left[ \|\Delta^k - \mathbb{E}_k [\Delta^k]\|^2 \right] \leq \frac{2}{mb} \sum_{j=1}^m L_j^2 (\|x^k - w^{k-1}\|^2 + \|x^k - x^{k-1}\|^2).$$

Applying the definition of  $\bar{L}$ , we obtain

$$\mathbb{E}_k \left[ \|\Delta^k - \mathbb{E}_k [\Delta^k]\|^2 \right] \leq \frac{2\bar{L}^2}{b} (\|x^k - w^{k-1}\|^2 + \|x^k - x^{k-1}\|^2).$$

Taking the full expectation concludes the proof.  $\square$

**Lemma 3.** *For iterates of Algorithm 1 with  $\gamma = p$  the following bound holds for any compact set  $\mathcal{C} \subseteq X$ :*

$$\mathbb{E} \left[ \max_{x \in \mathcal{C}} \sum_{k=0}^{K-1} e_1(x, k) \right] \leq 2 \max_{x \in \mathcal{C}} \{ \|x - x^0\|^2 \} + \frac{\gamma(1-\gamma)}{2} \sum_{k=0}^{K-1} \mathbb{E} \|x^{k+1} - \omega^k\|^2,$$

where  $e_1(k, x) = \|w^{k+1} - x\|^2 - \|w^k - x\|^2 + (1-\gamma)\|x^{k+1} - x\|^2$ .

*Proof.* For shortness we introduce

$$u^{k+1} = \gamma x^{k+1} + (1-\gamma)\omega^k - \omega^{k+1}.$$

With new notation, we can rewrite  $e_1(k, x)$  as

$$e_1(k, x) = 2 \langle u^{k+1}, x \rangle - \gamma \|x^{k+1}\|^2 - (1-\gamma)\|w^k\|^2 + \|w^k\|^2.$$

From line 8 of Algorithm 1 and using that  $\gamma = p$ , one can obtain

$$\mathbb{E} [\mathbb{E}_k [\|\omega^{k+1}\|^2 - \gamma \|x^{k+1}\|^2 - (1-\gamma)\|\omega^k\|^2]] = 0.$$

Using two properties above, we reach for any compact set  $\mathcal{C} \subseteq X$

$$\begin{aligned} \mathbb{E} \left[ \max_{x \in \mathcal{C}} \sum_{k=0}^{K-1} e_1(k, x) \right] &= 2 \mathbb{E} \left[ \max_{x \in \mathcal{C}} \sum_{k=0}^{K-1} \langle u^{k+1}, x \rangle \right] + \mathbb{E} [ -\gamma \|x^{k+1}\|^2 \\ &\quad - (1-\gamma)\|w^k\|^2 + \|w^k\|^2 ] \\ &= 2 \mathbb{E} \left[ \max_{x \in \mathcal{C}} \sum_{k=0}^{K-1} \langle u^{k+1}, x \rangle \right]. \end{aligned}$$

With  $\mathcal{F}_k = \sigma(\xi_0, \dots, \xi_k, x^k)$  we have that  $\mathbb{E}[u^{k+1} | \mathcal{F}_k] = 0$ . It means that we can apply Lemma 1. Thus, we get

$$\mathbb{E} \left[ \max_{x \in \mathcal{C}} \sum_{k=0}^{K-1} e_1(k, x) \right] \leq 2 \max_{x \in \mathcal{C}} \|x_0 - x\|^2 + \frac{1}{2} \sum_{k=0}^{K-1} \mathbb{E} \|u^{k+1}\|^2. \quad (10)$$

We estimate  $\|u_{k+1}\|^2$  using the fact that  $\mathbb{E}\|X - \mathbb{E}X\|^2 = \mathbb{E}\|X\|^2 - \|\mathbb{E}X\|^2$  and line 8 of Algorithm 1:

$$\begin{aligned} \mathbb{E}\|u_{k+1}\|^2 &= \mathbb{E}\left[\mathbb{E}_k\|u_{k+1}\|^2\right] = \mathbb{E}\left[\mathbb{E}_k\|\mathbb{E}_k[\omega_{k+1}] - \omega_{k+1}\|^2\right] \\ &= \mathbb{E}\left[\mathbb{E}_k\|\omega_{k+1}\|^2 - \|\mathbb{E}_k[\omega_{k+1}]\|^2\right] \\ &= \mathbb{E}\left[\gamma\|x_{k+1}\|^2 + (1-\gamma)\|\omega_k\|^2 - \|\gamma x_{k+1} + (1-\gamma)\omega_k\|^2\right] \\ &= \gamma(1-\gamma)\mathbb{E}\|x_{k+1} - \omega_k\|^2. \end{aligned}$$

Applying this result to (10), we get

$$\mathbb{E}\left[\max_{x \in \mathcal{C}} \sum_{k=0}^{K-1} e_1(k, x)\right] = 2 \max_{x \in \mathcal{C}} \|x_0 - x\|^2 + \frac{\gamma(1-\gamma)}{2} \sum_{k=0}^{K-1} \mathbb{E}\|x_{k+1} - \omega_k\|^2.$$

*Proof of Theorem 1.* We start from

$$\begin{aligned} \|x^{k+1} - x\|^2 &= \|x^k - x\|^2 + 2\langle x^{k+1} - x^k, x^{k+1} - x \rangle - \|x^{k+1} - x^k\|^2 \\ &= \|x^k - x\|^2 + 2\gamma\langle w^k - x^k, x^{k+1} - x \rangle - 2\eta\langle \Delta^k, x^{k+1} - x \rangle \\ &\quad - \|x^{k+1} - x^k\|^2 - 2[\langle x^k + \gamma(w^k - x^k) - \eta\Delta^k - x^{k+1}, x^{k+1} - x \rangle]. \end{aligned}$$

From line 7 of Algorithm 1 and according to the property (5) of proximal operator, it follows, that

$$x^k + \gamma(w^k - x^k) - \eta\Delta^k - x^{k+1} \in \partial(\eta g)(x^{k+1}).$$

From convexity of  $g(\cdot)$ , we obtain

$$\begin{aligned} \|x^{k+1} - x\|^2 &\leq \|x^k - x\|^2 + 2\gamma\langle w^k - x^k, x^{k+1} - x \rangle - 2\eta\langle \Delta^k, x^{k+1} - x \rangle \\ &\quad - \|x^{k+1} - x^k\|^2 + 2\eta g(x) - 2\eta g(x^{k+1}). \end{aligned}$$

Using  $2\gamma\langle w^k - x^k, x^{k+1} - x \rangle = 2\gamma\langle w^k - x, x^{k+1} - x \rangle - 2\gamma\langle x^k - x, x^{k+1} - x \rangle$  and the following property of scalar product:  $2\langle a, b \rangle = \|a + b\|^2 - \|a\|^2 - \|b\|^2$ , we get

$$\begin{aligned} \|x^{k+1} - x\|^2 &\leq \|x^k - x\|^2 + \gamma(\|w^k - x\|^2 + \|x^{k+1} - x\|^2 - \|x^{k+1} - w^k\|^2) \\ &\quad - 2\eta\langle \Delta^k, x^{k+1} - x \rangle - \gamma\|x^{k+1} - x\|^2 - \gamma\|x^k - x\|^2 \\ &\quad + \gamma\|x^{k+1} - x^k\|^2 - \|x^{k+1} - x^k\|^2 + 2\eta g(x) - 2\eta g(x^{k+1}) \\ &= \|x^k - x\|^2 + \gamma\|w^k - x\|^2 - \gamma\|x^k - x\|^2 - \gamma\|x^{k+1} - w^k\|^2 \\ &\quad - 2\eta\langle \Delta^k, x^{k+1} - x \rangle - (1-\gamma)\|x^{k+1} - x^k\|^2 + 2\eta g(x) - 2\eta g(x^{k+1}). \end{aligned}$$

Applying the properties of  $\mathbb{E}_k[\Delta^k]$  specified in Lemma 2, we obtain

$$\begin{aligned} \|x^{k+1} - x\|^2 &\leq \|x^k - x\|^2 + \gamma\|w^k - x\|^2 - \gamma\|x^k - x\|^2 - \gamma\|w^k - x^{k+1}\|^2 \\ &\quad - 2\eta\langle \mathbb{E}_k[\Delta^k], x^{k+1} - x \rangle - (1-\gamma)\|x^{k+1} - x^k\|^2 \end{aligned}$$

$$\begin{aligned}
& + 2\eta \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x^{k+1} - x^k \rangle + 2\eta \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x^k - x \rangle \\
& + 2\eta g(x) - 2\eta g(x^{k+1}) \\
= & \|x^k - x\|^2 + \gamma \|w^k - x\|^2 - \gamma \|x^k - x\|^2 - \gamma \|w^k - x^{k+1}\|^2 \\
& - 2\eta \langle F(x^k) + F(x^k) - F(x^{k-1}), x^{k+1} - x \rangle \\
& - (1 - \gamma) \|x^{k+1} - x^k\|^2 + 2\eta \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x^{k+1} - x^k \rangle \\
& + 2\eta \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x^k - x \rangle \\
& + 2\eta g(x) - 2\eta g(x^{k+1}) \\
= & \|x^k - x\|^2 + \gamma \|w^k - x\|^2 - \gamma \|x^k - x\|^2 - \gamma \|w^k - x^{k+1}\|^2 \\
& - 2\eta \langle F(x^k) - F(x^{k+1}) + F(x^k) - F(x^{k-1}), x^{k+1} - x \rangle \\
& - 2\eta \langle F(x^{k+1}), x^{k+1} - x \rangle \\
& - (1 - \gamma) \|x^{k+1} - x^k\|^2 + 2\eta \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x^{k+1} - x^k \rangle \\
& + 2\eta \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x^k - x \rangle + 2\eta g(x) - 2\eta g(x^{k+1}).
\end{aligned}$$

By a simple rearrangements, we obtain

$$\begin{aligned}
& 2\eta(g(x^{k+1}) - g(x)) + 2\eta \langle F(x^{k+1}), x^{k+1} - x \rangle \\
& \leq \|x^k - x\|^2 + \gamma \|w^k - x\|^2 - \gamma \|x^k - x\|^2 \\
& \quad - \gamma \|w^k - x^{k+1}\|^2 - \|x^{k+1} - x\|^2 \\
& \quad - 2\eta \langle F(x^k) - F(x^{k+1}) + F(x^k) - F(x^{k-1}), x^{k+1} - x \rangle \\
& \quad - (1 - \gamma) \|x^{k+1} - x^k\|^2 + 2\eta \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x^{k+1} - x^k \rangle \\
& \quad + 2\eta \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x^k - x \rangle \\
= & (1 - \gamma) \|x^k - x\|^2 + \|w^k - x\|^2 \\
& \quad - (1 - \gamma) \|x^{k+1} - x\|^2 - \|w^{k+1} - x\|^2 \\
& \quad + \|w^{k+1} - x\|^2 - \gamma \|x^{k+1} - x\|^2 \\
& \quad - (1 - \gamma) \|w^k - x\|^2 - \gamma \|w^k - x^{k+1}\|^2 \\
& \quad - 2\eta \langle F(x^k) - F(x^{k+1}), x^{k+1} - x \rangle \\
& \quad + 2\eta \langle F(x^{k-1}) - F(x^k), x^{k+1} - x \rangle \\
& \quad - (1 - \gamma) \|x^{k+1} - x^k\|^2 + 2\eta \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x^{k+1} - x^k \rangle \\
& \quad + 2\eta \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x^k - x \rangle \\
= & (1 - \gamma) \|x^k - x\|^2 + \|w^k - x\|^2 \\
& \quad - (1 - \gamma) \|x^{k+1} - x\|^2 - \|w^{k+1} - x\|^2 \\
& \quad + \|w^{k+1} - x\|^2 - \gamma \|x^{k+1} - x\|^2 \\
& \quad - (1 - \gamma) \|w^k - x\|^2 - \gamma \|w^k - x^{k+1}\|^2
\end{aligned}$$

$$\begin{aligned}
 & -2\eta \langle F(x^k) - F(x^{k+1}), x^{k+1} - x \rangle \\
 & + 2\eta \langle F(x^{k-1}) - F(x^k), x^k - x \rangle \\
 & + 2\eta \langle F(x^{k-1}) - F(x^k), x^{k+1} - x^k \rangle \\
 & - (1 - \gamma) \|x^{k+1} - x^k\|^2 + 2\eta \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x^{k+1} - x^k \rangle \\
 & + 2\eta \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x^k - x \rangle.
 \end{aligned}$$

After taking sum and then averaging, one can get

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

$$\begin{aligned}
& + 2\eta \cdot \frac{1}{K} \sum_{k=0}^{K-1} \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x^k - x \rangle \\
& - \frac{\gamma}{K} \sum_{k=0}^{K-1} \|w^k - x^{k+1}\|^2 - \frac{1-\gamma}{K} \sum_{k=0}^{K-1} \|x^{k+1} - x^k\|^2 \\
& + 2\eta \cdot \frac{1}{K} \sum_{k=0}^{K-1} \langle F(x^{k-1}) - F(x^k), x^{k+1} - x^k \rangle \\
& + 2\eta \cdot \frac{1}{K} \sum_{k=0}^{K-1} \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x^{k+1} - x^k \rangle. \tag{11}
\end{aligned}$$

Here we also used the initialization of Algorithm 1 with  $w^0 = x^{-1} = x^0$ . Applying Young's inequality, using the  $L$ -Lipshetzness of  $F$ , and taking into account the definition of  $\eta \leq \frac{1}{8L}$  from conditions of the theorem for any  $k$ , one can obtain

$$\begin{aligned}
-2\eta \langle F(x^{K-1}) - F(x^K), x^K - x \rangle & \leq 2\eta^2 \|F(x^{K-1}) - F(x^K)\|^2 + \frac{1}{2} \|x^K - x\|^2 \\
& \leq 2\eta^2 L^2 \|x^{K-1} - x^K\|^2 + \frac{1}{2} \|x^K - x\|^2 \\
& \leq \frac{1}{32} \|x^{K-1} - x^K\|^2 + \frac{1}{2} \|x^K - x\|^2. \tag{12}
\end{aligned}$$

Combining (11) and (12), we get

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



$$\begin{aligned}
 & + \frac{1}{K} \sum_{k=0}^{K-1} \left[ \|w^{k+1} - x\|^2 - \gamma \|x^{k+1} - x\|^2 - (1 - \gamma) \|w^k - x\|^2 \right] \\
 & + 2\eta \cdot \frac{1}{K} \sum_{k=0}^{K-1} \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x^k - x \rangle - \frac{\gamma}{K} \sum_{k=0}^{K-1} \|w^k - x^{k+1}\|^2 \\
 & - \frac{1 - \gamma}{K} \sum_{k=0}^{K-1} \|x^{k+1} - x^k\|^2 + \frac{1}{32K} \|x^{K-1} - x^K\|^2 \\
 & + 2\eta \cdot \frac{1}{K} \sum_{k=0}^{K-1} \langle F(x^{k-1}) - F(x^k), x^{k+1} - x^k \rangle \\
 & + 2\eta \cdot \frac{1}{K} \sum_{k=0}^{K-1} \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x^{k+1} - x^k \rangle.
 \end{aligned}$$

Next, we use monotonicity of  $F$ , apply Jensen's inequality for the convex function  $g$  and obtain

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

Using new notation  $\bar{x}^K = \frac{1}{K} \sum_{k=0}^{K-1} x^{k+1}$  and taking maximum on  $\mathcal{C}$ , we achieve

$$\begin{aligned}
 2\eta \text{Gap}(\bar{x}^K) & \leq \max_{x \in \mathcal{C}} \left\{ \frac{2 - \gamma}{K} \|x^0 - x\|^2 \right. \\
 & \left. + \frac{1}{K} \sum_{k=0}^{K-1} \left[ \|w^{k+1} - x\|^2 - \gamma \|x^{k+1} - x\|^2 - (1 - \gamma) \|w^k - x\|^2 \right] \right\}
 \end{aligned}$$

$$\begin{aligned}
& + 2\eta \cdot \frac{1}{K} \sum_{k=0}^{K-1} \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x^k - x \rangle \Big\} + \frac{1}{32K} \|x^{K-1} - x^K\|^2 \\
& - \frac{\gamma}{K} \sum_{k=0}^{K-1} \|w^k - x^{k+1}\|^2 - \frac{1-\gamma}{K} \sum_{k=0}^{K-1} \|x^{k+1} - x^k\|^2 \\
& + 2\eta \cdot \frac{1}{K} \sum_{k=0}^{K-1} \langle F(x^{k-1}) - F(x^k), x^{k+1} - x^k \rangle \\
& + 2\eta \cdot \frac{1}{K} \sum_{k=0}^{K-1} \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x^{k+1} - x^k \rangle \\
\leq & \max_{x \in \mathcal{C}} \left\{ \frac{2-\gamma}{K} \|x^0 - x\|^2 \right\} \\
& + \max_{x \in \mathcal{C}} \left\{ \frac{1}{K} \sum_{k=0}^{K-1} [\|w^{k+1} - x\|^2 - \gamma \|x^{k+1} - x\|^2 - (1-\gamma) \|w^k - x\|^2] \right\} \\
& + 2\eta \max_{x \in \mathcal{C}} \left\{ \frac{1}{K} \sum_{k=0}^{K-1} \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x^k - x \rangle \right\} + \frac{1}{32K} \|x^{K-1} - x^K\|^2 \\
& - \frac{\gamma}{K} \sum_{k=0}^{K-1} \|w^k - x^{k+1}\|^2 - \frac{1-\gamma}{K} \sum_{k=0}^{K-1} \|x^{k+1} - x^k\|^2 \\
& + 2\eta \cdot \frac{1}{K} \sum_{k=0}^{K-1} \langle F(x^{k-1}) - F(x^k), x^{k+1} - x^k \rangle \\
& + 2\eta \cdot \frac{1}{K} \sum_{k=0}^{K-1} \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x^{k+1} - x^k \rangle.
\end{aligned}$$

Here we also used that maximum of the sum not greater than the sum of the maximums. After that we take the an expectation and get

$$\begin{aligned}
2\eta \mathbb{E} [\text{Gap}(\bar{x}^K)] & \leq \mathbb{E} \left[ \max_{x \in \mathcal{C}} \left\{ \frac{2-\gamma}{K} \|x^0 - x\|^2 \right\} \right] \\
& + \mathbb{E} \left[ \max_{x \in \mathcal{C}} \left\{ \frac{1}{K} \sum_{k=0}^{K-1} [\|w^{k+1} - x\|^2 - \gamma \|x^{k+1} - x\|^2 \right. \right. \\
& \quad \left. \left. - (1-\gamma) \|w^k - x\|^2] \right\} \right] + \frac{1}{32K} \mathbb{E} [\|x^{K-1} - x^K\|^2] \\
& + 2\eta \mathbb{E} \left[ \max_{x \in \mathcal{C}} \left\{ \frac{1}{K} \sum_{k=0}^{K-1} \langle \mathbb{E}_k [\Delta^k] - \Delta^k, x^k - x \rangle \right\} \right] \\
& - \mathbb{E} \left[ \frac{\gamma}{K} \sum_{k=0}^{K-1} \|w^k - x^{k+1}\|^2 + \frac{1-\gamma}{K} \sum_{k=0}^{K-1} \|x^{k+1} - x^k\|^2 \right]
\end{aligned}$$

$$\begin{aligned}
 & + 2\eta\mathbb{E}\left[\frac{1}{K}\sum_{k=0}^{K-1}\langle F(x^{k-1}) - F(x^k), x^{k+1} - x^k \rangle\right] \\
 & + 2\eta\mathbb{E}\left[\frac{1}{K}\sum_{k=0}^{K-1}\langle \mathbb{E}_k[\Delta^k] - \Delta^k, x^{k+1} - x^k \rangle\right].
 \end{aligned}$$

With Lemma 3 for the second line of the previous estimate and Lemma 1 for the third line, we get

$$\begin{aligned}
 2\eta\mathbb{E}[\text{Gap}(\bar{x}^K)] & \leq \mathbb{E}\left[\max_{x \in \mathcal{C}} \left\{ \frac{2-\gamma}{K} \|x^0 - x\|^2 \right\}\right] \\
 & + \max_{x \in \mathcal{C}} \left\{ \frac{2}{K} \|x - x^0\|^2 \right\} + \frac{\gamma(1-\gamma)}{2K} \sum_{k=0}^{K-1} \mathbb{E}[\|x^{k+1} - \omega^k\|^2] \\
 & + \max_{x \in \mathcal{C}} \left\{ \frac{1}{K} \|x - x^0\|^2 \right\} + \frac{\eta^2}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\mathbb{E}_k[\Delta^k] - \Delta^k\|^2] \\
 & - \mathbb{E}\left[\frac{\gamma}{K} \sum_{k=0}^{K-1} \|w^k - x^{k+1}\|^2 + \frac{1-\gamma}{K} \sum_{k=0}^{K-1} \|x^{k+1} - x^k\|^2\right] \\
 & + \frac{1}{32K} \mathbb{E}[\|x^{K-1} - x^K\|^2] \\
 & + 2\eta\mathbb{E}\left[\frac{1}{K} \sum_{k=0}^{K-1} \langle F(x^{k-1}) - F(x^k), x^{k+1} - x^k \rangle\right] \\
 & + 2\eta\mathbb{E}\left[\frac{1}{K} \sum_{k=0}^{K-1} \langle \mathbb{E}_k[\Delta^k] - \Delta^k, x^{k+1} - x^k \rangle\right] \\
 & \leq \frac{4}{K} \mathbb{E}\left[\max_{x \in \mathcal{C}} \left\{ \|x^0 - x\|^2 \right\}\right] + \frac{\eta^2}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\mathbb{E}_k[\Delta^k] - \Delta^k\|^2] \\
 & - \mathbb{E}\left[\frac{\gamma}{2K} \sum_{k=0}^{K-1} \|w^k - x^{k+1}\|^2 + \frac{1-\gamma}{K} \sum_{k=0}^{K-1} \|x^{k+1} - x^k\|^2\right] \\
 & + \frac{1}{32K} \mathbb{E}[\|x^{K-1} - x^K\|^2] \\
 & + 2\eta\mathbb{E}\left[\frac{1}{K} \sum_{k=0}^{K-1} \langle F(x^{k-1}) - F(x^k), x^{k+1} - x^k \rangle\right] \\
 & + 2\eta\mathbb{E}\left[\frac{1}{K} \sum_{k=0}^{K-1} \langle \mathbb{E}_k[\Delta^k] - \Delta^k, x^{k+1} - x^k \rangle\right]. \tag{13}
 \end{aligned}$$

According to Young's inequality,

$$\mathbb{E}[2\eta\langle \mathbb{E}_k[\Delta^k] - \Delta^k, x^{k+1} - x^k \rangle] \leq 4\eta^2 \mathbb{E}[\|\mathbb{E}_k[\Delta^k] - \Delta^k\|^2] + \frac{1}{4} \mathbb{E}[\|x^{k+1} - x^k\|^2], \tag{14}$$

and

$$\mathbb{E}[2\eta\langle F(x^{k-1}) - F(x^k), x^{k+1} - x^k \rangle] \quad (15)$$

$$\leq 4\eta^2 \mathbb{E}[\|F(x^{k-1}) - F(x^k)\|^2] + \frac{1}{4} \mathbb{E}[\|x^{k+1} - x^k\|^2]. \quad (16)$$

Combining (14), (16) with (13), we obtain

$$\begin{aligned} 2\eta \mathbb{E}[\text{Gap}(\bar{x}^K)] &\leq \frac{4}{K} \mathbb{E} \left[ \max_{x \in \mathcal{C}} \{ \|x^0 - x\|^2 \} \right] + \frac{\eta^2}{K} \sum_{k=0}^{K-1} \mathbb{E} [\|\mathbb{E}_k[\Delta^k] - \Delta^k\|^2] \\ &\quad - \mathbb{E} \left[ \frac{\gamma}{2K} \sum_{k=0}^{K-1} \|w^k - x^{k+1}\|^2 + \frac{1-\gamma}{K} \sum_{k=0}^{K-1} \|x^{k+1} - x^k\|^2 \right] \\ &\quad + \frac{1}{32K} \mathbb{E} [\|x^{K-1} - x^K\|^2] + \frac{4\eta^2}{K} \sum_{k=0}^{K-1} \mathbb{E} [\|F(x^{k-1}) - F(x^k)\|^2] \\ &\quad + \frac{1}{4K} \sum_{k=0}^{K-1} \mathbb{E} [\|x^{k+1} - x^k\|^2] + \frac{4\eta^2}{K} \mathbb{E} \sum_{k=0}^{K-1} [\|\mathbb{E}_k[\Delta^k] - \Delta^k\|^2] \\ &\quad + \frac{1}{4K} \sum_{k=0}^{K-1} \mathbb{E} [\|x^{k+1} - x^k\|^2] \\ &\leq \frac{4}{K} \mathbb{E} \left[ \max_{x \in \mathcal{C}} \{ \|x^0 - x\|^2 \} \right] + \frac{5\eta^2}{K} \sum_{k=0}^{K-1} \mathbb{E} [\|\mathbb{E}_k[\Delta^k] - \Delta^k\|^2] \\ &\quad - \mathbb{E} \left[ \frac{\gamma}{2K} \sum_{k=0}^{K-1} \|w^k - x^{k+1}\|^2 + \frac{1/2 - \gamma}{K} \sum_{k=0}^{K-1} \|x^{k+1} - x^k\|^2 \right] \\ &\quad + \frac{1}{32K} \mathbb{E} [\|x^{K-1} - x^K\|^2] + \frac{4\eta^2}{K} \sum_{k=0}^{K-1} \mathbb{E} [\|F(x^{k-1}) - F(x^k)\|^2]. \end{aligned}$$

$L$  - Lipschitzness of  $F$  (Assumption 1) and the choice of  $\gamma \leq \frac{1}{L}$  give

$$\begin{aligned} 2\eta \mathbb{E}[\text{Gap}(\bar{x}^K)] &\leq \frac{4}{K} \mathbb{E} \left[ \max_{x \in \mathcal{C}} \{ \|x^0 - x\|^2 \} \right] + \frac{5\eta^2}{K} \sum_{k=0}^{K-1} \mathbb{E} [\|\mathbb{E}_k[\Delta^k] - \Delta^k\|^2] \\ &\quad - \mathbb{E} \left[ \frac{\gamma}{2K} \sum_{k=0}^{K-1} \|w^k - x^{k+1}\|^2 + \frac{1/2 - \gamma}{K} \sum_{k=0}^{K-1} \|x^{k+1} - x^k\|^2 \right] \\ &\quad + \frac{1}{32K} \mathbb{E} [\|x^{K-1} - x^K\|^2] + \frac{4\eta^2 L^2}{K} \sum_{k=0}^{K-1} \mathbb{E} [\|x^{k-1} - x^k\|^2] \\ &\leq \frac{4}{K} \mathbb{E} \left[ \max_{x \in \mathcal{C}} \{ \|x^0 - x\|^2 \} \right] + \frac{5\eta^2}{K} \sum_{k=0}^{K-1} \mathbb{E} [\|\mathbb{E}_k[\Delta^k] - \Delta^k\|^2] \end{aligned}$$

$$\begin{aligned}
& - \mathbb{E} \left[ \frac{\gamma}{2K} \sum_{k=0}^{K-1} \|w^k - x^{k+1}\|^2 + \frac{1/2 - \gamma}{K} \sum_{k=0}^{K-1} \|x^{k+1} - x^k\|^2 \right] \\
& + \frac{1}{32K} \mathbb{E} [\|x^{K-1} - x^K\|^2] + \frac{1}{4K} \sum_{k=0}^{K-1} \mathbb{E} [\|x^{k+1} - x^k\|^2] \\
& \leq \frac{4}{K} \mathbb{E} \left[ \max_{x \in \mathcal{C}} \{\|x^0 - x\|^2\} \right] + \frac{5\eta^2}{K} \sum_{k=0}^{K-1} \mathbb{E} [\|\mathbb{E}_k [\Delta^k] - \Delta^k\|^2] \\
& - \mathbb{E} \left[ \frac{\gamma}{2K} \sum_{k=0}^{K-1} \|w^k - x^{k+1}\|^2 + \frac{1/2 - \gamma}{K} \sum_{k=0}^{K-1} \|x^{k+1} - x^k\|^2 \right] \\
& + \frac{1}{4K} \sum_{k=0}^{K-1} \mathbb{E} [\|x^{k+1} - x^k\|^2] \\
& = \frac{4}{K} \mathbb{E} \left[ \max_{x \in \mathcal{C}} \{\|x^0 - x\|^2\} \right] + \frac{5\eta^2}{K} \sum_{k=0}^{K-1} \mathbb{E} [\|\mathbb{E}_k [\Delta^k] - \Delta^k\|^2] \\
& - \mathbb{E} \left[ \frac{\gamma}{2K} \sum_{k=0}^{K-1} \|w^k - x^{k+1}\|^2 + \frac{1/4 - \gamma}{K} \sum_{k=0}^{K-1} \|x^{k+1} - x^k\|^2 \right].
\end{aligned}$$

Here we also used the initialization of Algorithm 1 with  $x^{-1} = x^0$ . Applying Lemma 2, we obtain

$$\begin{aligned}
2\eta \mathbb{E} [\text{Gap}(\bar{x}^K)] & \leq \frac{4}{K} \mathbb{E} \left[ \max_{x \in \mathcal{C}} \{\|x^0 - x\|^2\} \right] \\
& + \frac{10\eta^2 \bar{L}^2}{bK} \sum_{k=0}^{K-1} (\|x^k - w^{k-1}\|^2 + \|x^k - x^{k-1}\|^2) \\
& - \mathbb{E} \left[ \frac{\gamma}{2K} \sum_{k=0}^{K-1} \|w^k - x^{k+1}\|^2 + \frac{1/4 - \gamma}{K} \sum_{k=0}^{K-1} \|x^{k+1} - x^k\|^2 \right] \\
& \leq \frac{4}{K} \mathbb{E} \left[ \max_{x \in \mathcal{C}} \{\|x^0 - x\|^2\} \right] \\
& + \frac{10\eta^2 \bar{L}^2}{bK} \sum_{k=0}^{K-1} (\|x^{k+1} - w^k\|^2 + \|x^k - x^{k+1}\|^2) \\
& - \mathbb{E} \left[ \frac{\gamma}{2K} \sum_{k=0}^{K-1} \|w^k - x^{k+1}\|^2 + \frac{1/4 - \gamma}{K} \sum_{k=0}^{K-1} \|x^{k+1} - x^k\|^2 \right] \\
& \leq \frac{4}{K} \mathbb{E} \left[ \max_{x \in \mathcal{C}} \{\|x^0 - x\|^2\} \right] \\
& - \mathbb{E} \left[ \left( \frac{\gamma}{2} - \frac{10\eta^2 \bar{L}^2}{b} \right) \frac{1}{K} \sum_{k=0}^{K-1} \|w^k - x^{k+1}\|^2 \right]
\end{aligned}$$

$$+ \left( \frac{1}{4} - \gamma - \frac{10\eta^2 \bar{L}^2}{b} \right) \frac{1}{K} \sum_{k=0}^{K-1} \|x^{k+1} - x^k\|^2 \Big].$$

Here we again used the initialization of Algorithm 1 with  $w^{-1} = x^{-1} = x^0$ . The choice of  $\eta \leq \frac{\sqrt{\gamma b}}{8L}$  and  $0 < \gamma \leq \frac{1}{16}$  gives

$$\begin{aligned} 2\eta \mathbb{E} [\text{Gap}(\bar{x}^K)] &\leq \frac{4}{K} \mathbb{E} \left[ \max_{x \in \mathcal{C}} \{ \|x^0 - x\|^2 \} \right] \\ &\quad - \mathbb{E} \left[ \left( \frac{1}{12} - \gamma \right) \frac{1}{K} \sum_{k=0}^{K-1} \|x^{k+1} - x^k\|^2 \right] \\ &\leq \frac{4}{K} \max_{x \in \mathcal{C}} \{ \|x^0 - x\|^2 \}. \end{aligned}$$

And we have

$$\mathbb{E} [\text{Gap}(\bar{x}^K)] \leq \frac{2}{\eta K} \max_{x \in \mathcal{C}} \{ \|x^0 - x\|^2 \}.$$

Substitution of  $\eta$  from the conditions of the theorem and  $\gamma = p$  finishes the proof.