Stochastic Halpern Iteration with Variance Reduction for Stochastic Monotone Inclusions

Anonymous Author(s) Affiliation Address email

Abstract

We study stochastic monotone inclusion problems, which widely appear in ma-1 chine learning applications, including robust regression and adversarial learning. 2 We propose novel variants of stochastic Halpern iteration with recursive variance З reduction. In the cocoercive-and more generally Lipschitz-monotone-setup, our 4 algorithm attains ϵ norm of the operator with $\mathcal{O}(\frac{1}{\epsilon^3})$ stochastic operator evalua-5 tions, which significantly improves over state of the art $\mathcal{O}(\frac{1}{c^4})$ stochastic operator 6 evaluations required for existing monotone inclusion solvers applied to the same 7 problem classes. We further show how to couple one of the proposed variants of 8 stochastic Halpern iteration with a scheduled restart scheme to solve stochastic 9 monotone inclusion problems with $\mathcal{O}(\frac{\log(1/\epsilon)}{\epsilon^2})$ stochastic operator evaluations un-10 der additional sharpness or strong monotonicity assumptions. Finally, we argue via 11 reductions between different problem classes that our stochastic oracle complexity 12 bounds are tight up to logarithmic factors in terms of their ϵ -dependence. 13

14 **1** Introduction

Recent trends in machine learning (ML) involve the study of models whose solutions do not reduce to 15 16 optimization but rather to equilibrium conditions. Standard examples include generative adversarial networks, adversarially robust training of ML models, and training of ML models under notions 17 of fairness. It turns out that several of these equilibrium conditions (including, but not limited to, 18 first-order stationary points, saddle-points, and Nash equilibria of minimax games) can be cast as 19 solutions to a monotone inclusion problem, which is defined as the problem of computing a zero of 20 a (maximal) monotone operator $\hat{F}: \mathbb{R}^d \to \mathbb{R}^d$ (see (MI) for a formal definition). In the context of 21 min-max optimization problems, monotone inclusion reduces to a stationarity condition, which for 22 unconstrained problems boils down to finding a point with small gradient norm. 23

Of particular interest to machine learning are stochastic versions of these problems, in which the operator F is not readily available, but can only be accessed through a stochastic oracle \hat{F} . Such are the settings mentioned above, where the definitions of equilibria involve expectations over continuous high-dimensional spaces. The corresponding problem, known as the *stochastic monotone inclusion*, has not been thoroughly studied, particularly in the context of its stochastic oracle complexity. Understanding stochastic oracle complexity of monotone inclusion in all standard settings with Lipschitz operators, from the algorithmic aspect, is the main motivation of this work.

31 1.1 Contributions

We study three main classes of stochastic monotone inclusion problems with Lipschitz operators, defined by the assumptions made about the operator itself: (i) cocoercive class, which is the most restricted class, but nevertheless fundamental for understanding monotone inclusion, as it relates to

Submitted to 36th Conference on Neural Information Processing Systems (NeurIPS 2022). Do not distribute.

the problem of finding a fixed point of a nonexpansive (1-Lipschitz) operator; (ii) Lipschitz monotone 35 class, which is perhaps the most basic class arising in the study of smooth convex-concave min-max 36 optimization problems; and (iii) Lipschitz monotone class with an additional sharpness property of 37 the operator. Sharpness is a widely studied property of optimization problems, often referred to as the 38 "local error bound" condition, which is weaker than strong convexity and roughly corresponds to the 39 problem landscape being curved outside of the solution set (see [38] for a survey of classical results). 40 From an algorithmic standpoint, we consider variants of classical Halpern iteration [19], which was 41 originally introduced for solving fixed point equations with nonexpansive operators. Variants of 42 this iteration have recently been shown to lead to (near-)optimal first-order oracle complexity for 43 all aforementioned standard problem classes in *deterministic* settings [10, 11, 50]. However, to the 44 best of our knowledge, stochastic variants of these methods have received very limited attention 45 prior to our work. The only results we are aware of are for a two-step extragradient-like variant 46 of Halpern iteration in co-monotone Lipschitz settings [26] and which show that when variance of 47 operator estimates is bounded by order- $\frac{\epsilon^2}{k}$ in iteration k, the method attains operator norm ϵ after $\mathcal{O}(\frac{1}{\epsilon})$ iterations. However, [26] does not discuss how such variance control would be obtained. 48 49 Simple mini-batching, as we show, only leads to $\mathcal{O}(\frac{1}{\epsilon^4})$ stochastic oracle complexity. 50

⁵¹ We show that existing variants of the Halpern iteration [10, 48] can be effectively combined with ⁵² recursive variance reduction [29] to obtain $\mathcal{O}(\frac{1}{\epsilon^3})$ stochastic oracle complexity in the cocoercive and ⁵³ Lipschitz monotone setups. We then show that the complexity can be further reduced to $\mathcal{O}(\frac{1}{\epsilon^2}\log(\frac{1}{\epsilon}))$ ⁵⁴ under an additional sharpness assumption about the operator. The latter two bounds can be certified ⁵⁵ to be near-optimal in ϵ , by reductions between different problem classes, as argued in Appendix A

To the best of our knowledge, our work is the first to use variance reduction to reduce stochastic oracle complexity of monotone inclusion (small gradient norm in min-max optimization settings), and the attained bounds are the best achieved to date.

59 1.2 Techniques

Inspired by the potential function originally used by $\boxed{10}$ and later used either in the same or slightly 60 modified form by [11, 26, 48, 50], we adapt this potential function-based argument to account for 61 stochastic error terms arising due to the stochastic oracle access to the operator. We first show that 62 in the cocoercive minibatch setting, this argument only leads to $\mathcal{O}(\frac{1}{\epsilon^4})$ stochastic oracle complexity, 63 and it is unclear how to improve it directly, as the analysis appears tight. We then combine the 64 cocoercive variant of Halpern iteration [10] with the PAGE estimator [29] to reduce the stochastic 65 oracle complexity to $\mathcal{O}(\frac{1}{\epsilon^3})$. The same variance reduced estimator is also used in conjunction with 66 the two-step extrapolated variant of Halpern iteration introduced by [48], as a direct application of 67 Halpern iteration is not known to converge on the class of Lipschitz monotone operators. 68

While the basic ideas in our arguments are simple, their realization requires addressing major technical 69 obstacles. First, the variance reduced estimator that we use [29] was originally devised for smooth 70 nonconvex optimization problems, where it was coupled with a stochastic variant of gradient descent. 71 This is significant, because the proof relies on a descent lemma, which allows cancelling the error 72 arising from the variance of the estimator by the "descent" part. Such an argument is not possible 73 in our setting, as there is no objective function to descend on. Instead, our analysis relies on an 74 intricate inductive argument that ensures that the expected norm of the operator is bounded in each 75 iteration, assuming a suitable bound on the variance of the estimator. To obtain our desired result for 76 the variance, we propose a data-dependent batch allocation in PAGE estimator [29] (see Corollary 77 2.2), which scales proportionally to the squared distance between successive iterates, similar to 3. 78 We inductively argue that the squared distance between successive iterates arising in the batch size of the estimator reduces at rate $\frac{1}{k^2}$ in expectation. This allows us to further certify that the estimators do 79 80 not only remain accurate, but their variance decreases as $\mathcal{O}(\epsilon^2/k)$, where k is the iteration count. 81

In the context of the potential function argument, unlike in the deterministic settings, we *do not* establish that the potential function is non-increasing, even in expectation. The stochastic error terms that arise due to the stochastic nature of the operator evaluations are controlled by taking slightly smaller step sizes than in the vanilla methods from [10, 48], which allows us to "leak" negative quadratic terms that are further used in controlling the stochastic error. The argument for controlling the value of the potential function is itself coupled with the inductive argument for ensuring that the expected operator norm remains bounded. ⁸⁹ Finally, while applying a restarting strategy is standard under sharpness conditions [43], obtaining

the claimed stochastic oracle complexity result of $\mathcal{O}\left(\frac{1}{\epsilon^2}\log(\frac{1}{\epsilon})\right)$ requires a rather technical argument

to bound the total number of stochastic queries to the operator.

92 1.3 Related Work

Monotone Inclusion and Variational Inequalities. Variational inequality problems were originally devised to deal with approximating equilibria. Their systematic study was initiated by [47].
The relationship between variational inequalities and min-max optimization was observed soon after [41], while one of the earliest papers to study solving monotone inclusion as a generalization of variational inequalities, convex and min-max optimization, and complementarity problems is [42].
For a historical overview of this area and an extensive review of classical results, see [15].

In the case of monotone operators, standard variants of variational inequality problems (see Section 2) 99 and monotone inclusion are equivalent—their solution sets coincide. This is a consequence of the 100 celebrated Minty Theorem [32]. However, there is a major difference between these problems when 101 it comes to solving them to a finite accuracy. In particular, on unbounded domains, approximating 102 103 variational inequalities is meaningless, whereas monotone inclusion remains well-defined. This is most readily seen from the observation that mapping from min-max optimization, variational 104 inequalities correspond to primal-dual gap guarantees, while monotone inclusion corresponds to 105 a guarantee in gradient norm. For a simple bilinear function f(x, y) = xy which has the unique 106 min-max solution at (x, y) = (0, 0), the primal-dual gap is infinite for any point other than (0, 0), 107 while the gradient remains finite and is a good proxy for measuring quality of a solution. Further, 108 even on bounded domains or using restricted gap functions on unbounded domains as in e.g., [34], 109 optimal oracle complexity guarantees for approximate monotone inclusion imply optimal complexity 110 guarantees for approximately satisfied variational inequalities (see, e.g., 10). The opposite does not 111 hold in general. In particular, in deterministic settings, standard algorithms such as the celebrated 112 extragradient [25, 33], dual extrapolation [34], or Popov's method [39] that have the optimal oracle complexity $O(\frac{1}{\epsilon})$ for approximating variational inequalities are suboptimal for monotone inclusion and attain oracle complexity of the order $O(\frac{1}{\epsilon^2})$ [11, 18]. 113 114 115

Halpern iteration. Halpern iteration is a classical fixed point iteration originally introduced by [19], and studied extensively in terms of both its asymptotic and non-asymptotic convergence guarantees [24, 28, 30, 49]. The first tight nonasymptotic convergence rate guarantee of 1/t was obtained in [30, 44]. This rate was also matched by an alternative method proposed by [23].

The usefulness of Halpern iteration for solving monotone inclusion problems was first observed by [10], who showed that its variants can be used to obtain near-optimal oracle complexity results for all standard classes of monotone inclusion problems with Lipschitz operators also studied in this work. The near-tightness (up to poly-logarithmic factors) of the results from [10] was certified using lower bound reductions from min-max optimization lower bounds introduced by [36]. These lower bounds were made tight for the cocoercive setup in [11].

The generalization of Halpern iteration from the cocoercive to Lipschitz monotone setup in [10]utilized approximating what is known as the resolvent operator, which led to a double-loop algorithm and an additional $\log(1/\epsilon)$ in the resulting complexity. This log factor was shaved off in [50], who introduced a two-step variant of Halpern iteration, inspired by the extragradient method of [25]. The results of [10, 50] were further extended to other classes of Lipschitz operators by [26, 48]. Except for [26] which considered controlled variance as discussed above, all of the existing results only targeted deterministic settings.

Stochastic Settings and Variance Reduction. Vanilla stochastic gradient methods have constant variance of stochastic gradients, which creates a bottleneck in the convergence rate. To improve the convergence rate, in the past decade, powerful variance reduction techniques have been proposed.

For strongly convex finite-sum problems, SAG [45], which used a biased stochastic estimator of the full gradient, was the first stochastic gradient method with a linear convergence rate. [22] and

¹³⁸ [9] improved [45] by proposing unbiased estimators of SVRG-type and SAGA-type, respectively.

¹Interestingly, the algorithm proposed by [23] for cocoercive inclusion coincides with the Halpern iteration for a related nonexpansive operator (see [7] Proposition 4.3]).

Such unbiased estimators were further combined with Nesterov acceleration [2, 46], or applied 139 to nonconvex finite-sum/infinite-sum problems [27, 40]. For nonconvex stochastic (infinite-sum) 140 problems, SARAH [35] and SPIDER [16, 51, 52] estimators were proposed to attain the optimal 141 oracle complexity of $\mathcal{O}(1/\epsilon^3)$ for finding an ϵ -approximate stationary point. Both estimators are 142 referred to as "recursive" variance reduction estimators, as they are biased when taking expectation 143 w.r.t. current randomness but unbiased w.r.t. all the randomness in history. PAGE [29] and STORM 144 **[8]** significantly simplified SARAH and SPIDER in terms of reducing the number of loops and 145 avoiding large minibatches, respectively. 3 further extended this line of work by incorporating 146 second-order information and dynamic batch sizes. 147

In the setting of min-max optimization and variational inequalities/monotone inclusion, variance 148 reduction has primarily been used for approximating variational inequalities, corresponding to the 149 primal-dual gap in min-max optimization; see, for example [1], 5, 6, 21, 31, 37]. Under strong 150 monotonicity (or sharpness in the case of [31]), such results generalize to monotone inclusion; 151 however, to the best of our knowledge, there have been no results that address monotone inclusion 152 under the weaker assumptions considered in this work. In the context of monotone inclusion with 153 Lipschitz operators, the tightest complexity result that we are aware of is $\mathcal{O}(\frac{1}{4})$, due to [12], and it 154 applies to a more general class of structured non-monotone Lipschitz operators, for the best iterate. 155 The same oracle complexity can be deduced for the last iterate of a two-step variant of Halpern from 156 [26] Theorem 6.1], using mini-batching. All the results in our work are also for the last iterate. 157

2 **Preliminaries** 158

176

We consider a real d-dimensional normed space $(\mathbb{R}^d, \|\cdot\|)$, where $\|\cdot\|$ is induced by an inner product 159 associated with the space, i.e., $\|\cdot\| = \sqrt{\langle \cdot, \cdot \rangle}$. Let $\mathcal{U} \subseteq \mathbb{R}^d$ be closed and convex; in the unconstrained case, $\mathcal{U} \equiv \mathbb{R}^d$. When \mathcal{U} is bounded, $D = \max_{\mathbf{u}, \mathbf{v} \in \mathcal{U}} \|\mathbf{u} - \mathbf{v}\|$ denotes its diameter. 160 161

Classes of Monotone Operators. We say that an operator $F : \mathbb{R}^d \to \mathbb{R}^d$ is 162

1. monotone, if $\forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^d$, $\langle F(\mathbf{u}) - F(\mathbf{v}), \mathbf{u} - \mathbf{v} \rangle > 0$. 163

2. L-Lipschitz continuous for some L > 0, if $\forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^d$, $||F(\mathbf{u}) - F(\mathbf{v})|| \le L ||\mathbf{u} - \mathbf{v}||$. 164

3. γ -cocoercive for some $\gamma > 0$, if $\forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^d$, $\langle F(\mathbf{u}) - F(\mathbf{v}), \mathbf{u} - \mathbf{v} \rangle \geq \gamma \|F(\mathbf{u}) - F(\mathbf{v})\|^2$. 165

4. μ -strongly monotone for some $\mu > 0$, if $\forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^d$, $\langle F(\mathbf{u}) - F(\mathbf{v}), \mathbf{u} - \mathbf{v} \rangle \ge \mu \|\mathbf{u} - \mathbf{v}\|^2$. 166

Note that we can easily specialize these definitions to the set \mathcal{U} by restricting **u**, **v** to be from \mathcal{U} . 167

Throughout the paper, the minimum assumption that we make about an operator F is that it is 168

monotone and Lipschitz. Observe that any γ -cocoercive operator is monotone and $\frac{1}{2}$ -Lipschitz. The 169 converse to this statement does not hold in general. 170

Monotone Inclusion and Variational Inequalities. Monotone inclusion asks for u* such that 171 $\mathbf{0} \in F(\mathbf{u}^*) + \partial I_{\mathcal{U}}(\mathbf{u}^*),$ (MI)

where $I_{\mathcal{U}}$ is the indicator function of the set \mathcal{U} and $\partial I_{\mathcal{U}}(\cdot)$ denotes the subdifferential of $I_{\mathcal{U}}$. 172

If F is continuous and monotone, the solution set to (\overline{MI}) is the same as the solution set of the 173 Stampacchia Variational Inequality (SVI) problem, which asks for $\mathbf{u}^* \in \mathcal{U}$ such that 174

$$\forall \mathbf{u} \in \mathcal{U}): \quad \langle F(\mathbf{u}^*), \mathbf{u} - \mathbf{u}^* \rangle \ge 0. \tag{SVI}$$

Further, when F is monotone, the solution set of (SVI) is equivalent to the solution set of the Minty 175 Variational Inequality (MVI) problem consisting in finding \mathbf{u}^* such that

$$(\forall \mathbf{u} \in \mathcal{U}): \langle F(\mathbf{u}), \mathbf{u}^* - \mathbf{u} \rangle \le 0.$$
 (MVI)

We assume throughout the paper that a solution to monotone inclusion (MI) exists, which implies that 177 solutions to both (SVI) and (MVI) exist as well. Existence of solutions follows from standard results 178 and is guaranteed whenever e.g., \mathcal{U} is compact, or, if there exists a compact set \mathcal{U}' such that $\mathrm{Id} - \frac{1}{L}F$ 179 maps \mathcal{U}' to itself, where Id is the identity map [15]. As remarked in the introduction, in unbounded 180 setups it is generally not possible to approximate (MVI) and (SVI), whereas approximating (MI) is 181 quite natural: we only need to find u such that $\mathbf{0} \in \overline{F(\mathbf{u})} + \partial I_{\mathcal{U}}(\mathbf{u}) + \mathcal{B}(\epsilon)$, where **0** denotes the zero 182 vector and $\mathcal{B}(\epsilon)$ denotes the centered ball of radius ϵ . 183

184 **Stochastic Access to the Operator.** We consider the stochastic setting for monotone inclusion 185 problems. More specifically, we make the following assumptions for stochastic queries to F. These

- assumptions are made throughout the paper, without being explicitly invoked.
- Assumption 1 (Unbiased samples with bounded variance). For each query point $\mathbf{x} \in \mathcal{U}$, we observe $\widehat{F}(\mathbf{x}, z)$ where $z \sim P_z$ is a random variable that satisfies the following assumptions:

$$\mathbb{E}_{z}[\widehat{F}(\mathbf{x},z)] = F(\mathbf{x}) \quad \text{ and } \quad \mathbb{E}_{z}[\|\widehat{F}(\mathbf{x},z) - F(\mathbf{x})\|^{2}] \leq \sigma^{2}.$$

Assumption 2 (Multi-point oracle). We can query a set of points (x_1, \ldots, x_n) and receive

$$\widehat{F}(\mathbf{x}_1, z), \dots, \widehat{F}(\mathbf{x}_n, z)$$
 where $z \sim P_z$.

190 Assumption 3 (Lipschitz in expectation). $\mathbb{E}_{z}\left[\left\|\widehat{F}(\mathbf{u},z) - \widehat{F}(\mathbf{v},z)\right\|^{2}\right] \leq L^{2} \|\mathbf{u} - \mathbf{v}\|^{2}, \forall \mathbf{u}, \mathbf{v} \in \mathcal{U}.$

We note that complexity results of the paper will bound the total number of queries made to this oracle. In particular, if multiple query points and/or multiple samples z are used in a single iteration, our complexity is given by the sum of all those queries throughout all iterations of the method. Also, Assumption 3 is primary with parameter L, by which F is also L-Lipschitz using Jensen's inequality.

PAGE Variance-Reduced Estimator. We now summarize a variant of the PAGE estimator, originally developed for smooth nonconvex optimization by [29], adapted to our setting. In particular, given queries to \hat{F} , we define the variance reduced estimator $\tilde{F}(\mathbf{u}_k)$ for $k \ge 1$ by

$$\widetilde{F}(\mathbf{u}_{k}) = \begin{cases} \frac{1}{S_{1}^{(k)}} \sum_{i=1}^{S_{1}^{(k)}} \widehat{F}(\mathbf{u}_{k}, z_{i}^{(k)}) & \text{w. p. } p_{k}, \\ \widetilde{F}(\mathbf{u}_{k-1}) + \frac{1}{S_{2}^{(k)}} \sum_{i=1}^{S_{2}^{(k)}} \left(\widehat{F}(\mathbf{u}_{k}, z_{i}^{(k)}) - \widehat{F}(\mathbf{u}_{k-1}, z_{i}^{(k)}) \right) & \text{w. p. } 1 - p_{k}, \end{cases}$$
(2.1)

where $p_0 = 1$, $z_i^{(k)} \stackrel{\text{i.i.d.}}{\sim} P_z$, and $S_1^{(k)}$ and $S_2^{(k)}$ are the sample sizes at iteration k. Observe that Assumption 2 guarantees that we can query \widehat{F} at \mathbf{u}_k and \mathbf{u}_{k-1} using the same random seed. Our analysis will make use of conditional expectations, and to that end, we define natural filtration \mathcal{F}_k by $\mathcal{F}_k := \sigma(\{\widetilde{F}(\mathbf{u}_j)\}_{j \le k})$; namely \mathcal{F}_k contains all the randomness that arises in the definitions of $\widetilde{F}(\mathbf{u}_j)$ for $j \le k$. Following a similar argument as in [29], we recursively bound the variance of the estimator \widetilde{F} , as summarized in the following lemma. The proof is provided in Appendix B.

Lemma 2.1. Let F be a monotone operator accessed via stochastic queries \widehat{F} , under Assumptions Then, the variance of \widetilde{F} defined by Eq. (2.1) satisfies the following recursive bound: for all $k \ge 1$, the variance of \widetilde{F} defined by Eq. (2.1) satisfies the following recursive bound: for all $k \ge 1$,

$$\mathbb{E}[\|\widetilde{F}(\mathbf{u}_k) - F(\mathbf{u}_k)\|^2] \le \frac{p_k \sigma^2}{S_1^{(k)}} + (1 - p_k) \Big(\mathbb{E}[\|\widetilde{F}(\mathbf{u}_{k-1}) - F(\mathbf{u}_{k-1})\|^2] + \mathbb{E}\Big[\frac{L^2 \|\mathbf{u}_k - \mathbf{u}_{k-1}\|^2}{S_2^{(k)}}\Big] \Big)$$

With the choices of $p_k, S_1^{(k)}, S_2^{(k)}$ specified in the following corollary and using induction with the inequality from Lemma 2.1, we obtain the following bound on the variance.

Corollary 2.2. Given a target error $\epsilon > 0$, if for all $k \ge 1$, $p_k = \frac{2}{k+1}, S_1^{(k)} \ge \left\lceil \frac{8\sigma^2}{p_k \epsilon^2} \right\rceil, S_2^{(k)} \ge \left\lceil \frac{8L^2 \|\mathbf{u}_k - \mathbf{u}_{k-1}\|^2}{p_k^2 \epsilon^2} \right\rceil$, then $\mathbb{E}\left[\left\| \widetilde{F}(\mathbf{u}_k) - F(\mathbf{u}_k) \right\|^2 \right] \le \frac{\epsilon^2}{k}$.

3 Stochastic Halpern Iteration for Cocoercive Operators

In this section, we consider the setting of $\frac{1}{L}$ -cocoercive operators *F*. While cocoercivity is a strong assumption that implies that an operator is both Lipschitz and monotone (as discussed in Section 2), it is nevertheless the most basic setup for studying the Halpern iteration. In particular, while Halpern iteration can be applied directly to the nonexpansive counterpart of a cocoercive operator *F* (i.e., to the linear transformation Id $-\frac{2}{L}F$, where $\frac{1}{L}$ is an upper bound on the cocoercivity parameter of *F*), convergence does not seem possible to establish for the more general class of Lipschitz monotone operators. We begin this section by providing a generic proof of stochastic oracle complexity, which we then use to briefly illustrate how to obtain $O(\frac{1}{\epsilon^4})$ oracle complexity with a simple minibatch

- stochastic estimator of F. We then show how to improve this bound to $\mathcal{O}(\frac{1}{\epsilon^3})$ by applying the proposed variant of the PAGE estimator from Eq. (2.1) to Halpern iteration.
- ²²² The stochastic variant of Halpern iteration that we consider is defined by

$$\mathbf{u}_{k+1} = \lambda_{k+1}\mathbf{u}_0 + (1 - \lambda_{k+1})\Big(\mathbf{u}_k - \frac{2}{L_{k+1}}\widetilde{F}(\mathbf{u}_k)\Big),\tag{3.1}$$

where \tilde{F} is a stochastic (possibly biased) estimator of F, $\lambda_{k+1} = \Theta(\frac{1}{k})$ is the step size, and $L_{k+1} \ge L$ is a parameter of the algorithm. Compared to the classical iteration $\mathbf{u}_{k+1} = \lambda_{k+1}\mathbf{u}_0 + (1 - \lambda_{k+1})T(\mathbf{u}_k)$, where $T : \mathbb{R}^d \to \mathbb{R}^d$ is a nonexpansive (1-Lipschitz) map [19], T is replaced by the mapping $\mathrm{Id} - \frac{2}{L_{k+1}}\tilde{F}$, which is stochastic and may not be nonexpansive (as the stochastic estimate \tilde{F} of F is not guaranteed to be cocoercive even when F is). Compared to the iteration variant considered by [10], the access to the monotone operator is stochastic and we also take slightly larger (by a factor of 2) values of L_{k+1} to bound the stochastic error terms.

Our argument for bounding the total number of stochastic queries to F is based on the use of the following potential function $C_k = \frac{A_k}{L_k} ||F(\mathbf{u}_k)||^2 + B_k \langle F(\mathbf{u}_k), \mathbf{u}_k - \mathbf{u}_0 \rangle$, where $\{A_k\}_{k \ge 1}$ and $\{B_k\}_{k \ge 1}$ are positive and non-decreasing sequences of real numbers, while the step size λ_k is defined by $\lambda_k := \frac{B_k}{A_k + B_k}$. Such potential function was previously used for the deterministic case of Halpern iteration in [10, 11]. Observe that even though we make oracle queries to \hat{F} , the potential function C_k and the final bound we obtain are in terms of the true operator value F.

²³⁶ Compared to the analysis of Halpern iteration in the deterministic case [10, 11], our analysis for the ²³⁷ stochastic case needs to account for the error terms caused by accessing F via stochastic queries and ²³⁸ is based on an intricate inductive argument. A generic bound on iteration complexity, under mild ²³⁹ assumptions about the estimator \tilde{F} , is summarized in Theorem 3.1. The proof is in Appendix C.

Theorem 3.1. Given an arbitrary $\mathbf{u}_0 \in \mathbb{R}^d$, suppose that iterates \mathbf{u}_k evolve according to Halpern iteration from Eq. (3.1) for $k \ge 1$, where $L_k = 2L$ and $\lambda_k = \frac{1}{k+1}$. Assume further that the stochastic estimate $\widetilde{F}(\mathbf{u})$ is unbiased for $\mathbf{u} = \mathbf{u}_0$ and $\mathbb{E}[||F(\mathbf{u}_0) - \widetilde{F}(\mathbf{u}_0)||^2] \le \frac{\epsilon^2}{8}$. Given $\epsilon > 0$, if for all

243 $k \ge 1$, we have that $\mathbb{E}\left[\left\|F(\mathbf{u}_k) - \widetilde{F}(\mathbf{u}_k)\right\|^2\right] \le \frac{\epsilon^2}{k}$, then for all $k \ge 1$,

$$\mathbb{E}[\|F(\mathbf{u}_k))\|] \le \frac{\Lambda_0}{k} + \Lambda_1 \epsilon, \tag{3.2}$$

where $\Lambda_0 = 76L \|\mathbf{u}_0 - \mathbf{u}^*\|$ and $\Lambda_1 = 4\sqrt{\frac{2}{3}}$. As a result, stochastic Halpern iteration from Eq. (3.1) returns a point \mathbf{u}_k such that $\mathbb{E}[\|F(\mathbf{u}_k)\|] \le 4\epsilon$ after at most $N = \lceil \frac{2\Lambda_0}{\epsilon} \rceil = \mathcal{O}(\frac{L\|\mathbf{u}_0 - \mathbf{u}^*\|}{\epsilon})$ iterations.

We remark that the previous result states an iteration complexity bound under a rather high accuracy assumption for the operator estimators at each iteration. In order to attain these accuracy requirements, we could either use a minibatch at every iteration, or use variance reduction. In what follows we explore both approaches. We further remark that we made no effort to optimize the constants in the bound above, and thus the constants are likely improvable.

Finally, observe that due to the required low error for the estimates $\mathbb{E}[\|F(\mathbf{u}_k) - \widetilde{F}(\mathbf{u}_k)\|^2] \leq \frac{\epsilon^2}{k}$, we can certify by Chebyshev bound that $\mathbb{P}[\|F(\mathbf{u}_k) - \widetilde{F}(\mathbf{u}_k)\| \geq \epsilon] \leq \frac{1}{k}$. In particular, after $O(\frac{1}{\epsilon})$ iterations, once we have $\|\widetilde{F}(\mathbf{u}_k)\| \leq \epsilon$, $\|F(\mathbf{u}_k)\|$ is also $O(\epsilon)$ with probability at least $1 - \epsilon$. This is particularly important for practical implementations, where a stopping criterion can be based on the value of $\|\widetilde{F}(\mathbf{u}_k)\|$, which, unlike $\|F(\mathbf{u}_k)\|$, can be efficiently evaluated.

256 3.1 Stochastic Oracle Complexity With a Simple Mini-batch Estimate

A direct consequence of Theorem 3.1 is that a simple estimator $\tilde{F}(\mathbf{u}_k) = \frac{1}{S_k} \sum_{i=1}^{S_k} \hat{F}(\mathbf{u}_k, z_i^{(k)})$ leads to the overall $\mathcal{O}(\frac{1}{\epsilon^4})$ oracle complexity, as stated below while the proof is deferred to Appendix C

Corollary 3.2. Under the assumptions of Theorem 3.1 if $\widetilde{F}(\mathbf{u}_k) = \frac{1}{S_k} \sum_{i=1}^{S_k} \widehat{F}(\mathbf{u}_k, z_i^{(k)})$, where $\widehat{F}(\mathbf{u}_k, z_i^{(k)})$ satisfies Assumption 1 and $z_i^{(k)} \stackrel{i.i.d.}{\sim} P_z$, then setting $S_k = \frac{\sigma^2(k+1)}{\epsilon^2}$ for all $k \ge 0$ guarantees that $\mathbb{E}[\|F(\mathbf{u}_k)\|] \le 4\epsilon$ after at most $\mathcal{O}(\frac{\sigma^2 L^2 \|\mathbf{u}_0 - \mathbf{u}^*\|^2}{\epsilon^4})$ queries to \widehat{F} .

3.2 Improved Oracle Complexity via Variance Reduction 262

We now consider using the recursive variance reduction method from Eq. (2.1) to obtain the variance 263 bound required in Theorem 3.1 as summarized in Algorithm 1. Of course, in practice, $||\mathbf{u}_0 - \mathbf{u}^*||$ is 264 not known, and instead of running the algorithm for a fixed number of iterations N, one could run it, 265 for example, until reaching a point with $\|\widetilde{F}(\mathbf{u}_k)\| < \epsilon$.

Algorithm 1: Stochastic Halpern-Cocoercive (Halpern)

Input: $\mathbf{u}_0 \in \mathbb{R}^d$, $\|\mathbf{u}_0 - \mathbf{u}^*\|$, $L, \epsilon > 0, \sigma$; Initialize: $\Lambda_0 = \frac{76L \|\mathbf{u}_0 - \mathbf{u}^*\|}{\epsilon}, N = \lceil \frac{2\Lambda_0}{\epsilon} \rceil, S_1^{(0)} = \lceil \frac{8\sigma^2}{\epsilon^2} \rceil, \widetilde{F}(\mathbf{u}_0) = \frac{1}{S_1^{(0)}} \sum_{i=1}^{S_1^{(0)}} \widehat{F}(\mathbf{u}_0, z_i^{(0)});$ for k = 1 : N do $\mathbf{u}_{k} = \frac{1}{k+1} \mathbf{u}_{0} + \frac{k}{k+1} \left(\mathbf{u}_{k-1} - \frac{1}{L} \widetilde{F}(\mathbf{u}_{k-1}) \right);$ $p_{k} = \frac{2}{k+1}, S_{1}^{(k)} = \left\lceil \frac{8\sigma^{2}}{p_{k}\epsilon^{2}} \right\rceil, S_{2}^{(k)} = \left\lceil \frac{8L^{2} \|\mathbf{u}_{k} - \mathbf{u}_{k-1}\|^{2}}{p_{k}^{2}\epsilon^{2}} \right\rceil;$ Compute $\widetilde{F}(\mathbf{u}_k)$ based on Eq. (2.1) **Return:** \mathbf{u}_N

266

Notice that convergence is guaranteed by Theorem 3.1 however it does not directly address the 267 problem of the oracle complexity (as batch sizes depend on successive iterate distances). To resolve 268 this issue, we first provide a bound on $\|\mathbf{u}_k - \mathbf{u}_{k-1}\|$, and making the appropriate parameter settings 269 for the estimator from Eq. (2.1), it is now possible to apply Theorem 3.1 to obtain the improved $\mathcal{O}(\frac{1}{\epsilon^3})$ stochastic oracle complexity bound, as stated below while the proof is deferred to Appendix C 270 271 **Corollary 3.3.** Given arbitrary $\mathbf{u}_0 \in \mathbb{R}^d$ and $\epsilon > 0$, consider \mathbf{u}_N returned by Algorithm I. Then, $\mathbb{E}[\|F(\mathbf{u}_N)\|] \leq 4\epsilon$ with expected $\mathcal{O}(\frac{\sigma^2 L \|\mathbf{u}_0 - \mathbf{u}^*\| + L^3 \|\mathbf{u}_0 - \mathbf{u}^*\|^3}{\epsilon^3})$ oracle queries to \widehat{F} . 272

273

We note in passing that the running time guarantee of this algorithm is of Las Vegas-type: despite 274 its iteration number being surely bounded by $\left\lceil \frac{2\Lambda_0}{\epsilon} \right\rceil = \mathcal{O}\left(\frac{L \|\mathbf{u}_0 - \mathbf{u}^*\|}{\epsilon}\right)$, the batch sizes (in particular 275 $S_2^{(k)}$) are random, and are not universally bounded. 276

We further argue that Algorithm 1 can be extended to **constrained settings** by defining the operator 277 mapping as in [10] and modifying the variance-reduced stochastic estimator accordingly based on 278 the projection of \vec{F} . We show that the newly defined operator mapping is also cocoercive while the 279 variance of the modified estimator is bounded by the variance of \tilde{F} , so arguments from Theorem 3.1 280 and Corollary 3.3 extend to this case. This modified estimator need not be unbiased (as neither is F); 281 however, this is irrelevant to our analysis as it does not require unbiasedness. For completeness, a 282 detailed extension to the constrained case is provided in Appendix C.2. 283

Monotone and Lipschitz Setup 4 284

Throughout this section, we assume that F is monotone and L-Lipschitz. While the previous section 285 addresses the cocoercive setup using the classical version of Halpern iteration adapted to cocoercive 286 operators, it is unclear how to directly generalize this result to the setting with monotone Lipschitz 287 operators. In the deterministic setting, generalization to monotone Lipschitz operators can be achieved 288 through the use of a resolvent operator (see 10). However, such an approach incurs an additional 289 $\log(1/\epsilon)$ factor in the iteration complexity coming from approximating the resolvent and it is further 290 unclear how to generalize it to stochastic settings, as the properties of the stochastic estimate F of F291 do not readily translate into the same or similar properties for the resolvent of \tilde{F} . Instead of taking 292 the approach based on the resolvent, we consider a recently proposed two-step variant of Halpern 293 iteration [48], adapted here to the stochastic setting. The variant uses extrapolation and is defined by 294

$$\begin{cases} \mathbf{v}_k &:= \lambda_k \mathbf{u}_0 + (1 - \lambda_k) \, \mathbf{u}_k - \eta_k F(\mathbf{v}_{k-1}), \\ \mathbf{u}_{k+1} &:= \lambda_k \mathbf{u}_0 + (1 - \lambda_k) \, \mathbf{u}_k - \eta_k \widetilde{F}(\mathbf{v}_k), \end{cases}$$
(4.1)

where $\lambda_k \in [0,1), \eta_k > 0$, and \widetilde{F} is defined by (2.1). The resulting algorithm with a complete 295 parameter setting is provided in Algorithm 2 296

Algorithm 2: Extrapolated Stochastic Halpern-Monotone (E-Halpern)

$$\begin{split} \overline{\mathbf{Input:} \ \mathbf{u}_{0} \in \mathbb{R}^{d}, \|\mathbf{u}_{0} - \mathbf{u}^{*}\|, 0 < \eta_{0} \leq \frac{1}{3\sqrt{3L}}, L, \epsilon > 0, \sigma; \\ \mathbf{Initialize:} \ \mathbf{v}_{-1} = \mathbf{u}_{0}, S_{1}^{(-1)} = S_{1}^{(0)} = \lceil \frac{8\sigma^{2}}{\epsilon^{2}} \rceil, M = 9L^{2}, \underline{\eta} = \frac{\eta_{0}(1-2M\eta_{0}^{2})}{1-M\eta_{0}^{2}}; \\ \mathbf{Set} \ \Lambda_{0} = \frac{4(L^{2}\eta_{0}\underline{\eta}+1)\|\mathbf{u}_{0} - \mathbf{u}^{*}\|^{2}}{\underline{\eta}^{2}}, \Lambda_{1} = \frac{5(1+M\underline{\eta}\eta_{0})}{M\underline{\eta}^{2}}, N = \lceil \frac{\sqrt{\Lambda_{0}}}{\sqrt{\Lambda_{1}\epsilon}} \rceil; \\ \widetilde{F}(\mathbf{v}_{-1}) = \frac{1}{S_{1}^{(-1)}} \sum_{i=1}^{S_{i}^{(-1)}} \widehat{F}(\mathbf{v}_{-1}, z_{i}^{(-1)}), \text{ where } z_{i}^{(-1)} \stackrel{\text{i.i.d.}}{\sim} \mathcal{P}_{z}; \\ \mathbf{for} \ k = 1: N \ \mathbf{do} \\ \mathbf{v}_{k-1} = \frac{1}{k+1} \mathbf{u}_{0} + \frac{k}{k+1} \mathbf{u}_{k-1} - \eta_{k-1} \widetilde{F}(\mathbf{v}_{k-2}); \\ p_{k-1} = \min(\frac{2}{k}, 1), S_{1}^{(k-1)} = \lceil \frac{8\sigma^{2}}{p_{k-1}\epsilon^{2}} \rceil, S_{2}^{(k-1)} = \lceil \frac{8L^{2}\|\mathbf{v}_{k-1} - \mathbf{v}_{k-2}\|^{2}}{p_{k-1}^{2}\epsilon^{2}} \rceil; \\ \mathbf{Compute} \ \widetilde{F}(\mathbf{v}_{k-1}) \ \mathbf{based on Eq.} \ (\underline{2.1}); \\ \mathbf{u}_{k} = \frac{1}{k+1} \mathbf{u}_{0} + \frac{k}{k+1} \mathbf{u}_{k-1} - \eta_{k-1} \widetilde{F}(\mathbf{v}_{k-1}); \\ \eta_{k} = \frac{(1-\frac{1}{(k+1)^{2}} - M\eta_{k-1}^{2})(k+1)^{2}}{(1-M\eta_{k-1}^{2})k(k+2)} \eta_{k-1} \\ \mathbf{Return:} \ \mathbf{u}_{N} \end{split}$$

To analyze the convergence of the extrapolated Halpern variant from Eq. (4.1), we use the potential 297 function $\mathcal{V}_k = A_k \|F(\mathbf{u}_k)\|^2 + B_k \langle F(\mathbf{u}_k), \mathbf{u}_k - \mathbf{u}_0 \rangle + c_k L^2 \|\mathbf{u}_k - \mathbf{v}_{k-1}\|^2$, previously used by [48], 298 where A_k , B_k and c_k are positive parameters to be determined later. Observe that this is essentially 299 the same potential function as C_k , corrected by the quadratic term $c_k L^2 ||\mathbf{u}_k - \mathbf{v}_{k-1}||^2$ to account 300 for error terms appearing in the analysis of the two-step variant from Eq. (4.1). Similarly as in 301 the cocoercive setup, the potential function is not monotonically non-increasing, due to the error 302 terms that arise due to the stochastic access to F. Bounding these error terms requires a careful 303 technical argument, and is the main technical contribution of this section. Due to space constraints, 304 the complete technical argument is deferred to Appendix \mathbf{D} , while the main results are stated below. 305 **Theorem 4.1.** Given an arbitrary initial point $\mathbf{u}_0 \in \mathbb{R}^d$ and target error $\epsilon > 0$, assume that the 306

iterates \mathbf{u}_k evolve according to Algorithm 2 for $k \ge 1$. Then, for all $k \ge 2$,

$$\mathbb{E}\left[\|F(\mathbf{u}_{k})\|^{2} + 2L^{2} \|\mathbf{u}_{k} - \mathbf{v}_{k-1}\|^{2}\right] \leq \frac{\Lambda_{0}}{(k+1)(k+2)} + \Lambda_{1}\epsilon^{2},$$
(4.2)

where $\Lambda_0 = \frac{4(L^2\eta_0\underline{\eta}+1)\|\mathbf{u}_0-\mathbf{u}^*\|^2}{\underline{\eta}^2}$ and $\Lambda_1 = \frac{5(1+M\underline{\eta}\eta_0)}{M\underline{\eta}^2}$. In particular, $\mathbb{E}[\|F(\mathbf{u}_N)\|^2 + 2L^2\|\mathbf{u}_N-\mathbf{v}_{N-1}\|^2] \leq 2\Lambda_1\epsilon^2 = \mathcal{O}(\epsilon^2)$ after at most $N = \lceil \frac{\sqrt{\Lambda_0}}{\sqrt{\Lambda_1\epsilon}} \rceil = \mathcal{O}(\frac{L\|\mathbf{u}_0-\mathbf{u}^*\|}{\epsilon})$ iterations. The total number of oracle queries to \widehat{F} is $\mathcal{O}(\frac{\sigma^2 L\|\mathbf{u}_0-\mathbf{u}^*\|+L^3\|\mathbf{u}_0-\mathbf{u}^*\|^3}{\epsilon^3})$ in expectation.

5 Faster Convergence Under a Sharpness Condition

We now show that by restarting Algorithm 2 we can achieve the $\mathcal{O}\left(\frac{1}{\epsilon^2}\log\frac{1}{\epsilon}\right)$ oracle complexity under a milder than strong monotonicity μ -sharpness condition: for all $\mathbf{u} \in \mathcal{U}$, $\langle F(\mathbf{u}) - F(\mathbf{u}^*), \mathbf{u} - \mathbf{u}^* \rangle \geq \mu \|\mathbf{u} - \mathbf{u}^*\|^2$. The scheme is summarized in Algorithm 3 and the proof is deferred to Appendix E. **Theorem 5.1.** Given L-Lipschitz and μ -sharp F and the precision parameter ϵ , Algorithm 3 outputs \mathbf{u}_N with $\mathbb{E}[\|\mathbf{u}_N - \mathbf{u}^*\|^2] \leq \epsilon^2$ as well as $\mathbb{E}[\|F(\mathbf{u}_N)\|^2] \leq L^2 \epsilon^2$ after $N = \mathcal{O}\left(\frac{L}{\mu}\log\frac{\|\mathbf{u}_0 - \mathbf{u}^*\|}{\epsilon}\right)$ iterations with at most $\mathcal{O}\left(\frac{\sigma^2(\mu+L)\log(\|\mathbf{u}_0 - \mathbf{u}^*\|/\epsilon) + L^3\|\mathbf{u}_0 - \mathbf{u}^*\|^2}{\mu^3 \epsilon^2}\right)$ queries to \widehat{F} in expectation.

318 6 Numerical Experiments and Discussion

We now illustrate the empirical performance of stochastic Halpern iteration on robust least square problems. Specifically, given data matrix $\mathbf{A} \in \mathbb{R}^{n \times d}$ and noisy observation vector $\mathbf{b} \in \mathbb{R}^{n}$ subject to bounded deterministic perturbation δ with $\|\delta\| \leq \rho$, robust least square (RLS) minimizes the worst-case residue as $\min_{\mathbf{x} \in \mathbb{R}^{d}} \max_{\delta: \|\delta\| \leq \rho} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_{2}^{2}$ with $\mathbf{y} = \mathbf{b} + \delta$ [I4]. We Algorithm 3: Restarted Extrapolated Stochastic Halpern-Sharp (Restarted E-Halpern) Input: $\mathbf{v}_{-1} = \mathbf{u}_0 \in \mathbb{R}^d$, $\|\mathbf{u}_0 - \mathbf{u}^*\|$, $0 < \eta_0 \le \frac{1}{3\sqrt{3L}}$, $L, \mu, \epsilon > 0, \sigma$; Initialize: $M = 9L^2$, $\underline{\eta} = \frac{\eta_0(1-2M\eta_0^2)}{1-M\eta_0^2}$, $N = \left\lceil \log\left(\frac{\sqrt{6}\|\mathbf{u}_0 - \mathbf{u}^*\|}{2\epsilon}\right) \right\rceil$; for k = 1 : N do Call Algorithm 2 with initialization $\mathbf{v}_{-1}^{(k)} = \mathbf{u}_0^{(k)} = \mathbf{u}_{k-1}$, $\epsilon_k = \frac{\mu\epsilon\sqrt{M\underline{\eta}^2}}{2\sqrt{5(1+M\underline{\eta}\eta_0)}}$, and $S_1^{(-1)} = S_1^{(0)} = \lceil \frac{8\sigma^2}{\epsilon_k^2} \rceil$, for $K = \left\lceil \frac{4\sqrt{L^2\eta_0\underline{\eta}+1}}{\mu\underline{\eta}} \right\rceil$ iterations, and return \mathbf{u}_k ; Return: \mathbf{u}_N



Figure 1: Empirical comparison of min-max algorithms on the robust least squares problem.

consider solving MI induced from RLS with Lagrangian relaxation where $\mathbf{u} = (\mathbf{x}, \mathbf{y})^T$ and 323 $F(\mathbf{u}) = \left(\nabla_{\mathbf{x}} L_{\lambda}(\mathbf{x}, \mathbf{y}), -\nabla_{\mathbf{y}} L_{\lambda}(\mathbf{x}, \mathbf{y})\right)^{T} \text{ for } L_{\lambda}(\mathbf{x}, \mathbf{y}) = \frac{1}{2n} \|A\mathbf{x} - \mathbf{y}\|_{2}^{2} - \frac{\lambda}{2n} \|\mathbf{y} - \mathbf{b}\|_{2}^{2}. \text{ We use a real-world superconductivity dataset [20] from UCI Machine Learning Repository [13] for our$ 324 325 experiment, which is of size 21263×81 . To ensure the problem is concave in y, we need that 326 $\lambda > 1$; in the experiments, we set $\lambda = 1.5$. For the experiment, we compare Halpern, E-Halpern, 327 and Restarted E-Halpern algorithms with gradient descent-ascent (GDA), extragradient (EG) [25], 328 and Popov's method [39] in stochastic settings. Even though our theoretical results for Restarted 329 E-Halpern require scheduled restarts based on known problem parameters, in the implementation, 330 to avoid complicated parameter tuning and illustrate empirical performance, we restart E-Halpern 331 whenever the norm of stochastic estimator \vec{F} used in E-Halpern halves. All Halpern variants are 332 implemented with PAGE estimator considered in our paper; all other algorithms are implemented 333 using minibatches. Additionally, we compare E-Halpern with the PAGE estimator against E-Halpern 334 with single-sample and mini-batch estimators. 335

We report and plot the (empirical) operator norm $||F(\mathbf{u})||$ against the number of stochastic operator 336 evaluations. Note that evaluations of $||F(\mathbf{u})||$ are only used for plotting but not for running any of 337 the algorithms. We use the same random initialization and tune the batch sizes and step sizes (to the 338 values achieving fastest convergence under noise) for each method by grid search. We use constant 339 batch sizes and constant step sizes for GDA, EG, and Popov. We also choose the batch sizes of PAGE 340 estimator to ensure $\mathbb{E}[\|F(\mathbf{u}_k) - F(\mathbf{u}_k)\|^2] \leq \mathcal{O}(\frac{1}{k})$, which handles error accumulation [26] and 341 early stagnation of stochastic Halpern iteration. We implement all the algorithms in Python and run 342 each algorithm using one CPU core on a macOS machine with Intel 2.3GHz Dual Core i5 Processor 343 and 8GB RAM. 344

We observe that (i) in Figure I(a) both Halpern and E-Halpern exhibit faster convergence to approximate stationary points (with much smaller gradient norm after same number of gradient evaluations) than other algorithms, and restarting E-Halpern provides additional speedup, validating our theoretical insights; (ii) in Figure I(b), E-Halpern with PAGE estimator displays faster convergence compared to other two estimators, in agreement with our theoretical analysis.

350 **References**

- [1] Ahmet Alacaoglu and Yura Malitsky. Stochastic variance reduction for variational inequality
 methods. *arXiv preprint arXiv:2102.08352*, 2021.
- [2] Zeyuan Allen-Zhu. Katyusha: The first direct acceleration of stochastic gradient methods. *The Journal of Machine Learning Research*, 18(1):8194–8244, 2017.
- [3] Yossi Arjevani, Yair Carmon, John C Duchi, Dylan J Foster, Ayush Sekhari, and Karthik Sridha ran. Second-order information in non-convex stochastic optimization: Power and limitations.
 In *Proc. COLT'20*, 2020.
- ³⁵⁸ [4] Amir Beck. *First-order methods in optimization*, volume 25. SIAM, 2017.
- [5] Yair Carmon, Yujia Jin, Aaron Sidford, and Kevin Tian. Variance reduction for matrix games.
 Proc. NeurIPS'19, 32, 2019.
- [6] Tatjana Chavdarova, Gauthier Gidel, François Fleuret, and Simon Lacoste-Julien. Reducing noise in gan training with variance reduced extragradient. *Proc. NeurIPS'19*, 32, 2019.
- [7] Juan Pablo Contreras and Roberto Cominetti. Optimal error bounds for nonexpansive fixed-point iterations in normed spaces. *arXiv preprint, arXiv:2108.10969*, 2021.
- [8] Ashok Cutkosky and Francesco Orabona. Momentum-based variance reduction in non-convex
 SGD. In *Proc. NeurIPS*'19, 2019.
- [9] Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. SAGA: A fast incremental gradient
 method with support for non-strongly convex composite objectives. In *Proc. NIPS'14*, 2014.
- [10] Jelena Diakonikolas. Halpern iteration for near-optimal and parameter-free monotone inclusion
 and strong solutions to variational inequalities. In *Proc. COLT'20*, 2020.
- [11] Jelena Diakonikolas and Puqian Wang. Potential function-based framework for making the gradients small in convex and min-max optimization. *arXiv preprint arXiv:2101.12101*, 2021.
- [12] Jelena Diakonikolas, Constantinos Daskalakis, and Michael Jordan. Efficient methods for
 structured nonconvex-nonconcave min-max optimization. In *Proc. AISTATS*'21, 2021.
- [13] Dheeru Dua and Casey Graff. UCI machine learning repository, 2017. URL http://
 archive.ics.uci.edu/ml.
- [14] Laurent El Ghaoui and Hervé Lebret. Robust solutions to least-squares problems with uncertain
 data. *SIAM Journal on matrix analysis and applications*, 18(4):1035–1064, 1997.
- [15] Francisco Facchinei and Jong-Shi Pang. *Finite-dimensional variational inequalities and com- plementarity problems*. Springer Science & Business Media, 2003.
- [16] Cong Fang, Chris Junchi Li, Zhouchen Lin, and Tong Zhang. Spider: Near-optimal non-convex
 optimization via stochastic path-integrated differential estimator. *Proc. NeurIPS'18*, 2018.
- [17] Saeed Ghadimi and Guanghui Lan. Accelerated gradient methods for nonconvex nonlinear and
 stochastic programming. *Mathematical Programming*, 156(1-2):59–99, 2016.
- [18] Noah Golowich, Sarath Pattathil, Constantinos Daskalakis, and Asuman Ozdaglar. Last iterate
 is slower than averaged iterate in smooth convex-concave saddle point problems. In *Conference on Learning Theory*, pages 1758–1784. PMLR, 2020.
- [19] Benjamin Halpern. Fixed points of nonexpanding maps. *Bulletin of the American Mathematical Society*, 73(6):957–961, 1967.
- [20] Kam Hamidieh. A data-driven statistical model for predicting the critical temperature of a superconductor. *Computational Materials Science*, 154:346–354, 2018. ISSN 0927-0256. doi: https://doi.org/10.1016/j.commatsci.2018.07.052. URL https://www.sciencedirect.
 ³⁹³ com/science/article/pii/S0927025618304877.
 - 10

- [21] Alfredo N Iusem, Alejandro Jofré, Roberto Imbuzeiro Oliveira, and Philip Thompson. Extragra dient method with variance reduction for stochastic variational inequalities. *SIAM Journal on Optimization*, 27(2):686–724, 2017.
- Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance
 reduction. In *Advances in Neural Information Processing Systems*, volume 26, 2013.
- [23] Donghwan Kim. Accelerated proximal point method and forward method for monotone
 inclusions. *arXiv preprint arXiv:1905.05149*, 2019.
- [24] Ulrich Kohlenbach. Applied proof theory: proof interpretations and their use in mathematics.
 Springer Science & Business Media, 2008.
- [25] GM Korpelevich. Extragradient method for finding saddle points and other problems. *Matekon*,
 13(4):35–49, 1977.
- [26] Sucheol Lee and Donghwan Kim. Fast extra gradient methods for smooth structured nonconvex nonconcave minimax problems. In *Proc. NeurIPS*'21, volume 34, 2021.
- Lihua Lei, Cheng Ju, Jianbo Chen, and Michael I Jordan. Non-convex finite-sum optimization
 via SCSG methods. In *Proc. NIPS'17*, 2017.
- [28] Laurentiu Leustean. Rates of asymptotic regularity for halpern iterations of nonexpansive
 mappings. *Journal of Universal Computer Science*, 13(11):1680–1691, 2007.
- [29] Zhize Li, Hongyan Bao, Xiangliang Zhang, and Peter Richtárik. PAGE: A simple and optimal
 probabilistic gradient estimator for nonconvex optimization. *arXiv preprint arXiv:2008.10898*,
 2020.
- [30] Felix Lieder. On the convergence rate of the Halpern-iteration. *Optimization Letters*, 15(2):
 405–418, 2021.
- [31] Nicolas Loizou, Hugo Berard, Gauthier Gidel, Ioannis Mitliagkas, and Simon Lacoste-Julien.
 Stochastic gradient descent-ascent and consensus optimization for smooth games: Convergence
 analysis under expected co-coercivity. *Proc. NeurIPS'21*, 2021.
- [32] George J Minty. Monotone (nonlinear) operators in hilbert space. *Duke Mathematical Journal*,
 29(3):341–346, 1962.
- [33] Arkadi Nemirovski. 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.
- 424 [34] Yurii Nesterov. Dual extrapolation and its applications to solving variational inequalities and 425 related problems. *Mathematical Programming*, 109(2-3):319–344, 2007.
- [35] Lam M Nguyen, Jie Liu, Katya Scheinberg, and Martin Takáč. Sarah: A novel method for
 machine learning problems using stochastic recursive gradient. In *Proc. ICML'17*, 2017.
- ⁴²⁸ [36] Yuyuan Ouyang and Yangyang Xu. Lower complexity bounds of first-order methods for ⁴²⁹ convex-concave bilinear saddle-point problems. *Mathematical Programming*, Aug 2019.
- [37] Balamurugan Palaniappan and Francis Bach. Stochastic variance reduction methods for saddle point problems. *Proc. NIPS'16*, 29, 2016.
- [38] Jong-Shi Pang. Error bounds in mathematical programming. *Mathematical Programming*, 79 (1):299–332, 1997.
- [39] L. D. Popov. A modification of the Arrow-Hurwicz method for search of saddle points.
 Mathematical notes of the Academy of Sciences of the USSR, 28(5):845–848, Nov 1980.
- [40] Sashank J Reddi, Ahmed Hefny, Suvrit Sra, Barnabas Poczos, and Alex Smola. Stochastic
 variance reduction for nonconvex optimization. In *Proc. ICML'16*, 2016.

- [41] R Tyrrell Rockafellar. Monotone operators associated with saddle-functions and minimax
 problems. *Nonlinear functional analysis*, 18(part 1):397–407, 1970.
- [42] R Tyrrell Rockafellar. Monotone operators and the proximal point algorithm. *SIAM journal on control and optimization*, 14(5):877–898, 1976.
- [43] Vincent Roulet and Alexandre d'Aspremont. Sharpness, restart, and acceleration. *SIAM Journal on Optimization*, 30(1):262–289, 2020.
- [44] Shoham Sabach and Shimrit Shtern. A first order method for solving convex bilevel optimization
 problems. *SIAM Journal on Optimization*, 27(2):640–660, 2017.
- [45] Mark Schmidt, Nicolas Le Roux, and Francis Bach. Minimizing finite sums with the stochastic
 average gradient. *Mathematical Programming*, 162(1):83–112, 2017.
- [46] Chaobing Song, Yong Jiang, and Yi Ma. Variance reduction via accelerated dual averaging for
 finite-sum optimization. In *Proc. NeurIPS'20*, 2020.
- [47] Guido Stampacchia. Formes bilineaires coercitives sur les ensembles convexes. *Académie des Sciences de Paris*, 258:4413–4416, 1964.
- [48] Quoc Tran-Dinh and Yang Luo. Halpern-type accelerated and splitting algorithms for monotone
 inclusions. *arXiv preprint arXiv:2110.08150*, 2021.
- [49] Rainer Wittmann. Approximation of fixed points of nonexpansive mappings. Archiv der
 Mathematik, 58(5):486–491, 1992.
- [50] Taeho Yoon and Ernest K Ryu. Accelerated algorithms for smooth convex-concave minimax problems with $O(1/k^2)$ rate on squared gradient norm. In *Proc. ICML*'21, 2021.
- [51] Dongruo Zhou, Pan Xu, and Quanquan Gu. Finding local minima via stochastic nested variance
 reduction. *arXiv preprint arXiv:1806.08782*, 2018.
- [52] Dongruo Zhou, Pan Xu, and Quanquan Gu. Stochastic nested variance reduction for nonconvex
 optimization. In *Proc. NeurIPS'18*, 2018.

462 Checklist

464

465

466

467

468

469

470 471

472

473

474

475 476

477

478

479

480

481

- (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes] The contributions are summarized in Introduction (Section []) and fully supported by detailed proofs in the main body and the appendix.
- (b) Did you describe the limitations of your work? [Yes] Some specific examples are: (i) the algorithm being of the Las Vegas-type, in Section 3.2 and (ii) the dependence on problem parameters other than ϵ and σ being likely suboptimal in Appendix [A].
 - (c) Did you discuss any potential negative societal impacts of your work? [N/A] This is a purely theoretical paper.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
- 2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Section 2]
 - (b) Did you include complete proofs of all theoretical results? [Yes] See Sections 3-5 in the main body and the appendix.
- If you ran experiments...
- (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] See supplemental materials.
- (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
 were chosen)? [Yes] See Section 6.

484 485	(c) Did you report error bars (e.g., with respect to the random seed after running experi- ments multiple times)? [Yes] We use the same initialization for all methods, and tune
486	the parameters of each method to be optimal individually.
487	(d) Did you include the total amount of compute and the type of resources used (e.g., type
488	of GPUs, internal cluster, or cloud provider)? [Yes] See Section 6.
489	4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets
490	(a) If your work uses existing assets, did you cite the creators? [Yes] See Section 6.
491	(b) Did you mention the license of the assets? [Yes] See Section 6.
492	(c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
493	Our code is included in the supplemental material.
494	(d) Did you discuss whether and how consent was obtained from people whose data you're
495	using/curating? [N/A] The data we used is from online public repository.
496	(e) Did you discuss whether the data you are using/curating contains personally identifiable
497	information or offensive content? [N/A] The data we used is from online public
498	repository.
499	5. If you used crowdsourcing or conducted research with human subjects
500	(a) Did you include the full text of instructions given to participants and screenshots, if
501	applicable? [N/A]
502	(b) Did you describe any potential participant risks, with links to Institutional Review
503	Board (IRB) approvals, if applicable? [N/A]
504	(c) Did you include the estimated hourly wage paid to participants and the total amount
505	spent on participant compensation? [N/A]