

SARAH-based Variance-reduced Algorithm for Stochastic Finite-sum Cocoercive Variational Inequalities

Aleksandr Beznosikov^{1,2} Alexander Gasnikov^{1,3,4}

¹ Moscow Institute of Physics and Technology, Dolgoprudny, Russia

² HSE University, Moscow, Russia

³ IITP RAS, Moscow, Russia

⁴ Caucasus Mathematical Center, Adyghe State University, Maikop, Russia

Abstract. Variational inequalities are a broad formalism that encompasses a vast number of applications. Motivated by applications in machine learning and beyond, stochastic methods are of great importance. In this paper we consider the problem of stochastic finite-sum cocoercive variational inequalities. For this class of problems, we investigate the convergence of the method based on the SARAH variance reduction technique. We show that for strongly monotone problems it is possible to achieve linear convergence to a solution using this method. Experiments confirm the importance and practical applicability of our approach.

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

1 Introduction

In this paper we focus on the following unconstrained variational inequality (VI) problem:

$$\text{Find } z^* \in \mathbb{R}^d \text{ such that } F(z^*) = 0, \quad (1)$$

where $F : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is some operator. This formulation is broad and encompasses many popular classes of tasks arising in practice. The simplest, however, widely encountered example of the VI is the minimization problem:

$$\min_{z \in \mathbb{R}^d} f(z).$$

To represent it in the form (1), it is sufficient to take $F(z) = \nabla f(z)$. As another also popular practical example, we can consider a saddle point or min-max problem:

$$\min_{x \in \mathbb{R}^{d_x}} \max_{y \in \mathbb{R}^{d_y}} g(x, y).$$

Here we need to take $F(z) = [\nabla_x g(x, y), -\nabla_y g(x, y)]$.

From a machine learning perspective, it is interesting not the deterministic formulation (1), but the stochastic one. More specifically, we want to consider

the setup with the operator $F(z) = \mathbb{E}_{\xi \sim \mathcal{D}} [F_\xi(z)]$, where ξ is a random variable, \mathcal{D} is some distribution, $F_\xi : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is a stochastic operator. But it is often the case (especially in practical problems) that the distribution \mathcal{D} is unknown, but we have some samples from \mathcal{D} . Then, one can replace with a finite-sum Monte Carlo approximation, i.e.

$$F(z) = \frac{1}{n} \sum_{i=1}^n F_i(z). \quad (2)$$

In the case of minimization problems, statements of the form (1)+(2) are also called empirical risk minimization [37]. These types of problems arise both in classical machine learning problems such as simple regressions and in complex, large-scale problems such as neural networks [23]. When it comes to saddle point problems, in recent times the so-called adversarial approach has become popular. Here one can highlight Generative Adversarial Networks (GANs) [13] and the adversarial training of models [40,25].

Based on the examples mentioned above, it can be noted that for operators of the form (2), computing the full value of F is a possible but not desirable operation, since it is typically very expensive compared to computing a single operator F_i . Therefore, when constructing an algorithm for the problem (1)+(2), one wants to avoid computing (or compute very rarely) the full F operator. This task can be solved by a stochastic gradient descent (SGD) framework. Currently, stochastic methods for minimization problems already have a huge background [15]. The first methods of this type were proposed back in the 1950s by [35]. For example, in the most classic variant, SGD could be written as follows:

$$z^{k+1} = z^k - \eta v^k, \quad (3)$$

where $\eta > 0$ is a predefined step-size and $v^k = \nabla f_i(z^k)$, where $i \in [n]$ is chosen randomly [38]. In this case, the variance of v_t is the main source of slower convergence or convergence only to the neighbourhood of the solution [7,31,21].

But for minimization problems of the finite-sum type, one can achieve stronger theoretical and practical results compared to the method (3). This requires the use of a variance reduction technique. Recently, many variance-reduced variants of SGD have been proposed, including SAG/SAGA [36,10,34], SVRG [19,3,39], MISO [27], SARAH [29,32,30,18], SPIDER [11], STORM [9], PAGE [24]. The essence of one of the earliest and best known variance-reduced methods SVRG is to use $v^k = \nabla f_i(z^k) - \nabla f_i(\tilde{z}) + \nabla f(\tilde{z})$, where $i \in [n]$ is picked at random, where $i \in [n]$ is picked at random and the point \tilde{z} is updated very rarely (hence we do not need to compute the full gradient often). With this type of methods it is possible to achieve a linear convergence to the solution. But for both convex and non-convex smooth minimization problems, the best theoretical guarantees of convergence are given by other variance-reduced technique SARAH (and its modifications: SPIDER, STORM, PAGE).

In turn, stochastic methods are also investigated for variational inequalities and saddle point problems [20,12,17,28,16,6,14,4,5], including methods based on variance reduction techniques [33,8,2,1,22,4,5]. Most of these methods are based

on the SVRG approach. At the same time, SARAH-based methods have not been explored for VIs. But as we noted earlier, these methods are the most attractive from the theoretical point of view for minimization problems. The purpose of this paper is to partially close the question of SARAH approach for stochastic finite-sum variational inequalities.

2 Problem setup and assumptions

Notation. We use $\langle x, y \rangle := \sum_{i=1}^n x_i y_i$ to denote standard inner product of $x, y \in \mathbb{R}^d$ where x_i corresponds to the i -th component of x in the standard basis in \mathbb{R}^d . It induces ℓ_2 -norm in \mathbb{R}^d in the following way $\|x\|_2 := \sqrt{\langle x, x \rangle}$.

Recall that we consider the problem (1), where the operator F has the form (2). Additionally, we assume

Assumption 1 (Cocoercivity) *Each operator F_i is ℓ -cocoercive, i.e. for all $u, v \in \mathbb{R}^d$ we have*

$$\|F_i(u) - F_i(v)\|^2 \leq \ell \langle F_i(u) - F_i(v), u - v \rangle. \quad (4)$$

This assumption is somehow a more restricted analogue of the Lipschitzness of F_i . For convex minimization problems, ℓ -Lipschitzness and ℓ -cocoercivity are equivalent. Regarding variational inequalities and saddle point problems, see [26].

Assumption 2 (Strong monotonicity) *The operator F is μ -strongly monotone, i.e. for all $u, v \in \mathbb{R}^d$ we have*

$$\langle F(u) - F(v); u - v \rangle \geq \mu \|u - v\|^2. \quad (5)$$

For minimization problems this property means strong convexity, and for saddle point problems strong convexity–strong concavity.

3 Main part

For general Lipschitzness variational inequalities, stochastic methods are usually based not on SGD, but on the Stochastic Extra Gradient method [20]. But due to the fact that we consider cocoercive VIs, it is sufficient to look at SGD like methods for this class of problems. For example, [26] considers SGD, [5] - SVRG. Following this reasoning, we base our method on the original SARAH [29].

Algorithm 1 SARAH [29] for Stochastic Cocoercive Variational Inequalities

```

1: Parameters: Step size  $\gamma > 0$ , number of iterations  $K, S$ .
2: Initialization: Choose  $\tilde{z}^0 \in \mathbb{R}^d$ .
3: for  $s = 1, 2, \dots, S$  do
4:    $z^0 = \tilde{z}^{s-1}$ 
5:    $v^0 = F(z^0)$ 
6:    $z^1 = z^0 - \gamma v^0$ 
7:   for  $k = 1, 2, \dots, K - 1$  do
8:     Sample  $i_k$  independently and uniformly from  $[n]$ 
9:      $v^k = F_{i_k}(z^k) - F_{i_k}(z^{k-1}) + v^{k-1}$ 
10:     $z^{k+1} = z^k - \gamma v^k$ 
11:   end for
12:    $\tilde{z}^s = z^K$ 
13: end for

```

Next, we analyse the convergence of this method. Note that we will use the vector v^K in the analysis, but in reality this vector is not calculated by the algorithm. Our proof are heavily based on the original work on SARAH [29]. Lemma 1 gives an understanding of how $\|v^k\|^2$ behaves during the internal loop of Algorithm 1.

Lemma 1. *Suppose that Assumptions 1 and 2 hold. Consider SARAH (Algorithm 1) with $\gamma \leq \frac{1}{\ell}$. Then, we have*

$$\mathbb{E}[\|v^K\|^2] \leq (1 - \gamma\mu)^K \mathbb{E}[\|F(z^0)\|^2].$$

Proof. We start the proof with an update for v^k :

$$\|v^k\|^2 = \|v^{k-1}\|^2 + \|F_{i_k}(z^k) - F_{i_k}(z^{k-1})\|^2 + 2\langle F_{i_k}(z^k) - F_{i_k}(z^{k-1}), v^{k-1} \rangle.$$

Next, we use an update for z^k and make a small rearrangement

$$\begin{aligned} \|v^k\|^2 &= \|v^{k-1}\|^2 + \|F_{i_k}(z^k) - F_{i_k}(z^{k-1})\|^2 - \frac{2}{\gamma} \langle F_{i_k}(z^k) - F_{i_k}(z^{k-1}), z^k - z^{k-1} \rangle \\ &= \|v^{k-1}\|^2 + \|F_{i_k}(z^k) - F_{i_k}(z^{k-1})\|^2 - \frac{1}{\gamma} \langle F_{i_k}(z^k) - F_{i_k}(z^{k-1}), z^k - z^{k-1} \rangle \\ &\quad - \frac{1}{\gamma} \langle F_{i_k}(z^k) - F_{i_k}(z^{k-1}), z^k - z^{k-1} \rangle. \end{aligned}$$

Taking the full mathematical expectation, we obtain

$$\begin{aligned} \mathbb{E}[\|v^k\|^2] &= \mathbb{E}[\|v^{k-1}\|^2] + \mathbb{E}[\|F_{i_k}(z^k) - F_{i_k}(z^{k-1})\|^2] \\ &\quad - \frac{1}{\gamma} \mathbb{E}[\langle F_{i_k}(z^k) - F_{i_k}(z^{k-1}), z^k - z^{k-1} \rangle] \\ &\quad - \frac{1}{\gamma} \mathbb{E}[\langle F_{i_k}(z^k) - F_{i_k}(z^{k-1}), z^k - z^{k-1} \rangle]. \end{aligned}$$

Independence of the i_k generation gives

$$\begin{aligned}
 \mathbb{E}[\|v^k\|^2] &= \mathbb{E}[\|v^{k-1}\|^2] + \mathbb{E}[\|F_{i_k}(z^k) - F_{i_k}(z^{k-1})\|^2] \\
 &\quad - \frac{1}{\gamma} \mathbb{E}[\langle F_{i_k}(z^k) - F_{i_k}(z^{k-1}), z^k - z^{k-1} \rangle] \\
 &\quad - \frac{1}{\gamma} \mathbb{E}[\langle \mathbb{E}_{i_k}[F_{i_k}(z^k) - F_{i_k}(z^{k-1})], z^k - z^{k-1} \rangle] \\
 &= \mathbb{E}[\|v^{k-1}\|^2] + \mathbb{E}[\|F_{i_k}(z^k) - F_{i_k}(z^{k-1})\|^2] \\
 &\quad - \frac{1}{\gamma} \mathbb{E}[\langle F_{i_k}(z^k) - F_{i_k}(z^{k-1}), z^k - z^{k-1} \rangle] \\
 &\quad - \frac{1}{\gamma} \mathbb{E}[\langle F(z^k) - F(z^{k-1}), z^k - z^{k-1} \rangle].
 \end{aligned}$$

With Assumptions 1 and 2, we get

$$\begin{aligned}
 \mathbb{E}[\|v^k\|^2] &\leq \mathbb{E}[\|v^{k-1}\|^2] + \mathbb{E}[\|F_{i_k}(z^k) - F_{i_k}(z^{k-1})\|^2] \\
 &\quad - \frac{1}{\gamma\ell} \mathbb{E}[\|F_{i_k}(z^k) - F_{i_k}(z^{k-1})\|^2] \\
 &\quad - \frac{\mu}{\gamma} \mathbb{E}[\|z^k - z^{k-1}\|^2] \\
 &= (1 - \gamma\mu) \mathbb{E}[\|v^{k-1}\|^2] + \left(\frac{\gamma\ell - 1}{\gamma\ell} \right) \mathbb{E}[\|F_{i_k}(z^k) - F_{i_k}(z^{k-1})\|^2].
 \end{aligned}$$

In the last step we substitute $z^{k-1} - z^k = \gamma v^k$. The choice of $0 < \gamma \leq \frac{1}{\ell}$ gives

$$\mathbb{E}[\|v^k\|^2] \leq (1 - \gamma\mu) \mathbb{E}[\|v^{k-1}\|^2].$$

Running recursion and using $v^0 = F(z^0)$, we finish the proof. \square

The following lemma gives how different v^K and $F(z^K)$ are in the inner loop of Algorithm 1.

Lemma 2. *Suppose that Assumption 1 holds. Consider SARAH (Algorithm 1). Then, we have*

$$\mathbb{E}[\|F(z^K) - v^K\|^2] \leq \frac{\gamma\ell}{2 - \gamma\ell} \mathbb{E}[\|F(z^0)\|^2].$$

Proof. Let us consider the following chain of reasoning:

$$\begin{aligned}
 \mathbb{E}[\|F(z^k) - v^k\|^2] &= \mathbb{E}[\|[F(z^{k-1}) - v^{k-1}] + [F(z^k) - F(z^{k-1})] - [v^k - v^{k-1}]\|^2] \\
 &= \mathbb{E}[\|F(z^{k-1}) - v^{k-1}\|^2] + \mathbb{E}[\|F(z^k) - F(z^{k-1})\|^2] \\
 &\quad + \mathbb{E}[\|v^k - v^{k-1}\|^2] \\
 &\quad + 2\mathbb{E}[\langle F(z^{k-1}) - v^{k-1}, F(z^k) - F(z^{k-1}) \rangle]
 \end{aligned}$$

$$\begin{aligned}
& -2\mathbb{E}[\langle F(z^{k-1}) - v^{k-1}, v^k - v^{k-1} \rangle] \\
& -2\mathbb{E}[\langle F(z^k) - F(z^{k-1}), v^k - v^{k-1} \rangle] \\
= & \mathbb{E}[\|F(z^{k-1}) - v^{k-1}\|^2] + \mathbb{E}[\|F(z^k) - F(z^{k-1})\|^2] \\
& + \mathbb{E}[\|v^k - v^{k-1}\|^2] \\
& + 2\mathbb{E}[\langle F(z^{k-1}) - v^{k-1}, F(z^k) - F(z^{k-1}) \rangle] \\
& - 2\mathbb{E}[\langle F(z^{k-1}) - v^{k-1}, \mathbb{E}_{i_k}[v^k - v^{k-1}] \rangle] \\
& - 2\mathbb{E}[\langle F(z^k) - F(z^{k-1}), \mathbb{E}_{i_k}[v^k - v^{k-1}] \rangle] \\
= & \mathbb{E}[\|F(z^{k-1}) - v^{k-1}\|^2] - \mathbb{E}[\|F(z^k) - F(z^{k-1})\|^2] \\
& + \mathbb{E}[\|v^k - v^{k-1}\|^2] \\
\leq & \mathbb{E}[\|F(z^{k-1}) - v^{k-1}\|^2] + \mathbb{E}[\|v^k - v^{k-1}\|^2].
\end{aligned}$$

Here we also use that

$$\mathbb{E}_{i_k}[v^k - v^{k-1}] = \mathbb{E}_{i_k}[F_{i_k}(z^k) - F_{i_k}(z^{k-1})] = F(z^k) - F(z^{k-1}).$$

Running recursion and using $v^0 = F(z^0)$, we have

$$\mathbb{E}[\|F(z^K) - v^K\|^2] \leq \sum_{k=1}^K \mathbb{E}[\|v^k - v^{k-1}\|^2]. \quad (6)$$

In the same way as in Lemma 1, we can derive

$$\begin{aligned}
\|v^k\|^2 &= \|v^{k-1}\|^2 + \|F_{i_k}(z^k) - F_{i_k}(z^{k-1})\|^2 + 2\langle F_{i_k}(z^k) - F_{i_k}(z^{k-1}), v^{k-1} \rangle \\
&= \|v^{k-1}\|^2 + \|F_{i_k}(z^k) - F_{i_k}(z^{k-1})\|^2 - \frac{2}{\gamma} \langle F_{i_k}(z^k) - F_{i_k}(z^{k-1}), z^k - z^{k-1} \rangle \\
&\leq \|v^{k-1}\|^2 + \|F_{i_k}(z^k) - F_{i_k}(z^{k-1})\|^2 - \frac{2}{\gamma\ell} \|F_{i_k}(z^k) - F_{i_k}(z^{k-1})\|^2 \\
&= \|v^{k-1}\|^2 + \left(\frac{\gamma\ell - 2}{\gamma\ell}\right) \|F_{i_k}(z^k) - F_{i_k}(z^{k-1})\|^2 \\
&= \|v^{k-1}\|^2 + \left(\frac{\gamma\ell - 2}{\gamma\ell}\right) \|v^k - v^{k-1}\|^2.
\end{aligned}$$

After a small rewriting and with the full expectation, we get

$$\mathbb{E}[\|v^k - v^{k-1}\|^2] \leq \frac{\gamma\ell}{2 - \gamma\ell} \mathbb{E}[\|v^{k-1}\|^2 - \|v^k\|^2].$$

By substituting this into the expression (6) and using $v^0 = F(z^0)$, we finish the proof. □

Let us combine Lemmas 1 and 2 into the main theorem of this paper.

Theorem 1. *Suppose that Assumptions 1 and 2 hold. Consider SARAH (Algorithm 1) with $\gamma = \frac{2}{9\ell}$ and $K = \frac{10\ell}{\mu}$. Then, we have*

$$\mathbb{E}[\|F(\tilde{z}^s)\|^2] \leq \frac{1}{2}\mathbb{E}[\|F(\tilde{z}^{s-1})\|^2].$$

Proof. We start from

$$\mathbb{E}[\|F(z^K)\|^2] \leq 2\mathbb{E}[\|F(z^K) - v^K\|^2] + 2[\mathbb{E}\|v^K\|^2].$$

Applying Lemma 1 and 2, we have

$$\begin{aligned} \mathbb{E}[\|F(z^K)\|^2] &\leq \left[\frac{2\gamma\ell}{2-\gamma\ell} + 2(1-\gamma\mu)^K \right] \mathbb{E}[\|F(z^0)\|^2] \\ &\leq \left[\frac{2\gamma\ell}{2-\gamma\ell} + 2\exp(-\gamma\mu K) \right] \mathbb{E}[\|F(z^0)\|^2]. \end{aligned}$$

Here we also use that $\gamma\mu \in (0; 1)$ (for $\gamma \leq \frac{2}{9\ell}$) and then $(1-\gamma\mu) \leq \exp(-\gamma\mu)$. The substitution γ and K gives

$$\mathbb{E}[\|F(z^K)\|^2] \leq \frac{1}{2}\mathbb{E}[\|F(z^0)\|^2].$$

We know that $z^0 = \tilde{z}^{s-1}$ and $z^K = \tilde{z}^s$ and have

$$\mathbb{E}[\|F(\tilde{z}^s)\|^2] \leq \frac{1}{2}\mathbb{E}[\|F(\tilde{z}^{s-1})\|^2].$$

□

Since we need to find a point z such that $F(z) \approx F(z^*) = 0$, we can easily get an estimate on the oracle complexity (number of F_i calls) to achieve precision ε .

Corollary 1. *Suppose that Assumptions 1 and 2 hold. Consider SARAH (Algorithm 1) with $\gamma = \frac{2}{9\ell}$ and $K = \frac{10\ell}{\mu}$. Then, to achieve ε -solution ($\mathbb{E}\|F(\tilde{z}^S)\|^2 \sim \varepsilon^2$), we need*

$$\mathcal{O}\left(\left[n + \frac{\ell}{\mu}\right] \log_2 \frac{\|F(z^0)\|^2}{\varepsilon^2}\right) \quad \text{oracle calls.}$$

Proof. From Theorem 1 we need the following number of outer iterations:

$$S = \mathcal{O}\left(\log_2 \frac{\|F(z^0)\|^2}{\varepsilon^2}\right).$$

At each outer iteration we compute the full operator one time, and at the remaining $K - 1$ iterations we call the single operator F_i two times per one inner iteration. Then, the total number of oracle calls is

$$S \times (2 \times (K - 1) + n) = \mathcal{O}\left(\left[n + \frac{\ell}{\mu}\right] \log_2 \frac{\|F(z^0)\|^2}{\varepsilon^2}\right).$$

□

Note that the obtained oracle complexity coincides with the similar complexity for SVRG from [5]. It is interesting to see how these methods behave in practice.

4 Experiments

The aim of our experiments is to compare the performance of different methods for stochastic finite-sum cocoercive variational inequalities. In particular, we use SGD from [26], SVRG from [5] and SARAH. We conduct our experiments on a finite-sum bilinear saddle point problem:

$$g(x, y) = \frac{1}{n} \sum_{i=1}^n \left[g_i(x, y) = x^\top A_i y + a_i^\top x + b_i^\top y + \frac{\lambda}{2} \|x\|^2 - \frac{\lambda}{2} \|y\|^2 \right], \quad (7)$$

where $A_i \in \mathbb{R}^{d \times d}$, $a_i, b_i \in \mathbb{R}^d$. This problem is λ -strongly convex–strongly concave and, moreover, L -smooth with $L = \|A\|_2$ for $A = \frac{1}{n} \sum_{i=1}^n A_i$. We take $n = 10$, $d = 100$ and generate matrix A and vectors a_i, b_i randomly, $\lambda = 1$. For this problem the cocoercivity constant $\ell = \frac{\|A\|_2^2}{\lambda}$. The steps of the methods are selected for best convergence. For SVRG and SARAH the number of iterations for the inner loops is taken as $\frac{\ell}{\lambda}$. We run three experiment setups: with small $\ell \approx 10^2$, medium $\ell \approx 10^3$ and big $\ell \approx 10^4$.

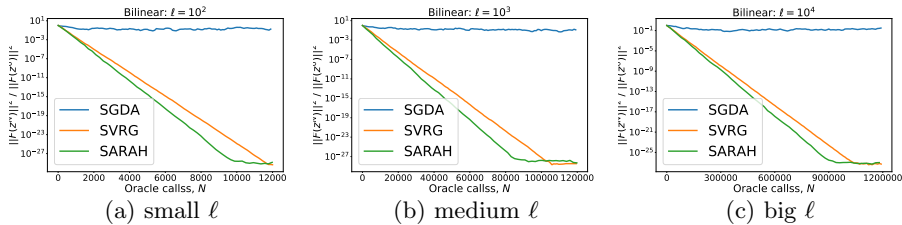


Fig. 1: Bilinear problem (7): Comparison of state-of-the-art SGD-based methods for stochastic cocoercive VIs.

See Figure 1 for the results. We see that SARAH converges better than SVRG, and SGD converges much slower.

Acknowledgments

The authors would like to congratulate Boris Mirkin on his jubilee and wish him good health and new scientific advances.

The work of A. Beznosikov was supported by the strategic academic leadership program 'Priority 2030' (Agreement 075-02-2021-1316 30.09.2021). The work of A. Gasnikov was supported by the Ministry of Science and Higher Education of the Russian Federation (Goszadaniye), No. 075-00337-20-03, project No. 0714-2020-0005.

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., Yuan, Y.: Improved svrg for non-strongly-convex or sum-of-non-convex objectives. In: *International conference on machine learning*. pp. 1080–1089. PMLR (2016)
4. Beznosikov, A., Gasnikov, A., Zainulina, K., Maslovskiy, A., Pasechnyuk, D.: A unified analysis of variational inequality methods: Variance reduction, sampling, quantization and coordinate descent. arXiv preprint arXiv:2201.12206 (2022)
5. Beznosikov, A., Gorbunov, E., Berard, H., Loizou, N.: Stochastic gradient descent-ascent: Unified theory and new efficient methods. arXiv preprint arXiv:2202.07262 (2022)
6. Beznosikov, A., Samokhin, V., Gasnikov, A.: Distributed saddle-point problems: Lower bounds, optimal and robust algorithms. arXiv preprint arXiv:2010.13112 (2020)
7. Bottou, L., Curtis, F.E., Nocedal, J.: Optimization methods for large-scale machine learning. *Siam Review* **60**(2), 223–311 (2018)
8. Chavdarova, T., Gidel, G., Fleuret, F., Lacoste-Julien, S.: Reducing noise in gan training with variance reduced extragradient. *Advances in Neural Information Processing Systems* **32** (2019)
9. Cutkosky, A., Orabona, F.: Momentum-based variance reduction in non-convex sgd. arXiv preprint arXiv:1905.10018 (2019)
10. 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)
11. Fang, C., Li, C.J., Lin, Z., Zhang, T.: Spider: Near-optimal non-convex optimization via stochastic path integrated differential estimator. arXiv preprint arXiv:1807.01695 (2018)
12. 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)
13. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial networks. *Communications of the ACM* **63**(11), 139–144 (2020)
14. Gorbunov, E., Berard, H., Gidel, G., Loizou, N.: Stochastic extragradient: General analysis and improved rates. In: *International Conference on Artificial Intelligence and Statistics*. pp. 7865–7901. PMLR (2022)
15. Gorbunov, E., Hanzely, F., Richtárik, P.: A unified theory of sgd: Variance reduction, sampling, quantization and coordinate descent. In: *International Conference on Artificial Intelligence and Statistics*. pp. 680–690. PMLR (2020)
16. Hsieh, Y.G., Iutzeler, F., Malick, J., Mertikopoulos, P.: Explore aggressively, update conservatively: Stochastic extragradient methods with variable stepsize scaling. *Advances in Neural Information Processing Systems* **33**, 16223–16234 (2020)
17. Hsieh, Y.G., Iutzeler, F., Malick, J., Mertikopoulos, P.: On the convergence of single-call stochastic extra-gradient methods (2019)

18. Hu, W., Li, C.J., Lian, X., Liu, J., Yuan, H.: Efficient smooth non-convex stochastic compositional optimization via stochastic recursive gradient descent (2019)
19. Johnson, R., Zhang, T.: Accelerating stochastic gradient descent using predictive variance reduction. *Advances in neural information processing systems* **26**, 315–323 (2013)
20. Juditsky, A., Nemirovski, A., Tauvel, C.: Solving variational inequalities with stochastic mirror-prox algorithm. *Stochastic Systems* **1**(1), 17–58 (2011)
21. Khaled, A., Richtárik, P.: Better theory for sgd in the nonconvex world. arXiv preprint arXiv:2002.03329 (2020)
22. 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)
23. LeCun, Y., Bengio, Y., Hinton, G.: Deep learning. *nature* **521**(7553), 436–444 (2015)
24. Li, Z., Bao, H., Zhang, X., Richtárik, P.: Page: A simple and optimal probabilistic gradient estimator for nonconvex optimization. In: *International Conference on Machine Learning*. pp. 6286–6295. PMLR (2021)
25. Liu, X., Cheng, H., He, P., Chen, W., Wang, Y., Poon, H., Gao, J.: Adversarial training for large neural language models. arXiv preprint arXiv:2004.08994 (2020)
26. Loizou, N., Berard, H., Gidel, G., Mitliagkas, I., Lacoste-Julien, S.: Stochastic gradient descent-ascent and consensus optimization for smooth games: Convergence analysis under expected co-coercivity. *Advances in Neural Information Processing Systems* **34**, 19095–19108 (2021)
27. Mairal, J.: Incremental majorization-minimization optimization with application to large-scale machine learning. *SIAM Journal on Optimization* **25**(2), 829–855 (2015)
28. Mishchenko, K., Kovalev, D., Shulgin, E., Richtárik, P., Malitsky, Y.: Revisiting stochastic extragradient. In: *International Conference on Artificial Intelligence and Statistics*. pp. 4573–4582. PMLR (2020)
29. 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)
30. Nguyen, L.M., Liu, J., Scheinberg, K., Takáč, M.: Stochastic recursive gradient algorithm for nonconvex optimization. arXiv preprint arXiv:1705.07261 (2017)
31. Nguyen, L.M., Nguyen, P.H., Richtárik, P., Scheinberg, K., Takáč, M., van Dijk, M.: New convergence aspects of stochastic gradient algorithms. *J. Mach. Learn. Res.* **20**, 176–1 (2019)
32. Nguyen, L.M., Scheinberg, K., Takáč, M.: Inexact SARAH algorithm for stochastic optimization. *Optimization Methods and Software* **36**(1), 237–258 (2021)
33. Palaniappan, B., Bach, F.: Stochastic variance reduction methods for saddle-point problems. In: *Advances in Neural Information Processing Systems*. pp. 1416–1424 (2016)
34. Qian, X., Qu, Z., Richtárik, P.: Saga with arbitrary sampling. In: *International Conference on Machine Learning*. pp. 5190–5199. PMLR (2019)
35. Robbins, H., Monro, S.: A stochastic approximation method. *The annals of mathematical statistics* pp. 400–407 (1951)
36. Schmidt, M., Le Roux, N., Bach, F.: Minimizing finite sums with the stochastic average gradient. *Mathematical Programming* **162**(1-2), 83–112 (2017)
37. Shalev-Shwartz, S., Ben-David, S.: *Understanding machine learning: From theory to algorithms*. Cambridge university press (2014)

38. Shalev-Shwartz, S., Singer, Y., Srebro, N., Cotter, A.: Pegasos: Primal estimated sub-gradient solver for svm. *Mathematical programming* **127**(1), 3–30 (2011)
39. Yang, Z., Chen, Z., Wang, C.: Accelerating mini-batch SARAH by step size rules. *Information Sciences* **558**, 157–173 (2021)
40. Zhu, C., Cheng, Y., Gan, Z., Sun, S., Goldstein, T., Liu, J.: Freelb: Enhanced adversarial training for natural language understanding. arXiv preprint arXiv:1909.11764 (2019)