

AN ADAPTIVE SAMPLING AUGMENTED LAGRANGIAN METHOD FOR STOCHASTIC OPTIMIZATION WITH DETERMINISTIC CONSTRAINTS

R. Bollapragada, C. Karamanli, B. Keith, B. Lazarov, S. Petrides, J. Wang

May 1, 2023

Computers & Mathematics with Applications

Disclaimer

This document was prepared as an account of work sponsored by an agency of the United States government. Neither the United States government nor Lawrence Livermore National Security, LLC, nor any of their employees makes any warranty, expressed or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States government or Lawrence Livermore National Security, LLC. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States government or Lawrence Livermore National Security, LLC, and shall not be used for advertising or product endorsement purposes.

AN ADAPTIVE SAMPLING AUGMENTED LAGRANGIAN METHOD FOR STOCHASTIC OPTIMIZATION WITH DETERMINISTIC CONSTRAINTS *

RAGHU BOLLAPRAGADA[†], CEM KARAMANLI[†], BRENDAN KEITH[‡], BOYAN LAZAROV[§], SOCRATIS PETRIDES[¶], AND JINGYI WANG[¶]

Dedicated with respect and admiration to Leszek Demkowicz on the occasion of his 70th birthday anniversary.

8 **Abstract.** The primary goal of this paper is to provide an efficient solution algorithm based on 9 the augmented Lagrangian framework for optimization problems with a stochastic objective func-10 tion and deterministic constraints. Our main contribution is combining the augmented Lagrangian framework with adaptive sampling, resulting in an efficient optimization methodology validated with 11 practical examples. To achieve the presented efficiency, we consider inexact solutions for the augmented Lagrangian subproblems, and through an adaptive sampling mechanism, we control the variance 14 in the gradient estimates. Furthermore, we analyze the theoretical performance of the proposed scheme by showing equivalence to a gradient descent algorithm on a Moreau envelope function, and we prove sublinear convergence for convex objectives and linear convergence for strongly convex ob-16 jectives with affine equality constraints. The worst-case sample complexity of the resulting algorithm, 17 for an arbitrary choice of penalty parameter in the augmented Lagrangian function, is $\mathcal{O}(\epsilon^{-3-\delta})$, 18 19 where $\epsilon > 0$ is the expected error of the solution and $\delta > 0$ is a user-defined parameter. If the penalty parameter is chosen to be $\mathcal{O}(\epsilon^{-1})$, we demonstrate that the result can be improved to $\mathcal{O}(\epsilon^{-2})$ 20 21 is competitive with the other methods employed in the literature. Moreover, if the objective function is strongly convex with affine equality constraints, we obtain $\mathcal{O}(\epsilon^{-1}\log(1/\epsilon))$ complexity. Finally, we empirically verify the performance of our adaptive sampling augmented Lagrangian framework in machine learning optimization and engineering design problems, including topology optimization 25 of a heat sink with environmental uncertainty.

1. Introduction. We consider constrained stochastic optimization problems of the form

28 (1.1)
$$\min_{x \in \mathcal{X}} f(x) \text{ subject to } c(x) = 0,$$

where the objective function $f: \mathbb{R}^n \to \mathbb{R}$ is the expected value $f(x) = \mathbb{E}_{\zeta}[F(x,\zeta)]$ of smooth random functions $F(\cdot,\zeta): \mathbb{R}^n \to \mathbb{R}$, the constraint set $\mathcal{X} \subset \mathbb{R}^n$ is compact and convex, and the constraint function $c: \mathbb{R}^n \to \mathbb{R}^m$,

$$c(x) \stackrel{\text{def}}{=} Ax - b,$$

2

3

4

5

6

7

26

2.7

is an affine map with $A \in \mathbb{R}^{m \times n}$ and $b \in \mathbb{R}^m$. Our primary motivation is to develop feasible strategies for solving optimal design problems with manufacturing and operational uncertainties [3, 27, 41, 73] (cf. Subsections 6.2 and 6.3) by efficiently solving

Funding: This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344 and the LLNL-LDRD Program under Project tracking No. 22-ERD-009. Release number LLNL-JRNL-848453.

 † Operations Research and Industrial Engineering, The University of Texas at Austin, Austin, TX 78712 (raghu.bollapragada@utexas.edu, cem.karamanli@utexas.edu)

[‡]Division of Applied Mathematics, Brown University, Providence, RI 02912 USA (brendan_keith@brown.edu).

§Center for Design Optimization, Lawrence Livermore National Laboratory, Livermore, CA 94550 (lazarov2@llnl.gov)

¶Center for Applied Scientific Computing, Lawrence Livermore National Laboratory, Livermore, CA 94550 (petrides1@llnl.gov, wang125@llnl.gov)

^{*}Submitted to the editors 29th March 2024.

optimization problems of the form (1.1). Due to the inherently high computational cost, current design problems are often limited to low-dimensional sources of uncertainty or involve smoothly-varying random fields, which can be parameterized by a truncated series expansion with a small number of discrete random variables [46]. The above limitations restrict the practical applicability of some optimization approaches and lead to simplified heuristic procedures requiring subsequent manual intervention and suboptimal design performance [53]. Therefore, designing an efficient and robust optimization framework to address these challenges is crucial. Moreover, constrained stochastic optimization problems like (1.1) also commonly arise in other applications such as machine learning and finance; see, e.g., [4, 28, 75, 78] and references therein.

One of the well-known techniques to solve constrained optimization problems is the augmented Lagrangian method [13, 29, 30, 39, 70]. This method transforms the original constrained optimization problem (1.1) into a sequence of subproblems where the constraint violation is penalized in the objective. The main advantage of this transformation is that it enables using efficient algorithms for solving the subproblems. On the other hand, the major drawback is that multiple subproblems must be solved sequentially. To mitigate the cost of solving the subproblems, inexact solution mechanisms are widely used [43,50,54,56,72,77,86]. Although these mechanisms are well-understood for deterministic problems, the literature on their usage in stochastic settings is limited [55,65]. Indeed, from our perspective, the main challenge in extending the augmented Lagrangian framework to stochastic approximation techniques lies in defining inexactness criteria for the stochastic methods used to solve the subproblems. In this work, we propose stochastic inexactness termination conditions that address this gap and guarantee convergence in expectation.

Adaptive sampling is a powerful technique that is used in stochastic optimization to control the accuracy of gradient estimates in a computationally efficient manner. The idea comes from the following observation, which is made mathematically precise later in the text: There is little need for an accurate gradient estimate in a stochastic solver when the iterates are far from the optimal solution. However, stochastic algorithms require increasingly accurate gradient estimates as the iterates get closer to the solution. To maintain accuracy, adaptive sampling methods dynamically increase the batch/sample size in response to an a posteriori estimate of the variance of the sampled gradients. Theoretical results from the adaptive sampling literature are promising. Indeed, in [23], the authors show that this methodology matches the best achievable complexity bound for unconstrained stochastic programs. Adaptive sampling is also known to be efficient in practice [18]. Recently, adaptive sampling methods have been used to develop efficient proximal/projected gradient algorithms for constrained optimization problems [7, 83]. Nevertheless, projecting gradients at every iteration can be challenging or inefficient, depending on the structure of the constraint set. Therefore, we go beyond the work in [7,83] and consider augmented Lagrangian techniques. In turn, we address a more general class of algorithms and provide greater flexibility for treating the constraint set.

1.1. Contributions. In this paper, we propose an adaptive sampling augmented Lagrangian (ASAL) method by combining the augmented Lagrangian framework with adaptive sampling techniques to solve constrained stochastic optimization problems. We use adaptive sampling to control the accuracy of the gradient estimates when solving the subproblems obtained by penalizing the linear equality constraints. Moreover, we employ inexact solution mechanisms by imposing stochastic inexactness conditions to terminate the inner (i.e., subproblem) iterations. In this way,

we maximize the overall computational efficiency of our approach without sacrificing accuracy. Another important aspect of the methodology is that it relies on proximal/projected gradients to achieve feasibility with respect to the constraint $x \in \mathcal{X}$. Since the method relies only on gradient information, we establish sublinear convergence in the outer iterations for convex objective functions. Furthermore, given a user-defined algorithm parameter $\delta > 0$ and an arbitrary penalty parameter $\alpha > 0$, we find the total expected number of gradient evaluations to achieve an ϵ -accurate solution to be $\mathcal{O}(\epsilon^{-3-\delta})$. Moreover, if the penalty parameter is chosen to be sufficiently large, i.e., $\mathcal{O}(\epsilon^{-1})$, then our result improves to $\mathcal{O}(\epsilon^{-2})$. Finally, the worst-case complexity becomes $\mathcal{O}(\epsilon^{-1}\log(1/\epsilon))$ for strongly convex objective functions and when $\mathcal{X} = \mathbb{R}^n$. Table 1.1 compares our setting and theoretical results with the relevant literature. To evaluate the efficacy of our framework, we compare its performance to baseline algorithms in a collection of model problems from machine learning (Subsection 6.1) and engineering (Subsections 6.2 and 6.3).

Table 1.1

Summary of the theoretical convergence rate and sample complexity results in the relevant literature under different problem settings. In all the works mentioned here, the constraints are deterministic. Here, K denotes the (outer) iteration number, and ϵ denotes the required accuracy. Convergence rates are deterministic for deterministic problems and are in expectation for stochastic problems. Finally, sample complexity for stochastic solvers denotes the total number of expected stochastic gradient evaluations required to get ϵ -accurate solutions.

paper	objective	set (X)	constraints	rate (outer iter)	sample complexity
[50]	convex deterministic	convex compact	linear	$\mathcal{O}(1/K)$	-
[86]	convex deterministic	convex closed	convex	$\mathcal{O}(1/K)$	-
[84]	strongly convex stochastic	\mathbb{R}^n	linear	linear	$\mathcal{O}(\epsilon^{-1})$
[85]	convex stochastic nonsmooth	convex	convex	$\mathcal{O}(1/\sqrt{K})$	$\mathcal{O}(\epsilon^{-2})$
[85]	strongly convex stochastic nonsmooth	convex	convex	$\mathcal{O}(\log(K)/K)$	$\mathcal{O}(\epsilon^{-1}\log(1/\epsilon))$
Theorem 4.5 (arbitrary penalty parameter)	convex stochastic	convex compact	linear	$\mathcal{O}(1/K)$	$\mathcal{O}(\epsilon^{-3-\delta})$
Corollary 4.6 $(\mathcal{O}(\epsilon^{-1}) \text{ penalty parameter})$	convex stochastic	convex compact	linear	$\mathcal{O}(1/K)$	$\mathcal{O}(\epsilon^{-2})$
Theorem 4.12	strongly convex stochastic	\mathbb{R}^n	linear	linear	$\mathcal{O}(\epsilon^{-1}\log(1/\epsilon))$

1.2. Literature Review. The augmented Lagrangian method, also known as the method of multipliers, was first proposed by Hestenes [39] and Powell [70]. In [13], its performance is analyzed and compared to other common approaches, such as penalty and Lagrangian methods; see also [11, 33, 71, 72]. Although there have been extensive research efforts to enhance the performance of the basic augmented Lagrangian method to solve deterministic optimization problems (see, e.g., [14,29,30,50,54,56,86]), the current literature on stochastic optimization problems is limited [42,55,85]. In [42], the authors apply a stochastic augmented Lagrangian method to the domain adaptation problem. In [85], Xu developed stochastic primal-dual methods using the augmented Lagrangian function for solving nonsmooth optimization problems with a large number of constraints. In the aforestated approach, a projected stochastic gradient method is employed for the primal updates, while a randomized coordinate method is used for the dual updates.

For structured optimization problems with linear constraints, the alternating dir-

ection method of multipliers (ADMM) framework is often preferred [21]. There has been significant work on stochastic versions of the ADMM method [65,81,84,88,89]. In [65], the authors consider stochastic ADMM and show a $\mathcal{O}(\log(K)/K)$ convergence rate for strongly convex and $\mathcal{O}(1/\sqrt{K})$ for general convex objective functions. In [84], the authors design an inexact solution mechanism for the subproblems in stochastic ADMM when $\mathcal{X} = \mathbb{R}^n$. There, the authors employ the stochastic gradient method to solve the subproblems and show a linear convergence rate for strongly convex functions. Although our approach also involves inexact solutions, we consider adaptive sampling techniques to solve the subproblems and analyze both general convex and strongly convex functions. Moreover, our formulation allows us to consider implicit constraint sets (i.e., $\mathcal{X} \subset \mathbb{R}^n$) and utilizes only projected (or proximal) stochastic gradients. Other works achieve improved convergence rates by introducing stochastic variance reduction techniques (see, e.g., [59,81,88,89]).

There are many articles on stochastic optimization methods with dynamic sample sizes [7,16-19,23,25,34-36,47,67,74,83]. Most of these works focus on unconstrained problems. Of note is the work by Friedlander and Schmidt [35], which shows linear convergence for finite-sum problems when the sample size increases at a geometric rate. Our work relates to the approach taken in Byrd et al. [23], which shows linear convergence of the expected risk minimization problem and calculates the worst-case complexity bounds for the number of gradient evaluations required to get ϵ -accurate solutions. Byrd et al. [23] also study the theoretical and practical aspects of the so-called *norm test*, which controls the sample sizes. Finally, in [7,83], the authors consider adaptive sampling mechanisms for constrained stochastic programs. In both works, the constraints are represented by an abstract convex set, and the authors propose generalizations of the norm test that utilize projected (reduced) gradients.

Another common methodology to approach (1.1) is using sample average approximation (SAA) techniques [48, 49, 68, 76, 78] which replace the expected value in the objective function with a fixed sample average or other empirical approximation. When it comes to alternative techniques to solve constrained stochastic programs, the sequential quadratic programming (SQP) framework [9, 10, 31, 32, 61, 62] is also often utilized.

- **1.3. Notation.** We denote the set of natural numbers by $\mathbb{N} \stackrel{\text{def}}{=} \{0,1,2,\dots\}$, and the set of positive natural numbers as $\mathbb{N}_+ \stackrel{\text{def}}{=} \{1,2,\dots\}$. Throughout this work, $\|\cdot\|$ denotes the ℓ_2 vector norm or matrix norm and $\langle\cdot,\cdot\rangle$ denotes the ℓ_2 -inner product. Finally, a matrix $A \in \mathbb{R}^{m \times n}$ is indicated to be positive definite by writing $A \succ 0$ and positive semi-definite by writing $A \succeq 0$. $A^T \in \mathbb{R}^{n \times m}$ denotes the transpose of a matrix A.
- 1.4. Organization. This paper is organized as follows. In Section 2, we introduce the preliminary material and assumptions used throughout the paper. The algorithmic framework and its components are given in Section 3. In Section 4, we analyze the convergence and complexity properties of our approach. Practical implementation of the algorithmic components is discussed in Section 5. We demonstrate the numerical performance of our methodology in Section 6. Finally, in Section 7, we provide concluding remarks and discuss avenues for future research.
- 2. Preliminaries and Assumptions. We provide preliminaries regarding the deterministic augmented Lagrangian method and its interpretation as a gradient descent method applied to the Moreau envelope of the dual function. We also state preliminary assumptions and recall results from the literature that are relied on later

- in the paper.
- 2.1. Deterministic Augmented Lagrangian Method. The Lagrangian function for the problem (1.1) is

164 (2.1)
$$\ell(x,\lambda) = f(x) - \langle \lambda, c(x) \rangle,$$

- where $\lambda \in \mathbb{R}^m$ is the Lagrangian (dual) parameter associated to the constraint function c(x). Using (2.1), we can define the saddle-point problem,
- 167 (2.2) $\min_{x \in \mathcal{X}} \sup_{\lambda \in \mathbb{R}^m} \ell(x, \lambda),$
- 168 and note that

$$\sup_{\lambda \in \mathbb{R}^m} \ell(x, \lambda) = \begin{cases} f(x) & \text{for } c(x) = 0, \\ \infty & \text{for } c(x) \neq 0. \end{cases}$$

- Hence, if there exists $x \in \mathcal{X} \cap \{x \in \mathbb{R}^n \mid c(x) = 0\}$, then (2.2) is equivalent to (1.1) in the sense that
- $\min_{x \in \mathcal{X}} \sup_{\lambda \in \mathbb{R}^m} \ell(x, \lambda) = \min_{\{x \in \mathcal{X} | c(x) = 0\}} f(x)$
- 175 and

$$\underset{177}{\operatorname{arg\,min}} \sup_{x \in \mathcal{X}} \ell(x, \lambda) = \underset{\{x \in \mathcal{X} | c(x) = 0\}}{\operatorname{arg\,min}} f(x).$$

A primal-dual iterate pair $(\hat{x}, \hat{\lambda})$ is said to be a stationary point of (2.2) if

179 (2.3)
$$(\hat{x}, \hat{\lambda}) \in \left\{ (x, \lambda) \middle| \frac{\operatorname{proj}_{\mathcal{X}}(x - \eta \nabla \ell_x(x, \lambda)) - x}{\eta} = 0 \text{ and } c(x) = 0 \right\},$$

where $\eta > 0$ and

181 (2.4)
$$\operatorname{proj}_{\mathcal{X}}(y) = \underset{x \in \mathcal{X}}{\operatorname{arg \, min}} \|x - y\|^2,$$

- is the projection of $y \in \mathbb{R}^n$ onto the set \mathcal{X} (see [29,50]). We also refer to the conditions
- 183 in (2.3) as the
- $\frac{184}{185}$ (2.5a) feasibility error: ||c(x)||,
- and the
- 187 (2.5b) stationarity error: $\left\| \frac{\operatorname{proj}_{\mathcal{X}}(x \eta \nabla \ell_x(x, \lambda)) x}{\eta} \right\|$.
- The augmented Lagrangian method is a class of iterative methods that produce stationary points satisfying (2.3) by solving a sequence of subproblems where the
- objective function is the sum of the Lagrangian function $\ell(x,\lambda)$ and a quadratic
- penalty term that penaltizes violation of the equality constraint c(x) = 0. Specifically,

at any given iteration $k \in \mathbb{N}$, the basic primal and dual update rules are given as 193 194 follows:

195 (2.6a)
$$x_k^* \in \underset{x \in \mathcal{X}}{\arg \min} \ \mathcal{L}(x, \lambda_k; \alpha_k),$$

$$\lambda_{k+1} = \lambda_k - \alpha_k c(x_k^*),$$

where $\alpha_k > 0$ is the penalty parameter and 198

199 (2.7)
$$\mathcal{L}(x,\lambda;\alpha) = f(x) - \langle \lambda, c(x) \rangle + \frac{\alpha}{2} ||c(x)||^2,$$

- is the augmented Lagrangian function. Without restrictions on the objective function 200 201 f(x), the subproblem in (2.6a) may be unbounded. In this paper, we invoke assumptions that ensure this is not the case (cf. Assumption 2.2 or Assumption 2.3) as well 202
- as some other basic assumptions of additional utility. 203

2.2. Assumptions. We make the following assumptions about the objective 204 function, the constraint function, and the existence of the solution. 205

Assumption 2.1. The objective function $f: \mathbb{R}^n \to \mathbb{R}$ is a convex continuously 206 differentiable function on \mathcal{X} . That is, $\nabla^2 f(x) \succeq 0$, for all $x \in \mathcal{X}$. In addition, the 207 gradient of the objective function $\nabla f: \mathbb{R}^n \to \mathbb{R}^n$ is Lipschitz continuous on \mathcal{X} with 208 Lipschitz constant $L < \infty$. That is, 209

$$\|\nabla f(x) - \nabla f(y)\| \le L\|x - y\| \quad \forall x, y \in \mathcal{X}.$$

212 Assumption 2.1 implies that the augmented Lagrangian function is also a convex 213 function with respect to x on \mathcal{X} . That is, for any $\alpha > 0$,

$$\nabla_{xx}^2 \mathcal{L}(x,\lambda;\alpha) = \nabla^2 f(x) + \alpha A^T A \succeq 0 \quad \forall x \in \mathcal{X}, \lambda \in \mathbb{R}^m.$$

Note that the affine constraint function $c: \mathbb{R}^n \to \mathbb{R}^m$ is Lipschitz continuous on 216 \mathcal{X} with Lipschitz constant ||A||. That is, for all $x, y \in \mathcal{X}$, 217

$$||c(x) - c(y)|| = ||A(x - y)|| \le ||A|| ||x - y||.$$

- Moreover, as a consequence of Assumption 2.1 and (2.8), we can show that the gradient 219
- of the augmented Lagrangian function is Lipschitz continuous with respect to x on \mathcal{X} 220
- with Lipschitz constant $L + \alpha ||A||^2$. That is, due to (2.7), 221

222
$$\nabla_x \mathcal{L}(x,\lambda;\alpha) - \nabla_y \mathcal{L}(y,\lambda;\alpha) = \nabla f(x) - \nabla f(y) + \alpha \langle A, Ax - Ay \rangle,$$

and so 223

$$224 \quad (2.9) \qquad \|\nabla_x \mathcal{L}(x,\lambda;\alpha) - \nabla_y \mathcal{L}(y,\lambda;\alpha)\| \le (L+\alpha||A||^2) \|x-y\|,$$

for all $x, y \in \mathcal{X}$. 225

Assumption 2.2. The set $\mathcal{X} \subset \mathbb{R}^n$ is nonempty, convex, and compact. Also, 226 there exists an optimal primal-dual pair (x^*, λ^*) that satisfies the optimality conditions 227 (2.3).228

The compactness of set \mathcal{X} implies that there exists a positive $D < \infty$ such that

230 (2.10)
$$||x - y|| \le D \quad \forall x, y \in \mathcal{X}.$$

Also, the existence of an optimal solution x^* implies that the problem in (2.6a) is bounded below. That is, for any $x \in \mathcal{X}$, $\lambda \in \mathbb{R}^m$, and $\alpha \geq 0$,

233
$$\mathcal{L}(x,\lambda;\alpha) \geq f(x) - \langle \lambda, c(x) \rangle = f(x) - \langle \lambda, c(x) - c(x^*) \rangle$$
234
$$\geq f(x) - \|\lambda\| \|c(x) - c(x^*)\|$$
235
$$\geq f(x^*) + \langle \nabla f(x^*), x - x^* \rangle - \|\lambda\| \|c(x) - c(x^*)\|$$
236
$$\geq f(x^*) - \|\nabla f(x^*)\| D - \|\lambda\| \|A\| D,$$

where the first inequality is due to $||c(x)||^2 \ge 0$, the equality is due to $c(x^*) = 0$, the third inequality is due to convexity of function f (Assumption 2.1) and (2.10), and the last inequality is due to (2.8) and (2.10). Therefore, (2.6a) is well-defined.

We also develop results for the special case where the augmented Lagrangian function is strongly convex with respect to $x \in \mathcal{X}$.

ASSUMPTION 2.3. The augmented Lagrangian is μ -strongly convex with respect to $x \in \mathcal{X}$. That is,

$$\nabla^2_{xx} \mathcal{L}(x,\lambda) \succeq \mu I \quad \forall x \in \mathcal{X}, \lambda \in \mathbb{R}^m,$$

where $I \in \mathbb{R}^{n \times n}$ is an identity matrix.

241242

266

267

268

269270

271

Note that if the objective function f(x) is μ -strongly convex or A has full column rank, then Assumption 2.3 is trivially satisfied. Moreover, if Assumption 2.3 holds, then (2.6a) is well-defined for any $\lambda_k \in \mathbb{R}^m$.

We also make a standard assumption about the stochastic gradient of $f(x) = \mathbb{E}_{\zeta}[F(x,\zeta)]$.

ASSUMPTION 2.4. The variance in the stochastic gradient of f(x) is bounded. That is, there exist constants $\omega_1, \omega_2 \geq 0$ such that

$$\mathbb{E}_{\zeta}[\|\nabla F(x,\zeta) - \nabla f(x)\|^2] \le \omega_1 \|\nabla f(x)\|^2 + \omega_2, \quad \forall x \in \mathcal{X}.$$

Using Assumptions 2.1, 2.2 and 2.4, it follows that

258
$$\mathbb{E}_{\zeta}[\|\nabla F(x,\zeta) - \nabla f(x)\|^{2}] \leq 2\omega_{1} \|\nabla f(x) - \nabla f(x^{*})\|^{2} + 2\omega_{1} \|\nabla f(x^{*})\|^{2} + \omega_{2}$$

$$\leq 2\omega_{1} L^{2} D^{2} + 2\omega_{1} \|\nabla f(x^{*})\|^{2} + \omega_{2} \stackrel{\text{def}}{=} \omega,$$

where the first inequality is due to the fact that $||a+b||^2 \le 2||a||^2 + 2||b||^2$ for any $a, b \in \mathbb{R}^n$. In turn, we note that combining the assumptions above implies the existence of $\omega \ge 0$ such that

$$\mathbb{E}_{\zeta}[\|\nabla F(x,\zeta) - \nabla f(x)\|^2] \le \omega, \quad \forall x \in \mathcal{X}.$$

2.3. Gradient Descent and the Moreau envelope. The convergence properties of the augmented Lagrangian method are often analyzed by showing its equivalence to a method (e.g., proximal point method) applied to dual problem (cf. [82]). We follow a similar approach in our analysis and show the equivalence of the augmented Lagrangian method and gradient descent method applied to the Moreau envelope [60] of the (negative) dual function. The negative of the dual function of (1.1) is denoted

272 (2.12)
$$q(\lambda) = -\min_{x \in \mathcal{X}} \ell(x, \lambda),$$

and is known to be a convex, proper and continuous function from \mathbb{R}^m to \mathbb{R} [20]. For 273 274 any given $\alpha > 0$, the Moreau envelope of $q(\lambda)$ is defined as follows [60]:

275 (2.13)
$$q_{\alpha}(u) \stackrel{\text{def}}{=} \min_{\lambda} \left[q(\lambda) + \frac{1}{2\alpha} \|\lambda - u\|^2 \right].$$

- 277 In the following lemma, we summarize the important properties of Moreau envelopes.
- LEMMA 2.1. The function $q_{\alpha}(u)$ given in (2.13) is called the Moreau envelope of 278 $q(\lambda)$ and satisfies the following properties. 279
 - (i) [66, Equation 3.2] The gradient of the Moreau envelope is

$$\nabla q_{\alpha}(u) = \frac{1}{\alpha}(u - \operatorname{prox}_{\alpha q}(u)),$$

where 283

280

295

$$\operatorname{prox}_{\alpha q}(u) = \arg\min_{\lambda} \left[q(\lambda) + \frac{1}{2\alpha} \|\lambda - u\|^2 \right].$$

- [6, Corollary 18.19] The gradients $\nabla q_{\alpha}(u)$ are Lipschitz continuous with 286 Lipschitz constant $L_{\alpha} = \alpha^{-1}$. That is, 287
- $\|\nabla q_{\alpha}(u) \nabla q_{\alpha}(v)\| < \alpha^{-1} \|u v\|, \quad \forall u, v \in \mathbb{R}^m.$ (2.15)289
- (iii) [66, Page 136] The Moreau envelope retains the optimal value and the set of 290 minimizers. That is, 291

292
$$\min_{\lambda} q(\lambda) = \min_{\lambda} q_{\alpha}(\lambda) \quad and \quad \underset{\lambda}{\arg\min} q(\lambda) = \underset{u}{\arg\min} q_{\alpha}(u),$$

- 294 where the unique common minimizer $\lambda^* \in \mathbb{R}^m$ satisfies the fixed point equation $\lambda^* = \operatorname{prox}_{\alpha q}(\lambda^*).$
 - (iv) [69, Lemma 2.23] $q(\lambda)$ is strongly convex with parameter $\mu_q > 0$ if and only if $q_{\alpha}(u)$ is strongly convex with parameter $\mu_{\alpha} = \frac{\mu_q}{\mu_q \alpha + 1} > 0$.
- Due to Assumption 2.2 and weak duality [20], we have that $q(\lambda)$ is bounded below. 298 That is, the optimal value q^* is finite. Indeed, 299

300 (2.17)
$$q^* = \min_{\lambda} q_{\alpha}(\lambda) = \min_{\lambda} q(\lambda) = -\max_{\lambda} [-q(\lambda)] \ge -\min_{x \in \mathcal{X}, c(x) = 0} f(x) = -f(x^*)$$
.

- Owing to this fact and the properties of $q_{\alpha}(\lambda)$ in Lemma 2.1, the dual variable $\lambda \to \lambda^*$ 301 will converge by iteratively minimizing $q_{\alpha}(\lambda)$ as in the gradient descent method. More 302 explicitly, we may form a convergent sequence of dual variables as follows: 303
- $\lambda_{k+1} = \lambda_k \alpha \nabla q_{\alpha}(\lambda_k)$ (2.18)304 $= \arg\min_{\lambda} \left[q(\lambda) + \frac{1}{2\alpha} \|\lambda - \lambda_k\|^2 \right]$ 305 $= \arg\min_{\lambda} \left[-\min_{x \in \mathcal{X}} [\ell(x, \lambda)] + \frac{1}{2\alpha} \|\lambda - \lambda_k\|^2 \right]$ 306 $= \arg\max_{\lambda} \left[\min_{x \in \mathcal{X}} [\ell(x, \lambda) - \frac{1}{2\alpha} \|\lambda - \lambda_k\|^2] \right] \,,$ 307 308

where the second equality is due to (2.14) and third equality is due to (2.12). The function $\ell(x,\lambda) - \frac{1}{2\alpha} ||\lambda - \lambda_k||^2$ is convex with respect to x on \mathcal{X} and strongly concave with respect to λ . By Sion's Minimax Theorem [80], we can interchange the min and max operations (cf. [82, Section 10.5.2]) and obtain an equivalent characterization. That is,

$$(2.19)$$

$$314 \quad \max_{\lambda} \left[\min_{x \in \mathcal{X}} [\ell(x, \lambda) - \frac{1}{2\alpha} \|\lambda - \lambda_k\|^2] \right] = \min_{x \in \mathcal{X}} \left[\max_{\lambda} [\ell(x, \lambda) - \frac{1}{2\alpha} \|\lambda - \lambda_k\|^2] \right]$$

$$315 \quad (2.20) \quad = \min_{x \in \mathcal{X}} \left[\max_{\lambda} [f(x) - \langle \lambda, c(x) \rangle - \frac{1}{2\alpha} \|\lambda - \lambda_k\|^2] \right].$$

Note that the optimal solution to the max problem (strongly concave in λ) in the second equality is $\lambda = \lambda_k - \alpha c(x)$. Substituting this expression into (2.20), we find

$$\max_{\lambda} \left[\min_{x \in \mathcal{X}} [\ell(x, \lambda) - \frac{1}{2\alpha} \|\lambda - \lambda_k\|^2] \right] = \min_{x \in \mathcal{X}} \mathcal{L}(x, \lambda_k; \alpha).$$

321 Hence, the dual update λ_{k+1} is given as follows:

322 (2.21a)
$$x_k^* \in \underset{x \in \mathcal{X}}{\operatorname{arg \, min}} \ \mathcal{L}(x, \lambda_k; \alpha)$$

$$\lambda_{k+1} = \lambda_k - \alpha c(x_k^*).$$

We now observe that the primal updates in (2.6a) and (2.21a) are both minimizers of the augmented Lagrangian function within the set \mathcal{X} . This optimization problem can have multiple optimal solutions when the augmented Lagrangian function $\mathcal{L}(x, \lambda_k; \alpha)$ is only a general convex function (not strongly convex). Hence, the updates (2.6a) and (2.21a) may not be the same. However, the dual updates are equivalent due to the following inequality [50, Equation 2.16]: For any $x \in \mathcal{X}$ and $x_k^* \in \arg\min_{x \in \mathcal{X}} \mathcal{L}(x, \lambda_k; \alpha)$,

332 (2.22)
$$||c(x_k^*) - c(x)||^2 \le \frac{2}{\alpha} \left(\mathcal{L}(x, \lambda_k; \alpha) - \mathcal{L}(x_k^*, \lambda_k; \alpha) \right).$$

Therefore, all solutions of $\min_{x \in \mathcal{X}} \mathcal{L}(x, \lambda_k; \alpha)$ have the same constraint function value c(x) and the augmented Lagrangian method is equivalent to the gradient descent method applied to the Moreau envelope (2.13). Finally, we conclude this section on preliminary material by noting that

$$\nabla q_{\alpha}(\lambda_k) = c(x_k^*),$$

339 due to (2.18) and (2.21b).

3. Algorithmic Framework. This section begins with a description of a generic inexact augmented Lagrangian framework for solving (1.1). We then provide a complete description of our algorithm, which employs the adaptive sampling proximal gradient method [7,83] to minimize the augmented Lagrangian function (2.7) defined at each iteration.

Each primal variable update (2.6a) in the augmented Lagrangian method involves solving a computationally expensive optimization problem, namely,

347 (3.1)
$$\min_{x \in \mathcal{X}} \mathcal{L}(x, \lambda_k; \alpha),$$

where $\alpha > 0$ is the penalty parameter and λ_k is the dual variable at iteration $k \in \mathbb{N}$. 348 349 Owing to the stochastic nature of sampling the objective function $f(x) = \mathbb{E}_{\zeta}[F(x,\zeta)]$, the exact solutions to these subproblems cannot be obtained efficiently. Therefore, we 350 work with the inexact augmented Lagrangian framework outlined in Algorithm 3.1. At each iteration of this meta-algorithm, the subproblem (3.1) is solved (inexactly) 352 by a given subproblem solver \mathcal{S} until certain as yet unspecified inexactness conditions 353 hold (cf. Subsection 3.2). Of course, the dual variable update incurs errors attributed 354 to the inexact primal solves. However, if appropriate inexactness conditions are used 355 to terminate the subproblem solver, then Algorithm 3.1 will still converge at the same 356 rate as the exact algorithm (2.6), albeit in expectation. 357

Algorithm 3.1 Inexact Augmented Lagrangian Framework

```
Require: x_{-1} \in \mathbb{R}^n, \lambda_0 \in \mathbb{R}^m, \alpha > 0, inexactness conditions, solver \mathcal{S}.
```

```
1: for k = 0, 1, ... do
```

- 2: Set starting point $x_{k,0} \leftarrow x_{k-1}$
- 3: Find an approximate minimizer x_k of (3.1) using solver S, starting with $x_{k,0}$ such that some inexactness conditions are satisfied
- 4: Update $\lambda_{k+1} \leftarrow \lambda_k \alpha c(x_k)$
- 5: end for

358

360

361

362

363

364

365 366

367

368

369

371

372

373

374

375

376

377

378

379

380 381

382

383

Remark 3.1. We make the following remarks about Algorithm 3.1.

- Solver and inexactness conditions: For the sake of generality, we leave the description of the solver and inexactness conditions arbitrary and specify them in Subsection 3.1 and Subsection 3.2 respectively. We assume that the solver S can compute an approximate minimizer x_k that satisfies the inexactness conditions. The sequences of primal and dual iterates obtained in the algorithm are random due to the stochastic nature of the objective function f(x). Therefore, this assumption is reasonable when the inexactness conditions are also stochastic.
- Penalty parameter ($\alpha > 0$): The algorithm employs a constant penalty parameter $\alpha > 0$. In Section 4, we show that the algorithm converges for any choice of this parameter and does not depend on problem characteristics or other algorithmic parameters.
- Starting points $(x_{k,0})$: At each iteration, the algorithm uses the previous primal iterate as starting point in the solver S to solve (3.1). This is meant to reduce the computational effort to solve (3.1). Since the successive augmented Lagrangian functions differ only in the dual variable λ_k , the approximate minimizer of the previous subproblem is an intuitive estimate of the solution to the current problem. In Section 4, we quantify the efficiency of this starting point rule in terms of total computational work.

We now describe the unspecified components of this algorithm: the solver S and the tolerance conditions.

3.1. Adaptive Sampling Proximal Gradient Method. Projected or proximal stochastic gradient methods are a popular class of methods for solving (3.1) when the projection or proximal operators are easy to compute [66]. The iterate update of a projected stochastic gradient method is given as follows:

384 (3.2)
$$x_{k,t+1} = x_{k,t} + \eta R_{S_{k,t}}(x_{k,t}, \lambda_k; \alpha, \eta),$$

where $\eta > 0$ is the step size parameter, $k \in \mathbb{N}$ denotes the *outer* augmented Lagrangian iteration counter, $t \in \mathbb{N}$ denotes the *inner* projected stochastic gradient iteration

counter, $S_{k,t}$ is a set consisting of i.i.d. samples of ζ ,

388 (3.3)
$$R_{S_{k,t}}(x_{k,t},\lambda_k;\alpha,\eta) \stackrel{\text{def}}{=} \frac{\operatorname{proj}_{\mathcal{X}}(x_{k,t} - \eta \nabla_x \mathcal{L}_{S_{k,t}}(x_{k,t},\lambda_k;\alpha)) - x_{k,t}}{\eta},$$

389 (3.4)
$$\nabla_x \mathcal{L}_{S_{k,t}}(x_{k,t}, \lambda_k; \alpha) \stackrel{\text{def}}{=} \frac{1}{|S_{k,t}|} \sum_{\zeta_i \in S_{k,t}} \nabla_x \mathcal{L}(x_{k,t}, \lambda_k, \zeta_i; \alpha),$$

391 and

396

397

398

399

400

401 402

403

404

405 406

392
$$\nabla_x \mathcal{L}(x_{k,t}, \lambda_k, \zeta_i; \alpha) = \nabla_x F(x_{k,t}, \zeta_i) - \langle \lambda_k, \nabla c(x_{k,t}) \rangle + \alpha \langle c(x_{k,t}), \nabla c(x_{k,t}) \rangle.$$

In what follows, it is helpful to note that $R_{S_{k,t}}(x_{k,t}, \lambda_k; \alpha, \eta)$ denotes a stochastic approximation of the true projected (reduced) gradient

395 (3.5)
$$R(x_{k,t}, \lambda_k; \alpha, \eta) \stackrel{\text{def}}{=} \frac{\text{proj}_{\mathcal{X}}(x_{k,t} - \eta \nabla_x \mathcal{L}(x_{k,t}, \lambda_k; \alpha)) - x_{k,t}}{\eta}.$$

Two adaptive sampling strategies have recently been proposed for the projected stochastic gradient method [7,83]. Both strategies employ a mechanism for improving the quality of the stochastic gradient approximation by updating the sample size $|S_{k,t}|$ on the fly at each (subproblem) iteration t. In turn, they overcome a significant limitation of fixed sample size strategies without compromising efficiency, while also maintaining the fast convergence of their deterministic counterparts. Indeed, fixed sample size strategies can only guarantee convergence to a neighborhood of the solution or must compromise on the convergence rate.

Adaptive sampling strategies aim to ensure that the variance in the stochastic gradient is controlled by the squared norm of the projected gradient. In [83], this is written as follows:

407 (3.6)
$$\mathbb{E}_{k,t} \left[\| \nabla_x \mathcal{L}_{S_{k,t}}(x_{k,t}, \lambda_k) - \nabla_x \mathcal{L}(x_{k,t}, \lambda_k) \|^2 \right] \le \theta_g^2 \| \mathbb{E}_{k,t} [R_{S_{k,t}}(x_{k,t}, \lambda_k; \alpha, \eta)] \|^2$$
,

408 where $\theta_q > 0$ is a given parameter, and

409 (3.7)
$$\mathbb{E}_{k,t}[\cdot] \stackrel{\text{def}}{=} \mathbb{E}[\cdot | \lambda_k, x_{k,t}],$$

denotes the expectation conditioned on the past iterates until $\lambda_k, x_{k,t}$. Specifically, $\mathbb{E}_{k,t}$ is the conditional expectation conditioned on the filtration

$$\mathbb{T}_{k,t} = \sigma(\lambda_0, x_{-1,0}, S_{0,0}, \dots, S_{0,T_0}, \dots, S_{k-1,0}, \dots, S_{k-1,T_{k-1}}, S_{k,T_0}, \dots, S_{k,t}),$$

410 where T_i denotes the number of inner iterations performed at the outer iteration i.

Using the definition of the gradient of the augmented Lagrangian function in (3.6) results in the following equivalent condition:

413 (3.8)
$$\mathbb{E}_{k,t} \left[\|\nabla F_{S_{k,t}}(x_{k,t}) - \nabla f(x_{k,t})\|^2 \right] \le \theta_a^2 \|\mathbb{E}_{k,t} [R_{S_{k,t}}(x_{k,t}, \lambda_k; \alpha, \eta)]\|^2,$$

414 where

415 (3.9)
$$\nabla F_{S_{k,t}}(x_{k,t}) \stackrel{\text{def}}{=} \frac{1}{|S_{k,t}|} \sum_{\zeta_i \in S_{k,t}} \nabla F(x_{k,t}, \zeta_i).$$

Since the samples of ζ are i.i.d., Bienaymé's identity may be used to simplify the left-hand side of (3.8). This results in the following equivalent condition:

CONDITION 3.1 (Theoretical Sampling Condition). For any given $\theta_g > 0$, the variance in the stochastic gradient of the objective function f is controlled by the squared norm of the expected projected gradient $R_{S_{k,t}}$. That is,

421 (3.10)
$$\frac{\mathbb{E}_{\zeta}[\|\nabla F(x_{k,t},\zeta) - \nabla f(x_{k,t})\|^2]}{|S_{k,t}|} \le \theta_g^2 \|\mathbb{E}_{k,t}[R_{S_{k,t}}(x_{k,t},\lambda_k;\alpha,\eta)]\|^2.$$

This condition involves computing true variances and exact projected gradients that are unavailable in practice. Therefore, in Section 5, we also propose a practical version of this condition to control the sample sizes.

We conclude this subsection with remark and the following well-known result (adapted to our setting) [64, Corollary 2.3.2] that is used in the coming analysis.

PROPOSITION 3.1. Suppose Assumptions 2.1 and 2.2 hold. Then, for any $0 < \frac{1}{L+\alpha ||A||^2}$ and for all $x \in \mathcal{X}, \lambda \in \mathbb{R}^m, \alpha > 0$,

429 (3.11)
$$||R(x,\lambda;\alpha,\eta)||^2 \leq \frac{2}{\eta} \left(\mathcal{L}(x,\lambda,\alpha) - \mathcal{L}(x_{\mathcal{L}}^*,\lambda,\alpha) \right),$$

430 where $x_{\mathcal{L}}^* \in \arg\min_{x \in \mathcal{X}} \mathcal{L}(x, \lambda, \alpha)$. Moreover, if Assumption 2.3 also holds, then,

431 (3.12)
$$\frac{\mu}{2} \|x - x_{\mathcal{L}}^*\|^2 + \frac{\eta}{2} \|R(x, \lambda; \alpha, \eta)\|^2 \le \langle R(x, \lambda; \alpha, \eta), x - x_{\mathcal{L}}^* \rangle.$$

REMARK 3.2 (Alternative Sampling Condition). An alternative sampling condition is proposed in [7] that would replace the right-hand side of (3.6) by a constant factor times the squared norm of the projected gradient (3.5). Following the procedure above, we would then arrive at a somewhat simpler inequality taking the place of (3.10), namely,

437 (3.13)
$$\frac{\mathbb{E}_{\zeta}[\|\nabla F(x_{k,t},\zeta) - \nabla f(x_{k,t})\|^{2}]}{|S_{k,t}|} \leq \tilde{\theta}_{g}^{2} \|R(x_{k,t},\lambda_{k};\alpha,\eta)\|^{2},$$

for some $\tilde{\theta}_g^2 > 0$. It turns out that the two conditions (3.10) and (3.13) are equivalent in the sense that their right-hand sides bound each other from above and below:

441 (3.14)
$$\frac{\|R(x_{k,t},\lambda_k;\alpha,\eta)\|}{1+\theta_q} \le \|\mathbb{E}_{k,t}[R_{S_{k,t}}(x_{k,t},\lambda_k;\alpha,\eta)\| \le \frac{\|R(x_{k,t},\lambda_k;\alpha,\eta)\|}{1-\theta_q}.$$

442 Indeed, note that

432

433

434

435 436

450

451 452

453

454

455

456

$$\begin{aligned}
443 & \|\mathbb{E}_{k,t}[R_{S_{k,t}}(x_{k,t},\lambda_{k};\alpha,\eta) - R(x_{k,t},\lambda_{k};\alpha,\eta)]\|^{2} \\
444 & \leq \mathbb{E}_{k,t}[\|R_{S_{k,t}}(x_{k,t},\lambda_{k};\alpha,\eta) - R(x_{k,t},\lambda_{k};\alpha,\eta)\|^{2}] \\
445 & = \eta^{-2}\mathbb{E}_{k,t}\Big[\|\operatorname{proj}_{\mathcal{X}}(x_{k,t} - \eta\nabla_{x}\mathcal{L}_{S_{k,t}}(x_{k,t},\lambda_{k};\alpha)) \\
& - \operatorname{proj}_{\mathcal{X}}(x_{k,t} - \eta\nabla_{x}\mathcal{L}(x_{k,t},\lambda_{k};\alpha))\|^{2}\Big] \\
446 & \leq \mathbb{E}_{k,t}\left[\|\nabla F_{S_{k,t}}(x_{k,t}) - \nabla f(x_{k,t})\|^{2}\right] \\
\leq \theta_{g}^{2}\|\mathbb{E}_{k,t}[R_{S_{k,t}}(x_{k,t},\lambda_{k};\alpha,\eta)]\|^{2},
\end{aligned}$$

where the first line follows from Jensen's inequality, the proceeding equality is due to (3.3) and (3.5), the second inequality follows from the non-expansiveness property of projections [64], and the last inequality is due to (3.8). Rearranging terms and using the reverse triangle inequality, $||a|| - ||b|| \le ||a - b||$, for all $a, b \in \mathbb{R}^n$, we arrive at (3.14). Moreover, both conditions (3.10) and (3.13) lead to identical practical algorithms; cf. Section 5. We choose to work with (3.10) instead of (3.13) because it leads to a simpler presentation of the complexity theory in Section 4.

3.2. Inexactness Conditions. The efficiency of an inexact augmented Lagrangian framework depends on its inexactness conditions. These conditions must balance the accuracy of the solution computed at each (outer) iteration and the overall computational efficiency. Due to the stochastic nature of the iterates obtained by our choice of the subproblem solver (cf. Subsection 3.1), these inexactness conditions must also be stochastic. We now propose inexactness conditions that meet these requirements based on the Moreau envelope perspective developed in Subsection 2.3.

Recall from Subsection 2.3 that the exact augmented Lagrangian method can be interpreted as a gradient descent method applied to the Moreau envelope of the (negative) dual function. Therefore, the inexact augmented Lagrangian method leads to inexact dual variable updates. That is, from the dual update (line 3 in Algorithm 3.1), (2.6b), and (2.23), we have,

469
$$\lambda_{k+1} = \lambda_k - \alpha c(x_k)$$

$$= \lambda_k - \alpha c(x_k^*) + \alpha c(x_k^*) - \alpha c(x_k)$$

$$= \lambda_k - \alpha \nabla q_\alpha(\lambda_k) + \alpha \epsilon_k,$$

$$= \lambda_k - \alpha \nabla q_\alpha(\lambda_k) + \alpha \epsilon_k,$$

where x_k^* is an exact minimizer of (3.1) and $\epsilon_k \stackrel{\text{def}}{=} c(x_k^*) - c(x_k)$. We note that we have not imposed any structure on the subproblems (3.1) that would give the update error ϵ_k zero mean; i.e., $\mathbb{E}_k[\epsilon_k] \neq 0$, where

476 (3.16)
$$\mathbb{E}_k[\cdot] = \mathbb{E}[\cdot|\lambda_k]$$

is the expected value operator conditioned on the iterates up until λ_k . Specifically, \mathbb{E}_k is the conditional expectation conditioned on the filtration

$$\mathbb{T}_k = \sigma(\lambda_0, x_{-1,0}, S_{0,0}, \dots, S_{0,T_0}, \dots, S_{k-1,0}, \dots, S_{k-1,T_{k-1}}),$$

where T_i denotes the number of inner iterations performed at the outer iteration i. In turn, we choose to view the additive update rule (3.15) as a *biased* stochastic gradient estimator update.

It is natural to consider an additional sampling condition when the sampling error can control the bias; cf. [7, Condition 2]. Such additional conditions are also common in trust-region methods [1, 5, 15]. In the present setting, however, the error is due to the subproblem solver. To address this, we aim to design a tolerance condition for terminating the inner loop. The following condition allows us to control the inexactness of the Moreau envelope gradient estimates in the dual update:

CONDITION 3.2 (Tolerance Condition I). For any given $\theta_e \in [0,1)$ and $\tau_k \geq 0$ with $\lim_{k\to\infty} \tau_k = 0$,

$$\mathbb{E}_k \left[\| c(x_k^*) - c(x_k) \|^2 \right] \le \theta_e^2 \| c(x_k^*) \|^2 + \tau_k,$$

490 where x_k^* is a minimizer of (3.1).

From the Moreau envelope perspective (3.15), this condition ensures that the expected squared norm of the error, $\mathbb{E}_k[\|\epsilon_k\|^2]$, is controlled by the squared norm of the gradient of $q_{\alpha}(\lambda_k)$ and a vanishing positive constant τ_k . That is, from (2.23), (3.15), and (3.17),

494 it follows that

458

$$\mathbb{E}_{k}[\|\epsilon_{k}\|^{2}] = \mathbb{E}_{k}\left[\|c(x_{k}^{*}) - c(x_{k})\|^{2}\right] \leq \theta_{e}^{2}\|c(x_{k}^{*})\|^{2} + \tau_{k} = \theta_{e}^{2}\|\nabla q_{\alpha}(\lambda_{k})\|^{2} + \tau_{k}.$$

Likewise, this condition ensures that inexact gradient information can be employed far away from the solution (i.e., when the gradient's norm is large). Meanwhile, it also ensures accurate gradient information near the solution (i.e., when the gradient norm is small). Although this condition is derived from controlling the inexactness in the dual update, it directly relates to inexactness in the minimization of (3.1) (cf. (2.22)). Therefore, we can replace Condition 3.2 by the following alternative condition:

Condition 3.3 (Tolerance Condition II). For any given $\theta_e \in [0,1)$ and $\tau_k \geq 0$ with $\lim_{k\to\infty} \tau_k = 0$,

$$\mathbb{E}_{k} \left[\mathcal{L}(x_{k}, \lambda_{k}; \alpha) - \mathcal{L}(x_{k}^{*}, \lambda_{k}; \alpha) \right] \leq \frac{\alpha \theta_{e}^{2} \|c(x_{k}^{*})\|^{2}}{2} + \frac{\alpha \tau_{k}}{2},$$

- 507 where x_k^* is a minimizer of (3.1).
- Condition 3.3 controls the error in the minimization of (3.1) and directly implies Condition 3.2. Indeed, using (2.22) and (3.18), we find

510 (3.19)
$$\mathbb{E}_k \left[\| c(x_k^*) - c(x_k) \|^2 \right] \leq \frac{2}{\alpha} \mathbb{E}_k \left[\mathcal{L}(x_k, \lambda_k; \alpha) - \mathcal{L}(x_k^*, \lambda_k; \alpha) \right] \leq \theta_e^2 \| c(x_k^*) \|^2 + \tau_k.$$

When the augmented Lagrangian functions are strongly convex, we can also control the norm of the projected gradient (3.5). That is, by Assumption 2.3, (2.8), and (3.12), we have

514 (3.20)
$$||c(x_k^*) - c(x_k)||^2 \le ||A||^2 ||x_k^* - x_k||^2 \le \frac{4||A||^2}{\mu^2} ||R(x_k, \lambda_k; \alpha, \eta)||^2.$$

- Therefore, we can impose the following alternate condition when f(x) is strongly convex.
- CONDITION 3.4 (Tolerance Condition III). For any given $\tilde{\theta}_e \in [0,1)$ and $\tilde{\tau}_k \geq 0$ with $\lim_{k\to\infty} \tilde{\tau}_k = 0$,

$$\mathbb{E}_{k} \left[\| R(x_{k}, \lambda_{k}; \alpha, \eta) \|^{2} \right] \leq \tilde{\theta}_{e}^{2} \| c(x_{k}^{*}) \|^{2} + \tilde{\tau}_{k},$$

521 where x_k^* is a minimizer of (3.1).

525

528 529

530

531

532

- 522 Condition 3.4 also controls the error in the subproblem (3.1) and implies Condition 3.2.
- Indeed, set $\tilde{\theta}_e \leq \frac{\mu \theta_e}{2\|A\|}$ and $\tilde{\tau}_k \leq \frac{\mu^2 \tau_k}{4\|A\|^2}$. Then, using (3.20) and (3.21), it holds that

524 (3.22)
$$\mathbb{E}_k \left[\|c(x_k^*) - c(x_k)\|^2 \right] \le \frac{4\|A\|^2}{\mu^2} \mathbb{E}_k \left[\|R(x_k, \lambda_k; \alpha, \eta)\|^2 \right] \le \theta_e^2 \|c(x_k^*)\|^2 + \tau_k.$$

REMARK 3.3. We observe that conditions similar to Conditions 3.2 through 3.4 have been proposed in the literature (cf. [50, 56, 72, 86]). The primary advantage of employing our conditions lies in their adaptive control over the subproblem error. Although verifying Conditions 3.2 through 3.4 for a stochastic subproblem solver can be challenging because they each require evaluating deterministic quantities, these conditions can still help us gain insight into the errors permitted in the algorithm while retaining desirable convergence properties. Furthermore, these conditions can guide the development of practical algorithms.

4. Theory. We now establish theoretical convergence guarantees and total sample complexity results for the proposed inexact augmented Lagrangian algorithmic framework when the inexactness conditions proposed in Subsection 3.2 are satisfied. We use the following notation for the full expectation:

538 (4.1)
$$\mathbb{E}[\cdot] = \mathbb{E}_0[\mathbb{E}_1[\cdots \mathbb{E}_k[\cdot]]].$$

- 4.1. Convergence Results. We start by establishing a technical lemma.
- LEMMA 4.1. Suppose Assumptions 2.1 and 2.2 hold. For any x_{-1} , λ_0 and $\alpha > 0$, let $\{x_k, \lambda_k\}$ be the sequence of primal-dual iterates generated by Algorithm 3.1. Then, for all $k \in \mathbb{N}$,

$$q_{\alpha}(\lambda_{k+1}) \le q_{\alpha}(\lambda_k) - \frac{\alpha}{2} \|\nabla q_{\alpha}(\lambda_k)\|^2 + \frac{\alpha}{2} \|\epsilon_k\|^2,$$

- 545 where $\epsilon_k = c(x_k^*) c(x_k)$ and x_k^* is a minimizer of (3.1).
- 546 Proof. From the dual update rule (line 3 in Algorithm 3.1), (2.6b), and (2.23), 547 it follows that

$$\lambda_{k+1} = \lambda_k - \alpha \nabla q_\alpha(\lambda_k) + \alpha \epsilon_k.$$

Using the Lipschitz continuity of $\nabla q_{\alpha}(\lambda)$ with Lipschitz constant $L_{\alpha} = \alpha^{-1}$ (cf. (2.15)) and the descent lemma [12], we have,

552
$$q_{\alpha}(\lambda_{k+1}) \leq q_{\alpha}(\lambda_{k}) - \alpha \langle \nabla q_{\alpha}(\lambda_{k}) - \epsilon_{k}, \nabla q_{\alpha}(\lambda_{k}) \rangle + \frac{\alpha^{2} L_{\alpha}}{2} \|\nabla q_{\alpha}(\lambda_{k}) - \epsilon_{k}\|^{2}$$
553
$$= q_{\alpha}(\lambda_{k}) - \frac{\alpha}{2} \|\nabla q_{\alpha}(\lambda_{k})\|^{2} + \frac{\alpha}{2} \|\epsilon_{k}\|^{2},$$

555 as necessary.

539

We are now ready to establish convergence results for the inexactness conditions developed in Subsection 3.2.

- THEOREM 4.2. Suppose Assumptions 2.1, 2.2 and 2.4 hold. For any x_{-1}, λ_0 and $\alpha > 0$, let $\{(x_k, \lambda_k)\}$ be the sequence of primal-dual iterates generated by Alson gorithm 3.1. Furthermore, let $\theta_e \in [0,1)$ and $\tau_k \geq 0$ such that $\tau_0^{-1} \sum_{k=0}^{\infty} \tau_k = a_{\infty} < \infty$. If any of the following three statements hold at each iteration $k \in N$:
 - (a) the primal iterates x_k satisfy Condition 3.2;
 - (b) the primal iterates x_k satisfy Condition 3.3; or
 - (c) Assumption 2.3 also holds and the primal iterates x_k satisfy Condition 3.4 with $\tilde{\theta}_e \leq \frac{\mu \theta_e}{2||A||}$ and $\tilde{\tau}_k \leq \frac{\mu^2 \tau_k}{4||A||^2}$;
- 566 then

562

563

564

565

567 568

$$\lim_{k \to \infty} \mathbb{E}[\|c(x_k)\|^2] = 0.$$

569 Moreover, for any $K \in \mathbb{N}$, we have that,

$$\min_{570} \quad (4.3) \qquad \min_{0 \le k \le K - 1} \mathbb{E}[\|c(x_k)\|^2] \le \frac{4(1 + \theta_e^2)}{\alpha (1 - \theta_e^2) K} [q_\alpha(\lambda_0) - q^*] + \frac{4}{1 - \theta_e^2} \frac{\tau_0 a_\infty}{K},$$

where $q^* > -\infty$ is defined in (2.17). In addition, if either (b) or (c) is satisfied, then

$$\lim_{k \to \infty} \mathbb{E} \left[\left\| \frac{\operatorname{proj}_{\mathcal{X}}(x_k - \eta \nabla \ell_x(x_k, \lambda_{k+1})) - x_k}{\eta} \right\|^2 \right] = 0,$$

574 for every $0 < \eta < \frac{1}{L + \alpha ||A||^2}$.

Proof. If (a), (b), or (c) holds, then (3.17) holds as well due to (3.19) and (3.22). By Lemma 4.1, we have,

$$q_{\alpha}(\lambda_{k+1}) \le q_{\alpha}(\lambda_k) - \frac{\alpha}{2} \|\nabla q_{\alpha}(\lambda_k)\|^2 + \frac{\alpha}{2} \|\epsilon_k\|^2.$$

Taking the conditional expectation (3.16) of both sides and invoking (2.23) and (3.17),

580 we arrive at

$$\mathbb{E}_{k}\left[q_{\alpha}(\lambda_{k+1})\right] \leq q_{\alpha}(\lambda_{k}) - \frac{\alpha}{2} \|\nabla q_{\alpha}(\lambda_{k})\|^{2} + \frac{\alpha}{2} \mathbb{E}_{k}\left[\|\epsilon_{k}\|^{2}\right]$$

$$\leq q_{\alpha}(\lambda_k) - \frac{\alpha(1 - \theta_e^2)}{2} ||c(x_k^*)||^2 + \frac{\alpha}{2} \tau_k.$$

584 Rearranging terms, we find

$$||c(x_k^*)||^2 \le \frac{2}{\alpha(1-\theta_e^2)} (q_\alpha(\lambda_k) - \mathbb{E}_k[q_\alpha(\lambda_{k+1})]) + \frac{1}{1-\theta_e^2} \tau_k.$$

587 Therefore,

588

589

590 591

596

597 598

602 603

$$\mathbb{E}_{k} \left[\| c(x_{k}) \|^{2} \right] \leq 2\mathbb{E}_{k} \left[\| c(x_{k}) - c(x_{k}^{*}) \|^{2} \right] + 2\| c(x_{k}^{*}) \|^{2}
\leq 2(1 + \theta_{e}^{2}) \| c(x_{k}^{*}) \|^{2} + 2\tau_{k}
\leq \frac{4(1 + \theta_{e}^{2})}{\alpha(1 - \theta_{e}^{2})} (q_{\alpha}(\lambda_{k}) - \mathbb{E}_{k}[q_{\alpha}(\lambda_{k+1})]) + \frac{4}{1 - \theta_{e}^{2}} \tau_{k},$$

where the first inequality is due to $||a+b||^2 \le 2||a||^2 + 2||b||^2$ for any $a, b \in \mathbb{R}^n$, the second inequality is due to (3.17) and the last inequality follows from (4.6). Taking the full expectation (4.1), and summing the above inequality from k=0 to K-1, delivers

$$\sum_{k=0}^{K-1} \mathbb{E}\left[\|c(x_k)\|^2\right] \le \frac{4(1+\theta_e^2)}{\alpha(1-\theta_e^2)} \mathbb{E}[q_\alpha(\lambda_0) - q_\alpha(\lambda_K)] + \frac{4}{1-\theta_e^2} \sum_{k=0}^{K-1} \tau_k$$

$$\le \frac{4(1+\theta_e^2)}{\alpha(1-\theta_e^2)} [q_\alpha(\lambda_0) - q^*] + \frac{4}{1-\theta_e^2} \tau_0 a_\infty,$$

where the second inequality follows from (2.16) and the assumption $\sum_{k=0}^{\infty} \tau_k = \tau_0 a_{\infty} < \infty$. Observe that $q^* > -\infty$ due to (2.17), which follows from Assumption 2.2. Therefore,

$$\sum_{k=0}^{K-1} \mathbb{E}\left[\|c(x_k)\|^2 \right] < \infty,$$

604 which implies that

$$\lim_{k \to \infty} \mathbb{E}\left[\|c(x_k)\|^2 \right] = 0.$$

607 Moreover,

608
$$\min_{0 \le k \le K-1} \mathbb{E}[\|c(x_k)\|^2] \le \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|c(x_k)\|^2]$$

$$\le \frac{4(1+\theta_e^2)}{\alpha(1-\theta_e^2)K} [q_\alpha(\lambda_0) - q^*] + \frac{4}{1-\theta_e^2} \frac{\tau_0 a_\infty}{K}.$$

We will now analyze the stationarity error (4.4). Using (2.1), $\lambda_{k+1} = \lambda_k - \alpha c(x_k)$, and (2.7), it follows that,

613
$$x_k - \eta \nabla \ell_x(x_k, \lambda_{k+1}) = x_k - \eta (\nabla f(x_k) - \langle \lambda_{k+1}, \nabla c(x_k) \rangle)$$
614
$$= x_k - \eta (\nabla f(x_k) - \langle \lambda_k - \alpha c(x_k), \nabla c(x_k) \rangle)$$
615
$$= x_k - \eta \nabla_x \mathcal{L}(x_k, \lambda_k; \alpha).$$

617 Therefore,

618 (4.7)
$$\left\| \frac{\operatorname{proj}_{\mathcal{X}}(x_k - \eta \nabla \ell_x(x_k, \lambda_{k+1})) - x_k}{\eta} \right\|^2 = \left\| R(x_k, \lambda_k; \alpha, \eta) \right\|^2.$$

620 If statement (b) holds, then it follows from (3.11) that

621
$$\mathbb{E}\left[\left\|R(x_{k},\lambda_{k};\alpha,\eta)\right\|^{2}\right] \leq \frac{2}{\eta}\left(\mathbb{E}\left[\mathcal{L}(x_{k},\lambda_{k};\alpha) - \mathcal{L}(x_{k}^{*},\lambda_{k};\alpha)\right]\right)$$
622
623
$$\leq \frac{\alpha\theta_{e}^{2}}{\eta}\mathbb{E}\left[\left\|c(x_{k}^{*})\right\|^{2}\right] + \frac{\alpha}{\eta}\tau_{k}.$$

624 If statement (c) holds, then

$$\mathbb{E}\left[\left\|R(x_k, \lambda_k; \alpha, \eta)\right\|^2\right] \le \tilde{\theta}_e^2 \mathbb{E}\left[\left\|c(x_k^*)\right\|^2\right] + \tilde{\tau}_k.$$

627 In turn, if either statements (b) or (c) holds, it follows that

628 (4.8)
$$\mathbb{E}\left[\left\|R(x_k, \lambda_k; \alpha, \eta)\right\|^2\right] \le \max\left\{\frac{\alpha \theta_e^2}{\eta}, \tilde{\theta}_e^2\right\} \mathbb{E}\left[\left\|c(x_k^*)\right\|^2\right] + \max\left\{\frac{\alpha}{\eta} \tau_k, \tilde{\tau}_k\right\}.$$

Taking the full expectation and summing the inequality (4.6) from k = 0 to K - 1, and observing that $q^* > -\infty$, we arrive at

632 (4.9)
$$\sum_{k=0}^{K-1} \mathbb{E}\left[\|c(x_k^*)\|^2\right] \le \frac{2}{\alpha(1-\theta_e^2)} [q_\alpha(\lambda_0) - q^*] + \frac{1}{1-\theta_e^2} \tau_0 a_\infty < \infty,$$

634 which implies that

$$\lim_{k \to \infty} \mathbb{E}\left[\|c(x_k^*)\|^2 \right] = 0.$$

Taking limits on both sides of (4.8) and using (4.10) completes the proof.

Theorem 4.2 establishes that the *expected* feasibility error vanishes as $k \to 0$, and meanwhile, the *smallest* feasibility error converges to zero at a sublinear rate. Moreover, the stationarity error also converges to zero *in expectation* when either Condition 3.3 or Condition 3.4 holds. However, the theorem does not guarantee any rate of convergence of the stationarity error. To establish such a result, we can perform one additional update at the iterate at which the expected feasibility error attains a

one additional update at the iterate at which the expected feasibility error attains a minimum.

COROLLARY 4.3. Suppose Assumptions 2.1, 2.2 and 2.4 hold. Let k_* be the iteration number at which $\min_{0 \le k \le K-1} \mathbb{E}[\|c(x_k)\|^2]$ is attained. That is,

$$\mathbb{E}[\|c(x_{k_*})\|^2] = \min_{0 \le k \le K - 1} \mathbb{E}[\|c(x_k)\|^2].$$

For any given $\tilde{\alpha} > 0$, $h \ge 0$, and $0 < \tilde{\eta} < \frac{1}{L + \tilde{\alpha} ||A||^2}$, let \tilde{x} be an approximate minimizer of $\min_{x \in \mathcal{X}} \mathcal{L}(x, \lambda_{k_*}, \tilde{\alpha})$ obtained with the starting point x_{k_*} that satisfies either of the following two statements:

(i)
$$\tilde{x}$$
 satisfies

$$\mathbb{E}_{k_*} \left[\mathcal{L}(\tilde{x}, \lambda_{k_*}; \tilde{\alpha}) - \mathcal{L}(x_{k_*}^*, \lambda_{k_*}; \tilde{\alpha}) \right] \le \frac{\tilde{\alpha}h}{2K}; \text{ or }$$

(ii) Assumption 2.3 also holds and \tilde{x} satisfies

$$\mathbb{E}_{k_*} \left[\| R(\tilde{x}, \lambda_{k_*}; \tilde{\alpha}, \tilde{\eta}) \|^2 \right] \le \frac{\mu^2 h}{4 \| A \|^2 K}.$$

654 Then, for $\tilde{\lambda} = \lambda_{k_*} - \tilde{\alpha}c(\tilde{x})$, we have

$$\mathbb{E}[\|c(\tilde{x})\|^2] \le \frac{2h}{K} + \frac{8(1+\theta_e^2)}{\alpha(1-\theta_e^2)K} [q_\alpha(\lambda_0) - q^*] + \frac{8}{1-\theta_e^2} \frac{\tau_0 a_\infty}{K}$$

657 and

652

653

658 (4.11b)
$$\mathbb{E}\left[\left\|\frac{\operatorname{proj}_{\mathcal{X}}(\tilde{x}-\tilde{\eta}\nabla\ell_{x}(\tilde{x},\tilde{\lambda}))-\tilde{x}}{\tilde{\eta}}\right\|^{2}\right] \leq \frac{b}{K},$$

where $b = \frac{\tilde{\alpha}h}{\tilde{\eta}}$ if (i) holds and $b = \frac{\mu^2 h}{4\|A\|^2}$ if (ii) holds.

Proof. If either (i) or (ii) is satisfied, then using $\tilde{\alpha}$ as the penalty parameter in (3.19) and (3.22), it follows that

$$\mathbb{E}_{k_*}[\|c(\tilde{x}) - c(x_{k_*})\|^2] \le \frac{h}{K}.$$

665 Therefore, taking the full expectation,

666
$$\mathbb{E}[\|c(\tilde{x})\|^{2}] = \mathbb{E}[\|c(\tilde{x}) - c(x_{k_{*}}) + c(x_{k_{*}})\|^{2}]$$
667
$$\leq 2\mathbb{E}[\|c(\tilde{x}) - c(x_{k_{*}})\|^{2}] + 2\mathbb{E}[\|c(x_{k_{*}})\|^{2}]$$
668
$$\leq \frac{2h}{K} + \frac{8(1 + \theta_{e}^{2})}{\alpha(1 - \theta_{e}^{2})K}[q_{\alpha}(\lambda_{0}) - q^{*}] + \frac{8}{1 - \theta_{e}^{2}} \frac{\tau_{0}a_{\infty}}{K},$$

where the first inequality is due to $||a+b||^2 \le 2||a||^2 + 2||b||^2$ for any $a, b \in \mathbb{R}^n$, and the last inequality follows from (4.3). Now, consider the stationarity error. Similar

to (4.7), we can show that

673 (4.12)
$$\mathbb{E}\left[\left\|\frac{\operatorname{proj}_{\mathcal{X}}(\tilde{x}-\tilde{\eta}\nabla\ell_{x}(\tilde{x},\tilde{\lambda}))-\tilde{x}}{\tilde{\eta}}\right\|^{2}\right] = \mathbb{E}\left[\left\|R(\tilde{x},\lambda_{k_{*}};\tilde{\alpha},\tilde{\eta})\right\|^{2}\right].$$

Therefore, if (i) holds, it follows from (3.11) with $\tilde{\alpha}$ as the penalty parameter that

$$\mathbb{E}\left[\|R(\tilde{x},\lambda_{k_*};\tilde{\alpha},\tilde{\eta})\|^2\right] \leq \frac{2}{\tilde{\eta}}\left(\mathbb{E}\left[\mathcal{L}(\tilde{x},\lambda_{k_*};\tilde{\alpha}) - \mathcal{L}(x_{k_*}^*,\lambda_{k_*};\tilde{\alpha}]\right) \leq \frac{\tilde{\alpha}h}{\tilde{\eta}K}\right)$$

Likewise, if (ii) holds, then

$$\mathbb{E}\left[\|R(\tilde{x}, \lambda_{k_*}; \tilde{\alpha}, \tilde{\eta})\|^2\right] \le \frac{\mu^2 h}{4\|A\|^2 K}.$$

Substituting these inequalities into (4.12) completes the proof.

4.2. Sample Complexity. We now establish the sample complexity for our inexact augmented Lagrangian algorithm, i.e., we estimate the worst-case expected total number of stochastic gradient evaluations to reach an ϵ -accurate solution. To define accuracy, we specifically consider the following metric:

683 (4.13)
$$\max \left\{ \mathbb{E}[\|c(\tilde{x})\|^2], \mathbb{E}\left[\left\|\frac{\operatorname{proj}_{\mathcal{X}}(\tilde{x} - \tilde{\eta}\nabla \ell_x(\tilde{x}, \tilde{\lambda})) - \tilde{x}}{\tilde{\eta}}\right\|^2\right] \right\} \leq \epsilon,$$

for some $\epsilon \in (0,1)$. For the sake of brevity in this analysis, we employ Condition 3.3 as the inexactness condition with $\theta_e = 0$. At any outer iteration k, x_{k-1} is used as the starting point in the adaptive sampling proximal gradient method to solve the inner subproblem (3.1) until Condition 3.3 is satisfied. Recall that we define the index for the inner iterations as t, and the iterates in the inner loop as $x_{k,t}$. Since x_{k-1} is used as the starting iterate, we set $x_{k,0} \stackrel{\text{def}}{=} x_{k-1}$.

The adaptive sampling projected gradient method used to solve the inner subproblems (see Subsection 3.1) converges at a sublinear rate (cf. [7, Theorem 2.11], [83, Theorem 3.7]). The following theorem reformulates this result for augmented Lagrangian subproblems (3.1).

Theorem 4.4. Suppose Assumptions 2.1 and 2.3 hold. If $\eta = \frac{(1-2\theta_g^2)}{L+\alpha||A||^2}$ with 695 $\theta_g \in [0, \frac{1}{\sqrt{2}})$ and Condition 3.1 is satisfied, then for any outer iteration $k \in \mathbb{N}$ and 696 inner iteration $t \in \mathbb{N}_+$, it holds that

698 (4.14)
$$\mathbb{E}_{k}[\mathcal{L}(x_{k,t},\lambda_{k};\alpha) - \mathcal{L}(x_{k}^{*},\lambda_{k};\alpha)] \leq \frac{(L+\alpha||A||^{2}) \min_{x_{k}^{*} \in \mathcal{X}_{k}^{*}} ||x_{k-1} - x_{k}^{*}||^{2}}{2(1-2\theta_{g}^{2})t},$$

where $\mathcal{X}_k^* = \arg\min_{x \in \mathcal{X}} \mathcal{L}(x, \lambda_k; \alpha)$. 700

679 680

681

685

687

689

690

691 692

693

694

701

703 704

705

706

707

708

709

710

711

Recall that \mathbb{E}_k denotes expectation conditioned on the filtration \mathbb{T}_k , and note that the initial distance to optimality is in this filtration, i.e., $\min_{x_k^* \in \mathcal{X}_k^*} ||x_{k-1} - x_k^*||^2 \in \mathbb{T}_k$. 702

Adaptive sampling methods are more efficient and robust in practice than methods that increase the sample sizes at predetermined rates. However, their sample complexity analysis has proven to be difficult, and establishing an upper bound on the sample sizes at each iteration poses significant challenges. Therefore, we make the following assumption based on the sample size growth rate over inner iterations t.

Assumption 4.1. At any given outer iteration $k \in \mathbb{N}$, the expected sample size required to satisfy Condition 3.1 increases at a polynomial rate over the inner iterations t. More specifically, there exists $c_0 \ge 0$ and $\delta_0 > 0$ arbitrarily close to zero, such that

712 (4.15)
$$\mathbb{E}_{k}[|S_{k,t}|] = \frac{c_0 \omega}{\theta_{q}^{2} \min_{x_{k}^{*} \in \mathcal{X}_{k}^{*}} ||x_{k-1} - x_{k}^{*}||^{2}} (t+1)^{1+\delta_{0}} \quad \forall k, t \in \mathbb{N},$$

where $\mathcal{X}_k^* = \arg\min_{x \in \mathcal{X}} \mathcal{L}(x, \lambda_k; \alpha)$ and ω is defined in (2.11). 713

Predetermined sample growth rates similar to Assumption 4.1 are employed in un-714 constrained and constrained stochastic optimization settings [9,67]. We acknowledge 715 that Assumption 4.1 pertains to algorithmic quantities and is, therefore, less than 716 ideal. Nevertheless, while we cannot rigorously prove this statement, we provide the 717 following set of supporting (heuristic) arguments.

Consider rewriting Condition 3.1 in the following way: 719

720 (4.16)
$$b_{k,t} \stackrel{\text{def}}{=} \frac{\mathbb{E}_{\zeta}[\|\nabla f(x_{k,t},\zeta) - \nabla F(x_{k,t})\|^2]}{\theta_q^2 \|\mathbb{E}_{k,t}[R_{S_{k,t}}(x_{k,t},\lambda_k;\alpha,\eta)]\|^2} \le |S_{k,t}|.$$

This inequality is tight, i.e., $|S_{k,t}| = b_{k,t}$ when Condition 3.1 is satisfied with equality. 722

In this case, due to (2.11), it follows that 723

724 (4.17)
$$b_{k,t} \le \frac{\omega}{\theta_g^2 \|\mathbb{E}_{k,t}[R_{S_{k,t}}(x_{k,t}, \lambda_k; \alpha, \eta)]\|^2}.$$

On the other hand, using (4.14) and taking the expected value of both sides of (3.11)726 727

728 (4.18)
$$\mathbb{E}_{k}[\|R(x_{k,t},\lambda_{k};\alpha,\eta)\|^{2}] \leq \frac{(L+\alpha\|A\|^{2})\min_{x_{k}^{*}\in\mathcal{X}_{k}^{*}}\|x_{k-1}-x_{k}^{*}\|^{2}}{\eta(1-2\theta_{g}^{2})t}, \quad \forall t \in \mathbb{N}_{+}.$$

This inequality implies that the expected squared norm of the reduced gradient goes 730 to zero at a sublinear rate. 731

Now, recall (3.14). In particular,

732

739

740

741

742 743

744

733
$$\|\mathbb{E}_{k,t}[R_{S_{k,t}}(x_{k,t},\lambda_k;\alpha,\eta)]\|^2 \le \frac{1}{(1-\theta_g)^2} \|R(x_{k,t},\lambda_k;\alpha,\eta)\|^2.$$

Taking the conditional expectation \mathbb{E}_k of both sides and invoking (4.18), it follows 735 736

737
$$\mathbb{E}_{k}[\|\mathbb{E}_{k,t}[R_{S_{k,t}}(x_{k,t},\lambda_{k};\alpha,\eta)]\|^{2}] \leq \frac{(L+\alpha\|A\|^{2})\min_{x_{k}^{*}\in\mathcal{X}_{k}^{*}}\|x_{k-1}-x_{k}^{*}\|^{2}}{(1-\theta_{g})^{2}\eta(1-2\theta_{g}^{2})t}, \quad \forall t \in \mathbb{N}_{+}.$$

This inequality implies that the expected squared norm of the stochastic reduced gradient goes to zero at a sublinear rate. Therefore, it is possible to replace the righthand-side of Condition 3.1 with a sublinearly convergent sequence $((t+1)^{-(1+\delta_0)})$ for any $t \in \mathbb{N}$ and achieve a similar sublinear convergence result as in Theorem 4.4. In such a scenario, the sample sizes satisfy Assumption 4.1. For the sake of brevity, in the rest of our analysis, we assume $\delta_0 = 0$ since δ_0 is arbitrarily close to zero.

We now state an equation that is useful to bound finite sum expressions in the 745 complexity analysis. For any $\delta > 0$ and $K \in \mathbb{N}_+$, we have 746

747 (4.19)
$$\sum_{k=0}^{K} k^{1+\delta} < \int_{t=0}^{K+1} t^{1+\delta} dt = \frac{(K+1)^{2+\delta}}{2+\delta}.$$

We are now ready to prove the main theorem about outer iteration and sample com-748 plexity. 749

Theorem 4.5. Suppose Assumptions 2.1, 2.2 and 4.1 hold with $\delta_0 = 0$. For any 750 x_{-1}, λ_0 and $\alpha > 0$, let $\{(x_k, \lambda_k)\}$ be the sequence of primal-dual iterates generated by 751 Algorithm 3.1 where x_k satisfies Condition 3.3 at each outer iteration k with $\theta_e = 0$, $\tau_k = \frac{\tau_0}{(k+1)^{1+\delta/2}}$, $a_\infty = \sum_{k=0}^\infty \frac{1}{(k+1)^{1+\delta/2}}$, with $\tau_0 \ge 0$, and $\delta > 0$. Suppose the sample 752 753

754

sizes $|S_{k,t}|$ satisfy Condition 3.1 with $\theta_g \in [0, \frac{1}{\sqrt{2}})$, $\eta = \frac{(1-2\theta_g^2)}{L+\alpha \|A\|^2}$. Under the conditions of Corollary 4.3 with \tilde{x} satisfying (i) with $h \geq 0$, $\tilde{\eta} = \frac{(1-2\theta_g^2)}{L+\tilde{\alpha}\|A\|^2}$ and $\tilde{\alpha} > 0$, the number 755

of outer iterations to get an ϵ -accurate solution $(\tilde{x}, \tilde{\lambda})$ satisfying (4.13) is 756

757 (4.20)
$$K_{\epsilon} = \left\lceil \frac{1}{\epsilon} \max \left\{ \frac{8}{\alpha} [q_{\alpha}(\lambda_0) - q^*] + 2(4\tau_0 a_{\infty} + h), \frac{\tilde{\alpha}h}{\tilde{\eta}} \right\} \right\rceil.$$

Moreover, if the gradients of the Lagrangian function are bounded, i.e.,

759 $\|\nabla \mathcal{L}(x_{k-1}, \lambda_k; \alpha)\|^2 \leq D_{\mathcal{L}}$ for all $k \in \mathbb{N}$, then the expected number of stochastic

760 gradient evaluations is

761 (4.21)
$$\mathbb{E}[\mathcal{W}] \le \frac{B(K_{\epsilon} + 1)^{3+\delta}}{\alpha^2 \tau_0^2 (3+\delta)} + \frac{\tilde{B}K_{\epsilon}^2}{\tilde{\alpha}^2 h^2},$$

762 where

763
$$B = \frac{2c_0\omega}{\theta_q^2} \left(\frac{(L+\alpha||A||^2)^2 D^2}{(1-2\theta_q^2)^2} + 4D_{\mathcal{L}} \right),$$

764 and

765
$$\tilde{B} = \frac{2c_0\omega}{\theta_g^2} \left(\frac{(L + \tilde{\alpha} ||A||^2)^2 D^2}{(1 - 2\theta_g^2)^2} + 4D_{\mathcal{L}} \right).$$

766 Proof. Substituting $\theta_e = 0$ into (4.11), we obtain

$$\mathbb{E}[\|c(\tilde{x})\|^2] \le \frac{8}{\alpha K} [q_{\alpha}(\lambda_0) - q^*] + \frac{2(4\tau_0 a_{\infty} + h)}{K},$$

769 and

773

774

775

776

777

$$\mathbb{E}\left[\left\|\frac{\operatorname{proj}_{\mathcal{X}}(\tilde{x}-\tilde{\eta}\nabla\ell_{x}(\tilde{x},\tilde{\lambda}))-\tilde{x}}{\tilde{\eta}}\right\|^{2}\right] \leq \frac{\tilde{\alpha}h}{\tilde{\eta}K}.$$

772 Therefore, for any $K \geq K_{\epsilon}$ defined in (4.20), $(\tilde{x}, \tilde{\lambda})$ satisfies (4.13).

Let T_k be the first inner iteration at which Condition 3.3 is satisfied with $\theta_e = 0$. If $T_k = 0$, then we would have a sufficiently accurate starting point for the algorithm to terminate before the first complete iteration. Therefore, without loss of generality, we assume that $T_k > 0$. By Theorem 4.4, the inner subproblem termination condition, Condition 3.3 with $\theta_e = 0$, is satisfied at a given inner iteration $t \in \mathbb{N}_+$ if

778
$$\frac{(L+\alpha||A||^2)\min_{x_k^* \in \mathcal{X}_k^*} ||x_{k-1} - x_k^*||^2}{2(1-2\theta_q^2)t} \le \frac{\alpha \tau_k}{2}.$$

Thus, we have a deterministic upper bound $\Omega_k > 0$ on the random variable T_k ; namely,

780 (4.22)
$$T_k \le \Omega_k \stackrel{\text{def}}{=} \left[\frac{(L + \alpha ||A||^2)}{(1 - 2\theta_q^2)\alpha \tau_k} \min_{x_k^* \in \mathcal{X}_k^*} ||x_{k-1} - x_k^*||^2 \right].$$

We now analyze the total number of expected stochastic gradient evaluations. First, consider the expected sample complexity at each outer iteration k:

783
$$\mathbb{E}_{k}[\mathcal{W}_{k}] = \mathbb{E}_{k} \left[\sum_{t=0}^{T_{k}-1} |S_{k,t}| \right] \leq \frac{c_{0}\omega}{\theta_{g}^{2} \min_{x_{k}^{*} \in \mathcal{X}_{k}^{*}} \|x_{k-1} - x^{*}\|^{2}} \sum_{t=0}^{\Omega_{k}-1} (t+1)$$
784 (4.23)
$$\leq \frac{c_{0}\omega \Omega_{k}^{2}}{\theta_{g}^{2} \min_{x_{k}^{*} \in \mathcal{X}_{k}^{*}} \|x_{k-1} - x^{*}\|^{2}},$$

where the first inequality is due to Assumption 4.1 with $\delta_0 = 0$. Substituting (4.22)

787 into (4.23), using $\lceil x \rceil \le x + 1$ and $||a + b||^2 \le 2||a||^2 + 2||b||^2$ for any $a, b \in \mathbb{R}^n$, we have

788 that

(4.24)

$$\mathbb{E}_{k}[\mathcal{W}_{k}] \leq \frac{2c_{0}\omega(L+\alpha\|A\|^{2})^{2}\min_{x_{k}^{*}\in\mathcal{X}_{k}^{*}}\|x_{k-1}-x_{k}^{*}\|^{2}}{\theta_{g}^{2}(1-2\theta_{g}^{2})^{2}\alpha^{2}\tau_{k}^{2}} + \frac{2c_{0}\omega}{\theta_{g}^{2}\min_{x_{k}^{*}\in\mathcal{X}_{k}^{*}}\|x_{k-1}-x^{*}\|^{2}}.$$

Now, $T_k > 0$ implies that Condition 3.3 is violated at $x_{k,0} = x_{k-1}$. That is,

$$\mathbb{E}_{k}\left[\mathcal{L}(x_{k-1},\lambda_{k};\alpha)-\mathcal{L}(x_{k}^{*},\lambda_{k};\alpha)\right] > \frac{\alpha\tau_{k}}{2}.$$

Recalling that x_{k-1} is in the filtration \mathbb{T}_k , and using convexity of \mathcal{L} , it follows that

795
$$\frac{\alpha \tau_k}{2} < \mathcal{L}(x_{k-1}, \lambda_k; \alpha) - \mathcal{L}(x_k^*, \lambda_k; \alpha)$$
796
$$\leq \|\nabla \mathcal{L}(x_{k-1}, \lambda_k; \alpha)\| \|x_{k-1} - x_k^*\|$$
797
$$\leq \sqrt{D_{\mathcal{L}}} \|x_{k-1} - x_k^*\|,$$

799 for all $x_k^* \in \mathcal{X}_k^*$. Therefore,

$$\min_{\substack{x_k^* \in \mathcal{X}_k^* \\ 801}} \|x_{k-1} - x_k^*\|^2 > \frac{\alpha^2 \tau_k^2}{4D_{\mathcal{L}}}$$

Now, summing the inequality (4.24) from k=0 to K-1, taking full expectation,

803 using (2.10) and (4.25), $\tau_k = \tau_0(k+1)^{-1-\delta/2}$, and (4.19), it follows that

$$\mathbb{E}\left[\sum_{k=0}^{K-1} \mathcal{W}_{k}\right] \leq \frac{2c_{0}\omega(L+\alpha\|A\|^{2})^{2}}{\theta_{g}^{2}(1-2\theta_{g}^{2})^{2}\alpha^{2}} \sum_{k=0}^{K-1} \frac{\mathbb{E}[\min_{x_{k}^{*} \in \mathcal{X}_{k}^{*}} \|x_{k-1} - x_{k}^{*}\|^{2}]}{\tau_{k}^{2}} + \frac{8c_{0}\omega D_{\mathcal{L}}}{\theta_{g}^{2}\alpha^{2}} \sum_{k=0}^{K-1} \frac{1}{\tau_{k}^{2}}$$

$$\leq \frac{B}{\alpha^2 \tau_0^2} \sum_{k=0}^{K-1} (k+1)^{2+\delta}$$

$$\underset{807}{806} \quad (4.26) \qquad \leq \frac{B}{\alpha^2 \tau_0^2 (3+\delta)} (K+1)^{3+\delta} \,.$$

808 We now consider the total number of stochastic gradients evaluated in the final step

described in Corollary 4.3 with \tilde{x} satisfying (i). Following a similar approach to the

derivation of (4.24) and (4.25), and using (2.10), we have that

$$\mathbb{E}[\tilde{W}] \leq \frac{2c_0\omega(L + \tilde{\alpha}||A||^2)^2K^2D^2}{\theta_q^2(1 - 2\theta_q^2)^2\tilde{\alpha}^2\tilde{h}^2} + \frac{8c_0\omega K^2D_{\mathcal{L}}}{\theta_q^2\tilde{\alpha}^2\tilde{h}^2} = \frac{\tilde{B}K^2}{\tilde{\alpha}^2h^2}.$$

813 Finally, we can define the expected total number of gradient evaluations as

814 (4.28)
$$\mathbb{E}[\mathcal{W}] = \mathbb{E}\left[\sum_{k=0}^{K-1} \mathcal{W}_k\right] + \mathbb{E}[\tilde{\mathcal{W}}]$$

$$\leq \frac{B(K+1)^{3+\delta}}{\alpha^2 \tau_0^2 (3+\delta)} + \frac{\tilde{B}K^2}{\tilde{\alpha}^2 h^2}.$$

Substituting $K = K_{\epsilon}$ into (4.29) completes the proof.

REMARK 4.1. In Theorem 4.5, we state an additional assumption related to the boundedness of the gradients of the augmented Lagrangian functions at the iterates computed by the algorithm. We note that this is a mild assumption and can be proven using Assumptions 2.1 and 2.2, and if the dual variables λ_k are bounded. Due to the convergence results established in Subsection 4.1, it is reasonable to assume that the dual variables are bounded.

4.3. Sample Complexity: $\alpha = \mathcal{O}(\epsilon^{-1})$. Theorem 4.5 establishes total outer iteration complexity, K_{ϵ} , and expected sample complexity, $\mathbb{E}[\mathcal{W}]$, for any choice of the penalty parameter α . If α and the other parameters (e.g., τ_0, h) given in Theorem 4.5 are chosen to be independent of the accuracy ϵ , then $K_{\epsilon} = \mathcal{O}(\epsilon^{-1})$ and $\mathbb{E}[\mathcal{W}] =$ $\mathcal{O}(\epsilon^{-3-\delta})$. However, this sample complexity bound is not tight as the optimal sample complexity for stochastic convex programs is $\mathcal{O}(\epsilon^{-2})$ [51, 85]. The next corollary establishes that this optimal sample complexity can be achieved when $\alpha = \mathcal{O}(\epsilon^{-1})$.

COROLLARY 4.6. Under the conditions of Theorem 4.5, if $\alpha = c_{\alpha} \epsilon^{-1}$, $\tau_0 = c_{\tau} \epsilon$, and $h = c_h \epsilon$ for some $c_{\alpha}, c_{\tau}, c_h \in (0, \infty)$. Then $K_{\epsilon} = \mathcal{O}(1)$ and 832

833 (4.30)
$$\mathbb{E}[\mathcal{W}] = \mathcal{O}(\epsilon^{-2}).$$

Proof. Substituting α, τ_0 , and h values into (4.20), it follows that 835

836
$$K_{\epsilon} \leq 1 + \frac{1}{\epsilon} \max \left\{ \frac{8\epsilon}{c_{\alpha}} [q_{\alpha}(\lambda_{0}) - q_{\alpha}(\lambda^{*})] + 2\epsilon (4c_{\tau}a_{\infty} + c_{h}), \frac{\tilde{\alpha}c_{h}\epsilon}{\tilde{\eta}} \right\}$$
837
$$= 1 + \max \left\{ \frac{8}{c_{\alpha}} [q_{\alpha}(\lambda_{0}) - q_{\alpha}(\lambda^{*})] + 2(4c_{\tau}a_{\infty} + c_{h}), \frac{\tilde{\alpha}c_{h}}{\tilde{\eta}} \right\}$$
838
$$= \mathcal{O}(1).$$

We now analyze the sample complexity. Using α, τ_0 , and h values, and $\epsilon < 1$, it 840 follows that

842
$$\frac{B(K_{\epsilon}+1)^{3+\delta}}{\alpha^{2}\tau_{0}^{2}(3+\delta)} = \frac{2c_{0}\omega}{\theta_{g}^{2}c_{\alpha}^{2}c_{\tau}^{2}(3+\delta)} \left(\frac{(L+c_{\alpha}\epsilon^{-1}||A||^{2})^{2}D^{2}}{(1-2\theta_{g}^{2})^{2}} + 4D_{\mathcal{L}}\right) (K_{\epsilon}+1)^{3+\delta}$$

843 (4.32)
$$\leq \frac{2c_0\omega}{\epsilon^2\theta_g^2c_\alpha^2c_\tau^2(3+\delta)} \left(\frac{(L+c_\alpha||A||^2)^2D^2}{(1-2\theta_g^2)^2} + 4D_\mathcal{L}\right) (K_\epsilon + 1)^{3+\delta},$$

and 845

818 819

820

822

823

824

825

826

828

829

830

831

846 (4.33)
$$\frac{\tilde{B}K_{\epsilon}^{2}}{\tilde{\alpha}^{2}h^{2}} = \frac{2c_{0}\omega}{\epsilon^{2}\theta_{g}^{2}\tilde{\alpha}^{2}c_{h}^{2}} \left(\frac{(L+\tilde{\alpha}\|A\|^{2})^{2}D^{2}}{(1-2\theta_{g}^{2})^{2}} + 4D_{\mathcal{L}}\right)K_{\epsilon}^{2}.$$

Substituting (4.32) and (4.33) in (4.21), we have that,

849
$$\mathbb{E}[\mathcal{W}] \leq \frac{2c_0\omega}{\epsilon^2 \theta_g^2 c_\alpha^2 c_\tau^2 (3+\delta)} \left(\frac{(L+c_\alpha ||A||^2)^2 D^2}{(1-2\theta_g^2)^2} + 4D_\mathcal{L} \right) (K_\epsilon + 1)^{3+\delta}$$
850
$$+ \frac{2c_0\omega}{\epsilon^2 \theta_g^2 \tilde{\alpha}^2 c_h^2} \left(\frac{(L+\tilde{\alpha}||A||^2)^2 D^2}{(1-2\theta_g^2)^2} + 4D_\mathcal{L} \right) K_\epsilon^2$$
851
$$= \mathcal{O}(\epsilon^{-2}),$$

where the last equality is due to the fact that all other constants in the inequality are 853 independent of the choice of ϵ . 854

REMARK 4.2. We observe that the complexity results given in Theorem 4.5 and Corollary 4.6 do not exploit the benefits of using the previous iterate $x_{k-1} = x_{k,0}$ as the starting point for solving the current subproblem. That is, the bound on $\mathbb{E}\left[\min_{x_k^* \in \mathcal{X}_k^*} \|x_{k-1} - x_k^*\|^2\right]$ is not tight. The difficulty in exploiting the benefits of this procedure is due to the fact that the augmented Lagrangian functions are only convex but not necessarily strongly convex. In Subsection 4.4, we consider strongly convex functions and establish the advantages of this procedure.

4.4. Sample Complexity: $\mathcal{X} = \mathbb{R}^n$. We provide improved convergence and complexity results when $\mathcal{X} = \mathbb{R}^n$ and the objective function f is μ -strongly convex.

Assumption 4.2. The objective function f is μ -strongly convex. That is,

$$\nabla^2 f(x) \succeq \mu I \quad \forall x \in \mathbb{R}^n$$

where $I \in \mathbb{R}^{n \times n}$ is an identity matrix.

855

856

857

858

859

860

861

862

863

864

882

883

884

We should note that Assumption 4.2 implies Assumption 2.3. In this case, the inner subproblems are unconstrained and have unique optimal solutions. Moreover, the optimality conditions given in (2.3) can be written as

$$\nabla \ell_x(x,\lambda) = 0$$
 and $c(x) = 0$.

It can also be shown that the negative dual function $q(\lambda)$ is strongly convex in this setting, as stated in the following proposition (cf. [37, Propositions 3.1 and 3.3] and the references therein, [90, Theorem 1], [38, Proposition 2.5]).

PROPOSITION 4.7. If Assumptions 2.1 and 4.2 hold with $\mathcal{X} = \mathbb{R}^n$, then $q(\lambda)$ 877 defined in (2.12) is strongly convex with the strong convexity parameter $\mu_q = \frac{\sigma}{\mu + L}$ 878 where $\sigma = \lambda_{\min}(AA^T)$.

For the sake of completeness, we include the proof of this proposition in Appendix A.
We also state the following well-known result for strongly convex functions with
Lipschitz continuous gradients (cf. [64, Theorem 2.1.5 and Theorem 2.1.10])

PROPOSITION 4.8. If the function $q_{\alpha}(\lambda)$ is strongly convex with parameter μ_{α} and has a Lipschitz continuous gradient with Lipschitz constant L_{α} , then for any $\lambda \in \mathbb{R}^m$, it holds that

$$2\mu_{\alpha}(q_{\alpha}(\lambda) - q_{\alpha}(\lambda^*)) \le \|\nabla q_{\alpha}(\lambda)\|^2 \le 2L_{\alpha}(q_{\alpha}(\lambda) - q_{\alpha}(\lambda^*)),$$

887 where $\lambda^* = \arg\min_{\lambda} q_{\alpha}(\lambda)$.

Note that $L_{\alpha} = \alpha^{-1}$ by Lemma 2.1. We now establish a linear rate of convergence of both feasibility error and stationarity error. For the sake of brevity, we only consider Condition 3.4.

THEOREM 4.9. Suppose Assumptions 2.1 and 4.2 hold and $\mathcal{X} = \mathbb{R}^n$. For any x_{-1}, λ_0 and $\alpha > 0$, let $\{(x_k, \lambda_k)\}$ be the sequence of primal-dual iterates generated by Algorithm 3.1. If the primal iterates x_k satisfy Condition 3.4 at each iteration $k \in \mathbb{N}$ with $\tilde{\theta}_e \leq \frac{\mu \theta_e}{2\|A\|}$, $\tilde{\tau}_k = \frac{\mu^2 \tau_k}{4\|A\|^2}$, $\theta_e \in [0, 1)$, and $\tau_k = \tau_0(1/a)^k$ for some $\tau_0 > 0$ and a > 1, then

$$\mathbb{E}[\|c(x_k)\|^2] \le A_1 \rho^k \quad and \quad \mathbb{E}\left[\|\nabla \ell_x(x_k, \lambda_{k+1}))\|^2\right] \le A_2 \rho^k,$$

898 where
$$A_1 = 4(1 + \theta_e^2)L_{\alpha}A_3 + 2\tau_0$$
, $A_2 = 2L_{\alpha}\tilde{\theta}_e^2A_3 + \frac{\mu^2\tau_0}{4\|A\|^2}$,
899 $A_3 = \max\left\{q_{\alpha}(\lambda_0) - q_{\alpha}(\lambda^*), \frac{\tau_0}{\mu_{\alpha}(1 - \theta_e^2)}\right\}$, $\rho = \max\left\{1 - \frac{\alpha\mu_{\alpha}(1 - \theta_e^2)}{2}, \frac{1}{a}\right\} < 1$, and $\mu_{\alpha} = \frac{\mu_q}{\mu_{\alpha}\alpha + 1}$.

Proof. Using Proposition 4.7 and Lemma 2.1, it follows that the $q_{\alpha}(\lambda)$ is a strongly convex function with strong convexity parameter $\frac{\mu_q}{\mu_q \alpha + 1}$. Therefore, substituting (4.34) into (4.5), using (2.23), subtracting $q_{\alpha}(\lambda^*)$ from both sides and taking full expectation we obtain

$$\mathbb{E}[q_{\alpha}(\lambda_{k+1}) - q_{\alpha}(\lambda^{*})] \leq (1 - \alpha\mu_{\alpha}(1 - \theta_{e}^{2})) \,\mathbb{E}[q_{\alpha}(\lambda_{k}) - q_{\alpha}(\lambda^{*})] + \frac{\alpha\tau_{k}}{2}$$

$$\leq (1 - \alpha\mu_{\alpha}(1 - \theta_{e}^{2})) \,\mathbb{E}[q_{\alpha}(\lambda_{k}) - q_{\alpha}(\lambda^{*})] + \frac{\alpha\tau_{0}}{2a^{k}}$$

where the second inequality is due to $\tau_k = \tau_0(1/a)^k$. It is now a straightforward exercise in mathematical induction to show that

910 (4.36)
$$\mathbb{E}[q_{\alpha}(\lambda_k) - q_{\alpha}(\lambda^*)] \le A_3 \rho^k \quad \forall k \in \mathbb{N}.$$

The statement is trivially true for k = 0. Let's assume it is true for iteration k. For iteration k + 1, it follows that

913
$$\mathbb{E}[q_{\alpha}(\lambda_{k+1}) - q_{\alpha}(\lambda^{*})] \leq \left(1 - \alpha\mu_{\alpha}(1 - \theta_{e}^{2})\right) \mathbb{E}[q_{\alpha}(\lambda_{k}) - q_{\alpha}(\lambda^{*})] + \frac{\alpha\tau_{0}}{2a^{k}}$$
914
$$\leq A_{3}\rho^{k} \left(1 - \alpha\mu_{\alpha}(1 - \theta_{e}^{2}) + \frac{\alpha\tau_{0}}{2A_{3}(a\rho)^{k}}\right)$$
915
$$\leq A_{3}\rho^{k} \left(1 - \alpha\mu_{\alpha}(1 - \theta_{e}^{2}) + \frac{\alpha\mu_{\alpha}(1 - \theta_{e}^{2})}{2}\right)$$
916
$$\leq A_{3}\rho^{k+1},$$

where the second inequality is due to the statement of the induction, third inequality is due to $\rho \geq 1/a$ and the definition of A_3 , and the last inequality is due to the definition of ρ . Hence, (4.36) is satisfied. Substituting $\lambda = \lambda_k$ in (4.34), by (2.23) and taking expectation of both sides it follows that,

922
$$\mathbb{E}[\|c(x_k^*)\|^2] \le 2L_\alpha \mathbb{E}[q_\alpha(\lambda_k) - q_\alpha(\lambda^*)]$$
933 (4.37)
$$\le 2L_\alpha A_3 \rho^k.$$

Therefore, using the definitions of A_1 and ρ , we have that

926
$$\mathbb{E}[\|c(x_k)\|^2] \leq 2\mathbb{E}[\|c(x_k^*) - c(x_k)\|^2] + 2\mathbb{E}[\|c(x_k^*)\|^2]$$
927
$$\leq 2(1 + \theta_e^2)\mathbb{E}[\|c(x_k^*)\|^2] + 2\tau_k$$
928
$$\leq \rho^k \left(4(1 + \theta_e^2)L_\alpha A_3 + \frac{2\tau_0}{(\rho a)^k}\right)$$
938
$$\leq A_1 \rho^k.$$

Using (4.37), Condition 3.4, and (4.7), it follows that

932
$$\mathbb{E}\left[\|\nabla \ell_x(x_k, \lambda_{k+1}))\|^2\right] \leq \tilde{\theta}_e^2 \mathbb{E}[\|c(x_k^*)\|^2] + \tilde{\tau}_k$$
933
$$\leq 2L_\alpha \tilde{\theta}_e^2 A_3 \rho^k + \frac{\mu^2 \tau_0}{4\|A\|^2 a^k}$$

$$\leq A_2 \rho^k,$$

- 936 where the last inequality is due to definitions of A_2 and $\rho a \geq 1$.
- 937 We now derive the sample complexity results. We will use the fact that the metric

938 (4.13) can be simplified to

939 (4.38)
$$\max \left\{ \mathbb{E}[\|c(x_k)\|^2], \mathbb{E}\left[\|\nabla \ell_x(x_k, \lambda_{k+1})\|^2\right] \right\} \le \epsilon \in (0, 1),$$

- since $\mathcal{X} = \mathbb{R}^n$. Moreover, the adaptive sampling projected gradient method employed
- 942 to solve the inner subproblems converges at a linear rate as stated below (cf. [7,
- 943 Theorem 2.10], [83, Theorem 3.7]).
- Theorem 4.10. Suppose Assumptions 2.1, 2.4 and 4.2 hold. If $\eta = \frac{(1-2\theta_g^2)}{L+\alpha||A||^2}$ with
- 945 $\theta_g \in [0, \frac{1}{2})$, and Condition 3.1 is satisfied. Then, for any outer iteration $k \in \mathbb{N}$ and
- 946 inner iteration $t \in \mathbb{N}$, it holds that

$$\mathbb{E}_{k}[\mathcal{L}(x_{k,t},\lambda_{k};\alpha) - \mathcal{L}(x_{k}^{*},\lambda_{k};\alpha)] \leq \rho_{\mathcal{L}}^{t}(\mathcal{L}(x_{k,0},\lambda_{k};\alpha) - \mathcal{L}(x_{k}^{*},\lambda_{k};\alpha)),$$

- 949 where $\rho_{\mathcal{L}} = 1 \frac{(1 2\theta_g^2)\mu}{L + \alpha ||A||^2} \in [0, 1).$
- 950 Using Proposition 4.8, it can be shown that the gradient of the augmented Lagrangian
- 951 function also converges to zero. That is, applying Proposition 4.8 to the augmented
- 952 Lagrangian function, we have that

953
$$2\mu(\mathcal{L}(x_{k,t},\lambda_k;\alpha) - \mathcal{L}(x_k^*,\lambda;\alpha)) \leq \|\nabla_x \mathcal{L}(x_{k,t},\lambda_k;\alpha)\|^2$$

$$\leq 2(L+\alpha\|A\|^2)(\mathcal{L}(x_{k,t},\lambda_k;\alpha) - \mathcal{L}(x_k^*,\lambda;\alpha)).$$

956 Combining (4.39) and (4.40), it follows that,

957 (4.41)
$$\mathbb{E}_{k}[\|\nabla_{x}\mathcal{L}(x_{k,t},\lambda_{k};\alpha)\|^{2}] \leq \frac{L+\alpha\|A\|^{2}}{\mu}\rho_{\mathcal{L}}^{t}\|\nabla_{x}\mathcal{L}(x_{k,0},\lambda_{k};\alpha)\|^{2}.$$

- 959 The next theorem establishes pessimistic upper bounds on the sample sizes employed
- at each outer iteration $k \in \mathbb{N}$ and each inner iteration $t \in \mathbb{N}$, and the number of
- 961 inner iterations T_k required to satisfy Condition 3.4. For the sake of brevity, in this
- complexity analysis, we employ Condition 3.4 with $\tilde{\theta}_e = 0$, and also assume that
- 963 Assumption 2.4 holds with $\omega_1 = 1$ and $\omega_2 = \omega$, i.e.,

$$\mathbb{E}_{\zeta}[\|\nabla F(x,\zeta) - \nabla f(x)\|^2] \le \omega.$$

THEOREM 4.11. Suppose the conditions of Theorem 4.10 are satisfied and (4.42)

967 holds. Then the number of inner iterations T_k required to satisfy Condition 3.4 with

968 $\theta_e = 0$ are bounded from above as follows:

(4.43)
$$T_{k} \leq \left\lceil \log_{1/\rho_{\mathcal{L}}} \left(\frac{2(L+\alpha||A||^{2}) \left(\alpha^{2}||A||^{2}||c(x_{k-1})||^{2} + ||R(x_{k-1},\lambda_{k-1};\alpha,\eta)||^{2}\right)}{\mu \tilde{\tau}_{k}} \right) \right\rceil.$$

Moreover, for any inner iteration $t < T_k$, the sample sizes $|S_{k,t}|$ are at most

972 (4.44)
$$|S_{k,t}| \le \frac{\omega}{\theta_g^2 \tilde{\tau}_k}.$$

Proof. At any outer iteration $k \in N$, let T_k denote the first inner iteration t at which the following condition holds:

$$\min \left\{ \| R(x_{k,t}, \lambda_k; \alpha, \eta) \|^2, \mathbb{E}_{k,t} \left[\| R(x_{k,t}, \lambda_k; \alpha, \eta) \|^2 \right] \right\} \le \tilde{\tau}_k.$$

Hence, at $t = T_k$, Condition 3.4 is satisfied with $\tilde{\theta}_e = 0$ at $x_k = x_{k,t}$. Therefore, for all $t < T_k$, it follows that

$$\|R(x_{k,t},\lambda_k;\alpha,\eta)\|^2 > \tilde{\tau}_k.$$

Using (3.14), (4.42), and (4.46), and choosing the smallest sample size $|S_{k,t}|$ satisfying (4.16), it follows that

984 (4.47)
$$|S_{k,t}| \le \frac{\omega(1+\theta_g)^2}{\theta_g^2 \tilde{\tau}_k}.$$

Now, let us bound the number of inner iterations required to satisfy (4.45). Using (3.5) and (4.41) and $\mathcal{X} = \mathbb{R}^n$ it follows that

988 (4.48)
$$\mathbb{E}_{k} \left[\| R(x_{k,t}, \lambda_{k}; \alpha, \eta) \|^{2} \right] \leq \frac{L + \alpha \|A\|^{2}}{\mu} \rho_{\mathcal{L}}^{t} \| R(x_{k,0}, \lambda_{k}; \alpha, \eta) \|^{2}.$$

990 Substituting $x_{k,0} = x_{k-1}$, it follows that

991
$$||R(x_{k-1}, \lambda_k; \alpha, \eta)||^2$$
992
$$= ||R(x_{k-1}, \lambda_k; \alpha, \eta) - R(x_{k-1}, \lambda_{k-1}; \alpha, \eta) + R(x_{k-1}, \lambda_{k-1}; \alpha, \eta)||^2$$
983
$$(4.49) \qquad \leq 2||R(x_{k-1}, \lambda_k; \alpha, \eta) - R(x_{k-1}, \lambda_{k-1}; \alpha, \eta)||^2 + 2||R(x_{k-1}, \lambda_{k-1}; \alpha, \eta)||^2$$

where the last inequality is due to $||a+b||^2 \le 2||a||^2 + 2||b||^2$ for any $a, b \in \mathbb{R}^n$. Consider

996
$$||R(x_{k-1}, \lambda_k; \alpha, \eta) - R(x_{k-1}, \lambda_{k-1}; \alpha, \eta)||$$
997
$$= ||\nabla_x \mathcal{L}(x_{k-1}, \lambda_k; \alpha) - \nabla_x \mathcal{L}(x_{k-1}, \lambda_{k-1}; \alpha)||$$
998
$$= ||\langle \lambda_k - \lambda_{k-1}, \nabla c(x_{k-1}) \rangle||$$
1888
$$(4.50) \leq \alpha ||A|| ||c(x_{k-1})||$$

where the first equality is due to (3.5) and the inequality is due to $\lambda_k = \lambda_{k-1} - \alpha c(x_{k-1})$. Using (4.48), (4.49), and (4.50), it follows that

1003
1004 (4.51)
$$\mathbb{E}_{k} \left[\| R(x_{k,t}, \lambda_{k}; \alpha, \eta) \|^{2} \right]$$

$$\leq \frac{L + \alpha ||A||^2}{\mu} \rho_{\mathcal{L}}^t \left(2\alpha^2 ||A||^2 ||c(x_{k-1})||^2 + 2||R(x_{k-1}, \lambda_{k-1}; \alpha, \eta)||^2 \right).$$

1007 Therefore, for any

$$1008 \qquad t \ge \left\lceil \log_{1/\rho_{\mathcal{L}}} \left(\frac{2(L + \alpha ||A||^2) \left(\alpha^2 ||A||^2 ||c(x_{k-1})||^2 + ||R(x_{k-1}, \lambda_{k-1}; \alpha, \eta)||^2 \right)}{\mu \tilde{\tau}_k} \right) \right\rceil$$

1010 we have,

$$\mathbb{E}_{k}\left[\left\|R(x_{k,t},\lambda_{k};\alpha,\eta)\right\|^{2}\right] \leq \tilde{\tau}_{k}.$$

1013 Using (4.45), it follows that,

- which completes the proof.
- 1017 We are now ready to provide the main complexity theorem for this subsection.

THEOREM 4.12. Suppose the conditions of Theorems 4.9 and 4.10, and (4.42) hold. Then the number of outer iterations to get an ϵ -accurate solution (x_k, λ_{k+1}) satisfying (4.38) is

$$K_{\epsilon} = \left\lceil \log_{1/\rho} \left(\frac{\max\{A_1, A_2\}}{\epsilon} \right) \right\rceil = \mathcal{O}\left(\log(1/\epsilon)\right),$$

where A_1 , A_2 are defined in Theorem 4.9, and the expected number of stochastic gradient evaluations is

$$\mathbb{E}[\mathcal{W}] = \mathcal{O}\left(\epsilon^{-1}\log(1/\epsilon)\right).$$

1027 Proof. Equation (4.53) directly follows from (4.35). Now, consider the sample 1028 complexity at each outer iteration k

1029 (4.55)
$$\mathcal{W}_k \stackrel{\text{def}}{=} \sum_{t=0}^{T_k - 1} |S_{k,t}| \le \frac{\omega (1 + \theta_g)^2}{\theta_g^2 \tilde{\tau}_k} T_k,$$

where the inequality is due to (4.44). Therefore, the expected total number of gradient evaluations is found to be

1033
$$\mathbb{E}\left[\sum_{k=0}^{K-1} \mathcal{W}_{k}\right] \leq \mathbb{E}\left[\sum_{k=0}^{K-1} \frac{\omega(1+\theta_{g})^{2}}{\theta_{g}^{2}\tilde{\tau}_{k}} T_{k}\right]$$

$$\leq \sum_{k=0}^{K-1} \frac{\omega(1+\theta_{g})^{2}}{\theta_{g}^{2}\tilde{\tau}_{k}} \mathbb{E}[T_{k}]$$

$$\leq \sum_{k=0}^{K-1} \frac{4\omega(1+\theta_{g})^{2} ||A||^{2} a^{k}}{\theta_{g}^{2}\tau_{0}\mu^{2}} \mathbb{E}[T_{k}],$$
1035 (4.56)

where the last inequality is due to $\tilde{\tau}_k = \frac{\mu^2 \tau_k}{4\|A\|^2}$ and $\tau_k = \tau_0 (1/a)^k$. Using (4.43) and

taking the full expectation of both sides, it follows that

$$\mathbb{E}[T_{k} - 1] \\
1040 \qquad \leq \mathbb{E}\left[\log_{1/\rho_{\mathcal{L}}} \left(\frac{2(L + \alpha \|A\|^{2}) \left(\alpha^{2} \|A\|^{2} \|c(x_{k-1})\|^{2} + \|R(x_{k-1}, \lambda_{k-1}; \alpha, \eta)\|^{2}\right)}{\mu\tilde{\tau}_{k}}\right)\right] \\
1041 \qquad = \mathbb{E}\left[\log_{1/\rho_{\mathcal{L}}} \left(\frac{2(L + \alpha \|A\|^{2}) \left(\alpha^{2} \|A\|^{2} \|c(x_{k-1})\|^{2} + \|\nabla_{x}\ell(x_{k-1}, \lambda_{k}; \alpha, \eta)\|^{2}\right)}{\mu\tilde{\tau}_{k}}\right)\right] \\
1042 \qquad \leq \log_{1/\rho_{\mathcal{L}}} \left(\frac{2(L + \alpha \|A\|^{2}) \left(\alpha^{2} \|A\|^{2} \mathbb{E}[\|c(x_{k-1})\|^{2}] + \mathbb{E}[\|\nabla_{x}\ell(x_{k-1}, \lambda_{k}; \alpha, \eta)\|^{2}]\right)}{\mu\tilde{\tau}_{k}}\right) \\
1043 \qquad \leq \log_{1/\rho_{\mathcal{L}}} \left(\frac{2(L + \alpha \|A\|^{2}) \left(\alpha^{2} \|A\|^{2} A_{1} + A_{2}\right) \rho^{k-1}}{\mu\tilde{\tau}_{k}}\right) \\
(4.57) \qquad (4.57) \\
1044 \qquad = \log_{1/\rho_{\mathcal{L}}} \left(\frac{8\|A\|^{2}(L + \alpha \|A\|^{2}) \left(\alpha^{2} \|A\|^{2} A_{1} + A_{2}\right) \rho^{k-1} a^{k}}{\mu^{3}\tau_{0}}\right) = \mathcal{O}(k),$$

where the second line is due to (4.7), third line due to Jensen's inequality, fourth line is due to Theorem 4.9, and the last line follows from $\tilde{\tau}_k = \frac{\mu^2 \tau_k}{4\|A\|^2}$ and $\tau_k = \tau_0 (1/a)^k$.

Therefore, (4.57) shows that there exist $s_1 > 0$ and $s_2 > 0$ such that

$$\mathbb{E}[T_k] \le s_1 + s_2 k.$$

Substituting (4.58) into (4.56), we have that

1052
$$\mathbb{E}\left[\sum_{k=0}^{K-1} \mathcal{W}_{k}\right] \leq \sum_{k=0}^{K-1} \frac{4\omega(1+\theta_{g})^{2} \|A\|^{2} a^{k}}{\theta_{g}^{2} \tau_{0} \mu^{2}} (s_{1}+s_{2}k)$$

$$1053 \qquad \leq \sum_{k=0}^{K-1} \frac{4\omega(1+\theta_{g})^{2} \|A\|^{2} a^{k}}{\theta_{g}^{2} \tau_{0} \mu^{2}} (s_{1}+s_{2}K)$$

$$1054 \qquad \leq \frac{4\omega(1+\theta_{g})^{2} \|A\|^{2} a^{K}}{\theta_{g}^{2} \tau_{0} \mu^{2} (a-1)} (s_{1}+s_{2}K)$$

$$1055 \qquad = \mathcal{O}\left(\epsilon^{-1} \log(1/\epsilon)\right),$$

where the last line is due to (4.53).

Remark 4.3. It is important to emphasize that performing sampling complexity analysis for adaptive sampling methods is quite challenging with present optimization techniques. However, these methods fall under a general class of increasing batch size methods where one can establish theoretical sample complexity analysis that shows stochastic gradient and increasing batch size mechanisms have similar total sample complexity results (see, e.g., [23]). We have established pessimistic (i.e., worst-case) complexity bounds where the sample sizes at each inner iteration are bounded above by the largest sample size employed across all inner iterations at any given outer iteration k (cf. (4.44)). Owing to this pessimistic bound on sample sizes, the overall complexity bound $\mathcal{O}\left(\epsilon^{-1}\log(1/\epsilon)\right)$ is slightly worse than the optimal sample complexity $\mathcal{O}\left(\epsilon^{-1}\right)$ for strongly convex stochastic programming problems [19,84].

5. Practical Algorithm. In this section, we present a complete and practical adaptive sampling augmented Lagrangian (ASAL) algorithm that uses an adaptive sampling proximal gradient method to inexactly solve the augmented Lagrangian subproblems. We describe the mechanism by which the sample size is selected at each inner iteration and the mechanism to terminate the subproblem solver.

The sample size selection and inexactness conditions described in Subsections 3.1 and 3.2 respectively are impractical as they require computing exact variances or deterministic quantities such as $\mathcal{L}(x_k, \lambda_k; \alpha)$ and $R(x_k, \lambda_k; \alpha, \eta)$. That being said, these quantities can be approximated using sample variances and sampled stochastic counterparts of the deterministic quantities following the ideas proposed in [7,16,83].

Sample Size Selection. We propose the following practical sampling test to approximate Condition 3.1 where the left-hand-side is the sample variance that approximates the exact variance and the right-hand-side is the stochastic projected (reduced) gradient that approximates the expectation of this quantity.

TEST 5.1 (Practical Sampling Test). For any given $\theta_g \geq 0$, the sample size $|S_{k,t}|$ satisfies

1085 (5.1)
$$\frac{\frac{1}{|S_{k,t}|-1} \sum_{\zeta_i \in S_{k,t}} \|\nabla F(x_{k,t},\zeta_i) - \nabla F_{S_{k,t}}(x_{k,t})\|^2}{|S_{k,t}|} \le \theta_g^2 \|R_{S_{k,t}}(x_{k,t},\lambda_k;\alpha,\eta)\|^2.$$

In our practical Algorithm 5.1, we aim to satisfy Test 5.1 at each inner iteration using the following procedure. Whenever (5.1) is *not* satisfied at the current inner iteration t, we attempt to ensure (5.1) will be satisfied at the next inner iteration t+1 by using the relative variance,

1090 (5.2)
$$\nu_t \stackrel{\text{def}}{=} \frac{\frac{1}{|S_{k,t}|-1} \sum_{\zeta_i \in S_{k,t}} \|\nabla F(x_{k,t}, \zeta_i) - \nabla F_{S_{k,t}}(x_{k,t})\|^2}{\theta_q^2 |S_{k,t}| \|R_{S_{k,t}}(x_{k,t}, \lambda_k; \alpha, \eta)\|^2},$$

to select the next sample size. More specifically, we set $|S_{k,t+1}| = \lceil \nu_t |S_{k,t}| \rceil$ whenever $\nu_t > 1$.

On the other hand, if (5.1) is satisfied at the current inner iteration t (i.e., $\nu_t \leq 1$), then keeping the sample size unchanged, $|S_{k,t+1}| = |S_{k,t}|$, is a simple rule to maintain control over the sample variance. However, if $\nu_t \ll 1$ is sufficiently small and the current sample size $|S_{k,t}| \gg 1$ is sufficiently large, then it may be beneficial to reduce cost by decreasing the sample size. We explore this possibility by providing an opportunity for the sample size to decrease like $|S_{k,t+1}| = \lceil \nu_t |S_{k,t} \rceil \rceil$ until $|S_{k,t}|$ reaches a minimum value. Lines 8 through 16 in Algorithm 5.1 encapsulate the sample size selection procedure.

Inexactness Conditions. We propose a practical test to terminate the inner subproblem solver. Owing to the difficulty in computing the optimal quantities $c(x_k^*)$ and $\mathcal{L}(x_k^*, \lambda_k; \alpha)$, and the equivalence of Conditions 3.2 through 3.4, we design the practical test based on Condition 3.4. Following a similar procedure employed in approximating the sample size test conditions, we approximate the projected (reduced) gradient with its stochastic counterpart and the optimal constraint violation with the current constraint violation. The resulting practical test is as follows:

¹Although the sample sizes are allowed to decrease, we do not observe sample size decreases in our numerical experiments; cf. Remark 6.1.

TEST 5.2 (Practical Tolerance Test). For any given $\tilde{\theta}_e \in [0,1)$ and $\tilde{\tau}_k \geq 0$ with $\lim_{k \to \infty} \tilde{\tau}_k = 0$,

```
\|R_{S_{k,t}}(x_{k,t},\lambda_k;\alpha,\eta)\|^2 \le \tilde{\theta}_e^2 \|c(x_{k,t})\|^2 + \tilde{\tau}_k.
```

We terminate the inner subproblem whenever (5.3) is violated. Algorithm 5.1 provides a complete description of the ASAL algorithm.

Algorithm 5.1 Adaptive Sampling Augmented Lagrangian (ASAL) Method

Input: $x_{-1} \in \mathbb{R}^n$, $\lambda_0 \in \mathbb{R}^m$, step size $\eta > 0$, penalty parameter $\alpha > 0$, initial sample size $|S_{0,0}|$, sample size test parameters $(\theta_g > 0, \nu_l \in (0,1), s_l > 0, s_{min} > 0)$, inexactness tolerance parameters $(\tilde{\theta}_e \in [0,1), \tilde{\tau}_k \geq 0)$

Initialization: Set $k \leftarrow 0$

1119 1120

1121

1122

```
1: loop
         Set t \leftarrow 0
 2:
 3:
         Set x_{k,0} \leftarrow x_{k-1}
         repeat
 4:
            Choose a set S_{k,t} consisting of |S_{k,t}| i.i.d. realizations of \zeta
 5:
 6:
             Compute R_{S_{k,t}}(x_{k,t},\lambda_k;\alpha,\eta) via (3.3) and (3.4)
 7:
             Update x_{k,t+1} \leftarrow x_{k,t} + \eta R_{S_{k,t}}(x_{k,t}, \lambda_k; \alpha, \eta)
             if Test 5.1 is not satisfied then
 8:
                Set |S_{k,t+1}| \leftarrow \lceil \nu_t |S_{k,t}| \rceil
 9:
10:
                if \nu_t < \nu_l and |S_{k,t}| > s_l then
11:
12:
                    Set |S_{k,t+1}| \leftarrow \max\{s_{min}, \lceil \nu_t | S_{k,t} \rceil \rceil \}
13:
                    Set |S_{k,t+1}| \leftarrow |S_{k,t}|
14:
                end if
15:
            end if
16:
17:
            Set t \leftarrow t + 1
         until Test 5.2 is satisfied
18:
         Set x_k \leftarrow x_{k,t}
19:
         Update \lambda_{k+1} \leftarrow \lambda_k - \alpha c(x_k)
20:
         Set |S_{k+1,0}| \leftarrow |S_{k,t}|
21:
22:
         Set k \leftarrow k+1
23: end loop
```

- 6. Numerical results. In this section, we study the performance of ASAL (Algorithm 5.1) using model problems from machine learning (Subsection 6.1) and engineering (Subsections 6.2 and 6.3). We implement Test 5.2 with $\tilde{\theta}_e = 0$ and $\tilde{\tau}_k = \tau_0/(k+1)$, treating τ_0 as a hyperparameter for this numerical study.
 - 6.1. Logistic regression with multiple disparate impact constraints. We first consider a constrained logistic regression problem. A decision-making system suffers from *disparate impact* if it provides outputs that affect a group of people sharing a value of a sensitive feature more frequently than other groups [4]. In [87, Section 4.4], it is shown that disparate impact can be controlled in binary classification problems by applying deterministic constraints. More explicitly, we consider the optimization

1125 problem

minimize
$$\frac{1}{N} \sum_{i=1}^{N} \left[\log(1 + \exp(-z_i \langle x, y_i \rangle)) \right] + \frac{\gamma}{2} ||x||^2$$
subject to $\langle a_1, x \rangle = b_1, \quad |\langle a_2, x \rangle| \le b_2,$

where $x \in \mathbb{R}^n$ is the optimization variable and $(y_i, z_i) \in \mathbb{R}^n \times \{-1, 1\}$ are input/output pairs from a classification data set. Here, $\gamma > 0$ is a fixed Tikhonov regularization parameter. Meanwhile, $a_1, a_2 \in \mathbb{R}^n$ and $b_1, b_2 \geq 0$ are constraint parameters. In [4], it is suggested to take, e.g., $a_1 = \mathbb{E}_{y,s}[(s - \mathbb{E}_s[s])y]$, where s is a secondary observable, in addition to y. However, for the purpose of demonstration, we arbitrarily set a_1 and a_2 from samples drawn for a standard multivariate normal distribution. Likewise, we set $b_1 = 0.1, b_2 = 0.02$. The initial x_{-1} and λ_0 variables are chosen to be zero vectors, and we set $\gamma = 1/N$.

In this experiment, we use the mushroom classification data set from the LIBSVM collection [26]. The size of this data set is N=8124, and the dimension of the problem is n=112. In order to evaluate the performance of ASAL, we record the feasibility and stationarity errors (2.5) until 200 training epochs (i.e., 200N cumulative gradient evaluations) have elapsed. We then compare ASAL to three separately-tuned fixed-batch-size implementations of ASAL using 10%, 20%, and 50% of the data set size at each iteration, respectively. In this experiment we use $\theta_g=0.99$, $\nu_l=0.5$, $s_l=0.1N$, and $s_{\min}=0.1N$. The value of θ_g is not tuned and is, instead, set at an arbitrary value close to the suggestion for unconstrained problems in [24,34]. The values of the other three fixed hyperparameters are also set arbitrarily. Yet, they appear to have little to no effect on performance; cf. Remark 6.1.

We treat τ_0 , α and the step size η as tunable hyperparameters. All of the hyperparameters are tuned using the following procedure: We run each augmented Lagrangian algorithm for all possible combinations of $\tau_0 = 10^4$, 10^3 , 10^2 , 10^1 , 10^0 , 10^{-1} , $\eta = 10^{-1}$, 10^{-2} , 10^{-3} , 10^{-4} , 10^{-5} , 10^{-6} , and $\alpha = 10^2$, 10^1 , 10^0 , 10^{-1} , 10^{-2} . Then, for each algorithm, we select the run with the smallest average objective function value in the final 5 inner iterations among all runs whose minimum feasibility error in the final 30 inner iterations is less than feasibility tolerance 10^{-4} .

The stationarity and feasibility errors corresponding to the best hyperparameters for each algorithm are overlaid in Figure 6.1. Because the hyperparameter tuning procedure we have used seeks the best stationarity error among runs reaching a feasibility error threshold, it is no surprise that ASAL and each of the three baseline algorithms achieve a similar *minimal* feasibility error (around feasibility tolerance 10^{-4}). Nevertheless, we observe that ASAL outperforms the three baseline algorithms with respect to stationarity error. We also present similar results for australian data set from the LIBSVM collection [26] in Appendix B.

REMARK 6.1. Notice from Figure 6.1 that the ASAL sample size never decreases. This is despite the safeguarding mechanism in line 17 of Algorithm 5.1. We have witnessed this non-decreasing sample size property in all of our experiments with ASAL after tuning the hyperparameters α , η , and τ_0 . Thus, we see little justification for allowing sample size decreases in future implementations of ASAL and do not report the hyperparameters ν_l , s_l , and s_{\min} in the remaining experiments.

Remark 6.2. The starting values for the cumulative gradient evaluations in Figure 6.1 represent the fact that we are recording errors only *after* advancing a single optimization step. Each algorithm began with the same initial guesses x_{-1} and λ_0 .

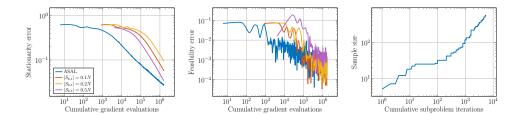


FIGURE 6.1. Results of running Algorithm 5.1 on the constrained logistic regression problem (6.1) with the mushroom classification data set. Notice that ASAL achieves the lowest average stationarity errors while matching the minimal feasibility error of the three baseline algorithms. To generate these results, we used the algorithm parameters $\theta_g = 0.99$ and individually tuned α , η , and τ_0 .

Remark 6.3. Observe that the expected feasibility error with ASAL steadily decreases in Figure 6.1. Meanwhile, the feasibility error in each of the other algorithms plateaus after around 5×10^5 cumulative gradient evaluations. This is due to the stable sample size growth provided by our adaptive sampling strategy and the fact that a fixed number of samples are used for each of the baselines; i.e., the baseline algorithms can only converge in expectation to a neighborhood of the solution. As a result, even though the slopes of the stationarity errors for the baseline algorithms are higher than ASAL after 200 epochs, we conclude that ASAL would remain the better practical algorithm even if a larger epoch threshold had been used.

Remark 6.4. The tuning procedure used in this experiment is expensive and impractical for more expensive problems. Owing to this fact, in the remaining sections, we only compare ASAL to baseline algorithms with a shared set of hyperparameters.

6.2. Optimal truss design. We consider optimizing the simply supported truss structure shown in Figure 6.2 in a problem inspired by an example presented in [73]. The truss elements are numbered as shown in Figure 6.2, and a random force F, pointing downwards, is applied in the middle of the bottom chord. The cross-sections of the truss elements are denoted by x^i , $i=1,2,\ldots,7$, and the yield stress associated with the members with σ_i , $i=1,2,\ldots,7$. The first two yield stress limits σ_i , i=1,2, are log-normal random variables with mean 100N/mm^2 and standard deviation 20N/mm^2 . The yield stresses for all other members are also log-normal, but with mean 200N/mm^2 and standard deviation 40N/mm^2 . The correlation coefficient between σ_1 and σ_2 is 0.8, and between σ_i , i=1,2, and σ_j , j=3,4,5,6,7, the correlation coefficients are each 0.5. The correlation coefficients between each $\sigma_i \neq \sigma_j$, $i,j \in \{3,4,5,6,7\}$, are set to 0.8. The applied force f is independent of the yield stresses and is distributed log-normally with mean 1000kN and standard deviation 400kN. The structure will fail if any member exceeds the associated yield stress, i.e., for each member, we can define the following random limit state function:

1197 (6.2)
$$g_i(x; f, \sigma) = \frac{f}{c_i x^i} - \sigma_i, \quad i \in \{1, 2, \dots, 7\},$$

where the fixed parameters c_i depend on the geometry and the loads. For this structure, $c_{\{i=1,2\}} = 1/(2\sqrt{3})$ and $c_{\{i=3,\dots,7\}} = 1/\sqrt{3}$.

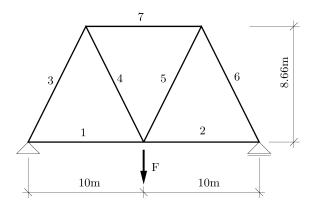


FIGURE 6.2. Definition of the geometry and the load applied to the truss.

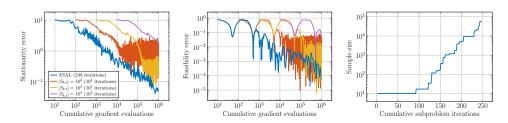


FIGURE 6.3. Results of running Algorithm 5.1 on the truss optimization problem (6.3). Notice that ASAL achieves the lowest average stationary and feasibility errors while simultaneously requiring the smallest number of iterations. To generate these results, we used the algorithm parameters $\theta_q = 0.99$, $\alpha = 0.01$, $\eta = 1.0$, and primal tolerance sequence $\tau_k = 10/k$.

We pose the following stochastic optimization problem:

minimize
$$\frac{1}{7\alpha}\mathbb{E}\left[\ln\left(\sum_{i=1}^{7}\exp(\alpha g_i(x))\right)\right],$$
 subject to $A \leq x \leq B, \quad \langle 1, x \rangle \leq C,$

1200

1206

1208

1209

1210

1211

1212

1213

1214

where $\alpha=1,\ A=1\times 10^4 \mathrm{mm^2},\ B=5\times 10^4 \mathrm{mm^2},\ \mathrm{and}\ C=15\times 10^4 \mathrm{mm^2}$ are user-defined parameters. The components of the optimal solution are estimated to be

1204 (6.4)
$$x_{\{i=1,2\}} = 4.342 \times 10^4 \text{mm}^2 \text{ and } x_{\{i=3,\dots,7\}} = 1.263 \times 10^4 \text{mm}^2.$$

To solve this problem, we use ASAL with $\theta_g = 0.99$ and compare its performance to the stochastic augmented Lagrangian method with fixed sample sizes under a 1 million cumulative sample budget. In each experiment, we use the penalty and step size values $\alpha = 0.01$ and $\eta = 1.0$. Figure 6.3 documents our findings. Notice that, even though it used less than 25% of the total iterations, the stationarity and feasibility errors from ASAL (248 iterations) are significantly lower after the sample budget expires than the best-performing fixed sample size algorithm (1000 iterations).

6.3. Optimal design of a heat sink. We close with a non-convex optimization problem of engineering interest. In this final experiment, we consider the optimal design of a heat sink within a hypothetical square domain $\Omega = (0,1)^2$ with a stochastic

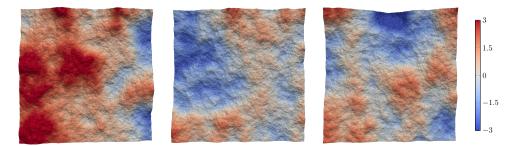


FIGURE 6.4. Three independent realizations of the Gaussian random field $f(\mathbf{x})$ generated by solving (6.5) on the square domain $\Omega = (0,1)^2$.

heat source f(x), $x \in \Omega$, described by a spatial Gaussian random field with Mátern covariance. More specifically, we follow [57,58] and define

1217 (6.5)
$$-\kappa^2 \Delta f + f = \mathcal{W} \text{ in } \Omega, \quad \nabla f \cdot \boldsymbol{n} = 0 \text{ on } \partial \Omega,$$

where W is spatial additive white Gaussian noise, $\kappa > 0$ is a correlation length parameter, and n denotes the outward-facing unit normal vector field on $\partial\Omega$. Mátern random fields can be used to model various random spatial phenomena [44, 45, 57], which makes them reasonable for modeling the heat source in this example. Figure 6.4 depicts three representative solutions to (6.5) for the reader's interest.

We use the two-field filtered density approach to topology optimization [79, Section 3.1.2] to formulate the optimal heat sink design problem. The goal is to find a material distribution $0 \le \rho \le 1$, where zero indicates no material, and one indicates the complete presence of material, that induces the smallest thermal compliance, $\int_{\Omega} uf \, dx$, in expectation. In the aforestated expression, the temperature distribution u is determined by ρ and f through the heat diffusion equation $-\operatorname{div} r(\widetilde{\rho}) \nabla u = f$, where $\widetilde{\rho}$ is a regularized (filtered) distribution function [22, 52] and $r(\widetilde{\rho}) > 0$ is a thermal conductivity model. In this work, we use the well-known (modified) solid isotropic material penalization (SIMP) model $r(\widetilde{\rho}) = \rho_- + \widetilde{\rho}^3 (1 - \rho_-)$, where $0 < \rho_- \ll 1$ is a nominal thermal diffusivity constant assigned to void regions in order to prevent the stiffness matrix from becoming singular [2].

The full problem formulation is written as follows:

1235 (6.6a)
$$\min_{\rho \in L^2(\Omega), \ u \in H^1(\Omega)} \left\{ \widehat{F}(\rho, u) := \mathbb{E}\left[\int_{\Omega} u f \, \mathrm{d}\boldsymbol{x} \right] \right\},$$

1236 subject to the constraints

$$\begin{cases}
-\epsilon^2 \Delta \widetilde{\rho} + \widetilde{\rho} = \rho & \text{in } \Omega, \quad \nabla \widetilde{\rho} \cdot \boldsymbol{n} = 0 \text{ on } \partial \Omega, \\
-\operatorname{div} \left(r(\widetilde{\rho}) \nabla u \right) = f & \text{in } \Omega, \quad u = 0 \text{ on } \Gamma_0, \quad \nabla u \cdot \boldsymbol{n} = 0 \text{ on } \partial \Omega \setminus \Gamma_0, \\
\int_{\Omega} \rho(\boldsymbol{x}) d\boldsymbol{x} \leq \gamma |\Omega|, \quad \text{and } 0 \leq \rho \leq 1 \text{ in } \Omega,
\end{cases}$$

where $0 < \gamma < 1$ is the *volume fraction*, which constrains the fraction of the domain occupied by design, and $\epsilon > 0$ is a *length scale* for the final design. The boundary conditions and solution to the optimization problem (6.6) with $\rho_{-} = 10^{-3}$, $\gamma = 0.5$, $\epsilon = 0.01$, $\kappa = 0.2$ are depicted in Figure 6.5.

To remove the PDE constraints from the optimization problem, we employ a reduced space formulation, often referred to in the literature as a nested formulation [8], which can be written as

1245 (6.7a)
$$\min_{\rho \in L^2(\Omega)} \left\{ F(\rho) := \mathbb{E} \left[\int_{\Omega} u\left(\widetilde{\rho}\left(\rho\right) \right) f \, \mathrm{d}\boldsymbol{x} \right] \right\},$$

1246 subject to the constraints

1255

1271

1282

1247 (6.7b)
$$\int_{\Omega} \rho(\boldsymbol{x}) d\boldsymbol{x} \leq \gamma |\Omega|, \text{ and } 0 \leq \rho \leq 1 \text{ in } \Omega.$$

In this formulation, it is understood that the temperature field $u = u(\widetilde{\rho}(\rho))$ solves the state equation

1250 (6.8)
$$-\operatorname{div}(r(\widetilde{\rho})\nabla u) = f \text{ in } \Omega, \quad u = 0 \text{ on } \Gamma_0, \quad \nabla u \cdot \boldsymbol{n} = 0 \text{ on } \partial \Omega \setminus \Gamma_0,$$

and the filtered density $\tilde{\rho} = \tilde{\rho}(\rho)$ solves the screened Poisson equation,

1252 (6.9)
$$-\epsilon^2 \Delta \widetilde{\rho} + \widetilde{\rho} = \rho \text{ in } \Omega, \quad \nabla \widetilde{\rho} \cdot \boldsymbol{n} = 0 \text{ on } \partial \Omega.$$

Since the inequality constraint in (6.7b) is always active, it is replaced by an equality constraint that our ASAL algorithm can handle. The gradients of the reduced objective function in (6.7a) are computed with FEM-discretized representations of the temperature u and filtered density $\tilde{\rho}$ using standard adjoint analysis techniques [8]. Finally, $L^2(\Omega)$ projections are used to enforce the box constraints found in (6.7b).

For comparison, Figure 6.5 also depicts a reference solution to (6.6) corresponding to the (deterministic) uniform heat field $f \equiv 1$. Close examination reveals significant differences between the designs with deterministic and stochastic inputs. The deterministic case results in an organic tree-like structure that aims to transfer the heat generated at any point in the computational domain using the shortest possible way to the Dirichlet boundary with zero temperature. The design does not depend on the magnitude of the heat source, and any constant input will result in the same material distribution if the initial material distribution is in the vicinity of the local solution. On the other hand, due to the oscillatory nature of the stochastic input, the heat source term can take positive and negative values. Such input distribution allows the optimization process to balance the heat transfer locally without linking the local subdomain directly to the boundary with a fixed temperature. Thus, the role of the closed loops of material appearing in the design with stochastic input is to establish a local heat equilibrium. In this case, the global tree-like structure transfers only the excess heat, which cannot be balanced locally.

In this experiment, we use ASAL with $\theta_g = 2$ and compare its performance to stochastic augmented Lagrangian with fixed sample sizes, $|S_{k,t}| = 10^j$, j = 1, 2, 3, under a 10^5 cumulative sample budget. In each execution, we use the step size values $\alpha = 0.1$ and $\eta = 2.0$. Figure 6.6 documents our findings. ASAL achieves the lowest combined average stationary and feasibility errors while requiring less than 20% of the iterations of the best-performing fixed sample size run ($|S_{k,t}| = 10^2$). Although the average feasibility errors with ASAL and this fixed sample size run are similar, the variance of the fixed sample size run is much greater. Finally, the average stationarity error for the best-performing fixed sample size run is significantly larger than the average stationarity error with ASAL.

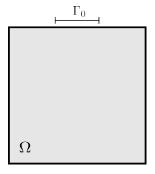
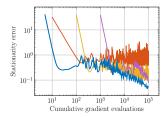






FIGURE 6.5. Left: Depiction of the subsets of the domain boundary $\partial\Omega$ for the boundary conditions for the heat field u in (6.5). We define u=0 on Γ_0 and $\nabla u \cdot \mathbf{n}=0$ on $\partial\Omega \setminus \Gamma_0$. Middle: Reference density field $\tilde{\rho}$ for the solution of the deterministic thermal compliance optimization problem (6.6) with $f\equiv 1$ everywhere in Ω . Right: The filtered density $\tilde{\rho}$ for the solution of the expected value thermal compliance optimization problem (6.6) with f given by (6.5). The presence of closed-loop branches in the optimal solution on the right indicates a preference for balancing the heat locally and transferring only the excess unbalanced heat to the external environment through Γ_0 .



1284

1285

1286

1287 1288

1289

1290

1292

1293

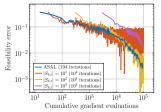
1294

1295

1296

1297

1298 1299



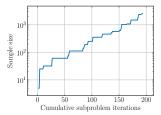


FIGURE 6.6. Results of running ASAL (Algorithm 5.1) on the topology optimization problem (6.6). Notice that ASAL achieves the lowest combined average stationary and feasibility errors while simultaneously requiring the smallest number of iterations. Indeed, the average stationarity error for the best-performing fixed sample size run, $|S_{k,t}| = 10^j$, j = 1, 2, 3, ends up larger than the average stationarity error with ASAL. Moreover, although the average feasibility errors of ASAL and the best-performing fixed sample size run are similar, the variance for the fixed sample size run is much greater. To generate these results, we used the algorithm parameters $\theta_g = 2$, $\alpha = 0.1$, $\eta = 2.0$, and primal tolerance sequence $\tau_k = 1/k$.

7. Final Remarks. Motivated by a growing interest in developing optimization algorithms for constrained stochastic optimization problems, we introduced a framework that combines augmented Lagrangian methods with adaptive sampling techniques. In our framework, we employed stochastic solvers for the subproblems and imposed stochastic tolerance criteria for the inexact solutions. We analyzed various theoretical tolerance conditions and designed a practical test. To establish convergence results, we first showed that our framework is equivalent to an inexact gradient descent algorithm on the Moreau envelope. Second, we showed sublinear convergence in the outer iterations when f is convex and linear convergence when f is strongly convex with $\mathcal{X} = \mathbb{R}^n$. We also analyzed the worst-case expected work complexity of our approach in terms of the number of gradient evaluations required to obtain an ϵ -accurate solution. For convex f and compact \mathcal{X} , we showed $\mathcal{O}(\epsilon^{-3-\delta})$ complexity where $\delta > 0$ is a user-defined parameter. This result improves to $\mathcal{O}(\epsilon^{-2})$ when the penalty parameter $\alpha = \mathcal{O}(\epsilon^{-1})$. If f is strongly convex and $\mathcal{X} = \mathbb{R}^n$, we proved $\mathcal{O}(\epsilon^{-1}\log(1/\epsilon))$ complexity.

To evaluate our framework's practical performance, we tested it on a constrained machine learning problem and in engineering applications. Here, we observed that our method minimizes the objective function more efficiently and reaches a feasible solution in a more stable manner than benchmark stochastic approximation algorithms.

Although our analysis holds for any penalty parameter $\alpha > 0$, this parameter should be tuned for optimal performance in practice. The other main hyperparameters are the step size $\eta > 0$ for the inner problems, and the subproblem tolerance values $\tau_k > 0$. Since tuning is computationally expensive, it would be helpful to develop methods that adaptively select these hyper-parameters in order to further improve the practical efficacy of our adaptive sampling framework. Two other natural extensions would be generalizing our methods to include nonlinear constraints and chance constraints.

Appendix A. Proof of Proposition 4.7.

Proof. Due to the strong convexity of f, (2.12) has a unique optimal solution, denoted as $x(\lambda)$. Using [40, Corollary 4.5.3], we can show that $q(\lambda)$ is differentiable and

1314 (A.1)
$$\nabla q(\lambda) = Ax(\lambda) - b.$$

Moreover, from the optimality conditions of (2.12), we have

1316 (A.2)
$$\nabla f(x(\lambda)) - A^T \lambda = 0.$$

1317 Let $\lambda_1, \lambda_2 \in \mathbb{R}^m$. Consider,

1318
$$\langle \nabla q(\lambda_2) - \nabla q(\lambda_1), \lambda_2 - \lambda_1 \rangle = \langle A(x(\lambda_2) - x(\lambda_1)), \lambda_2 - \lambda_1 \rangle$$
1319
$$= \langle x(\lambda_2) - x(\lambda_1), A^T(\lambda_2 - \lambda_1) \rangle$$
1320
$$= \langle x(\lambda_2) - x(\lambda_1), \nabla f(x(\lambda_2)) - \nabla f(x(\lambda_1)) \rangle$$
1321
$$\geq \frac{1}{\mu + L} \|\nabla f(x(\lambda_2)) - \nabla f(x(\lambda_1))\|^2$$
1322
$$= \frac{1}{\mu + L} \|A^T(\lambda_2 - \lambda_1)\|^2$$
1323
$$\geq \frac{\lambda_{\min}(AA^T)}{\mu + L} \|\lambda_2 - \lambda_1\|^2,$$

where the first equality is due to (A.1), the second and the third equalities are due to (A.2), and the first inequality is due to [63, Theorem 2.1.11]. Therefore, using [63, Theorem 2.1.9], we can claim that $q(\lambda)$ is strongly convex with parameter $\frac{\sigma}{L+\mu}$.

Appendix B. Logistic regression with multiple disparate impact constraints, australian dataset.

We consider problem (6.1) with australian classification data set from the LIBSVM collection [26]. The data set has N=690 rows, and the dimension of the problem is n=14. Considering the budget of cumulative gradient evaluations as 200N, and the fixed hyperparameters as $\theta_g=0.99$, $\nu_1=0.5$, $s_1=0.1N$, $s_{\min}=0.1N$, we compare three separately-tuned fixed-batch-size implementations of ASAL using 10%, 20%, and 50% of the data set size. We tune τ_0 , α and the step size η using the same procedure described in Subsection 6.1 with the sets of $\tau_0=10^{i-1}$, $\eta=10^{j-5}$, and $\alpha=10^{j-4}$, where i=0,1,2,3,4,5 and j=0,1,2,3,4,5,6.

For each algorithm, we select the run with the smallest average objective function value in the final 10 inner iterations among all runs whose minimum feasibility error in the final 50 inner iterations is less than the feasibility tolerance 10^{-3} . These

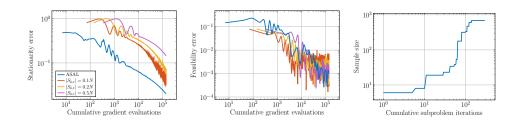


FIGURE B.1. Results of running Algorithm 5.1 on the constrained logistic regression problem (6.1) with the australian classification data set. Notice that ASAL achieves the lowest average stationarity errors while matching the minimal feasibility error of the three baseline algorithms. To generate these results, we used the algorithm parameters $\theta_g = 0.99$ and individually tuned α , η , and τ_0 .

values (i.e., 10, 50, and 10^{-3} , respectively) are slightly different than the values given in Subsection 6.1 to ensure that the best combinations of hyperparameter values correspond to a more stable set of runs. Because of the same reason, we restrict $\alpha = 10^{-1}$ for the ASAL algorithm while tuning, as we observe this value results in choosing runs that show a good balance between stationarity and feasibility errors. The comparison of the algorithms is given in Figure B.1. Similar to Subsection 6.1, we observe that ASAL and each of the three baseline algorithms achieve a similar minimal feasibility error (around feasibility tolerance 10^{-3}) and that ASAL performs better than the three baseline algorithms when it comes to stationarity error.

1350 REFERENCES

1341

1342

1344

1346

1347

1348 1349

1351

1352

1353 1354

1356

1357

1358

1359

 $1360 \\ 1361$

1362

 $1363 \\ 1364$

1365 1366

1367

1368 1369

1370

1371

1372

1373

1374

1375

1376

- [1] Stochastic optimization using a trust-region method and random models, Mathematical Programming, 169 (2018), pp. 447–487.
- [2] E. ANDREASSEN, A. CLAUSEN, M. SCHEVENELS, B. S. LAZAROV, AND O. SIGMUND, Efficient topology optimization in MATLAB using 88 lines of code, Structural and Multidisciplinary Optimization, 43 (2011), pp. 1–16.
- [3] E. Andreassen, B. S. Lazarov, and O. Sigmund, Design of manufacturable 3D extremal elastic microstructure, Mechanics of Materials, 69 (2014), pp. 1 10, https://doi.org/http://dx.doi.org/10.1016/j.mechmat.2013.09.018, http://www.sciencedirect.com/science/article/pii/S0167663613002093.
- [4] S. BAROCAS AND A. D. SELBST, Big data's disparate impact, Calif. L. Rev., 104 (2016), p. 671. [5] F. BASTIN, C. CIRILLO, AND P. L. TOINT, An adaptive monte carlo algorithm for computing
- [5] F. Bastin, C. Cirillo, and P. L. Toint, An adaptive monte carlo algorithm for computing mixed logit estimators, Computational Management Science, 3 (2006), pp. 55–79.
- [6] H. H. BAUSCHKE AND P. L. COMBETTES, Convex Analysis and Monotone Operator Theory in Hilbert Spaces, Springer, 2nd ed., 2017.
- [7] F. Beiser, B. Keith, S. Urbainczyk, and B. Wohlmuth, Adaptive sampling strategies for risk-averse stochastic optimization with constraints, IMA Journal of Numerical Analysis, (2023), https://doi.org/10.1093/imanum/drac083. drac083.
- [8] M. P. Bendsoe and O. Sigmund, Topology optimization: theory, methods, and applications, Springer Science & Business Media, 2003.
- [9] A. S. Berahas, R. Bollapragada, and B. Zhou, An adaptive sampling sequential quadratic programming method for equality constrained stochastic optimization, arXiv preprint arXiv:2206.00712, (2022).
- [10] A. S. Berahas, F. E. Curtis, D. Robinson, and B. Zhou, Sequential quadratic optimization for nonlinear equality constrained stochastic optimization, SIAM Journal on Optimization, 31 (2021), pp. 1352–1379.
- [11] D. Bertsekas, Convex optimization algorithms, Athena Scientific, 2015.
- [12] D. BERTSEKAS, A. NEDIC, AND A. OZDAGLAR, Convex analysis and optimization, vol. 1, Athena Scientific, 2003.
- 1379 [13] D. P. Bertsekas, Constrained optimization and Lagrange multiplier methods, Academic press, 1380 2014.

- 1381 [14] E. G. Birgin, Augmented lagrangian method with nonmonotone penalty parameters for con-1382 strained optimization, Computational Optimization and Applications, 51 (2012), pp. 941– 1383 965.
- 1384 [15] J. Blanchet, C. Cartis, M. Menickelly, and K. Scheinberg, Convergence rate analysis of 1385 a stochastic trust-region method via supermartingales, INFORMS journal on optimization, 1386 1 (2019), pp. 92–119.
- 1387 [16] R. BOLLAPRAGADA, R. BYRD, AND J. NOCEDAL, Adaptive sampling strategies for stochastic optimization, SIAM Journal on Optimization, 28 (2018), pp. 3312–3343.
- 1389 [17] R. BOLLAPRAGADA, R. H. BYRD, AND J. NOCEDAL, Exact and inexact subsampled newton methods for optimization, IMA Journal of Numerical Analysis, 39 (2019), pp. 545–578.
- 1391 [18] R. BOLLAPRAGADA, J. NOCEDAL, D. MUDIGERE, H.-J. SHI, AND P. T. P. TANG, A progress-1392 ive batching l-bfgs method for machine learning, in International Conference on Machine 1393 Learning, PMLR, 2018, pp. 620–629.
- 1394 [19] L. BOTTOU, F. E. CURTIS, AND J. NOCEDAL, Optimization methods for large-scale machine learning, SIAM review, 60 (2018), pp. 223–311.
- 1396 [20] S. BOYD, S. P. BOYD, AND L. VANDENBERGHE, Convex optimization, Cambridge university 1397 press, 2004.
- 1398 [21] S. BOYD, N. PARIKH, E. CHU, B. PELEATO, J. ECKSTEIN, ET AL., Distributed optimization and statistical learning via the alternating direction method of multipliers, Foundations and Trends® in Machine learning, 3 (2011), pp. 1–122.

1402

1403

1414

 $\begin{array}{c} 1415 \\ 1416 \end{array}$

1417

 $1420 \\ 1421$

- [22] T. E. Bruns and D. A. Tortorelli, Topology optimization of non-linear elastic structures and compliant mechanisms, Computer methods in applied mechanics and engineering, 190 (2001), pp. 3443–3459.
- 1404 [23] R. H. Byrd, G. M. Chin, J. Nocedal, and Y. Wu, Sample size selection in optimization methods for machine learning, Mathematical programming, 134 (2012), pp. 127–155.
- 1406 [24] A. CARLON, L. ESPATH, R. LOPEZ, AND R. TEMPONE, Multi-iteration stochastic optimizers, arXiv preprint arXiv:2011.01718, (2020).
- [25] C. Cartis and K. Scheinberg, Global convergence rate analysis of unconstrained optimization
 methods based on probabilistic models, Mathematical Programming, 169 (2018), pp. 337–
 375.
- [26] C.-C. CHANG AND C.-J. LIN, Libsum: A library for support vector machines, ACM Trans.
 Intell. Syst. Technol., 2 (2011), https://doi.org/10.1145/1961189.1961199, https://doi.org/10.1145/1961189.1961199.
 - [27] S. CHEN, W. CHEN, AND S. LEE, Level set based robust shape and topology optimization under random field uncertainties, Structural and Multidisciplinary Optimization, 41 (2010), pp. 507–524, https://doi.org/10.1007/s00158-009-0449-2, http://dx.doi.org/10.1007/s00158-009-0449-2.
- 1418 [28] G. CORNUEJOLS AND R. TÜTÜNCÜ, Optimization methods in finance, vol. 5, Cambridge Uni-1419 versity Press, 2006.
 - [29] F. E. CURTIS, N. I. GOULD, H. JIANG, AND D. P. ROBINSON, Adaptive augmented lagrangian methods: algorithms and practical numerical experience, Optimization Methods and Software, 31 (2016), pp. 157–186.
- 1423 [30] F. E. CURTIS, H. JIANG, AND D. P. ROBINSON, An adaptive augmented lagrangian method for large-scale constrained optimization, Mathematical Programming, 152 (2015), pp. 201–245.
- 1425 [31] F. E. Curtis, M. J. O'Neill, and D. P. Robinson, Worst-case complexity of an sqp method for nonlinear equality constrained stochastic optimization, arXiv preprint arXiv:2112.14799, (2021).
- [32] F. E. Curtis, D. P. Robinson, and B. Zhou, Inexact sequential quadratic optimization for minimizing a stochastic objective function subject to deterministic nonlinear equality constraints, arXiv preprint arXiv:2107.03512, (2021).
- 1431 [33] J. ECKSTEIN AND W. YAO, Augmented lagrangian and alternating direction methods for convex optimization: A tutorial and some illustrative computational results, RUTCOR Research Reports, 32 (2012), p. 44.
- 1434 [34] L. ESPATH, S. KRUMSCHEID, R. TEMPONE, AND P. VILANOVA, On the equivalence of different adaptive batch size selection strategies for stochastic gradient descent methods, 2021, https://arxiv.org/abs/2109.10933.
- 1437 [35] M. P. FRIEDLANDER AND M. SCHMIDT, Hybrid deterministic-stochastic methods for data fitting, 1438 SIAM Journal on Scientific Computing, 34 (2012), pp. A1380–A1405.
- 1439 [36] S. GANESH AND F. NOBILE, Gradient-based optimisation of the conditional-value-at-risk using the multi-level Monte Carlo method, arXiv preprint arXiv:2210.03485, (2022).
- 1441 [37] V. Guigues, On the strong concavity of the dual function of an optimization problem, arXiv preprint arXiv:2006.16781, (2020).

- 1443 [38] V. Guigues, Inexact stochastic mirror descent for two-stage nonlinear stochastic programs, 1444 Mathematical Programming, 187 (2021), pp. 533–577.
- 1445 [39] M. R. Hestenes, Multiplier and gradient methods, Journal of optimization theory and applications, 4 (1969), pp. 303–320.
- 1447 [40] J.-B. HIRIART-URRUTY AND C. LEMARÉCHAL, Convex analysis and minimization algorithms I:
 1448 Fundamentals, vol. 305, Springer science & business media, 2013.
- [41] M. Jansen, B. Lazarov, M. Schevenels, and O. Sigmund, On the similarities between micro/nano lithography and topology optimization projection methods, Structural and Multidisciplinary Optimization, 48 (2013), pp. 717-730, https://doi.org/10.1007/s00158-013-0941-6.
- [42] Z. Jiang, C. Liu, Y. M. Lee, C. Hegde, S. Sarkar, and D. Jiang, The stochastic augmented lagrangian method for domain adaptation, Knowledge-Based Systems, 235 (2022),
 p. 107593.
- [43] M. Kang, M. Kang, and M. Jung, Inexact accelerated augmented lagrangian methods, Computational Optimization and Applications, 62 (2015), pp. 373-404.
 [44] B. Keith, U. Khristenko, and B. Wohlmuth, A fractional PDE model for turbulent velocity
 - [44] B. KEITH, U. KHRISTENKO, AND B. WOHLMUTH, A fractional PDE model for turbulent velocity fields near solid walls, Journal of Fluid Mechanics, 916 (2021), p. A21.

1460

1461

1462

1470

1471

1474

1475

1476

1477

1481

1482

1483

1484

1485

- [45] U. Khristenko, A. Constantinescu, P. L. Tallec, J. T. Oden, and B. Wohlmuth, A statistical framework for generating microstructures of two-phase random materials: application to fatigue analysis, Multiscale Modeling & Simulation, 18 (2020), pp. 21–43.
- [46] U. Khristenko, L. Scarabosio, P. Swierczynski, E. Ullmann, and B. Wohlmuth, Analysis
 of boundary effects on PDE-based sampling of Whittle-Matérn random fields, SIAM/ASA
 Journal on Uncertainty Quantification, 7 (2019), pp. 948-974.
- [47] A. KODAKKAL, B. KEITH, U. KHRISTENKO, A. APOSTOLATOS, K.-U. BLETZINGER,
 B. WOHLMUTH, AND R. WÜCHNER, Risk-averse design of tall buildings for uncertain
 wind conditions, Computer Methods in Applied Mechanics and Engineering, 402 (2022),
 p. 115371.
 - [48] D. P. KOURI AND A. SHAPIRO, Optimization of PDEs with uncertain inputs, Frontiers in PDE-Constrained Optimization, (2018), pp. 41–81.
- 1472 [49] D. P. KOURI AND T. M. SUROWIEC, A primal-dual algorithm for risk minimization, Mathem-1473 atical Programming, 193 (2022), pp. 337–363.
 - [50] G. LAN AND R. D. MONTEIRO, Iteration-complexity of first-order augmented lagrangian methods for convex programming, Mathematical Programming, 155 (2016), pp. 511–547.
 - [51] G. LAN AND Z. ZHOU, Algorithms for stochastic optimization with function or expectation constraints, Computational Optimization and Applications, 76 (2020), pp. 461–498.
- 1478 [52] B. S. LAZAROV AND O. SIGMUND, Filters in topology optimization based on Helmholtz-type differential equations, International Journal for Numerical Methods in Engineering, 86 (2011), pp. 765–781.
 - [53] B. S. LAZAROV, F. WANG, AND O. SIGMUND, Length scale and manufacturability in density-based topology optimization, Archive of Applied Mechanics, 86 (2016), pp. 189–218, https://doi.org/10.1007/s00419-015-1106-4, http://dx.doi.org/10.1007/s00419-015-1106-4.
 - [54] Z. LI, P.-Y. CHEN, S. LIU, S. LU, AND Y. XU, Rate-improved inexact augmented lagrangian method for constrained nonconvex optimization, in International Conference on Artificial Intelligence and Statistics, PMLR, 2021, pp. 2170–2178.
- 1487 [55] Z. LI, P.-Y. CHEN, S. LIU, S. LU, AND Y. XU, Stochastic inexact augmented lagrangian method 1488 for nonconvex expectation constrained optimization, arXiv preprint arXiv:2212.09513, 1489 (2022).
- [56] Z. Li and Y. Xu, Augmented lagrangian-based first-order methods for convex-constrained programs with weakly convex objective, INFORMS Journal on Optimization, 3 (2021), pp. 373–397.
- 1493 [57] F. LINDGREN, D. BOLIN, AND H. RUE, The SPDE approach for Gaussian and non-Gaussian 1494 fields: 10 years and still running, Spatial Statistics, (2022), p. 100599.
- [58] F. LINDGREN, H. RUE, AND J. LINDSTRÖM, An explicit link between Gaussian fields and gaussian markov random fields: The stochastic partial differential equation approach, Journal of the Royal Statistical Society: Series B (Statistical Methodology), 73 (2011), pp. 423–498.
- [59] Y. Liu, F. Shang, H. Liu, L. Kong, L. Jiao, and Z. Lin, Accelerated variance reduction stochastic admm for large-scale machine learning, IEEE Transactions on Pattern Analysis and Machine Intelligence, 43 (2021), pp. 4242–4255, https://doi.org/10.1109/TPAMI.2020.
 3000512.
- 1502 [60] J.-J. MOREAU, Proximité et dualité dans un espace hilbertien, Bulletin de la Société mathématique de France, 93 (1965), pp. 273–299.
- 1504 [61] S. NA, M. ANITESCU, AND M. KOLAR, An adaptive stochastic sequential quadratic program-

- $\begin{array}{ll} 1505 & \textit{ming with differentiable exact augmented lagrangians}, \text{Mathematical Programming, (2022)}, \\ 1506 & \text{pp. 1-71}. \end{array}$
- [62] S. NA, M. ANITESCU, AND M. KOLAR, Inequality constrained stochastic nonlinear optimization via active-set sequential quadratic programming, Mathematical Programming, (2023), pp. 1–75.
- [63] Y. Nesterov, Introductory lectures on convex optimization: A basic course, vol. 87, Springer
 Science & Business Media, 2003.
 - [64] Y. Nesterov, Lectures on convex optimization, vol. 137, Springer, 2018.

1519

1520

1522

1523

1530

1531

1532

1533

1534 1535

1536

1542

1543

1544

1545

1546

1547

- 1513 [65] H. OUYANG, N. HE, L. TRAN, AND A. GRAY, Stochastic alternating direction method of mul-1514 tipliers, in International conference on machine learning, PMLR, 2013, pp. 80–88.
- 1515 [66] N. Parikh, S. Boyd, et al., *Proximal algorithms*, Foundations and trends® in Optimization, 1516 1 (2014), pp. 127–239.
- 1517 [67] R. PASUPATHY, P. GLYNN, S. GHOSH, AND F. S. HASHEMI, On sampling rates in simulation-1518 based recursions, SIAM Journal on Optimization, 28 (2018), pp. 45–73.
 - [68] C. PHELPS, J. O. ROYSET, AND Q. GONG, Optimal control of uncertain systems using sample average approximations, SIAM Journal on Control and Optimization, 54 (2016), pp. 1–29.
 - [69] C. Planiden and X. Wang, Strongly convex functions, moreau envelopes, and the generic nature of convex functions with strong minimizers, SIAM Journal on Optimization, 26 (2016), pp. 1341–1364.
- 1524 [70] M. J. D. POWELL, A method for nonlinear constraints in minimization problems, Optimization, 1525 (1969), pp. 283–298.
- 1526 [71] R. T. ROCKAFELLAR, Augmented lagrange multiplier functions and duality in nonconvex pro-1527 gramming, SIAM Journal on Control, 12 (1974), pp. 268–285.
- 1528 [72] R. T. ROCKAFELLAR, Augmented lagrangians and applications of the proximal point algorithm
 1529 in convex programming, Mathematics of operations research, 1 (1976), pp. 97–116.
 - [73] R. T. ROCKAFELLAR AND J. O. ROYSET, On buffered failure probability in design and optimization of structures, Reliability Engineering & System Safety, 95 (2010), pp. 499 510, https://doi.org/http://dx.doi.org/10.1016/j.ress.2010.01.001.
 - [74] F. ROOSTA-KHORASANI AND M. W. MAHONEY, Sub-sampled newton methods, Mathematical Programming, 174 (2019), pp. 293–326.
 - [75] J. O. ROYSET, Risk-adaptive approaches to learning and decision making: A survey, 2022, https://doi.org/10.48550/ARXIV.2212.00856, https://arxiv.org/abs/2212.00856.
- 1537 [76] J. O. ROYSET AND R. SZECHTMAN, Optimal budget allocation for sample average approximation, 1538 Operations Research, 61 (2013), pp. 762–776.
- 1539 [77] M. F. Sahin, A. Alacaoglu, F. Latorre, V. Cevher, et al., An inexact augmented lag-1540 rangian framework for nonconvex optimization with nonlinear constraints, Advances in 1541 Neural Information Processing Systems, 32 (2019).
 - [78] A. SHAPIRO, D. DENTCHEVA, AND A. RUSZCZYNSKI, Lectures on stochastic programming: Modeling and theory, SIAM, 2021.
 - [79] O. SIGMUND AND K. MAUTE, Topology optimization approaches: A comparative review, Structural and Multidisciplinary Optimization, 48 (2013), pp. 1031–1055.
 - [80] M. Sion, On general minimax theorems., Pacific Journal of mathematics, 8 (1958), pp. 171–176.
 - [81] T. Suzuki, Stochastic dual coordinate ascent with alternating direction method of multipliers, in International Conference on Machine Learning, PMLR, 2014, pp. 736–744.
- 1549 [82] S. WRIGHT AND B. RECHT, Optimization for Data Analysis, Cambridge University Press, 2022.
- [83] Y. XIE, R. BOLLAPRAGADA, R. BYRD, AND J. NOCEDAL, Constrained and composite optimization via adaptive sampling methods, IMA Journal of Numerical Analysis, (2023), p. drad020.
- 1553 [84] Y. XIE AND U. V. SHANBHAG, Si-admm: A stochastic inexact admm framework for stochastic tonvex programs, IEEE Transactions on Automatic Control, 65 (2019), pp. 2355–2370.
- 1555 [85] Y. Xu, Primal-dual stochastic gradient method for convex programs with many functional constraints, SIAM Journal on Optimization, 30 (2020), pp. 1664–1692.
- 1557 [86] Y. Xu, Iteration complexity of inexact augmented lagrangian methods for constrained convex programming, Mathematical Programming, 185 (2021), pp. 199–244.
- 1559 [87] M. B. Zafar, I. Valera, M. Gomez-Rodriguez, and K. P. Gummadi, Fairness constraints: A flexible approach for fair classification, J. Mach. Learn. Res., 20 (2019), pp. 1–42.
- 1561 [88] S. ZHENG AND J. T. KWOK, Fast-and-light stochastic ADMM, in International Joint Conference on Artificial Intelligence, 2016.
- 1563 [89] W. Zhong and J. T. Kwok, Fast stochastic alternating direction method of multipliers, in 1564 International conference on machine learning, PMLR, 2014, pp. 46–54.
- [90] X. Zhou, On the fenchel duality between strong convexity and lipschitz continuous gradient,
 2018, https://arxiv.org/abs/1803.06573.