
~~Projection-free~~ Constrained Stochastic Nonconvex Optimization with State-dependent Markov Data

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We study a projection-free conditional gradient-type algorithm for constrained
2 nonconvex stochastic optimization problems with Markovian data. In particular, we
3 focus on the case when the transition kernel of the Markov chain is state-dependent.
4 Such stochastic optimization problems arise in various machine learning problems
5 including strategic classification and reinforcement learning. For this problem, we
6 establish that the number of calls to the stochastic first-order oracle and the linear
7 minimization oracle to obtain an appropriately defined ϵ -stationary point, are of
8 the order $\mathcal{O}(1/\epsilon^{2.5})$ and $\mathcal{O}(1/\epsilon^{5.5})$ respectively. We also empirically demonstrate
9 the performance of our algorithm on the problem of strategic classification with
10 neural networks.

11 1 Introduction

12 We consider the following stochastic optimization problem

$$\operatorname{argmin}_{\theta \in \Theta} f(\theta) = \operatorname{argmin}_{\theta \in \Theta} \mathbb{E} [F(\theta; x)], \quad (1)$$

13 where (i) the expectation is taken over the stationary distribution, π_θ , of the random vector x , (ii) F
14 (and hence f) is a potentially non-convex function in θ , and (iii) Θ is a compact and convex constraint
15 set. Stochastic approximation algorithms for solving problem (1), given an independent and identically
16 distributed (iid) data stream $\{x_k\}_k$ drawn from π , are well-studied. Such iid assumptions are
17 commonly made in various machine learning and statistical problems including empirical risk
18 minimization [SSBD14], sparse recovery [BJMO12] and compressed sensing [FR13, Lan20]. We
19 refer to [MB11, ABRW12, RSS12, GL13, SZ13, LZ16, ACD⁺19] for a partial list of non-asymptotic
20 upper and lower bounds on the oracle complexity of widely-used stochastic approximation algorithms
21 like the Stochastic Gradient Descent (SGD) and the Stochastic Conditional Gradient Algorithm.

22 Our focus in this work is on the case when the data sequence $\{x_k\}_k$ is drawn from a Markov
23 chain with a state-dependent transition kernel P_θ . Such a setting arises in several machine learning
24 applications including but not limited to strategic classification [HMPW16, CDP15, MDPZH20,
25 LW22] and reinforcement learning [Bar92, GSK13, ZJM21, KMMW19, QW20]. Despite their
26 prevalence in practice, a deeper understanding of the non-asymptotic oracle complexity of stochastic
27 approximation for Markovian data is only now starting to emerge. We establish non-asymptotic
28 oracle complexity results for the stochastic conditional gradient algorithm for non-convex constrained
29 stochastic optimization with Markovian data. To establish our results, from a methodological
30 point-of-view, we leverage the moving-average stochastic gradient estimation technique recently
31 used in [ZSM⁺20, GRW20, XBG22] in the context of constrained optimization with iid data. This
32 technique avoids having to use a mini-batch of samples in each iteration, which turns out to be crucial
33 in the non-iid setup we consider. From a theoretical point-of-view, we assume the so-called drift
34 conditions, a classical assumption in Markov Chain literature [AMP05]. This ensures the existence

35 of a solution to the Poisson equation associated with the underlying Markov chain [DMPS18] which
 36 enables one to decompose the noise present in the stochastic gradient into three components: a
 37 martingale difference sequence, a time-decaying sequence, and a telescopic sum type sequence. The
 38 key idea of our paper is to use this decomposition to construct an auxiliary sequence of iterates with a
 39 time-decaying noise-variance and show that these sequence of iterates are *close* to the iterates of the
 40 original sequence produced by our algorithm. This novel technique is then used in combination with
 41 a merit-function based analysis to establish the oracle complexity results.

42 1.1 Motivating Example

43 Problems of the form in (1) arise in various important applications, e.g., strategic classification, and
 44 reinforcement learning as mentioned above. Below we illustrate the motivation of this work through
 45 the example of strategic classification with adapted best response [LW22]. In strategic classification,
 46 there is a *learner* whose task is to classify a given dataset which is collected from a set of *agents*.
 47 Given the knowledge of the classifier, the agents can distort some of their personal features, in order to
 48 get classified in a predetermined target class. This scenario arises in various applications, e.g., spam
 49 email filtering, and credit score classification. Optimizing the classifier to classify such strategically
 50 modified data where the agents modify the data iteratively can be formulated as problem (1).

51 Formally, let the classifier be $h(x, \theta)$ where $x \in \mathbb{R}^d$ is the feature and θ is the parameter to be
 52 optimized. $h(x; \cdot) : \Theta \rightarrow \mathbb{R}$ is potentially nonconvex. Let the loss function be logistic loss which for
 53 a sample (x, y) , where $y \in \{-1, 1\}$ denotes the corresponding class, is given by,

$$L(\theta; x, y) = \log(1 + \exp(-h(x; \theta))) + (1 - y)h(x; \theta)/2. \quad (2)$$

54 We use x_S , and x_{-S} to denote the subset of feature x which are respectively strategically modifiable,
 55 and non-modifiable by the agents. Then the modified feature (the best response) x'_S reported by the
 56 agent is the solution to the following optimization problem:

$$x'_S = \operatorname{argmax}_{x_S} (h(x; \theta) - c(x_S, x'_S)), \quad (3)$$

57 where $c(x, x')$ is the cost of modifying x_S to x'_S . Let the agents iteratively learn x'_S similar to
 58 [LW22]. Note that unlike [LW22], where the authors deploy a logistic regression classifier and the
 59 closed form solution of the best response is readily known to the agents, it may not be the case in
 60 general. In that case the agents have to possibly learn the best response x'_S using some iterative
 61 optimization algorithm. For example, if the agents use Gradient Ascent then, at every iteration k , a
 62 set \mathcal{I}_k of $n_1 \leq M$ randomly chosen agents out of M agents modify their features as:

$$x_{S,i}^k = \begin{cases} x_{S,i}^{k-1} + \alpha \left(\nabla h(x_{S,i}^{k-1}; \theta_k) - \nabla c(x_{S,i}^{k-1}, x_{S,i}^0) \right) & i \in \mathcal{I}_k \\ x_{S,i}^{k-1} & i \notin \mathcal{I}_k \end{cases} \quad (4)$$

where α is the stepsize. With a little abuse of notation, we use $\nabla h(x_{S,i}^{k-1}; \theta)$ in (4) to denote the
 fact that the gradient is with respect to $x_{S,i}^{k-1}$ while $x_{-S,i}$ remains unchanged. This introduces the
 state-dependent Markov chain dynamics in the training data. The objective function, analogous to
 $f(\theta)$ in (1), is

$$\min_{\theta \in \Theta} \mathbb{E}_{\pi_\theta} [L(\theta; x, y)],$$

63 where π_θ is the stationary joint distribution of (x, y) , and Θ is a convex and compact set, e.g., sparsity
 64 inducing constraint $\|\theta\|_1 \leq R$ from some $R > 0$. The loss evaluated at a single data point (x, y) ,
 65 $L(\theta; x, y)$, is analogous to $F(\theta; x)$ in (1). [DX20], and [LW22] study this problem theoretically and
 66 empirically respectively in an unconstrained strongly convex setting. Our results takes a step towards
 67 analyzing this problem in constrained nonconvex setting. We empirically show the performance of
 68 the stochastic conditional gradient algorithm on a strategic classification problem in Section 4.

69 1.2 Preliminaries and Main Contributions

70 Before we present our main contributions, we introduce our convergence criterion. In constrained
 71 optimization literature, most commonly used convergence criteria are: (i) *Gradient Mapping* (GM),
 72 and (ii) *Frank-Wolfe Gap* (FW-gap). The *Gradient Mapping* at a point $\bar{\theta} \in \Theta$ is defined as

$$\mathcal{G}_\Theta(\bar{\theta}, \nabla f(\bar{\theta}), \beta) := \beta \left(\bar{\theta} - \Pi_\Theta \left(\bar{\theta} - \frac{1}{\beta} \nabla f(\bar{\theta}) \right) \right), \quad (5)$$

73 where $\Pi_{\Theta}(x)$ denotes the orthogonal projection of the vector x onto the set Θ , i.e.,

$$\Pi_{\Theta}\left(\bar{\theta} - \frac{1}{\beta}\nabla f(\bar{\theta})\right) = \operatorname{argmin}_{y \in \Theta} \left\{ \langle \nabla f(\bar{\theta}), y - \bar{\theta} \rangle + \frac{\beta}{2} \|y - \bar{\theta}\|_2^2 \right\}.$$

74 We will use $\Pi_{\Theta}(\bar{\theta}, \nabla f(\bar{\theta}), \beta)$ to denote $\Pi_{\Theta}(\bar{\theta} - \nabla f(\bar{\theta})/\beta)$ when there is no confusion. Note that
 75 when $\Theta \equiv \mathbb{R}^d$ we have $\mathcal{G}_{\Theta}(\bar{\theta}, \nabla f(\bar{\theta}), \beta) = \nabla f(\bar{\theta})$. In other words, for constrained optimization
 76 gradient mapping plays an analogous role of the gradient for unconstrained optimization. The gradient
 77 mapping is a frequently used measure in the literature as a convergence criterion for nonconvex
 78 constrained optimization [Nes18]. We should emphasize here that although the gradient mapping
 79 cannot be computed in the stochastic setting, one can still use it as a convergence measure.

80 [BG22] shows that the above notion of convergence criterion is closely related to the so-called
 81 *Frank-Wolfe Gap*. The FW-gap is defined as

$$g_{\Theta}(\bar{\theta}, \nabla f(\bar{\theta})) := \min_{y \in \Theta} \langle \nabla f(\bar{\theta}), y - \bar{\theta} \rangle. \quad (6)$$

82 The following proposition from [BG22] establishes the relation between the gradient mapping
 83 criterion and the Frank-Wolfe gap:

84 **Proposition 1.1** [BG22] *Let $g_{\Theta}(\cdot)$ be the Frank-Wolfe gap defined in (6) and $\mathcal{G}_{\Theta}(\cdot)$ be the gradient*
 85 *mapping defined in (5). Then, we have*

$$\|\mathcal{G}_{\Theta}(\bar{\theta}, \nabla f(\bar{\theta}), \beta)\|_2^2 \leq g_{\Theta}(\bar{\theta}, \nabla f(\bar{\theta})), \quad \forall \bar{\theta} \in \Theta.$$

86 *Moreover, under standard regularity assumption in smooth optimization (specifically, Assumption 2.1,*
 87 *and 2.2), we have*

$$g_{\Theta}(\bar{\theta}, \nabla f(\bar{\theta})) \leq L \|\mathcal{G}_{\Theta}(\bar{\theta}, \nabla f(\bar{\theta}), \beta)\|_2 / \beta. \quad (7)$$

88 In this work we use a suboptimality measure, closely related to both GM and the FW-gap. At point
 89 $\bar{\theta} \in \Theta$, we define the suboptimality measure $V(\bar{\theta}, z) : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ as [GRW20]

$$V(\bar{\theta}, z) := \|\Pi_{\Theta}(\bar{\theta} - z/\beta) - \bar{\theta}\|_2^2 + \|z - \nabla f(\bar{\theta})\|_2^2, \quad (8)$$

90 where z , formally defined in Algorithm 1, is the moving-average estimate of $\nabla f(\bar{\theta})$. We show the
 91 relation among $V(\bar{\theta}, z)$, GM $\mathcal{G}_{\Theta}(\bar{\theta}, z, \beta)$, and FW-gap $g(\bar{\theta}_k, \nabla f(\bar{\theta}_k))$ in the following proposition.

92 **Proposition 1.2** *Let $\{z_k\}$ be the sequence generated in Algorithm 1. Then, for $k = 1, 2, \dots, N$,*

$$\max\{\|\mathcal{G}_{\Theta}(\bar{\theta}_k, z_k, \beta)\|_2^2, g(\bar{\theta}_k, \nabla f(\bar{\theta}_k))\} \leq V(\bar{\theta}_k, z_k).$$

93 The proof is provided in Appendix A. Observe that by establishing ϵ -stationarity in terms of
 94 $\mathbb{E}[V(\bar{\theta}_k, z_k)]$ one can have a tighter bound on FW-gap than gradient mapping ϵ -stationarity as
 95 described in Definition 1. Indeed $\mathbb{E}[V(\bar{\theta}_k, z_k)] \leq \epsilon$ implies $\mathbb{E}[g_{\Theta}(\bar{\theta}_k, \nabla f(\bar{\theta}_k))] \leq \epsilon$ unlike gradient
 96 mapping criterion where, according to (7), $\mathbb{E}[\|\mathcal{G}_{\Theta}(\bar{\theta}_k, z_k, \beta)\|_2^2] \leq \epsilon$ implies $\mathbb{E}[g_{\Theta}(\bar{\theta}_k, \nabla f(\bar{\theta}_k))] \leq \sqrt{\epsilon}$.

97 The main objective of this work is to find an ϵ -stationary solution to (1), where an ϵ -stationary solution
 98 is defined as follows:

99 **Definition 1** *A point $\bar{\theta}$ is said to be an ϵ -stationary solution to (1), if $\mathbb{E}[V(\bar{\theta}, z)] \leq \epsilon$, where the*
 100 *expectation is taken over all the randomness involved in the problem.*

101 For stochastic Frank-Wolfe-type algorithms, the oracle complexity is measured in terms of number of
 102 calls to the Stochastic First-order Oracle (SFO) and the Linear Minimization Oracle (LMO) used to
 103 solve the sub-problems of the algorithm which involves minimizing a linear function over the
 104 convex constraint set. Formally, we have the following definition.

105 **Definition 2** *For a given point $\theta \in \Theta$, SFO returns the stochastic gradient $\nabla F(\theta, x)$. Given a vector*
 106 *z , LMO returns a vector $v := \operatorname{argmin}_{y \in \Theta} \langle z, y \rangle$.*

107 Hence, in this work, the oracle complexity is measured in terms of the number of calls to SFO and
 108 LMO required by the proposed algorithm to obtain an ϵ -stationary solution as in Definition 1. With
 109 the above preliminaries, we now list our **main contributions**:

Algorithm	Criterion	iid		non-iid			
		SFO	LMO	State-independent MC		State-dependent MC	
				SFO	LMO	SFO	LMO
1-SFW [ZSM ⁺ 20]	FW-gap	$\mathcal{O}(\epsilon^{-3})$	$\mathcal{O}(\epsilon^{-3})$	\times	\times	\times	\times
(ASA+ICG) [XBG22]	GM	$\mathcal{O}(\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-3})$	\times	\times	\times	\times
(ASA+ICG) [This paper]	GM/FW-Gap	\times	\times	$\tilde{\mathcal{O}}(\epsilon^{-2})$	$\tilde{\mathcal{O}}(\epsilon^{-3})$	$\mathcal{O}(\epsilon^{-2.5})$	$\mathcal{O}(\epsilon^{-5.5})$

Table 1: Oracle complexity of projection-free one-sample stochastic conditional gradient algorithms for constrained non-convex optimization, to find an ϵ -stationary point.

- In Theorem 3.1, we show that the number of calls to the SFO and LMO required by the stochastic conditional gradient-type method in Algorithm 1, with *state-dependent* Markovian data, is of order $\mathcal{O}(\epsilon^{-2.5})$ and $\mathcal{O}(\epsilon^{-5.5})$ respectively. To the best of our knowledge, these are the first oracle complexity results for projection-free one-sample stochastic optimization algorithm for constrained nonconvex optimization in the Markovian setting.
- In Theorem 3.2, for the sake of completion, we also show that the number of calls to the SFO and LMO required for the case of *state-independent* Markovian data is of the order $\tilde{\mathcal{O}}(\epsilon^{-2})$ and $\tilde{\mathcal{O}}(\epsilon^{-3})$ respectively. In particular, this turns out to be of the same order as that of iid data ignoring the logarithmic factors.

A summary of the our contributions is provided in Table 1. We also empirically evaluate our algorithm on a strategic classification problem with 2-layer neural network classifier and show that the proposed method obtains encouraging results. We provide an experiment on single-index model regression with sparsity-inducing nuclear-norm ball constraint in Appendix C.

1.3 Related Work

Stochastic Optimization with Dependent Data. Understanding stochastic approximation algorithms like SGD with dependent data in the asymptotic setting has been well-explored in the optimization literature. We refer to [KY03, Bor09, BMP12] for a text-book introduction to such classical results. A few recent results include [AMP05, TD17]. In the unconstrained non-asymptotic setting, [DAJJ12] studies convex optimization with ergodic data sequence. [DL22] uses multi-level gradient estimator and analyze AdaGrad for nonconvex optimization with Markovian Data. Block coordinate descent with homogeneous Markov chain has been analyzed in [SSXY20] for nonconvex unconstrained optimization. [DX20] studies stochastic optimization with decision-dependent data distribution for strongly convex functions in the context of strategic classification.

Sample-average approximation algorithms for constrained convex optimization with ϕ -mixing data was considered in [WPT⁺21]. [SSY18], and [AL22] analyze projected SGD for constrained non-convex optimization with time-homogeneous Markov chain. None of these works consider state-dependent data distribution except [DX20]. But unlike [DX20], we consider constrained nonconvex optimization. There also exists work in the reinforcement learning literature on understanding stochastic optimization with Markovian data; see, for example [XXLZ20, BRS18, DNPR20]. However, such works are invariably focused on specific objective functions arising in the reinforcement learning setup, while our focus is on obtaining results for a general class of functions.

Conditional Gradient-Type Method. There has been significant recent advancements in the conditional gradient algorithm literature although it was developed long back [FW56, LP66]; see [Mig94, Jag13, LJJ15, LJJ15, HJN15, GKS21, BS17], for a non-exhaustive list of recent works. [HK12, HL16] provided expected oracle complexity results for stochastic conditional gradient algorithm in the stochastic convex setup. Better rates were provided by a sliding procedure in [LZ16]. In the non-convex setting, [RSPS16, YSC19, HL16] considered variance reduced stochastic conditional gradient algorithms, and provided expected oracle complexities. [QLX18] analyzed the sliding algorithm in the non-convex setting and provided results for the gradient mapping criterion. All of the above works use increasing orders of mini-batch based gradient-estimate.

To avoid mini-batches, a moving-average gradient estimator based on only one-sample in each iteration for a stochastic conditional gradient-type algorithm was proposed in [MHK20] and [ZSM⁺20] for the convex and non-convex setting. However, several restrictive assumptions have been made in [MHK20] and [ZSM⁺20]. Specifically, [ZSM⁺20] requires that the stochastic gradient $G_1(x, \xi_1)$ has uniformly bounded function value, gradient-norm, and Hessian spectral-norm, and the distribution of the random vector ξ_1 has an absolutely continuous density p such that the norm of the gradient of

156 $\log p$ and spectral norm of the Hessian of $\log p$ has finite fourth and second-moments respectively. In
 157 contrasts, we do not require such stringent assumptions.

158 2 Assumptions

159 We now introduce the precise assumptions we make in this work. Let \mathcal{F}_k be the filtration generated
 160 by $\{\theta_0, \dots, \theta_k, z_0, \dots, z_k, x_1, \dots, x_k\}$. For any mapping $g : \mathbb{R}^d \rightarrow \mathbb{R}^d$ define the following norm
 161 with respect to a function $\mathcal{V} : \mathbb{R}^d \rightarrow [1, \infty)$:

$$\|g\|_{\mathcal{V}} = \sup_{x \in \mathcal{X}} (\|g(x)\|_2 / \mathcal{V}(x)),$$

162 and let $L_{\mathcal{V}} = \{g : \mathbb{R}^d \rightarrow \mathbb{R}^d, \sup_{x \in \mathcal{X}} \|g\|_{\mathcal{V}} < \infty\}$.

163 **Assumption 2.1 (Constraint set)** *The set $\Theta \subset \mathbb{R}^d$ is convex and closed with $\max_{x, y \in \Theta} \|x - y\|_2 \leq D_{\Theta}$,
 164 form some $D_{\Theta} > 0$.*

165 **Assumption 2.2** *Let f be a continuously differentiable function.*

166 **Assumption 2.3** *Let $\xi_{k+1}(\theta_k, x_{k+1}) := \nabla F(\theta_k, x_{k+1}) - \nabla f(\theta_k)$. Then,*

$$\mathbb{E} \left[\|\xi_{k+1}(\theta_k, x_{k+1})\|_2^2 | \mathcal{F}_k \right] \leq \sigma_1^2 \quad \mathbb{E} \left[\|\nabla F(\theta_k, x_{k+1})\|_2^2 | \mathcal{F}_k \right] \leq \sigma_2^2 \quad \sigma^2 := \max(\sigma_1^2, \sigma_2^2).$$

167 **Assumption 2.4** *Let $\{x_k\}_k$ be a Markov chain with transition kernel P_{θ} . For any $\theta \in \Theta$, P_{θ} is
 168 irreducible and aperiodic. Additionally, there exists a function $\mathcal{V} : \mathbb{R}^d \rightarrow [1, \infty)$ and a constant
 169 $\alpha \geq 2$ such that for any compact set $\Theta' \subset \Theta$:*

170 (a) *There exist a set $C \subset \mathbb{R}^d$, an integer I , constants $0 < \lambda < 1$, $b, \kappa, \delta > 0$, and a probability
 171 measure ν such that,*

$$\sup_{\theta \in \Theta'} P_{\theta}^I \mathcal{V}^{\alpha}(x) \leq \lambda \mathcal{V}^{\alpha}(x) + bI(x \in C) \quad \forall x \in \mathbb{R}^d, \quad (9)$$

$$\sup_{\theta \in \Theta'} P_{\theta} \mathcal{V}^{\alpha}(x) \leq \kappa \mathcal{V}^{\alpha}(x) \quad \forall x \in \mathbb{R}^d, \quad (10)$$

$$\inf_{\theta \in \Theta'} P_{\theta}^I(x, A) \geq \delta \nu(A) \quad \forall x \in C, \forall A \in \mathcal{B}_{\mathbb{R}^d}. \quad (11)$$

172 where $\mathcal{B}_{\mathbb{R}^d}$ is the Borel σ -algebra over \mathbb{R}^d .

173 (b) *There exists a constant $c > 0$, such that, for all $x \in \mathbb{R}^d$,*

$$\sup_{\theta \in \Theta'} \|\nabla F(\theta, x)\|_{\mathcal{V}} \leq c, \quad (12)$$

$$\sup_{(\theta, \theta') \in \Theta'} \|\nabla F(\theta, x) - \nabla F(\theta', x)\|_{\mathcal{V}} \leq c \|\theta - \theta'\|_2. \quad (13)$$

174 (c) *There exists a constant $c > 0$, such that, for all $(\theta, \theta') \in \Theta' \times \Theta'$,*

$$\|P_{\theta} g - P_{\theta'} g\|_{\mathcal{V}} \leq c \|g\|_{\mathcal{V}} \|\theta - \theta'\|_2 \quad \forall g \in L_{\mathcal{V}} \quad (14)$$

$$\|P_{\theta} g - P_{\theta'} g\|_{\mathcal{V}^{\alpha}} \leq c \|g\|_{\mathcal{V}^{\alpha}} \|\theta - \theta'\|_2 \quad \forall g \in L_{\mathcal{V}^{\alpha}}. \quad (15)$$

175 Some comments regarding the assumptions are in order. Assumption 2.1, and Assumption 2.2 are
 176 common for constrained optimization [GRW20, XBG22, AL22, ZSM⁺20]. Assumption 2.1, and
 177 Assumption 2.2 together imply the Lipschitz continuity of $f(\cdot)$, i.e., there is a constant $L > 0$ such
 178 that for any $\theta_1, \theta_2 \in \Theta$, we have $|f(\theta_1) - f(\theta_2)| \leq L \|\theta_1 - \theta_2\|_2$. Assumption 2.3 is common in
 179 stochastic optimization literature. Assumption 2.4(a) is a frequently used assumption in Markov
 180 chain literature. It implies that for every $\theta \in \Theta$, there exists a stationary distribution $\pi_{\theta}(x)$, and the
 181 chain is \mathcal{V}^{α} -uniformly ergodic [AMP05]. Assumption 2.4(c) provides smoothness guarantee on the
 182 function $f(\cdot)$. More formally, we have the following proposition.

183 **Proposition 2.1 (Lipschitz continuous gradient [AMP05])** *Let Assumption 2.4 be true. Then $f(\cdot)$
 184 has Lipschitz continuous gradient, i.e., there is a constant $L_G > 0$ such that for any $\theta_1, \theta_2 \in \Theta$:*

$$\|\nabla f(\theta_1) - \nabla f(\theta_2)\|_2 \leq L_G \|\theta_1 - \theta_2\|_2. \quad (16)$$

185 Finally, the most important implication of Assumption 2.4 is that it ensures the existence and regularity
 186 of a solution $u(\theta, x)$ to Poisson equation of the transition kernel P_θ given by $u(\theta, x) - P_\theta u(\theta, x) =$
 187 $\nabla F(\theta, x) - \nabla f(\theta)$. Solution of Poisson equation has been crucial in analyzing additive functionals
 188 of Markov chain (see [AMPO5] for details). In this work, the Poisson equation solution facilitates a
 189 decomposition of the noise as presented in Lemma 3.1 which is a key component of our analysis.

190 3 Main Result

191 In this section we present our main result on the oracle complexity to establish a bound on
 $\mathbb{E}[V(\theta_k, z_k)]$. In order to do so we use Algorithm 1, and 2 similar to [XBG22]. If an exact

Algorithm 1 Inexact Averaged Stochastic Approximation (I-ASA)

Input: $z_0, \theta_0 \in \mathbb{R}^d, \eta_k = (N + k)^{-a}, 1/2 < a < 1, \beta$.

for $k = 1, 2, \dots, N$ **do**

$$y_k = \begin{cases} \min_{y \in \Theta} \left\{ \langle z_k, y - \theta_k \rangle + \frac{\beta}{2} \|y - \theta_k\|_2^2 \right\} & \text{(Projection)} \\ \text{ICG}(z_k, \theta_k, \beta, t_k, \omega) & \text{(No Projection)} \end{cases}$$

$$\theta_{k+1} = \theta_k + \eta_{k+1}(y_k - \theta_k)$$

$$z_{k+1} = (1 - \eta_{k+1})z_k + \eta_{k+1}\nabla F(\theta_k, x_{k+1})$$

end for

Output: θ_k

Algorithm 2 Inexact Conditional Gradient (ICG)

Input: $z, \theta, \beta, t, \omega$.

Set $w_0 = \theta$

for $i = 1, 2, \dots, t - 1$ **do**

Find v_i such that

$$\langle v_i, z + \beta(w_i - \theta) \rangle \leq \operatorname{argmin}_{v \in \Theta} \langle v, z + \beta(w_i - \theta) \rangle + \beta\omega\mathcal{D}_\Theta^2/(i + 2)$$

$$w_{i+1} = (1 - \mu_i)w_i + \mu_i v_i \text{ where } \mu_i = \frac{2}{i+2}$$

end for

Output: w_t

192
 193 minimizer of the following subproblem, which is the projection of $\theta_k - z_k/\beta$ on to Θ , is available,
 194 then Algorithm 1 is same as ASA algorithm introduced in [GRW20].

$$\min_{y \in \Theta} \left\{ \langle z_k, y - \theta_k \rangle + \frac{\beta}{2} \|y - \theta_k\|_2^2 \right\}. \quad (17)$$

195 When a projection operator is unavailable or computationally costly, we use Algorithm 2 instead to
 196 solve (17). At iteration k , Algorithm 2 finds an approximate solution to (17) based on the conditional
 197 gradient algorithm. Algorithm 2 needs access to LMO which is often much cheaper and simpler
 198 to compute than projection operator. We should emphasize that our results are not limited to ICG
 199 method but are valid for any method which can solve (17) within an error of the order of $\{\eta_k\}$.

200 **Theorem 3.1** *Let Assumption 2.1-2.4 be true. Then, for Algorithm 1,*

201 (a) *when a projection operator is available, choosing*

$$\eta_k = (N + k)^{-3/5}, \quad \beta = 1 \quad (18)$$

202 *for $k = 1, 2, \dots, N$ we have*

$$\mathbb{E}[V(\theta_R, z_R)] = \mathcal{O}\left(N^{-\frac{2}{5}}\right),$$

203 (b) *when Algorithm 2 is used to solve (17), choosing*

$$\eta_k = (N + k)^{-3/5}, \quad t_k = \eta_k^{-2}, \quad \beta = 1, \quad \omega = 1, \quad \mu_i = 2/(i + 2) \quad (19)$$

204

for $k = 1, 2, \dots, N$ we have

$$\mathbb{E} [V(\theta_R, z_R)] = \mathcal{O} \left(N^{-\frac{2}{5}} \right),$$

205 where the expectations are taken with respect to all the randomness of the algorithm, and an
206 independent integer random variable $R \in \{1, 2, \dots, N\}$ with probability mass function,

$$P(R = k) = \eta_k / \sum_{k=1}^N \eta_k \quad k \in \{1, 2, \dots, N\}.$$

207 **Remark 1** Note that total number of LMO calls are $\sum_{k=1}^N t_k = \sum_{k=1}^N t_k = \sum_{k=1}^N (N+k)^{2a} =$
208 $\mathcal{O}(N^{11/5})$. In other words, to achieve $\|\mathcal{G}_\Theta(\theta_R, \nabla f(\theta_R), \beta)\|_2^2 \leq \mathbb{E} [V(\theta_R, z_R)] \leq \epsilon$, SFO and LMO
209 complexities are respectively $\epsilon^{-2.5}$, and $\epsilon^{-5.5}$. Note that the SFO complexity will be $\epsilon^{-2.5}$ as long as
210 one has an approximation of the projection operator with approximation error $\mathcal{O}(\eta_k)$.

211 **Remark 2** In Theorem 3.1, one obtains sublinear rate $\max(N^{a-1}, N^{2-4a})$ with $\eta_k = (N+k)^{-a}$
212 for $1/2 < a < 1$. Choosing $a = 3/5$ provides the fastest rate of convergence.

213 Before sketching the outline of the proof, we present the following lemma which provides a decom-
214 position of the noise $\xi_k(\theta_{k-1}, x_k)$ – one of the key result used in the proof of the main theorem.
215 The lemma and its proof are almost same as Lemma A.5 in [Lia10] with the only difference that
216 unlike [Lia10], where the iterates are of SGD, we need to prove it for the iterates of Algorithm 1. We
217 provide the proof in Appendix A.

218 **Lemma 3.1** Let Assumption 2.1-2.4 be true. Then the following decomposition takes place:

$$\xi_k(\theta_{k-1}, x_k) = e_k + \nu_k + \zeta_k,$$

219 where, $\{e_k\}$ is martingale difference sequence, $\mathbb{E} [\|\nu_k\|_2] \leq \eta_k$, and $\zeta_k = (\tilde{\zeta}_k - \tilde{\zeta}_{k+1})/\eta_k$, where
220 $\mathbb{E} [\|\tilde{\zeta}_k\|_2] \leq \eta_k$.

221 **Outline of the proof:** A key step in the analysis of Algorithm 1 involves controlling the expectation
222 of interaction with noise of the form $\langle \nabla f(\theta_k) - \nabla f(\theta_{k-1}), \xi_{k+1}(\theta_k, x_{k+1}) \rangle$. For iid or martingale
223 difference data it is easy to control because $\mathbb{E} [\langle \nabla f(\theta_k) - \nabla f(\theta_{k-1}), \xi_{k+1}(\theta_k, x_{k+1}) \rangle | \mathcal{F}_k] = 0$.
224 But this is no longer true for Markov chain data. To resolve the issue, first notice that under our
225 assumptions, the noise sequence ξ_k can be decomposed into the sum of a martingale difference
226 sequence $\{e_k\}$ and some residual terms $\{\nu_k\}$, and $\{\zeta_k\}$ as shown in Lemma 3.1. Then the key step
227 is to introduce a different sequence of hypothetical iterates $(\tilde{\theta}_k, \tilde{y}_k, \tilde{z}_k)$ for which the noise is small
228 enough so that we can bound $\mathbb{E} [V(\tilde{\theta}_k, \tilde{z}_k)]$, and then show that these hypothetical iterates and the
229 original sequence generated by Algorithm 1 are close enough so that $\mathbb{E} [V(\theta_k, z_k)]$ is of the same
230 order as $\mathbb{E} [V(\tilde{\theta}_k, \tilde{z}_k)]$. This step is the main novelty of the proof.

231 Specifically, consider the following sequence:

$$\tilde{\theta}_0 = \theta_0 \quad \tilde{z}_0 = z_0 \tag{20}$$

$$\tilde{y}_k = \operatorname{argmin}_{y \in \Theta} \left\{ \langle \tilde{z}_k, y - \tilde{\theta}_k \rangle + \frac{\beta}{2} \|y - \tilde{\theta}_k\|_2^2 \right\} \tag{21}$$

$$\tilde{\theta}_{k+1} = \tilde{\theta}_k + \eta_{k+1}(\tilde{y}_k - \tilde{\theta}_k) \tag{22}$$

$$\tilde{z}_{k+1} = z_{k+1} + \tilde{\zeta}_{k+2} \tag{23}$$

232 This also means,

$$\tilde{z}_{k+1} = (1 - a\eta_{k+1})\tilde{z}_k + a\eta_{k+1}(\nabla f(\theta_k) + \tilde{\epsilon}_{k+1}), \tag{24}$$

233 where, $\tilde{\epsilon}_k = e_k + \nu_k + \tilde{\zeta}_k$. Note that by Lemma 3.1, $\mathbb{E} [e_k] = 0$, and $\mathbb{E} [\|\nu_k + \tilde{\zeta}_k\|_2] \leq \eta_k$. First we
234 show that by choosing $\eta_k = (N+k)^{-a}$, $1/2 < a < 1$, and $t_k = 1/\eta_k^2$ one has $\mathbb{E} [\|\tilde{\theta}_k - \theta_k\|_2^2] =$
235 $\mathcal{O}(N^{2-4a})$, and $\mathbb{E} [V(\theta_k, z_k)] \leq 2\mathbb{E} [V(\tilde{\theta}_k, \tilde{z}_k)] + \mathcal{O}(N^{2-4a})$. Then we establish the bound on
236 $V(\tilde{\theta}_k, \tilde{z}_k)$. Combining the above two facts proves Theorem 3.1. We defer the detailed proof to
237 Appendix A.1.

238 **3.1 State-independent Markov Chain**

239 While our main goal in this work is to analyze Algorithm 1 for constrained nonconvex optimization
 240 with state-dependent Markov chain data, we provide the following result on the complexity of
 241 Algorithm 1 for Markov chain data with state-independent transition kernel for the sake of completion.
 242 Here we use P to denote the transition kernel (as opposed to P_θ for state-dependent kernel). Note
 243 that under Assumption 2.4(a), for each θ , the chain is \mathcal{V} -uniformly ergodic, and hence, exponentially
 244 mixing [MT12] in the following sense:

245 **Definition 3** A Markov chain is said to be exponentially mixing, if there exists $C, r > 0$ such that,
 246 for any initial state x ,

$$\|P^n(x, \cdot) - \pi\|_{\mathcal{V}} \leq C \exp(-rn), \quad (25)$$

247 where $P^n(x, \cdot)$ is the distribution of X_n with initial state $X_0 = x$.

248 Now we present our result on the complexity of Algorithm 1 to find an ϵ -stationary solution to (1) for
 249 exponentially-mixing Markov chain data with state-independent transition kernel.

250 **Theorem 3.2** Let Assumption 2.1-2.3 be true. Let Assumption 2.4(a)-(b) be true with P_θ replaced by
 251 P . Then, for Algorithm 1,

252 (a) when the projection operator is available, choosing

$$\eta_k = 1/\sqrt{N}, \quad \beta = 1 \quad (26)$$

253 for $k = 1, 2, \dots, N$ we have

$$\mathbb{E}[V(\theta_R, z_R)] = \mathcal{O}\left(\log N/\sqrt{N}\right),$$

254 (b) when Algorithm 2 is used, choosing

$$\eta_k = 1/\sqrt{N}, \quad t_k = \lceil \sqrt{k} \rceil, \quad \beta = 1, \quad \omega = 1, \quad \mu_i = 2/(i+2) \quad (27)$$

255 for $k = 1, 2, \dots, N$ we have

$$\mathbb{E}[V(\theta_R, z_R)] = \mathcal{O}\left(\log N/\sqrt{N}\right),$$

256 where the expectation is taken with respect to all the randomness of the algorithm, and an independent
 257 integer random variable $R \in \{1, 2, \dots, N\}$ whose probability mass function is given by,

$$P(R = k) = \eta_k / \sum_{k=1}^N \eta_k \quad k \in \{1, 2, \dots, N\}.$$

258 We defer the proof to the Appendix.

259 **Remark 3** To find an ϵ -stationary point, the total number of calls to SFO and LMO are $\tilde{\mathcal{O}}(\epsilon^{-2})$,
 260 and $\tilde{\mathcal{O}}(\epsilon^{-3})$, where $\tilde{\mathcal{O}}(\cdot)$ denotes the order ignoring logarithmic factors.

261 **Remark 4** The authors of [AL22] obtain the same rate as in Theorem 3.2 for constrained (but
 262 projection-based) nonconvex optimization with state-independent exponentially mixing data. In the
 263 state-dependent case, since the transition kernel of the Markov chain is controlled by θ_k , and the
 264 transition kernel is assumed to be only Lipschitz smooth in θ (15), the chain does not necessarily
 265 exponentially mix. In the state-independent case, since the chain mixes exponentially we obtain the
 266 same rate as well. While their results are for projection-based algorithms, we analyze a projection-
 267 free LMO-based algorithm since LMO is often computationally cheaper than projection.

268 **4 Experimental Evaluation**

269 In this section we illustrate our algorithm on the strategic classification problem as described in
 270 Section 1.1 with the GiveMeSomeCredit¹ dataset. The main task is a credit score classification

¹Available at <https://www.kaggle.com/c/GiveMeSomeCredit/data>

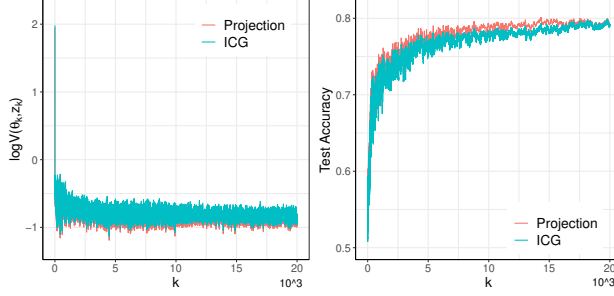


Figure 1: Strategic Classification: (Left): Performance of Algorithm 1 with and without the projection operator. (Right): Test Accuracy with Algorithm 1 with and without the projection operator.

271 problem where the bank (learner) has to decide whether a loan should be granted to a client. Given
 272 the knowledge of the classifier the clients (agents) can distort some of their personal traits in order to
 273 get approved for a loan. Here we use a 2-layer neural network with width m as the classifier, given by

$$h(x; \mathcal{W}, \mathcal{A}, \mathcal{B}) = \sum_{i=1}^m \mathcal{A}_i v(\mathcal{W}_i^\top x + \mathcal{B}_i),$$

274 where $v(\cdot)$ is the activation function, $\mathcal{W}_i \in \mathbb{R}^d$, $\mathcal{W} = [\mathcal{W}_1, \mathcal{W}_2, \dots, \mathcal{W}_d]^\top \in \mathbb{R}^{m \times d}$, $\mathcal{A} =$
 275 $(\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_m) \in \mathbb{R}^m$, $\mathcal{B} = (\mathcal{B}_1, \mathcal{B}_2, \dots, \mathcal{B}_m) \in \mathbb{R}^m$. We will use θ to collectively denote
 276 $(\mathcal{W}, \mathcal{A}, \mathcal{B})$. We impose the constraint of sparsity on the classifier given by $\|\theta\|_1 \leq R$ for some $R > 0$.
 277 As loss function we consider logistic loss as shown in (2). We consider a quadratic cost given by
 278 $c(x, x') = \|x_S - x'_S\|_2^2 / (2\lambda)$ where λ is the sensitivity of the underlying distribution on θ . We
 279 assume that the agents iteratively learn x'_S similar to [LW22]. Note that unlike [LW22], the closed
 280 form of best response is not known here. So we assume that the agents use Gradient Ascent (GA) to
 281 learn the best response. For $\|\theta\|_1 \leq R$ constraint, the LMO in Algorithm 2 at iteration k is given by
 282 $-R \text{sign}(q_i)$, where $i = \text{argmax}_{j=1, \dots, d} |q_j|$, $q = z + \beta(w_k - \theta)$, and q_j is the j -th coordinate of q .
 283 We select a subset of randomly chosen $M = 2000$ samples (agents) such that the dataset is balanced.
 284 Each agent has 10 features. Note that since Algorithm 1 computes the gradient on one sample at
 285 every iterate, the computation time is independent of the total number of agents. We assume that
 286 the agents can modify Revolving Utilization, Number of Open Credit Lines, and Number of Real
 287 Estate Loans or Lines. In this experiment we set $n_1 = 200$. Similar to [LW22], we set $\alpha = 0.5\lambda$,
 288 and $\lambda = 0.01$. For the classifier, the activation function is chosen as *sigmoidal*, and $m = 400$. We
 289 set $N = 20000$, and $R = 4000$. All the parameters of Algorithm 1 are chosen as described in (19).
 290 Figure 4 shows that Algorithm 1 finds an ϵ -stationary point of the strategic classification problem.
 291 We show that Algorithm 1 performs comparably with Averaged Stochastic Approximation with the
 292 projection operator. Each curve in Figure 4 is an average of 50 repetitions.

293 5 Discussion

294 In this work we provide oracle complexity results for the stochastic conditional gradient algorithm
 295 to find an ϵ -stationary point of a constrained nonconvex optimization problem with state-dependent
 296 Markovian data. In Theorem 3.1, we show that the number of calls to the SFO and LMO required
 297 by the stochastic conditional gradient-type method in Algorithm 1, with *state-dependent* Markovian
 298 data, is $\mathcal{O}(\epsilon^{-2.5})$ and $\mathcal{O}(\epsilon^{-5.5})$ respectively. To the best of our knowledge, these are the first oracle
 299 complexity results in this setting. In Theorem 3.2, we show that SFO and LMO complexity in the
 300 case of state-independent Markovian data is $\tilde{\mathcal{O}}(\epsilon^{-2})$ and $\tilde{\mathcal{O}}(\epsilon^{-3})$ respectively, which matches the
 301 corresponding results in the iid setting.

302 There are various avenues for further extensions. Establishing lower bounds on the oracle complexity
 303 of projection-free algorithms in the Markovian data setting is extremely interesting. It is also
 304 intriguing to establish upper and lower bounds on the oracle complexity for more general types of
 305 dependent data sequences arising in applications, including ϕ and α mixing sequences. Yet another
 306 exciting direction is that of designing algorithms adaptive to the dependency in the data that achieve
 307 potentially better oracle complexity bounds.

308 **References**

- 309 [ABRW12] Alekh Agarwal, Peter Bartlett, Pradeep Ravikumar, and Martin Wainwright.
310 Information-theoretic lower bounds on the oracle complexity of stochastic convex
311 optimization. *IEEE Transactions on Information Theory*, 58(5):3235–3249, 2012.
312 (Cited on page 1.)
- 313 [ACD⁺19] Yossi Arjevani, Yair Carmon, John C Duchi, Dylan J Foster, Nathan Srebro, and Blake
314 Woodworth. Lower bounds for non-convex stochastic optimization. *arXiv preprint*
315 *arXiv:1912.02365*, 2019. (Cited on page 1.)
- 316 [AL22] Ahmet Alacaoglu and Hanbaek Lyu. Convergence and complexity of stochastic
317 subgradient methods with dependent data for nonconvex optimization. *arXiv preprint*
318 *arXiv:2203.15797*, 2022. (Cited on pages 4, 5, and 8.)
- 319 [AMP05] Christophe Andrieu, Éric Moulines, and Pierre Priouret. Stability of stochastic ap-
320 proximation under verifiable conditions. *SIAM Journal on control and optimization*,
321 44(1):283–312, 2005. (Cited on pages 1, 4, 5, 6, 15, and 24.)
- 322 [Bar92] Peter L Bartlett. Learning with a slowly changing distribution. In *Proceedings of the*
323 *fifth annual workshop on Computational learning theory*, pages 243–252, 1992. (Cited
324 on page 1.)
- 325 [BG22] Krishnakumar Balasubramanian and Saeed Ghadimi. Zeroth-order nonconvex stochas-
326 tic optimization: Handling constraints, high dimensionality, and saddle points. *Founda-*
327 *tions of Computational Mathematics*, 22(1):35–76, 2022. (Cited on page 3.)
- 328 [BJMO12] Francis Bach, Rodolphe Jenatton, Julien Mairal, and Guillaume Obozinski. Structured
329 sparsity through convex optimization. *Statistical Science*, 27(4):450–468, 2012. (Cited
330 on page 1.)
- 331 [BMP12] Albert Benveniste, Michel Métivier, and Pierre Priouret. *Adaptive algorithms and*
332 *stochastic approximations*, volume 22. Springer Science & Business Media, 2012.
333 (Cited on page 4.)
- 334 [Bor09] Vivek S Borkar. *Stochastic approximation: a dynamical systems viewpoint*, volume 48.
335 Springer, 2009. (Cited on page 4.)
- 336 [BRS18] Jalaj Bhandari, Daniel Russo, and Raghav Singal. A finite time analysis of temporal
337 difference learning with linear function approximation. In *Conference on learning*
338 *theory*, pages 1691–1692. PMLR, 2018. (Cited on page 4.)
- 339 [BS17] Amir Beck and Shimrit Shtern. Linearly convergent away-step conditional gradient
340 for non-strongly convex functions. *Mathematical Programming*, 164(1-2):1–27, 2017.
341 (Cited on page 4.)
- 342 [CDP15] Yang Cai, Constantinos Daskalakis, and Christos Papadimitriou. Optimum statistical
343 estimation with strategic data sources. In *Conference on Learning Theory*, pages
344 280–296. PMLR, 2015. (Cited on page 1.)
- 345 [CXY21] Rong Chen, Han Xiao, and Dan Yang. Autoregressive models for matrix-valued time
346 series. *Journal of Econometrics*, 222(1):539–560, 2021. (Cited on page 27.)
- 347 [DAJJ12] John C Duchi, Alekh Agarwal, Mikael Johansson, and Michael I Jordan. Ergodic
348 mirror descent. *SIAM Journal on Optimization*, 22(4):1549–1578, 2012. (Cited on
349 page 4.)
- 350 [DL22] Ron Dorfman and Kfir Y Levy. Adapting to mixing time in stochastic optimization
351 with markovian data. *arXiv preprint arXiv:2202.04428*, 2022. (Cited on page 4.)
- 352 [DMPS18] Randal Douc, Eric Moulines, Pierre Priouret, and Philippe Soulier. *Markov chains*.
353 Springer, 2018. (Cited on page 2.)
- 354 [DNPR20] Think T Doan, Lam M Nguyen, Nhan H Pham, and Justin Romberg. Convergence rates
355 of accelerated markov gradient descent with applications in reinforcement learning.
356 *arXiv preprint arXiv:2002.02873*, 2020. (Cited on page 4.)
- 357 [DX20] Dmitry Drusvyatskiy and Lin Xiao. Stochastic optimization with decision-dependent
358 distributions. *arXiv preprint arXiv:2011.11173*, 2020. (Cited on pages 2 and 4.)

- 359 [FR13] Simon Foucart and Holger Rauhut. An invitation to compressive sensing. In *A*
360 *mathematical introduction to compressive sensing*, pages 1–39. Springer, 2013. (Cited
361 on page 1.)
- 362 [FW56] Marguerite Frank and Philip Wolfe. An algorithm for quadratic programming. *Naval*
363 *research logistics quarterly*, 3(1-2):95–110, 1956. (Cited on page 4.)
- 364 [GKS21] Dan Garber, Atara Kaplan, and Shoham Sabach. Improved complexities of condi-
365 tional gradient-type methods with applications to robust matrix recovery problems.
366 *Mathematical Programming*, 186(1):185–208, 2021. (Cited on page 4.)
- 367 [GL13] Saeed Ghadimi and Guanghui Lan. Stochastic first-and zeroth-order methods for
368 nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368,
369 2013. (Cited on page 1.)
- 370 [GRW20] Saeed Ghadimi, Andrzej Ruszczyński, and Mengdi Wang. A single timescale stochastic
371 approximation method for nested stochastic optimization. *SIAM Journal on Optimiza-*
372 *tion*, 30(1):960–979, 2020. (Cited on pages 1, 3, 5, 6, 15, 22, and 23.)
- 373 [GSK13] Yair Goldberg, Rui Song, and Michael R Kosorok. Adaptive q-learning. In *From Prob-*
374 *ability to Statistics and Back: High-Dimensional Models and Processes—A Festschrift*
375 *in Honor of Jon A. Wellner*, pages 150–162. Institute of Mathematical Statistics, 2013.
376 (Cited on page 1.)
- 377 [HJN15] Zaid Harchaoui, Anatoli Juditsky, and Arkadi Nemirovski. Conditional gradient algo-
378 rithms for norm-regularized smooth convex optimization. *Mathematical Programming*,
379 152(1-2):75–112, 2015. (Cited on pages 4 and 28.)
- 380 [HK12] Elad Hazan and Satyen Kale. Projection-free online learning. In *29th International*
381 *Conference on Machine Learning, ICML 2012*, pages 521–528, 2012. (Cited on page 4.)
- 382 [HL16] Elad Hazan and Haipeng Luo. Variance-reduced and projection-free stochastic opti-
383 mization. In *International Conference on Machine Learning*, pages 1263–1271, 2016.
384 (Cited on page 4.)
- 385 [HMPW16] Moritz Hardt, Nimrod Megiddo, Christos Papadimitriou, and Mary Wootters. Strategic
386 classification. In *Proceedings of the 2016 ACM conference on innovations in theoretical*
387 *computer science*, pages 111–122, 2016. (Cited on page 1.)
- 388 [Jag13] Martin Jaggi. Revisiting frank-wolfe: Projection-free sparse convex optimization. In
389 *International Conference on Machine Learning*, pages 427–435. PMLR, 2013. (Cited
390 on pages 4, 16, 24, and 28.)
- 391 [JS10] Martin Jaggi and Marek Sulovský. A simple algorithm for nuclear norm regularized
392 problems. In *ICML*, 2010. (Cited on page 28.)
- 393 [KMMW19] Belhal Karimi, Blazej Miasojedow, Eric Moulines, and Hoi-To Wai. Non-asymptotic
394 analysis of biased stochastic approximation scheme. In *Conference on Learning Theory*,
395 pages 1944–1974. PMLR, 2019. (Cited on page 1.)
- 396 [KY03] Harold Kushner and George Yin. *Stochastic approximation and recursive algorithms*
397 *and applications*, volume 35. Springer Science & Business Media, 2003. (Cited on
398 page 4.)
- 399 [Lan20] Guanghui Lan. *First-order and stochastic optimization methods for machine learning*.
400 Springer, 2020. (Cited on page 1.)
- 401 [Lia10] Faming Liang. Trajectory averaging for stochastic approximation mcmc algorithms.
402 *The Annals of Statistics*, 38(5):2823–2856, 2010. (Cited on pages 7 and 15.)
- 403 [LJJ15] Simon Lacoste-Julien and Martin Jaggi. On the global linear convergence of Frank-
404 Wolfe optimization variants. In *Advances in Neural Information Processing Systems*,
405 pages 496–504, 2015. (Cited on page 4.)
- 406 [LP66] Evgeny Levitin and Boris Polyak. Constrained minimization methods. *USSR Compu-*
407 *tational mathematics and mathematical physics*, 6(5):1–50, 1966. (Cited on page 4.)
- 408 [LW22] Qiang Li and Hoi-To Wai. State dependent performative prediction with stochastic
409 approximation. In *International Conference on Artificial Intelligence and Statistics*,
410 pages 3164–3186. PMLR, 2022. (Cited on pages 1, 2, and 9.)

- 411 [LZ16] Guanghui Lan and Yi Zhou. Conditional gradient sliding for convex optimization.
412 *SIAM Journal on Optimization*, 26(2):1379–1409, 2016. (Cited on pages 1 and 4.)
- 413 [MB11] Eric Moulines and Francis Bach. Non-asymptotic analysis of stochastic approximation
414 algorithms for machine learning. *Advances in neural information processing systems*,
415 24, 2011. (Cited on page 1.)
- 416 [MDPZH20] Celestine Mendler-Dünger, Juan Perdomo, Tijana Zrnic, and Moritz Hardt. Stochastic
417 optimization for performative prediction. *Advances in Neural Information Processing
418 Systems*, 33:4929–4939, 2020. (Cited on page 1.)
- 419 [MHK20] Aryan Mokhtari, Hamed Hassani, and Amin Karbasi. Stochastic conditional gradient
420 methods: From convex minimization to submodular maximization. *Journal of machine
421 learning research*, 2020. (Cited on page 4.)
- 422 [Mig94] Athanasios Migdalas. A regularization of the frank—wolfe method and unification of
423 certain nonlinear programming methods. *Mathematical Programming*, 65(1):331–345,
424 1994. (Cited on page 4.)
- 425 [MT12] Sean P Meyn and Richard L Tweedie. *Markov chains and stochastic stability*. Springer
426 Science & Business Media, 2012. (Cited on page 8.)
- 427 [Nes18] Yurii Nesterov. *Lectures on convex optimization*, volume 137. Springer, 2018. (Cited
428 on page 3.)
- 429 [QLX18] Chao Qu, Yan Li, and Huan Xu. Non-convex conditional gradient sliding. In *International
430 Conference on Machine Learning*, pages 4208–4217. PMLR, 2018. (Cited on
431 page 4.)
- 432 [QW20] Guannan Qu and Adam Wierman. Finite-time analysis of asynchronous stochastic
433 approximation and q -learning. In *Conference on Learning Theory*, pages 3185–3205.
434 PMLR, 2020. (Cited on page 1.)
- 435 [RSPS16] Sashank J Reddi, Suvit Sra, Barnabás Póczos, and Alex Smola. Stochastic frank-wolfe
436 methods for nonconvex optimization. In *2016 54th Annual Allerton Conference on
437 Communication, Control, and Computing (Allerton)*, pages 1244–1251. IEEE, 2016.
438 (Cited on page 4.)
- 439 [RSS12] Alexander Rakhlin, Ohad Shamir, and Karthik Sridharan. Making gradient descent
440 optimal for strongly convex stochastic optimization. In *Proceedings of the 29th
441 International Conference on Machine Learning*, pages 1571–1578, 2012. (Cited on page 1.)
- 443 [SSBD14] Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From
444 theory to algorithms*. Cambridge university press, 2014. (Cited on page 1.)
- 445 [SSXY20] Tao Sun, Yuejiao Sun, Yangyang Xu, and Wotao Yin. Markov chain block coordinate
446 descent. *Computational Optimization and Applications*, 75(1):35–61, 2020. (Cited on
447 page 4.)
- 448 [SSY18] Tao Sun, Yuejiao Sun, and Wotao Yin. On markov chain gradient descent. *Advances
449 in neural information processing systems*, 31, 2018. (Cited on page 4.)
- 450 [SZ13] Ohad Shamir and Tong Zhang. Stochastic gradient descent for non-smooth optimiza-
451 tion: Convergence results and optimal averaging schemes. In *International conference
452 on machine learning*, pages 71–79. PMLR, 2013. (Cited on page 1.)
- 453 [TD17] Vladislav B Tadić and Arnaud Doucet. Asymptotic bias of stochastic gradient search.
454 *The Annals of Applied Probability*, 27(6):3255–3304, 2017. (Cited on page 4.)
- 455 [WPT⁺21] Yafei Wang, Bo Pan, Wei Tu, Peng Liu, Bei Jiang, Chao Gao, Wei Lu, Shangling
456 Jui, and Linglong Kong. Sample average approximation for stochastic optimization
457 with dependent data: Performance guarantees and tractability. *arXiv preprint
458 arXiv:2112.05368*, 2021. (Cited on page 4.)
- 459 [XBG22] Tesi Xiao, Krishnakumar Balasubramanian, and Saeed Ghadimi. A projection-free
460 algorithm for constrained stochastic multi-level composition optimization. *arXiv
461 preprint arXiv:2202.04296*, 2022. (Cited on pages 1, 4, 5, and 6.)
- 462 [XXLZ20] Huaqing Xiong, Tengyu Xu, Yingbin Liang, and Wei Zhang. Non-asymptotic con-
463 vergence of adam-type reinforcement learning algorithms under markovian sampling.
464 *arXiv preprint arXiv:2002.06286*, 2020. (Cited on page 4.)

- 465 [YSC19] Alp Yurtsever, Suvrit Sra, and Volkan Cevher. Conditional gradient methods via
 466 stochastic path-integrated differential estimator. In *International Conference on Ma-*
 467 *chine Learning*, pages 7282–7291. PMLR, 2019. (Cited on page 4.)
- 468 [ZJM21] Kelly Zhang, Lucas Janson, and Susan Murphy. Statistical inference with m-estimators
 469 on adaptively collected data. *Advances in Neural Information Processing Systems*, 34,
 470 2021. (Cited on page 1.)
- 471 [ZSM⁺20] Mingrui Zhang, Zebang Shen, Aryan Mokhtari, Hamed Hassani, and Amin Karbasi.
 472 One-sample Stochastic Frank-Wolfe. In *International Conference on Artificial Intelli-*
 473 *gence and Statistics*, pages 4012–4023. PMLR, 2020. (Cited on pages 1, 4, and 5.)

474 Checklist

475 The checklist follows the references. Please read the checklist guidelines carefully for information on
 476 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or
 477 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing
 478 the appropriate section of your paper or providing a brief inline description. For example:

- 479 • Did you include the license to the code and datasets? **[Yes]** See Section ??.
- 480 • Did you include the license to the code and datasets? **[No]** The code and the data are
 481 proprietary.
- 482 • Did you include the license to the code and datasets? **[N/A]**

483 Please do not modify the questions and only use the provided macros for your answers. Note that the
 484 Checklist section does not count towards the page limit. In your paper, please delete this instructions
 485 block and only keep the Checklist section heading above along with the questions/answers below.

- 486 1. For all authors...
 - 487 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
 488 contributions and scope? **[Yes]**
 - 489 (b) Did you describe the limitations of your work? **[Yes]**
 - 490 (c) Did you discuss any potential negative societal impacts of your work? **[N/A]**
 - 491 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
 492 them? **[Yes]**
- 493 2. If you are including theoretical results...
 - 494 (a) Did you state the full set of assumptions of all theoretical results? **[Yes]**
 - 495 (b) Did you include complete proofs of all theoretical results? **[Yes]**
- 496 3. If you ran experiments...
 - 497 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
 498 mental results (either in the supplemental material or as a URL)? **[Yes]**
 - 499 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
 500 were chosen)? **[Yes]**
 - 501 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
 502 ments multiple times)? **[Yes]** We report the average of 50 trails
 - 503 (d) Did you include the total amount of compute and the type of resources used (e.g., type
 504 of GPUs, internal cluster, or cloud provider)? **[No]** We did not calculate the exact
 505 timings. However, our experiments are fairly small-scale ones run on a personal laptop
 506 computer, and our main contributions are theoretical.
- 507 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - 508 (a) If your work uses existing assets, did you cite the creators? **[Yes]**
 - 509 (b) Did you mention the license of the assets? **[Yes]**
 - 510 (c) Did you include any new assets either in the supplemental material or as a URL? **[N/A]**
 511
 - 512 (d) Did you discuss whether and how consent was obtained from people whose data you’re
 513 using/curating? **[N/A]**

- 514 (e) Did you discuss whether the data you are using/curating contains personally identifiable
515 information or offensive content? [N/A]
- 516 5. If you used crowdsourcing or conducted research with human subjects...
- 517 (a) Did you include the full text of instructions given to participants and screenshots, if
518 applicable? [N/A]
- 519 (b) Did you describe any potential participant risks, with links to Institutional Review
520 Board (IRB) approvals, if applicable? [N/A]
- 521 (c) Did you include the estimated hourly wage paid to participants and the total amount
522 spent on participant compensation? [N/A]