

An Adaptive Sampling Sequential Quadratic Programming Method for Equality Constrained Stochastic Optimization

Albert S. Berahas^{*†} Raghu Bollapragada[‡] Baoyu Zhou[§]

June 30, 2022

Abstract

This paper presents a methodology for using varying sample sizes in sequential quadratic programming (SQP) methods for solving equality constrained stochastic optimization problems. The first part of the paper deals with the delicate issue of dynamic sample selection in the evaluation of the gradient in conjunction with inexact solutions to the SQP subproblems. Under reasonable assumptions on the quality of the gradient approximations employed and the accuracy of the solutions to the SQP subproblems, we establish global convergence results for the proposed method. Motivated by these results, the second part of the paper describes a practical adaptive inexact stochastic sequential quadratic programming (PAIS-SQP) method. We propose criteria for controlling the sample size and the accuracy in the solutions of the SQP subproblems based on the variance estimates obtained as the optimization progresses. Finally, we demonstrate the practical performance of the method on a subset of the CUTE problems and constrained classification tasks.

1 Introduction

We consider stochastic optimization problems with deterministic equality constraints, and design, analyze and implement a stochastic adaptive sampling algorithm based on the sequential quadratic programming (SQP) paradigm. Optimization problems of this form arise in a plethora of real-world applications, including but not limited to computer vision [32], optimal control [8], network optimization [7], partial differential equation optimization [33], deep learning [36], and reinforcement learning [1].

The majority of the methods developed for solving constrained stochastic optimization problems, i.e., optimization problems with deterministic constraints and stochastic objective functions, are based on the penalty method approach [19, 26, 29, 32]. Such methods

^{*}Corresponding author.

[†]Dept. of Industrial and Operations Engineering, University of Michigan. (albertberahas@gmail.com)

[‡]Operations Research and Industrial Engineering Prog., UT Austin. (raghu.bollapragada@utexas.edu)

[§]Dept. of Industrial and Systems Engineering, Lehigh University. (baoyu.zhou@lehigh.edu)

transform the given constrained problem into an unconstrained problem by penalizing constraint violation in the objective function, and then apply classical unconstrained stochastic optimization methods on the modified objective. While these methods are well studied, recently a number of algorithms endowed with sound theoretical guarantees and superior empirical performance have been proposed [5, 28]. In [5], a stochastic SQP method with adaptive step size selection for the fully stochastic regime is proposed, and, in [28], a stochastic line search SQP method with adaptive gradient accuracy is proposed. Several other extensions of these methods have been developed, e.g., relaxing constraint assumptions [4], employing inexact computations [22] and employing variance reduction [6].

Adaptive sampling is a powerful technique used in stochastic optimization to control the variance in the approximations employed as the optimization progresses. The idea is simple yet powerful; far from the solution(s) inaccurate and cheap (gradient) information can be employed while near the solution(s) accurate (gradient) information is required for both theory and practice. Of course, the key to such methods is the mechanism by which the sample size (or accuracy of the approximation) is selected. In [23, 12], algorithms that increase the samples sizes employed with prescribed (geometric) rules are proposed, and in [12] the authors showed that these methods achieve optimal worst-case first-order complexity for unconstrained problems. Other algorithms utilize gradient approximation tests to control the accuracy in the approximations; e.g., norm test [15, 12], inner product test [9, 10], and others [16, 3, 25]. Adaptive sampling algorithms have also been developed for simulation-based optimization problems [24, 31]. Finally, [2, 35] apply adaptive sampling methods to constrained stochastic optimization problems with convex feasible sets.

1.1 Contributions

In this paper, motivated by the successes of adaptive sampling methods (unconstrained stochastic optimization problems) and SQP methods (constrained deterministic optimization problems), we propose an adaptive sampling stochastic sequential quadratic programming algorithm for solving optimization problems with deterministic constraints and stochastic objective functions. This is far from a trivial task as complication arises when incorporating existing adaptive sampling strategies into the SQP paradigm primarily due to the multi-objective nature of constrained optimization. We propose a novel mechanism to control the accuracy in gradient approximations employed and the search directions calculated that balances the goals of achieving feasibility and reducing the objective function value. The main ingredients of our proposed method are: (1) the extension of the norm condition [15, 12] for equality constrained stochastic optimization problems, and (2) the adaptation of the stochastic SQP method proposed in [5].

We prove convergence guarantees for two different inexactness conditions, predetermined and adaptive. Specifically, under standard stochastic condition on the accuracy of the gradient approximations and deterministic conditions on the quality of the solutions of the linear system, we prove that a measure of first-order stationarity evaluated at the

iterates generated by our proposed algorithm converges to zero in expectation from arbitrary starting points. We also established worst-case sampling complexity results when the sample sizes are controlled using predetermined rates instead of adaptive tests showing that even though sampling work per iteration is increasing, the overall sampling complexity is still $O(n^{2(1+\epsilon)})$, where $\epsilon > 1$. Moreover, we then propose a practical variant of our algorithm, inspired by [12, 9]. Finally, we present numerical results on binary classification tasks with equality constraints and on equality constrained problems from the CUTE collection that demonstrate the efficiency and efficacy of our proposed practical method.

1.2 Notation

The set of natural numbers is denoted by $\mathbb{N} := \{0; 1; 2; \dots; g\}$. The set of real numbers (i.e., scalars) is denoted by \mathbb{R} , and $\mathbb{R}_{>r}$ ($\mathbb{R}_{\geq r}$) denotes the set of real numbers greater than (greater than or equal to) $r \in \mathbb{R}$. The set of n -dimensional vectors is denoted by \mathbb{R}^n , the set of m -by- n matrices is denoted by $\mathbb{R}^{m \times n}$, and the set of n -by- n symmetric matrices is denoted by \mathbb{S}^n . Our proposed algorithms are iterative, and generate a sequence of iterates $\{x_k\}$ with $x_k \in \mathbb{R}^n$. Let $f_k := f(x_k)$, $g_k := \nabla f(x_k)$, $c_k := c(x_k)$, and $J_k := \nabla c(x_k)^T$ for all $k \in \mathbb{N}$.

1.3 Organization

The paper is organized as follows. In Section 2 we formalize the problem statement and main assumptions. The general algorithmic framework of our proposed method is presented in Section 3. Convergence and complexity guarantees are stated and proven in Section 4. We present a practical adaptive sampling SQP method in Section 5. In Section 6, we demonstrate the empirical performance of the proposed method on classification tasks. Finally, in Section 7, we make some final remarks and discuss avenues for future research.

2 Problem Statement and Assumptions

We consider the following potentially nonlinear and/or nonconvex equality constrained optimization problem

$$\min_{x \in \mathbb{R}^n} f(x) \quad \text{s.t.: } c(x) = 0; \quad \text{with } f(x) = E[F(x; \xi)]; \quad (2.1)$$

where the objective function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and constraint function $c : \mathbb{R}^n \rightarrow \mathbb{R}^m$ are smooth, ξ is a random variable with associated probability space $(\Omega; \mathcal{F}; P)$, $F : \mathbb{R}^n \times \Omega \rightarrow \mathbb{R}$, and $E[\cdot]$ denotes the expectation taken with respect to P . Throughout the paper, we assume that the constraint function and its associated derivatives can be computed exactly. With regards to the objective function and its associated derivatives, we assume that these quantities are expensive to compute, but that accurate evaluations can be obtained as required. We formalize the notions of accuracy in subsequent sections of the paper.

We make the following main assumption with regards to (2.1) and the iterates x_k generated by our proposed algorithms.

Assumption 2.1. Let $X \subseteq \mathbb{R}^n$ be an open convex set containing the sequence $\{x_k\}$ generated by any run of the algorithm. The objective function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is continuously differentiable and bounded below over X and its gradient function $\nabla f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is Lipschitz continuous with constant $L \in \mathbb{R}_{>0}$ (with respect to $\|\cdot\|_2$) and bounded over X . The constraint function $c : \mathbb{R}^n \rightarrow \mathbb{R}^m$ (with $m \leq n$) is continuously differentiable and bounded over X and its Jacobian function $J := \nabla c^T : \mathbb{R}^n \rightarrow \mathbb{R}^{m \times n}$ is Lipschitz continuous with constant $\mu \in \mathbb{R}_{>0}$ (with respect to $\|\cdot\|_2$) and bounded over X . In addition, for all $x \in X$, the Jacobian $J(x)$ has singular values that are bounded uniformly below by $\nu \in \mathbb{R}_{>0}$.

Remark 2.2. We make the following remarks about Assumption 2.1. Under Assumption 2.1, there exist constants $(f_{\text{inf}}; g; c; \mu; J) \in \mathbb{R} \times \mathbb{R}_{>0} \times \mathbb{R}_{>0} \times \mathbb{R}_{>0} \times \mathbb{R}_{>0}$ such that, for all $k \in \mathbb{N}$,

$$f_{\text{inf}} \leq f_k; \|\nabla f(x_k)\|_2 \leq g; \|c_k\|_1 \leq c; \|J_k\|_2 \leq J; \text{ and } \kappa(J_k J_k^T)^{-1} \leq \mu^{-2}.$$

The components pertaining to the objective and constraint functions are standard assumptions in the equality constrained optimization literature. With regards to the algorithmic components, we do not assume that the iterate sequence itself is bounded, however, we do assume that the objective function and constraints, function values and derivatives, are bounded over the set X containing the iterates. While this assumption is reasonable in the deterministic setting, it is not ideal in the stochastic setting. That being said, it is a common assumption in the equality constrained stochastic optimization literature [5, 28]. The justification is that in this constrained setting, iterates are presumably converging towards a deterministic feasible region. Furthermore, we assume that the accuracy in the gradient approximations can be controlled, and as such claim that the assumption above is reasonable.

Let $\mathcal{L} : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$ be the Lagrangian function corresponding to (2.1), $\mathcal{L}(x; y) = f(x) + c(x)^T y$, where $y \in \mathbb{R}^m$ is the vector of Lagrange multipliers. Under Assumption 2.1, the necessary conditions for first-order stationarity for (2.1) are

$$0 = \begin{pmatrix} \nabla_x \mathcal{L}(x; y) \\ \nabla_y \mathcal{L}(x; y) \end{pmatrix} = \begin{pmatrix} \nabla f(x) + J(x)^T y \\ c(x) \end{pmatrix}.$$

3 General Algorithmic Framework

Our proposed algorithms can be characterized as adaptive (selection of the step size/accuracy of approximations), inexact (linear system solves), and stochastic (gradient approximation employed) sequential quadratic optimization (SQP) methods. Specifically, given x_k for all

$k \in \mathbb{N}$, a search direction $\bar{d}_k \in \mathbb{R}^n$ is computed by inexactly solving a quadratic optimization subproblem based on a local quadratic model of the objective function constructed using

$$\min_{d \in \mathbb{R}^n} f_k + \bar{g}_k^T d + \frac{1}{2} d^T H_k d \quad \text{s.t.:} \quad c_k + J_k d = 0; \quad (3.1)$$

where the matrix $H_k \in \mathbb{R}^{n \times n}$ satisfies the assumption below (Assumption 3.1).

Assumption 3.1. The sequence of symmetric matrices $\{H_k\}_{k \in \mathbb{N}}$ is bounded in norm by $H \in \mathbb{R}_{>0}$ such that $\|H_k\|_2 \leq H$ for all $k \in \mathbb{N}$. In addition, there exists a constant $\mu \in \mathbb{R}_{>0}$ such that, for all $k \in \mathbb{N}$, the matrix H_k has the property that $u^T H_k u \geq \mu \|u\|_2^2$ for all $u \in \mathbb{R}^n$ such that $J_k u = 0$.

Under Assumptions 2.1 and 3.1, the optimal solution $\bar{c}_k \in \mathbb{R}^n$ of the subproblem (3.1), and an associated displacement in the Lagrange multiplier $\bar{e}_k \in \mathbb{R}^m$, can be obtained by solving the linear system of equations given by

$$\begin{bmatrix} H_k & J_k^T \\ J_k & 0 \end{bmatrix} \begin{bmatrix} \bar{c}_k \\ \bar{e}_k \end{bmatrix} = \begin{bmatrix} \bar{g}_k + J_k^T y_k \\ c_k \end{bmatrix}; \quad (3.2)$$

Solving such linear systems can be expensive, and as such, our algorithms employ inexact solutions to the above linear systems given by $(\bar{d}_k; \bar{c}_k) = (\bar{c}_k; \bar{e}_k)$, i.e.,

$$\begin{bmatrix} H_k & J_k^T \\ J_k & 0 \end{bmatrix} \begin{bmatrix} \bar{d}_k \\ \bar{c}_k \end{bmatrix} = \begin{bmatrix} \bar{g}_k + J_k^T y_k \\ c_k \end{bmatrix} + \begin{bmatrix} r_k \\ r_k \end{bmatrix}; \quad (3.3)$$

where the tuple $(r_k; r_k)$ denote the residuals, and $\|(r_k; r_k)\|$ can be controlled as required. Our algorithms impose conditions on the norm of the residuals.

Given a pair $(\bar{d}_k; \bar{c}_k)$, our algorithms proceed to compute a positive step size in order to update the primal variables x_k and Lagrange multipliers y_k . The step size selection strategy is similar to that in [5]. To this end, the algorithms employ a merit function $l : \mathbb{R}^n \times \mathbb{R}_{>0} \rightarrow \mathbb{R}$, parameterized by a merit parameter $\kappa \in \mathbb{R}_{>0}$, defined as

$$l(x_k; \kappa) = \kappa f_k + \kappa \alpha_k \|r_k\|; \quad (3.4)$$

The merit parameter is dynamically adjusted as the optimization progresses. We employ an local model $l : \mathbb{R}^n \times \mathbb{R}_{>0} \rightarrow \mathbb{R}$ of the merit function defined by

$$l(x_k; \kappa; \bar{g}_k; \bar{d}_k) = \kappa (f_k + \bar{g}_k^T \bar{d}_k) + \kappa \alpha_k + \kappa \|J_k \bar{d}_k\| = \kappa (f_k + \bar{g}_k^T \bar{d}_k) + \kappa \|r_k\|;$$

and its associated reduction function $l : \mathbb{R}^n \times \mathbb{R}_{>0} \rightarrow \mathbb{R}$ defined by

$$\begin{aligned} l(x_k; \kappa; \bar{g}_k; \bar{d}_k) &= l(x_k; \kappa; \bar{g}_k; 0) - l(x_k; \kappa; \bar{g}_k; \bar{d}_k) \\ &= \kappa \bar{g}_k^T \bar{d}_k + \kappa \alpha_k - \kappa (c_k + J_k \bar{d}_k)^T \bar{d}_k \\ &= \kappa \bar{g}_k^T \bar{d}_k + \kappa \alpha_k - \kappa \|r_k\|; \end{aligned} \quad (3.5)$$

to guide the selection of the merit parameter. Note, $l(x_k; g_k; d_k) = g_k^T d_k + \alpha_k k_1$, where $(g_k; e_k)$ is the solution of (3.2). The mechanism for updating the merit parameter α_k is similar to that proposed in [5], and is motivated by state-of-the-art SQP methods [20, 14]. First, a trial merit parameter α_k^{trial} is computed. Given user-defined parameters $(\beta_1; \beta_2; \beta) \in (0; 1) \times (0; 1) \times \mathbb{R}_{>0}$, if $\beta_1 \alpha_k k_1 \leq (1 - \beta_1) \beta_2 \alpha_k k_1$ and/or $\beta_2 \alpha_k k_1 \leq \beta \alpha_k k_1$, we set $\alpha_k^{\text{trial}} = 1$. Otherwise, we set

$$\alpha_k^{\text{trial}} = \begin{cases} \beta < 1 & \text{if } g_k^T d_k + \max_{\|d_k\|_2 \leq g} d_k^T H_k d_k; \beta_2 \alpha_k k_1 \geq 0 \\ \frac{(1 - \beta_1)(1 - \beta_2) \alpha_k k_1}{g_k^T d_k + \max_{\|d_k\|_2 \leq g} d_k^T H_k d_k; \beta_2 \alpha_k k_1} & \text{otherwise,} \end{cases} \quad (3.6)$$

where $\beta \in (0; 1)$ is user-defined. Then, the merit parameter value α_k is updated via

$$\alpha_k = \begin{cases} \alpha_{k-1} & \text{if } \alpha_{k-1} \leq (1 - \beta) \alpha_k^{\text{trial}} \\ (1 - \beta) \alpha_k^{\text{trial}} & \text{otherwise,} \end{cases} \quad (3.7)$$

where $\beta \in (0; 1)$ is user-defined. This rule ensures that $\alpha_k g$ is a monotonically non-increasing positive sequence with $\alpha_k \leq (1 - \beta) \alpha_k^{\text{trial}}$ for all $k \geq N$. Moreover, this rule aims to ensure that the reduction function (3.5) is non-negative and satisfies,

$$l(x_k; g_k; d_k) \leq \alpha_k \max_{\|d_k\|_2 \leq g} d_k^T H_k d_k; \beta_2 \alpha_k k_1 + \beta_1 \max_{\|d_k\|_2 \leq g} f_k(x_k; d_k); \beta_1 \alpha_k k_1. \quad (3.8)$$

We show (see Lemma 3.6 below) that this bound indeed holds for the iterates generated by our proposed algorithm.

Finally, the step size $\alpha_k \in \mathbb{R}_{>0}$ is set as follows. At the k th iteration, we set via

$$\alpha_k = \min \left\{ \frac{\beta_1 (1 - \beta) \alpha_k^{\text{trial}}}{(\beta_1 L_k + \beta_2) \alpha_k k_1}; \alpha_k^{\text{opt}}; \beta_2 \alpha_k k_1; 1 \right\}; \quad (3.9)$$

where $\beta_1 \in (0; 1)$, $\beta_2 \in (0; 1)$, $\beta \in [2; 4]$ and $\alpha_k^{\text{opt}} \in \mathbb{R}_{>0}$ are user-defined parameters, and

$$\alpha_k^{\text{opt}} = \max \left\{ \min \left\{ \frac{l(x_k; g_k; d_k)}{(\beta_1 L_k + \beta_2) \alpha_k k_1}; 1 \right\}; \frac{l(x_k; g_k; d_k) - 2\alpha_k k_1}{(\beta_1 L_k + \beta_2) \alpha_k k_1} \right\}; \quad (3.10)$$

The step-size choice is motivated by that proposed in [5], and adapted for our specific setting. The user-defined parameters β_1 and β_2 that influence the step size are also related to the accuracy of gradient approximations employed and the accuracy of the solutions of the linear system (3.2), which are introduced and discussed in Section 4.

Our algorithmic framework, Adaptive, Inexact, Stochastic SQP (AIS-SQP), is presented in Algorithm 1.

Remark 3.2. We make the following remarks about Algorithm 1.

- Lines 2 and 3: For ease of exposition, we leave these two steps arbitrary and specify them later in the paper (Section 4). We assume that gradient approximations of arbitrary accuracy can be computed, and that the linear system can be solved iteratively and to arbitrary accuracy. In Section 5, we present a practical adaptive sampling strategy and a mechanism for solving the linear system inexactly.

Algorithm 1 (AIS-SQP) Adaptive, Inexact, Stochastic SQP Algorithm

Require: $x_0 \in \mathbb{R}^n$; $y_0 \in \mathbb{R}^m$; $f, H_k \in \mathbb{S}^n$; $\alpha_1 \in \mathbb{R}_{>0}$; $\beta_1, \beta_2 \in \mathbb{R}_{>0}$; $\gamma \in (0, 1)$; $\alpha \in \mathbb{R}_{>0}$; $\beta \in \mathbb{R}_{>0}$; $\delta \in (0, 1]$; $\mu \in \mathbb{R}_{>0}$; $\nu \in (0, \infty)$; $\eta \in [2, 4]$

- 1: for all $k \in \mathbb{N}$ do
- 2: Compute some gradient approximation $g_k \in \mathbb{R}^n$
- 3: Solve (3.2) iteratively; compute a step $(d_k; \kappa_k)$ that satisfies (a) or (b):
 - (a) (3.8) and $\kappa_k \kappa_{k+1} \leq \alpha \frac{\|L(x_k; \kappa_k; g_k; d_k)\|}{\|g_k\|}$, with $\kappa_k = \kappa_{k-1}$,
and, additionally, $\|g_k^T d_k + \max_{\|d\| \leq \kappa_k} d^T H_k d\| \leq \beta \|d_k\| \kappa_k^2 \|g_k\|$ if $\kappa_k \kappa_{k+1} > 0$,
 - (b) $\kappa_k \kappa_{k+1} < \min\{\beta_1 \alpha_1, \beta_2 \alpha\} \frac{\|g_k\|}{\|d_k\|}$ and $\kappa_k \kappa_{k+1} < \beta \kappa_k \kappa_{k+1}$
- 4: Update κ_k via (3.6)–(3.7) only if (b) is satisfied
- 5: Compute a step size κ_k via (3.9)–(3.10)
- 6: Update $x_{k+1} = x_k + \kappa_k d_k$, and $y_{k+1} = y_k + \kappa_k \nu$
- 7: end for

^ Lines 3(a) and 3(b): Lemma 3.5 (below) shows that the algorithm is well defined, and that Line 3 will terminate in either (a) or (b). One could enforce a simpler condition on the solution to the linear system, however, this simpler condition would come at the cost of having to solve the linear system exactly if/when the iterates are feasible. In the special case where $\kappa_k \kappa_{k+1} = 0$, by Assumptions 3.1 and 3.3 (presented below; pertaining to the linear system solutions) and (3.3), Line 3 of Algorithm 1 is guaranteed to terminate in case (a). The additional condition in case (a) is added for technical reasons discussed in Section 4. Intuitively, if the additional condition holds, then sufficient progress can be made in terms of minimizing the objective function, and thus progress can be made without modifying the merit parameter.

^ Merit parameter update: If Line 3 terminates due to the condition (a), then Line 4 is skipped. If Line 4 is triggered, then the merit parameter is potentially updated to ensure (3.8) is satisfied.

^ Step size selection: The step size selection strategy (3.9)–(3.10) depends on Lipschitz constants (L and ν), or estimates of these quantities. If one knows the Lipschitz constants, one could simply set $L_k = L$ and $\nu_k = \nu$ for all $k \in \mathbb{N}$. If such Lipschitz constants are unknown, as is the case more often than not, one can approximate these constants following the approaches proposed in [5, 23, 9]. See Section 5 for details of our practical implementation.

^ Comparison to other algorithms: The step computation, merit parameter update, and step size selection mechanisms are similar to those proposed in [5, 22, 4]. However, there are some key differences, primarily due to the fact that in this work we assume that the accuracy in the gradient approximations can be controlled as the optimization progresses. Similar to [22], the linear system is solved inexactly, but with a simpler approach that

does not require an explicit step decomposition.

Before we proceed, we state and prove a few results that hold throughout the paper. For the purposes of the presentation of the analysis in subsequent sections, we introduce the following notation. (Note, these quantities are never explicitly computed in our algorithms.) Given the iterate x_k and multiplier y_k , let the tuple $(d_k; \kappa_k)$ denote the solution to

$$\begin{pmatrix} H_k & J_k^T \\ J_k & 0 \end{pmatrix} \begin{pmatrix} d_k \\ \kappa_k \end{pmatrix} = \begin{pmatrix} g_k + J_k^T y_k \\ c_k \end{pmatrix}; \quad (3.11)$$

the deterministic counter-part of (3.2), where \bar{g}_k is replaced by the true gradient of the objective function. Moreover, let $f_{\kappa_k} g$ and $f_{\kappa_k} g$ be the sequences of merit parameters and step sizes, respectively, computed at x_k for all $k \geq 1$ by the deterministic variant of the algorithm with $\kappa_{k+1} = \kappa_k$. Finally, we introduce $\epsilon_{\min} \in \mathbb{R}_{>0}$, a lower bound on the deterministic merit parameter value, whose existence is guaranteed by Assumptions 2.1 and 3.1; see [5, 13]. We make the following additional assumption.

Assumption 3.3. There exist constants $\epsilon_r; \epsilon_g \in \mathbb{R}_{>0}$ such that $\|r_k\| \leq \epsilon_r$ and $\|g_k - \bar{g}_k\| \leq \epsilon_g$ for all $k \geq 1$. Moreover, for any $k \geq 1$, a sequence of inexact solutions $(d_{k;t}; \kappa_{k;t})_{t \geq 1}$ is generated by some iterative linear system solver (denotes the iteration counter of the linear system solver), where $\lim_{t \rightarrow \infty} f(d_{k;t}; \kappa_{k;t}) g = (d_k; \kappa_k)$ and $(d_k; \kappa_k) := (d_{k;t}; \kappa_{k;t})$ for some $t \geq 1$. Furthermore, for technical reasons, we also assume that either $\kappa_k \leq \epsilon_{\min}$ or $g_k \notin \text{Range}(J_k^T)$ for all $k \geq 1$.

Remark 3.4. We make the following remarks about Assumption 3.3. Assumption 3.3 pertains to properties of three main components: (i) the iterative solver, (ii) acceptable inexact solutions, and, (iii) the stochastic gradient estimates. First, we require that the iterative solver is able to return the exact solution of (3.2) in the limit. Second, we assume that all acceptable inexact solutions $(d_k; \kappa_k)$ have norms of residuals $\|r_k\|$ and $\|g_k - \bar{g}_k\|$ that are reasonably small and bounded uniformly. Third, we assume that the stochastic gradients (g_k) computed do not lie exactly in the range space of the Jacobian of the constraint (J_k^T) for iterates that are feasible. In general, this is not a strong assumption in the stochastic setting. For details about practical linear system solvers see [34] and references therein, and, for details about our implemented linear system solver see Sections 5.2 and 6.

The first result shows that Algorithm 1 is well-defined.

Lemma 3.5. Line 3 of Algorithm 1 terminates finitely.

Proof. We consider two cases: (i) $\kappa_k \leq \epsilon_{\min}$ and (ii) $\kappa_k > \epsilon_{\min}$. If $\kappa_k \leq \epsilon_{\min}$, by Assumption 3.3 it follows that $f_{\kappa_k} g \leq \epsilon_{\min}$. Therefore, for sufficiently large $t \geq 1$, $\|r_{k;t}\| \leq \epsilon_r$ and $\|g_{k;t} - \bar{g}_{k;t}\| \leq \epsilon_g$ are satisfied, and Line 3 of Algorithm 1 terminates finitely in case (i).

On the other hand, if $\kappa_k k_1 = 0$, by (3.2) and (3.5), it follows that

$$l(x_k; k-1; g_k; d_k) = \frac{1}{k} g_k^T d_k + \kappa_k k_1 \frac{1}{k} c_k + J_k d_k k_1 = \frac{1}{k} g_k^T d_k.$$

By (3.2) and Assumption 3.3 it follows that $g_k \notin \text{Range}(J_k^T)$ and $\kappa_k k_2 > 0$. By Assumption 3.1, (3.2), $\kappa_k k_2 > 0$ and $\frac{1}{d} \geq 2$ ($0; \Rightarrow$), it follows that $d_k^T H_k d_k > \frac{1}{d} \kappa_k k_2^2 > 0$. Moreover, by (3.2), $g_k^T d_k + d_k^T H_k d_k = 0$. Combining the above and $\frac{1}{k} g_k^T d_k > 0$,

$$\begin{aligned} l(x_k; k-1; g_k; d_k) &= \frac{1}{k} \max\{d_k^T H_k d_k; \frac{1}{d} \kappa_k k_2^2\} \\ &= \frac{1}{k} g_k^T d_k + \frac{1}{k} d_k^T H_k d_k \\ &= \frac{1}{k} (g_k^T d_k + d_k^T H_k d_k) + \frac{1}{k} (1 - \frac{1}{d}) d_k^T H_k d_k \\ &= \frac{1}{k} (1 - \frac{1}{d}) d_k^T H_k d_k > 0: \end{aligned}$$

Therefore, by $\frac{1}{k} g_k^T d_k > 0$ and $\kappa_k k_2 > 0$ (by Assumption 3.3), for sufficiently large $t \geq N$ and sufficiently small $\frac{1}{k} \kappa_k k_1$, it follows that $l(x_k; k-1; g_k; d_{k,t}) > 0$ and

$$\begin{aligned} l(x_k; k-1; g_k; d_{k,t}) &= \frac{1}{k} \max\{d_{k,t}^T H_k d_{k,t}; \frac{1}{d} \kappa_k k_2^2\} \\ &\quad + \frac{1}{k} \max\{c_k k_1; \frac{1}{k} \kappa_k k_1\}; \end{aligned}$$

with $l(x_k; k-1; g_k; d_{k,t}) \geq \frac{1}{k} \frac{\kappa_k k_2^2}{2}$. So Line 3 (Algorithm 1) terminates finitely. \square

Next, we prove that (3.8) is satisfied at every iteration of Algorithm 1. As mentioned above, this component of the algorithm is central in the analysis.

Lemma 3.6. The sequence of iterates generated by Algorithm 1 satisfies (3.8).

Proof. If the first condition on Line 3 of Algorithm 1 is triggered, i.e., case (a), then the result holds trivially. Thus, we focus on the case where the second condition is triggered, i.e., case(b). In this case, the residual vectors satisfy $\frac{1}{k} \kappa_k k_1 < \min\{(1 - \frac{1}{d}) \frac{1}{2}; \frac{1}{d} \frac{1}{a} \frac{1}{2} g_k c_k k_1$ and $\frac{1}{k} \kappa_k k_1 < \frac{1}{d} \kappa_k k_1$, and, potentially, the merit parameter is updated via (3.6) and (3.7). By the residual conditions, it follows that $\kappa_k k_1 > 0$. (Note, Lemma 3.5 showed that if $\kappa_k k_1 = 0$, Line 3 of Algorithm 1 terminates in case (a).) By (3.5), (3.8) and $\frac{1}{k} \kappa_k k_1 < \frac{1}{d} \kappa_k k_1$, to complete the proof, it is equivalent to show that

$$\frac{1}{k} g_k^T d_k + \frac{1}{k} \max\{d_k^T H_k d_k; \frac{1}{d} \kappa_k k_2^2\} > (1 - \frac{1}{d}) \kappa_k k_1 - \frac{1}{k} r_k k_1. \quad (3.12)$$

Since $\frac{1}{k} r_k k_1 < \frac{1}{d} \kappa_k k_1 < (1 - \frac{1}{d}) \kappa_k k_1$, (3.12) directly holds if $\frac{1}{k} g_k^T d_k + \frac{1}{k} \max\{d_k^T H_k d_k; \frac{1}{d} \kappa_k k_2^2\} > 0$. On the other hand, if $\frac{1}{k} g_k^T d_k + \frac{1}{k} \max\{d_k^T H_k d_k; \frac{1}{d} \kappa_k k_2^2\} > 0$, which also implies $\frac{1}{k} g_k^T d_k + \frac{1}{k} \max\{d_k^T H_k d_k; \frac{1}{d} \kappa_k k_2^2\} > 0$, by (3.6) and (3.7),

$$\frac{1}{k} (1 - \frac{1}{d}) \kappa_k k_1^{\text{trial}} = \frac{(1 - \frac{1}{d})(1 - \frac{1}{d}) \frac{1}{2} \kappa_k k_1}{\frac{1}{k} g_k^T d_k + \frac{1}{k} \max\{d_k^T H_k d_k; \frac{1}{d} \kappa_k k_2^2\}} < \frac{(1 - \frac{1}{d}) \kappa_k k_1 - \frac{1}{k} r_k k_1}{\frac{1}{k} g_k^T d_k + \frac{1}{k} \max\{d_k^T H_k d_k; \frac{1}{d} \kappa_k k_2^2\}};$$

which implies that (3.12) and (3.8) both hold. \square

The next lemma provides an upper bound on the primal residuals (3.3).

Lemma 3.7. For all $k \in \mathbb{N}$, the residual vector $r_k \in \mathbb{R}^m$ (3.3) satisfies $\|r_k\|_1 \leq \alpha^{-2} l(x_k; \kappa; g_k; d_k)$, where $\alpha \in \mathbb{R}_{>0}$, $\beta \in (0; 1]$ and $\gamma \in [2; 4]$.

Proof. There are two cases to consider with regards to Line 3 of Algorithm 1. If condition (a) is triggered, the result holds trivially. If condition (b) is triggered, by Lemma 3.6 and (3.8), it follows that

$$\|r_k\|_1 < \min\left(\frac{1}{\beta}, \frac{1}{\gamma}\right) \alpha^{-2} \beta \|g_k\|_1 + \frac{1}{\beta} \alpha^{-2} \gamma \|c_k\|_1 + \frac{1}{\alpha} \alpha^{-2} l(x_k; \kappa; g_k; d_k):$$

Combining the two cases yields the desired result. \square

The next lemma provides an upper bound on the deterministic dual variable update, $\|y_k + \kappa k_1\|$, required for the analysis in Section 4.1. We should note that we never require to compute κ in Algorithm 1.

Lemma 3.8. There exists some constant $\gamma \in \mathbb{R}_{>0}$, such that for all $k \in \mathbb{N}$, $\|y_k + \kappa k_1\| \leq \gamma$.

Proof. For all $k \in \mathbb{N}$, let $Z_k \in \mathbb{R}^{n \times (n-m)}$ be an orthonormal basis for the null space of the Jacobian of the constraints J_k , i.e., $J_k Z_k = 0$. Then under Assumption 3.1, it follows that $Z_k^T H_k Z_k = I$. By (3.11), it follows that

$$\begin{aligned} d_k &= -J_k^T (J_k J_k^T)^{-1} c_k - Z_k (Z_k^T H_k Z_k)^{-1} Z_k^T (g_k - H_k J_k^T (J_k J_k^T)^{-1} c_k); \\ \text{and } y_k + \kappa &= (J_k J_k^T)^{-1} J_k (g_k + H_k d_k) \\ &= (J_k J_k^T)^{-1} J_k (I - H_k Z_k (Z_k^T H_k Z_k)^{-1} Z_k^T) g_k \\ &\quad + (J_k J_k^T)^{-1} J_k H_k (I - Z_k (Z_k^T H_k Z_k)^{-1} Z_k^T H_k) J_k^T (J_k J_k^T)^{-1} c_k; \end{aligned}$$

By the Cauchy-Schwarz inequality, and Assumptions 2.1 and 3.1, it follows that

$$\|y_k + \kappa k_1\| \leq \frac{4}{\beta} \frac{\gamma}{J} + \frac{2}{\beta} \frac{\gamma}{J} \|g\|;$$

Thus, selecting a sufficiently large constant $\gamma \in \mathbb{R}_{>0}$, completes the proof. \square

Finally, we show that the model reduction function based on deterministic quantities, i.e., $l(x_k; \kappa; g_k; d_k)$, is non-negative and bounded above. For non-optimal points, one can further show the model reduction function is strictly positive.

Lemma 3.9. There exists some fixed constant $\beta \in \mathbb{R}_{>0}$ such that for all $k \in \mathbb{N}$, $l(x_k; \kappa; g_k; d_k) \in [0; \beta)$.

Proof. Notice, that in this lemma, we consider only deterministic quantities. First, we show that $l(x_k; \kappa; g_k; d_k) \geq 0$ for all $k \geq 2$. We consider two cases for the outcome of Line 3 of Algorithm 1. If (3.8) is satisfied with $\kappa = \kappa_{k-1}$, then by (3.5) and (3.11), we have

$$l(x_k; \kappa; g_k; d_k) \geq \kappa_{k-1} \max \{ d_k^T H_k d_k; \kappa d_k k_2^2 \} + \kappa_{k-1} \kappa_k k_1 \geq 0.$$

Otherwise, we have $0 = \kappa_k + J_k d_k k_1 < (1 - \beta_1) \kappa_{k-1} \kappa_k k_1$. We consider two subcases. If $g_k^T d_k + \max \{ d_k^T H_k d_k; \kappa d_k k_2^2 \} g \geq 0$, then it follows that $g_k^T d_k \geq 0$, and by $\kappa > 0$, (3.5) and (3.11), we have

$$l(x_k; \kappa; g_k; d_k) = \kappa g_k^T d_k + \kappa_k k_1 \geq 0.$$

On the other hand, if $g_k^T d_k + \max \{ d_k^T H_k d_k; \kappa d_k k_2^2 \} g < 0$, by (3.6) and (3.7) and the fact that $\kappa_k (1 - \beta_1) \kappa_k^{\text{trial}} < \kappa_k^{\text{trial}}$, we have

$$\kappa g_k^T d_k < (1 - \beta_1)(1 - \beta_2) \kappa_k k_1 \leq \kappa_k \max \{ d_k^T H_k d_k; \kappa d_k k_2^2 \} g. \quad (3.13)$$

Combining (3.5), (3.11) and (3.13), it follows that

$$\begin{aligned} l(x_k; \kappa; g_k; d_k) &= \kappa g_k^T d_k + \kappa_k k_1 \\ &> \kappa \max \{ d_k^T H_k d_k; \kappa d_k k_2^2 \} g + (1 - (1 - \beta_1)(1 - \beta_2)) \kappa_k k_1 \geq 0. \end{aligned}$$

Thus, we have shown that $l(x_k; \kappa; g_k; d_k) \geq 0$ for all $k \geq 2$.

Next, we show that $l(x_k; \kappa; g_k; d_k) \geq \beta_1$ for all $k \geq 2$. For all $k \geq 2$, let $Z_k \in \mathbb{R}^{n \times (n-m)}$ be an orthonormal basis for the null space of the Jacobian of the constraints f_k , i.e., $J_k Z_k = 0$, and Assumption 3.1 further implies that $Z_k^T H_k Z_k \succeq I$. From (3.11) we have

$$g_k^T d_k = g_k^T (I - Z_k(Z_k^T H_k Z_k)^{-1} Z_k^T H_k) J_k^T (J_k J_k^T)^{-1} \kappa_k - g_k^T Z_k (Z_k^T H_k Z_k)^{-1} Z_k^T g_k. \quad (3.14)$$

By Assumption 2.1 and (3.14), we have

$$\kappa g_k^T d_k k_2 \geq \kappa g_k^T J_k^T (J_k J_k^T)^{-1} \kappa_k k_2 + \kappa g_k^T Z_k (Z_k^T H_k Z_k)^{-1} Z_k^T g_k k_2 \geq \beta_1 \kappa_k k_2 + \frac{2}{g} \kappa_k k_2.$$

Moreover, by (3.5) and (3.11), it follows that for all $k \geq 2$,

$$l(x_k; \kappa; g_k; d_k) = \kappa g_k^T d_k + \kappa_k k_1 \geq \beta_1 \kappa_k k_2 + \frac{2}{g} \kappa_k k_2 + \kappa_k k_1.$$

Selecting a sufficiently large constant $\beta_1 \geq \beta_0$ completes the proof. \square

4 Adaptive, Inexact and Stochastic Sequential Quadratic Programming

In this section, we prove under different conditions on the gradient and linear system solution accuracies, that Algorithm 1 has convergence properties that match those from the deterministic setting in expectation. First, we consider adaptive error bounds (Section 4.1) and then consider predetermined sublinear errors (Section 4.2).

4.1 Adaptive Iteration-Dependent Errors

In this section, we provide a comprehensive convergence analysis for Algorithm 1 under stochastic conditions on the error in the gradient approximations (Assumption 4.1), and inexact solutions to the SQP subproblems (3.3) (Assumptions 4.2 and 4.9). The following two assumptions are central to the analysis presented in this section.

Assumption 4.1. For all $k \in \mathbb{N}$, the stochastic gradient estimate $\bar{g}_k \in \mathbb{R}^n$ satisfies

$$\mathbb{E}_k \|\bar{g}_k - g_k\|_2^2 \leq \frac{1}{\alpha} l(x_k; \kappa; g_k; d_k); \quad (4.1)$$

where $\alpha \in \mathbb{R}_{>0}$, $\beta \in (0; 1)$, and $\gamma \in [2; 4]$. Additionally, for all $k \in \mathbb{N}$, the stochastic gradient estimate $\bar{g}_k \in \mathbb{R}^n$ is an unbiased estimator of the gradient of l at x_k , i.e., $\mathbb{E}_k [\bar{g}_k] = g_k$, where $\mathbb{E}_k[\cdot]$ denotes the expectation with respect to the distribution of \bar{g}_k conditioned on the event that the algorithm has reached $x_k \in \mathbb{R}^n$ in iteration $k \in \mathbb{N}$.

Assumption 4.2. For all $k \in \mathbb{N}$, the search directions $(\bar{d}_k; \kappa) \in \mathbb{R}^n \times \mathbb{R}^m$ in (3.3) (inexact solutions to (3.2)) satisfy

$$\|\bar{d}_k - d_k\|_2^2 \leq \frac{1}{\beta} l(x_k; \kappa; g_k; d_k); \quad (4.2)$$

where $\beta \in \mathbb{R}_{>0}$, $\beta \in (0; 1)$, and $\gamma \in [2; 4]$. Note, $(d_k; \kappa)$ and $(\bar{d}_k; \kappa)$ are the exact and inexact solutions of (3.2), respectively.

Remark 4.3. We note upfront that the conditions given in Assumptions 4.1 and 4.2 are not implementable in our stochastic setting, as the right-hand-side of the inequalities depend on deterministic quantities. That being said, we use these assumptions as this allows us to gain insights into the errors permitted in the algorithm while still retaining strong convergence guarantees, and will guide the development of our practical algorithm. An important choice in conditions (4.1)-(4.2) is the (deterministic) model reduction function (3.5) on the right-hand-side of the inequality. As we show in the analysis, and similar to [5, 4, 22], we use this quantity as a proxy of convergence, and as such it is an appropriate measure of the accuracy in the gradient approximations. Another interesting question pertains to the analogue of (4.1) in the unconstrained setting, i.e., no equality constraints. One can show that for appropriately chosen constants α and β , in the unconstrained setting (4.1) is the well-known "norm" condition (in expectation) [15, 12]. With regards to (4.2), under Assumption 3.3 the inequality is well-defined. Finally, we emphasize that the constants α and β are the same constants that appear in Algorithm 1 and that are used in the step size selection. Thus, the gradient accuracy, linear system solves and the step size selection are inherently connected. The precise permissible ranges of the constants in Assumptions 4.1 and 4.2 are made explicit in subsequent lemmas and theorems.

We prove convergence guarantees for Algorithm 1, where the stochastic gradients employed satisfy Assumption 4.1 and the search directions employed satisfy Assumption 4.2. Before we delve into the analysis, we discuss the behavior of the merit parameter sequence $\{f_k\}$, a key component of our algorithmic framework. In the deterministic setting (e.g., Algorithm 1 with $\bar{g}_k = g_k$ and $\bar{d}_k = d_k$ for all $k \geq N$), under Assumption 2.1 the merit parameter sequence is bounded away from zero $f_k \geq \min f_k > 0$. In the stochastic setting, where the gradients approximations employed satisfy Assumption 4.1, boundedness (away from zero) of the merit parameter cannot be guaranteed. However, under the assumption that the iterates generated by Algorithm 1 converge to a stationary point of (2.1), the gradient approximation eventually become sufficient accurate, as do the solutions to the linear system, and thus it is reasonable to assume that the stochastic merit parameter value remains bounded away from zero. That being said, this is not sufficient to prove strong convergence and complexity guarantees for Algorithm 1 for the whole merit parameter sequence. To this end, we impose an additional technical condition on the gradient approximation and the search direction employed (see Assumption 4.9 later in the text), and, to the best of our knowledge, prove the first complexity guarantees in this setting.

For ease of exposition, we make the following assumption throughout this section that unifies all prior assumptions, and for brevity, we do not remind the reader of this fact within the statement of each result.

Assumption 4.4. There exist universal quantities such that Assumptions 2.1, 3.1, 3.3, 4.1 and 4.2 hold for any realization of Algorithm 1.

Let us now commence our analysis of Algorithm 1 under Assumption 4.4. We build up to our main results through a series of lemmas. Our first set of lemmas show that the stochastic search directions computed by Algorithm 1 are well-behaved. To this end, we invoke the following orthogonal decomposition of the (stochastic) search direction: $d_k = u_k + v_k$ where $u_k \in \text{Null}(J_k)$ and $v_k \in \text{Range}(J_k^T)$ for all $k \geq N$.

Lemma 4.5. There exists $\nu \in \mathbb{R}_{>0}$ such that, for all $k \geq N$, the normal component v_k satisfies $\max\{k v_k\|_2; k v_k\|_2^2\} \leq \nu \max\{k \alpha_k\|_2; k r_k\|_2\}$.

Proof. Since $u_k \in \text{Null}(J_k)$ and $v_k \in \text{Range}(J_k^T)$,

$$v_k = J_k^T (J_k J_k^T)^{-1} J_k v_k = J_k^T (J_k J_k^T)^{-1} J_k d_k = J_k^T (J_k J_k^T)^{-1} (r_k - \alpha_k):$$

Thus, by the Cauchy inequality

$$k v_k\|_2 \leq k J_k^T (J_k J_k^T)^{-1} k_2 (k r_k\|_2 + k \alpha_k\|_2) \leq 2k J_k^T (J_k J_k^T)^{-1} k_2 \max\{k \alpha_k\|_2; k r_k\|_2\}:$$

Moreover, it follows that

$$k v_k\|_2^2 \leq 4k J_k^T (J_k J_k^T)^{-1} k_2^2 \max\{k \alpha_k\|_2; k r_k\|_2\} \max\{k \alpha_k\|_2; k r_k\|_2\}:$$

By Assumption 4.4, we have that $k \alpha_k\|_2$, $k r_k\|_2$ and $k J_k^T (J_k J_k^T)^{-1} k_2$ are uniformly bounded above for all $k \geq N$, which in turn completes the proof. \square

The next lemma shows that if the step d_k is tangentially dominated, i.e., $\|u_k\|_2$ dominates $\|v_k\|_2$, then H_k is sufficiently positive definite along the computed stochastic search direction.

Lemma 4.6. There exists $\bar{\mu} \in \mathbb{R}_{>0}$ such that, for any $k \in \mathbb{N}$, if $\|u_k\|_2 \geq \bar{\mu} \|v_k\|_2$, then $d_k^T H_k d_k \geq \frac{1}{2} \|u_k\|_2^2$ and $d_k^T H_k d_k \geq \frac{1}{2} \|d_k\|_2^2$, where $\bar{\mu} \in \mathbb{R}_{>0}$ is defined in Assumption 3.1 and $\bar{\mu} \geq 0$ (\Rightarrow).

Proof. When $\|u_k\|_2 = 0$, to satisfy the condition in the statement, we require $\|u_k\|_2 = \|v_k\|_2 = 0$ which implies $\|d_k\|_2 = 0$, and the statement holds trivially.

When $\|u_k\|_2 > 0$, by Assumption 3.1 it follows that

$$\begin{aligned} d_k^T H_k d_k &= u_k^T H_k u_k + 2 u_k^T H_k v_k + v_k^T H_k v_k \\ &\geq \frac{1}{2} \|u_k\|_2^2 - \frac{2\mu}{\bar{\mu}} \|u_k\|_2 \|v_k\|_2 + \frac{1}{2} \|v_k\|_2^2 \\ &\geq \frac{2\mu - \frac{2\mu}{\bar{\mu}}}{\bar{\mu}} \|u_k\|_2 \|v_k\|_2 + \frac{1}{2} \|u_k\|_2^2 \end{aligned}$$

for sufficiently large $\bar{\mu}$. Meanwhile, when $\|u_k\|_2 > 0$, it follows that

$$d_k^T H_k d_k \geq \frac{1}{2} \|u_k\|_2^2 \geq \frac{1}{2} (1 + \frac{1}{\bar{\mu}}) \|u_k\|_2^2 \geq \frac{1}{2} (\|u_k\|_2^2 + \|v_k\|_2^2) = \frac{1}{2} \|d_k\|_2^2$$

for sufficiently large $\bar{\mu}$. □

For $\bar{\mu} \in \mathbb{R}_{>0}$ (defined in Lemma 4.6), let $K_{\bar{\mu}} := \{k \in \mathbb{N} : \|u_k\|_2 \geq \bar{\mu} \|v_k\|_2\}$ and $K_{\bar{\mu}}^c := \{k \in \mathbb{N} : \|u_k\|_2 < \bar{\mu} \|v_k\|_2\}$ denote index sets that form a partition of \mathbb{N} , and let

$$\beta_k := \begin{cases} \|u_k\|_2^2 + \|v_k\|_2^2 & \text{if } k \in K_{\bar{\mu}}; \\ \max\{\|u_k\|_2^2, \|v_k\|_2^2\} & \text{if } k \in K_{\bar{\mu}}^c. \end{cases}$$

Our next result shows that the squared norms of the stochastic search directions and the constraint violations for all $k \in \mathbb{N}$ are bounded above by β_k .

Lemma 4.7. There exists $\bar{\mu} \in \mathbb{R}_{>0}$ such that, for all $k \in \mathbb{N}$, the search direction and constraint violation satisfy $\|d_k\|_2^2 \leq \beta_k$ and $\|c_k\|_2^2 + \|r_k\|_2^2 \leq (1 + \frac{1}{\bar{\mu}}) \beta_k$.

Proof. For all $k \in K_{\bar{\mu}}$,

$$\|d_k\|_2^2 = \|u_k\|_2^2 + \|v_k\|_2^2 \leq (1 + \frac{1}{\bar{\mu}}) \|u_k\|_2^2 \leq (1 + \frac{1}{\bar{\mu}}) \beta_k;$$

For all $k \in K_{\bar{\mu}}^c$, by Lemma 4.5,

$$\|d_k\|_2^2 = \|u_k\|_2^2 + \|v_k\|_2^2 < (1 + \frac{1}{\bar{\mu}}) \|v_k\|_2^2 \leq (1 + \frac{1}{\bar{\mu}}) \beta_k;$$

Therefore, we set $\bar{\mu} := \max\{1 + \frac{1}{\bar{\mu}}; (1 + \frac{1}{\bar{\mu}}) \bar{\mu}\}$ to satisfy $\|d_k\|_2^2 \leq \beta_k$. Finally, since $\|c_k\|_2^2 \leq \beta_k$ trivially, this concludes the proof. □

The next lemma shows that the stochastic model reduction, $l(x_k; g_k; d_k)$, is bounded below by a non-negative quantity.

Lemma 4.8. There exists $\bar{\gamma} \in \mathbb{R}_{>0}$ such that, for all $k \in \mathbb{N}$, the reduction in the model of the merit function satisfies, $l(x_k; g_k; d_k) \geq \bar{\gamma} \|k\|$.

Proof. For all $k \in K_{\bar{\gamma}}$, by Lemmas 3.6 and 4.6,

$$l(x_k; g_k; d_k) \geq \|k\| \min \left\{ \frac{1}{2} \max_{\|d_k\| \leq \frac{1}{2}} d_k^T H_k d_k; \frac{1}{2} \max_{\|d_k\| \leq \frac{1}{2}} \left(\frac{1}{2} \|d_k\| \right) \right\} + \frac{1}{2} \max_{\|d_k\| \leq \frac{1}{2}} \left(\frac{1}{2} \|d_k\| \right) \geq \frac{1}{2} \min \left\{ \frac{1}{2}; \frac{1}{2} \right\} \|k\|$$

For all $k \in K_{\bar{\gamma}}$, by Lemma 3.6,

$$l(x_k; g_k; d_k) \geq \|k\| \min \left\{ \frac{1}{2} \max_{\|d_k\| \leq \frac{1}{2}} d_k^T H_k d_k; \frac{1}{2} \max_{\|d_k\| \leq \frac{1}{2}} \left(\frac{1}{2} \|d_k\| \right) \right\} + \frac{1}{2} \max_{\|d_k\| \leq \frac{1}{2}} \left(\frac{1}{2} \|d_k\| \right) \geq \frac{1}{2} \min \left\{ \frac{1}{2}; \frac{1}{2} \right\} \|k\|$$

Setting $\bar{\gamma} := \frac{1}{2} \min \left\{ \frac{1}{2}; \frac{1}{2} \right\} \geq \frac{1}{4} \geq \frac{1}{2} \geq 0$ completes the proof. \square

Next we introduce a technical condition required for the analysis. Specifically, the condition allows us to establish an integral and useful upper bound for the difference between the stochastic and deterministic merit parameters values, and as a result is of vital importance in establishing complexity result.

Assumption 4.9. For all $k \in \mathbb{N}$, the stochastic gradient estimate $g_k \in \mathbb{R}^n$ and the search direction $d_k \in \mathbb{R}^n$ (inexact solution of (3.2)) satisfy

$$j(g_k^T d_k + \max_{\|d_k\| \leq \frac{1}{2}} d_k^T H_k d_k; \frac{1}{2} \|d_k\|) \leq (g_k^T d_k + \max_{\|d_k\| \leq \frac{1}{2}} d_k^T H_k d_k; \frac{1}{2} \|d_k\|) + \beta \|g_k^T d_k + \max_{\|d_k\| \leq \frac{1}{2}} d_k^T H_k d_k; \frac{1}{2} \|d_k\|\|g_k\|$$

where $\beta \in \mathbb{R}_{>0}$, $\alpha \in (0; 1)$, $\gamma \in [2; 4]$, and $\delta \in \mathbb{R}_{>0}$.

Remark 4.10. We make a few remarks about Assumption 4.9. This is a strong assumption, nonetheless, it is necessary in order to establish the strong non-asymptotic convergence and complexity results presented in this paper. In the unconstrained setting, Assumption 4.9 does not add any additional restrictions. That is, in the unconstrained setting, consider the gradient method, where H_k is the identity matrix, $d_k = -g_k$, $d_k = -g_k$. Clearly, $(g_k^T d_k + \max_{\|d_k\| \leq \frac{1}{2}} d_k^T H_k d_k; \frac{1}{2} \|d_k\|) = 0$ and $g_k^T d_k + \max_{\|d_k\| \leq \frac{1}{2}} d_k^T H_k d_k; \frac{1}{2} \|d_k\| = 0$, thus no additional restrictions are imposed. Thus, in the unconstrained setting, Assumptions 4.1, 4.2 and 4.9 reduce to the well known "norm-condition". Several difficulties arise in the setting with constraints. The primary difficulty pertains to the fact that the merit parameter is possibly adjusted across iterations making the merit function a moving target. Thus, in order to establish convergence and complexity results for all iterations, as compared to other

papers that only consider the iterations after the merit parameter has stabilized at a sufficiently small constant value, additional control on the permissible differences is required. We should note that if one happened to know a sufficiently small merit parameter value, then Assumption 4.9 would no longer be required for the theory. Finally, again we point out the connection between the accuracy in the gradient approximations, the quality of the solution to the linear system, the step size through user-defined parameters and .

Our next lemma shows that the stochastic merit parameter sequence $\epsilon_k g$ generated by Algorithm 1 is bounded away from zero.

Lemma 4.11. Under Assumptions 4.4 and 4.9, there exists a constant $\epsilon_{\min} \geq \epsilon_{\min} > 0$ that following Algorithm 1, $\epsilon_k \geq \epsilon_{\min}$ for all $k \geq N$.

Proof. Under Assumption 4.4 (specifically Assumptions 2.1 and 3.1), it is well established that $f_k^{\text{trial}} g$, the sequence of deterministic variant of (3.6), is always positive and bounded away from zero. By Assumption 4.9, it follows that either $(g_k^T d_k + \max f d_k^T H_k d_k; \epsilon_k d_k^T k^2 g)$ and $(g_k^T d_k + \max f d_k^T H_k d_k; \epsilon_k d_k^T k^2 g)$ are both positive or both negative.

If $g_k^T d_k + \max f d_k^T H_k d_k; \epsilon_k d_k^T k^2 g \leq 0$, by (3.6), $\epsilon_k^{\text{trial}} = 1$. By Assumption 4.9 it follows that $g_k^T d_k + \max f d_k^T H_k d_k; \epsilon_k d_k^T k^2 g \leq 0$ and $\epsilon_k^{\text{trial}} = 1$.

If $g_k^T d_k + \max f d_k^T H_k d_k; \epsilon_k d_k^T k^2 g > 0$, by Assumption 4.9, it follows that

$$0 < (1 - \epsilon_k^{\text{trial}}) \leq \frac{g_k^T d_k + \max f d_k^T H_k d_k; \epsilon_k d_k^T k^2 g}{g_k^T d_k + \max f d_k^T H_k d_k; \epsilon_k d_k^T k^2 g} (1 + \epsilon_k^{\text{trial}}):$$

By (3.6), the above inequality implies $\frac{\epsilon_k^{\text{trial}}}{(1 - \epsilon_k^{\text{trial}})} \geq \frac{\epsilon_k^{\text{trial}}}{(1 + \epsilon_k^{\text{trial}})}$.

In both cases considered above, given the fact that $f_k^{\text{trial}} g$ is positive and bounded away from zero, one can conclude that there exists a constant $\epsilon_{\min}^{\text{trial}} \geq \epsilon_{\min}^{\text{trial}} > 0$ such that $\epsilon_k^{\text{trial}} \geq \epsilon_{\min}^{\text{trial}}$ for all $k \geq N$. Finally, invoking (3.7), completes the proof since $\epsilon_k \geq \min\{\epsilon_{\min}^{\text{trial}}; (1 - \epsilon_{\min}^{\text{trial}})\} \epsilon_k^{\text{trial}} g$. \square

The next lemma provides a useful lower bound for the reduction in the merit function that is proportional to the stochastic search direction computed.

Lemma 4.12. There exist constants $\epsilon_{\min}^{\text{trial}}; \epsilon_{\min}^{\text{trial}} > 0$ such that for all $k \geq N$, $l(x_k; \epsilon_k; g_k; d_k) \geq (\epsilon_{\min}^{\text{trial}} L_k + \epsilon_{\min}^{\text{trial}}) \epsilon_k d_k^T k^2 g \geq \epsilon_{\min}^{\text{trial}} \epsilon_k d_k^T k^2 g$.

Proof. By Lemmas 4.7{4.11, it follows that

$$l(x_k; \epsilon_k; g_k; d_k) \geq \epsilon_k d_k^T k^2 g = \frac{\epsilon_k d_k^T k^2 g}{(\epsilon_{\min}^{\text{trial}} L_k + \epsilon_{\min}^{\text{trial}})} (\epsilon_{\min}^{\text{trial}} L_k + \epsilon_{\min}^{\text{trial}}) \epsilon_k d_k^T k^2 g \geq \frac{\epsilon_{\min}^{\text{trial}}}{(\epsilon_{\min}^{\text{trial}} L_k + \epsilon_{\min}^{\text{trial}})} (\epsilon_{\min}^{\text{trial}} L_k + \epsilon_{\min}^{\text{trial}}) \epsilon_k d_k^T k^2 g \geq \epsilon_{\min}^{\text{trial}} \epsilon_k d_k^T k^2 g;$$

where the result follows by choosing appropriate constants $\epsilon_{\min}^{\text{trial}}$ and $\epsilon_{\min}^{\text{trial}}$. \square

Remark 4.13. Corollaries of the above lemmas can be derived for the special case under which deterministic quantities are used instead of stochastic, i.e. (3.11) is solved exactly. Specifically, under the same logic as in Lemmas 4.7-4.12, it follows that there exist constants $f_{l;d}g_{R>0}$ such that for all $k \geq N$

$$l(x_k; k; g_k; d_k) \leq (L_k + \epsilon_k)kd_kk_2^2 + l;d_kd_kk_2^2 \quad (4.3)$$

The next lemma provides upper and lower bounds on the adaptive step sizes employed by Algorithm 1.

Lemma 4.14. Let ϵ_k be defined via (3.9)-(3.10). For all $k \geq N$, there exists a constant $\underline{\epsilon} \geq R_{>0}$ such that $\epsilon_k \geq \underline{\epsilon} \min_{u \in \{2^{-2}, 1\}}$.

Proof. The upper bound follows from (3.9). To derive the lower bound, we consider different cases. If $\epsilon_k = \min_{u \in \{2^{-2}, 1\}} \min_{k^{\text{opt}}; 1}$, under the conditions $\epsilon \in (0; 1]$ and $\epsilon \in [2; 4]$, it follows that $\epsilon_k \geq \min_{u \in \{2^{-2}, 1\}} \min_{k^{\text{opt}}; 1} g$. Otherwise, if $\epsilon_k = \frac{2(1 - \epsilon)^{(\epsilon = 2^{-1})} l(x_k; k; g_k; d_k)}{(L_k + \epsilon_k)kd_kk_2^2}$, by Lemma 4.12, $\epsilon \in (0; 1]$ and $\epsilon \in [2; 4]$, it follows that $\epsilon_k \geq 2(1 - \epsilon)^{(\epsilon = 2^{-1})} g$, where $\epsilon \geq R_{>0}$ is defined in Lemma 4.12. Combining (3.10) and Lemma 4.12, it follows that $\epsilon_k \geq \min_{k^{\text{opt}}; 1} \frac{l(x_k; k; g_k; d_k)}{(L_k + \epsilon_k)kd_kk_2^2} \geq \min_{u \in \{2^{-2}, 1\}} \min_{k^{\text{opt}}; 1} g$. Setting $\underline{\epsilon} := \min_{u \in \{2^{-2}, 1\}} \min_{k^{\text{opt}}; 1} g$ completes the proof. \square

The next result provides an upper bound on the change of merit function value after a step. Central to the proof of this lemma is our step size strategy (3.9)-(3.10).

Lemma 4.15. For all $k \geq N$, it follows that

$$\begin{aligned} & l(x_{k+1}; k; g_k; d_k) - l(x_k; k; g_k; d_k) \\ & \leq \epsilon_k l(x_k; k; g_k; d_k) + (1 - \epsilon_k)^{(\epsilon = 2^{-1})} l(x_k; k; g_k; d_k) \\ & \quad + \epsilon_k g_k^T (d_k - \hat{d}_k) + \epsilon_k (L_k + \epsilon_k) g_k^T d_k + \epsilon_k k J_k(d_k - \hat{d}_k) k_1; \end{aligned}$$

where \hat{d}_k and d_k are defined in (3.3) and (3.11), respectively, ϵ_k is updated via (3.6)-(3.7) and \hat{d}_k is its deterministic counterpart, and ϵ_k is computed via (3.9)-(3.10).

Proof. By the step size selection strategy (3.9)-(3.10), for all $k \geq N$, we have that $\epsilon_k \leq 1$.

By the triangle inequality, Assumption 2.1, and (3.2), it follows that

$$\begin{aligned}
& \langle x_k + \alpha_k d_k; x_k \rangle - \langle x_k; x_k \rangle \\
&= \alpha_k (f(x_k + \alpha_k d_k) - f(x_k)) + \langle \alpha_k \nabla f(x_k + \alpha_k d_k); x_k \rangle - \alpha_k \langle \nabla f(x_k); x_k \rangle \\
&= \alpha_k \left(\langle \alpha_k g_k^T d_k + \frac{L}{2} \alpha_k^2 \|d_k\|_2^2 \right) + \alpha_k \langle \alpha_k \nabla f(x_k + \alpha_k d_k); x_k \rangle + \frac{L}{2} \alpha_k^2 \|d_k\|_2^2 - \alpha_k \langle \nabla f(x_k); x_k \rangle \\
&= \alpha_k \left(\langle \alpha_k g_k^T d_k + \frac{L}{2} \alpha_k^2 \|d_k\|_2^2 \right) + \frac{1}{2} \alpha_k (L + 1) \alpha_k^2 \|d_k\|_2^2 + \alpha_k \langle \nabla f(x_k + \alpha_k d_k); x_k \rangle - \alpha_k \langle \nabla f(x_k); x_k \rangle \\
&= \alpha_k \left(\langle \alpha_k g_k^T d_k + \frac{L}{2} \alpha_k^2 \|d_k\|_2^2 \right) + \frac{1}{2} \alpha_k (L + 1) \alpha_k^2 \|d_k\|_2^2 \\
&\quad + \alpha_k \langle \alpha_k g_k^T (d_k - \tilde{d}_k) \rangle + \alpha_k \left(\langle \alpha_k \nabla f(x_k + \alpha_k d_k); x_k \rangle - \langle \alpha_k \nabla f(x_k + \alpha_k \tilde{d}_k); x_k \rangle \right) \\
&= \alpha_k \left(\langle \alpha_k g_k^T d_k + \frac{L}{2} \alpha_k^2 \|d_k\|_2^2 \right) + \frac{1}{2} \alpha_k (L + 1) \alpha_k^2 \|d_k\|_2^2 \\
&\quad + \alpha_k \langle \alpha_k g_k^T (d_k - \tilde{d}_k) \rangle + \alpha_k \left(\langle \alpha_k \nabla f(x_k + \alpha_k d_k); x_k \rangle - \langle \alpha_k \nabla f(x_k + \alpha_k \tilde{d}_k); x_k \rangle \right)
\end{aligned}$$

By (3.9), we have for all $k \geq N$, that $\alpha_k \leq \frac{2(1 - \epsilon) \alpha_k^2}{(L + 1) \|d_k\|_2^2} \frac{\langle \alpha_k \nabla f(x_k + \alpha_k d_k); x_k \rangle}{\|d_k\|_2^2}$, and

$$\begin{aligned}
& \langle x_k + \alpha_k d_k; x_k \rangle - \langle x_k; x_k \rangle \\
&= \alpha_k \left(\langle \alpha_k g_k^T d_k + \frac{L}{2} \alpha_k^2 \|d_k\|_2^2 \right) + \frac{1}{2} \alpha_k (L + 1) \alpha_k^2 \|d_k\|_2^2 \\
&\quad + \alpha_k \langle \alpha_k g_k^T (d_k - \tilde{d}_k) \rangle + \alpha_k \left(\langle \alpha_k \nabla f(x_k + \alpha_k d_k); x_k \rangle - \langle \alpha_k \nabla f(x_k + \alpha_k \tilde{d}_k); x_k \rangle \right) \\
&= \alpha_k \left(\langle \alpha_k g_k^T d_k + \frac{L}{2} \alpha_k^2 \|d_k\|_2^2 \right) + \frac{1}{2} \alpha_k (L + 1) \alpha_k^2 \|d_k\|_2^2 \\
&\quad + \alpha_k \langle \alpha_k g_k^T (d_k - \tilde{d}_k) \rangle + \alpha_k \left(\langle \alpha_k \nabla f(x_k + \alpha_k d_k); x_k \rangle - \langle \alpha_k \nabla f(x_k + \alpha_k \tilde{d}_k); x_k \rangle \right)
\end{aligned}$$

which is the desired result. \square

We now proceed to state and prove a series of lemmas (Lemmas 4.16-4.19) that provide bounds for the differences between the deterministic and stochastic gradients, search directions, and their inner products.

Lemma 4.16. For all $k \geq N$, $\mathbb{E}_k[\langle \nabla f(x_k); g_k \rangle] \leq \frac{2^p}{1 - 2^p} \frac{\langle \nabla f(x_k); g_k \rangle}{\|g_k\|_2}$.

Proof. By Assumption 4.1 and Jensen's inequality, it follows that

$$\mathbb{E}_k[\langle \nabla f(x_k); g_k \rangle] \leq \frac{\mathbb{E}_k[\langle \nabla f(x_k); g_k \rangle^2]}{\mathbb{E}_k[\|g_k\|_2^2]} \leq \frac{2^p}{1 - 2^p} \frac{\langle \nabla f(x_k); g_k \rangle}{\|g_k\|_2}. \quad \square$$

Lemma 4.17. For all $k \geq N$, $\mathbb{E}_k[d_k] = \tilde{d}_k$, $\mathbb{E}_k[u_k] = u_k$, and $\mathbb{E}_k[\tilde{d}_k] = \tilde{d}_k$. Moreover, there exists $\sigma \geq 2$, such that $\mathbb{E}_k[\|d_k - \tilde{d}_k\|_2] \leq \frac{2^p}{\sigma} \frac{\langle \nabla f(x_k); g_k \rangle}{\|g_k\|_2}$.

Proof. The first statement follows from the facts that, the matrix on the left-hand side of (3.2) is invertible and deterministic under Assumption 4.4 and conditioned on Algorithm 1

having reached iterate x_k at iteration k , and due to the fact that expectation is a linear operator. For the second statement, for any realization of g_k ,

$$\begin{pmatrix} d_k \\ \tilde{d}_k \end{pmatrix} = \begin{pmatrix} H_k & J_k^T \\ J_k & 0 \end{pmatrix}^{-1} \begin{pmatrix} g_k \\ 0 \end{pmatrix} \Rightarrow \| \begin{pmatrix} d_k \\ \tilde{d}_k \end{pmatrix} \|_2 \leq L \| g_k \|_2;$$

where (under Assumption 4.4) $L \geq 2R_{>0}$ is an upper bound on the norm of the matrix shown above. Thus, it follows by Lemma 4.16 that

$$E_k [\| \begin{pmatrix} d_k \\ \tilde{d}_k \end{pmatrix} \|_2^p] \leq E_k [L^p \| g_k \|_2^p] = L^p \frac{1}{p} \frac{1}{I(x_k; k; g_k; d_k)};$$

□

Lemma 4.18. For all $k \geq N$,

$$\begin{aligned} |g_k^T(d_k - \tilde{d}_k)| &\leq \frac{L \|g_k\|_2 \|g_k\|_2^p}{I(x_k; k; g_k; d_k)}; \\ |g_k^T \tilde{d}_k - g_k^T d_k| &\leq \frac{L \|g_k\|_2 \|g_k\|_2^p}{I(x_k; k; g_k; d_k)} + L \|g_k\|_2 \|g_k\|_2^2; \\ |g_k^T(d_k - \tilde{d}_k)| &\leq \frac{L \|g_k\|_2 \|g_k\|_2^2}{I(x_k; k; g_k; d_k)} + \frac{L \|g_k\|_2 \|g_k\|_2^2}{I(x_k; k; g_k; d_k)} \\ &\quad + \frac{L \|g_k\|_2 \|g_k\|_2^2}{I(x_k; k; g_k; d_k)}; \\ |g_k^T(d_k - d_k)| &\leq \frac{L \|g_k\|_2 \|g_k\|_2^p}{I(x_k; k; g_k; d_k)} \\ &\quad + \frac{L \|g_k\|_2 \|g_k\|_2^2}{I(x_k; k; g_k; d_k)} + \frac{L \|g_k\|_2 \|g_k\|_2^2}{I(x_k; k; g_k; d_k)}; \\ |g_k^T d_k - g_k^T d_k| &\leq \frac{L \|g_k\|_2 \|g_k\|_2^p}{I(x_k; k; g_k; d_k)} + 2 \frac{L \|g_k\|_2 \|g_k\|_2^2}{I(x_k; k; g_k; d_k)} \\ &\quad + \frac{L \|g_k\|_2 \|g_k\|_2^2}{I(x_k; k; g_k; d_k)} + \frac{L \|g_k\|_2 \|g_k\|_2^2}{I(x_k; k; g_k; d_k)}; \end{aligned}$$

and $\|J_k(d_k - \tilde{d}_k)\|_1 \leq \frac{L \|g_k\|_2 \|g_k\|_2^p}{I(x_k; k; g_k; d_k)}$;

where $\frac{L}{I} \geq 2R_{>0}$, $\frac{L}{I} \geq \frac{1}{I} \geq 2R_{>0}$, $L \geq 2R_{>0}$, $\frac{L}{I} \geq \frac{1}{I} \geq 2R_{>0}$, $\frac{L}{I} \geq \frac{1}{I} \geq 2R_{>0}$, $\frac{L}{I} \geq \frac{1}{I} \geq 2R_{>0}$, $\frac{L}{I} \geq \frac{1}{I} \geq 2R_{>0}$, $\frac{L}{I} \geq \frac{1}{I} \geq 2R_{>0}$, and $\frac{L}{I} \geq \frac{1}{I} \geq 2R_{>0}$.

Proof. (First inequality) By Assumption 3.1, and (3.2), (3.11), (4.3),

$$\begin{aligned} |g_k^T(d_k - \tilde{d}_k)| &= |j(g_k + J_k^T(y_k - x_k))^T(d_k - \tilde{d}_k)| \\ &= |j(H_k d_k)^T(d_k - \tilde{d}_k)| \leq \|H_k\|_F \|d_k - \tilde{d}_k\|_2 \|d_k - \tilde{d}_k\|_2 \\ &\leq \frac{L \|g_k\|_2 \|g_k\|_2^p}{I(x_k; k; g_k; d_k)} + L \|g_k\|_2 \|g_k\|_2^2; \end{aligned} \tag{4.4}$$

where the result follows using the definition of $\frac{L}{I}$.

(Second inequality) By Lemma 4.17, and (3.2), (3.11), (4.3), (4.4),

$$\begin{aligned}
& j g_k^T d_k - g_k^T d_k j \\
& j (g_k - g_k)^T d_k j + j g_k^T (d_k - d_k) j + j (g_k - g_k)^T (d_k - d_k) j \\
& \leq \frac{L}{r} \|g_k - g_k\|_2^2 + \frac{L}{g;dd} \|g_k - g_k\|_2^2 + \frac{L}{g;dd} \|g_k - g_k\|_2^2 + \frac{L}{g;dd} \|g_k - g_k\|_2^2 + \frac{L}{g;dd} \|g_k - g_k\|_2^2 \\
& \leq \frac{L}{g;dd} \|g_k - g_k\|_2^2 + \frac{L}{g;dd} \|g_k - g_k\|_2^2 + \frac{L}{g;dd} \|g_k - g_k\|_2^2 + \frac{L}{g;dd} \|g_k - g_k\|_2^2 + \frac{L}{g;dd} \|g_k - g_k\|_2^2;
\end{aligned} \tag{4.5}$$

where the result follows using the definition of $g;dd$.

(Third Inequality) The third inequality is proven as follows. By Assumptions 3.1 and 4.2, Lemmas 3.7, 3.8 and 4.17, and (3.2), (3.3), (3.11), (4.3),

$$\begin{aligned}
j g_k^T (d_k - d_k) j &= j (g_k + J_k^T (y_k + \kappa) - J_k^T (y_k + \kappa))^T (d_k - d_k) j \\
& \leq j (g_k + J_k^T (y_k + \kappa))^T (d_k - d_k) j + j (J_k^T (y_k + \kappa))^T (d_k - d_k) j \\
& \leq j (g_k + J_k^T (y_k + \kappa))^T (d_k - d_k) j + j (g_k - g_k)^T (d_k - d_k) j \\
& \quad + j (J_k^T (y_k + \kappa))^T (d_k - d_k) j \\
& \leq j (H_k d_k)^T (d_k - d_k) j + \frac{L}{g;dd} \|g_k - g_k\|_2^2 + \frac{L}{g;dd} \|y_k + \kappa\|_2^2 + \frac{L}{g;dd} \|y_k + \kappa\|_2^2 \\
& \leq \frac{L}{g;dd} \|g_k - g_k\|_2^2 + \frac{L}{g;dd} \|y_k + \kappa\|_2^2 + \frac{L}{g;dd} \|y_k + \kappa\|_2^2 + \frac{L}{g;dd} \|y_k + \kappa\|_2^2 \\
& \quad + \frac{L}{g;dd} \|g_k - g_k\|_2^2;
\end{aligned} \tag{4.6}$$

where the result follows using the definitions of $g;dd$ and $g;dd$.

(Fourth inequality) By Assumptions 3.1 and 4.2, Lemmas 3.7 and 3.8, and (3.2), (3.3), (3.11), (4.3), it follows that

$$\begin{aligned}
j g_k^T (d_k - d_k) j &= j (g_k + J_k^T (y_k + \kappa) - J_k^T (y_k + \kappa))^T (d_k - d_k) j \\
& \leq j (g_k + J_k^T (y_k + \kappa))^T (d_k - d_k) j + j (J_k^T (y_k + \kappa))^T (d_k - d_k) j \\
& \leq j (H_k d_k)^T (d_k - d_k) j + \frac{L}{g;dd} \|y_k + \kappa\|_2^2 + \frac{L}{g;dd} \|y_k + \kappa\|_2^2 \\
& \leq \frac{L}{g;dd} \|y_k + \kappa\|_2^2 + \frac{L}{g;dd} \|y_k + \kappa\|_2^2 + \frac{L}{g;dd} \|y_k + \kappa\|_2^2 + \frac{L}{g;dd} \|y_k + \kappa\|_2^2;
\end{aligned}$$

By the triangle inequality and (4.4),

$$\begin{aligned}
j g_k^T (d_k - d_k) j & \leq j g_k^T (d_k - d_k) j + j g_k^T (d_k - d_k) j \\
& \leq \frac{L}{g;dd} \|g_k - g_k\|_2^2 + \frac{L}{g;dd} \|y_k + \kappa\|_2^2 + \frac{L}{g;dd} \|y_k + \kappa\|_2^2 \\
& \quad + \frac{L}{g;dd} \|y_k + \kappa\|_2^2 + \frac{L}{g;dd} \|y_k + \kappa\|_2^2;
\end{aligned}$$

where the result follows using the definition of $g;dd$ and $g;dd$.

(Fifth inequality) By (4.5), (4.6),

$$\begin{aligned} & |jg_k^T d_k - g_k^T d_{k,j}| \leq |g_k^T d_k - g_k^T d_{k,j}| + |jg_k^T (d_k - d_{k,j})| \\ & \leq \rho_{gg;dd} k g_k + g_k k_2 \sqrt{I(x_k; k; g_k; d_k)} + 2L k g_k + g_k k_2^2 \\ & + \rho_{g;dd} =^2 I(x_k; k; g_k; d_k) + \rho_{g;dd} =^2 I(x_k; k; g_k; d_k); \end{aligned}$$

where the result follows using the definitions of $\rho_{gg;dd}$ and $\rho_{g;dd}$.

(Sixth inequality) By Lemma 3.7, and (3.2), (3.3),

$$k J_k(d_k - d_{k,j}) k_1 = k r_k k_1 \leq a =^2 I(x_k; k; g_k; d_k);$$

where the result follows using the definition of $J;dd$. □

The next lemma provides an upper bound for $|jg_k^T d_{k,j}|$ with respect to reduction in the model of the merit function, i.e., $I(x_k; k; g_k; d_k)$.

Lemma 4.19. For all $k \geq N$, $|jg_k^T d_{k,j}| \leq \rho_{gd;l} I(x_k; k; g_k; d_k)$, where $\rho_{gd;l} = \frac{H}{l;d} + \frac{\rho_{m;y}}{m;l \min} \geq 2R > 0$.

Proof. By Assumption 3.1, Lemma 4.8 (deterministic variant), and (3.11), (4.3),

$$\begin{aligned} |jg_k^T d_{k,j}| &= |j d_k^T H_k d_k + d_k^T J_k^T (y_k + \eta_k)| \leq H k d_k k_2^2 + \eta_k k \alpha_k k_1 \\ &\leq \frac{H}{l;d} + \frac{\rho_{m;y}}{m;l \min} I(x_k; k; g_k; d_k); \end{aligned}$$

where the result follows using the definition of $\rho_{gd;l}$. □

Next we state and prove an upper bound on the difference between the deterministic and stochastic merit parameters. When $k \alpha_k k = 0$, by (3.6)-(3.7), $\eta_k = \eta_k = \eta_{k-1}$. When $k \alpha_k k > 0$, for the ease of exposition, we equivalently reformulate the merit parameter sequence of η_k update (3.6)-(3.7) as

$$\eta_k = \begin{cases} \eta_{k-1} & \text{if } g_k^T d_k + \max f d_k^T H_k d_k; \eta_k d_k k_2^2 g \geq \frac{(1-\beta)(1-\beta_1)(1-\beta_2)k\alpha_k k_1}{k-1}; \\ \frac{(1-\beta)(1-\beta_1)(1-\beta_2)k\alpha_k k_1}{g_k^T d_k + \max f d_k^T H_k d_k; \eta_k d_k k_2^2 g} & \text{otherwise.} \end{cases} \quad (4.7)$$

(The update formula for the deterministic merit parameter η_k can be defined as above with the stochastic quantities replaced by their deterministic counterparts.) It is clear that if the merit parameter is updated from its previous value, i.e., $\eta_k \neq \eta_{k-1}$, then

$$\frac{(1-\beta)(1-\beta_1)(1-\beta_2)k\alpha_k k_1}{g_k^T d_k + \max f d_k^T H_k d_k; \eta_k d_k k_2^2 g} < \eta_{k-1}. \quad (4.8)$$

Similarly, given η_{k-1} if the deterministic merit parameter is updated,

$$\frac{(1-\beta)(1-\beta_1)(1-\beta_2)k\alpha_k k_1}{g_k^T d_k + \max f d_k^T H_k d_k; \eta_k d_k k_2^2 g} < \eta_{k-1}. \quad (4.9)$$

Lemma 4.20. For all $k \geq N$, $j(x_k; g_k; d_k) \leq \frac{2}{3} \rho(x_k; g_k; d_k)$, where $\rho = \frac{2}{3} \rho(x_k; g_k; d_k)$ and $\rho \in (0, \frac{2}{3}]$.

Proof. By the merit parameter updating mechanism (3.6)-(3.7), the merit parameter values (α_k and β_k) are only potentially updated if $\alpha_k \beta_k > 0$. We divide the proof into two cases based the outcome of Line 3 in Algorithm 1 (a) or case (b)). Let $h_k = g_k^T d_k + \max\{d_k^T H_k d_k; \alpha_k d_k^T d_k\}$ and $h_k = g_k^T d_k + \max\{d_k^T H_k d_k; \alpha_k d_k^T d_k\}$.

If Line 3 in Algorithm 1 terminates in case (a), and $\alpha_k \beta_k > 0$, it follows that $h_k > 0$, and by Assumption 4.9, $h_k > 0$; otherwise, when $\alpha_k \beta_k = 0$, by (3.11) and Assumptions 3.1, it follows that $d_k \in \text{Null}(J_k)$ and $h_k = g_k^T d_k + d_k^T H_k d_k = d_k^T J_k^T (y_k + \beta_k) = c_k^T (y_k + \beta_k) = 0$. Thus, neither the stochastic nor the deterministic merit parameters update, i.e., $\alpha_k = \alpha_{k-1}$. The result holds since $\rho(x_k; g_k; d_k)$ is non-negative.

Next, we consider the case where Line 3 in Algorithm 1 terminates with case(b). We divide merit parameter values ($\alpha_k; \beta_k$) into four cases.

Case (i) : Neither α_k or β_k are updated, i.e., $\alpha_k = \alpha_{k-1}$. In this case, the result holds since $\rho(x_k; g_k; d_k)$ is non-negative.

Case (ii) : Both α_k and β_k are updated, i.e., (4.8) and (4.9) hold. By (4.7), (4.8) and (4.9), and Assumption 4.9, it follows that

$$j(x_k; g_k; d_k) = \frac{(1 - \alpha_k)(1 - \beta_k)(1 - \alpha_{k-1})(1 - \beta_{k-1})\alpha_k \beta_k}{h_k} \frac{j(x_k; g_k; d_k)}{h_k} \leq \frac{2}{3} \rho(x_k; g_k; d_k)$$

Case (iii) : The stochastic merit parameter β_k is updated and the deterministic merit parameter α_k is not, i.e., $\alpha_k = \alpha_{k-1}$. Since the stochastic merit parameter is updated, $\beta_k^{\text{trial}} < 1$ and $h_k > 0$, and by Assumption 4.9 it follows that $h_k > 0$, and $\beta_k^{\text{trial}} < 1$. Moreover, since the deterministic merit parameter is not updated,

$$\frac{(1 - \alpha_k)(1 - \beta_k)(1 - \alpha_{k-1})\alpha_k \beta_k}{h_k} \leq \beta_k^{\text{trial}} \leq \frac{2}{3} \rho(x_k; g_k; d_k) \tag{4.10}$$

By (4.7), (4.8) and (4.10), and Assumption 4.9, it follows that

$$j(x_k; g_k; d_k) = \frac{(1 - \alpha_k)(1 - \beta_k)(1 - \alpha_{k-1})\alpha_k \beta_k}{h_k} \leq \frac{(1 - \alpha_k)(1 - \beta_k)(1 - \alpha_{k-1})\alpha_k \beta_k}{h_k} \frac{j(x_k; g_k; d_k)}{h_k} \leq \frac{2}{3} \rho(x_k; g_k; d_k)$$

Case (iv) : The deterministic merit parameter α_k is updated, while the stochastic merit parameter β_k is not, i.e., $\beta_k = \beta_{k-1}$. Since the deterministic merit parameter

is updated, $\alpha_k^{\text{trial}} < 1$ and $h_k > 0$, and by Assumption 4.9 it follows that $h_k > 0$ and $\alpha_k^{\text{trial}} < 1$. Moreover, since the stochastic merit parameter is not updated,

$$\frac{(1 - \alpha_k)(1 - \alpha_1)(1 - \alpha_2)k\alpha_k k_1}{h_k} \leq \alpha_k \quad (4.11)$$

By (4.7), (4.9) and (4.11), and Assumption 4.9, it follows that

$$\begin{aligned} j_k - \alpha_k j &= \alpha_k \frac{(1 - \alpha_k)(1 - \alpha_1)(1 - \alpha_2)k\alpha_k k_1}{h_k} \\ &< \frac{(1 - \alpha_k)(1 - \alpha_1)(1 - \alpha_2)k\alpha_k k_1}{h_k} - \frac{(1 - \alpha_k)(1 - \alpha_1)(1 - \alpha_2)k\alpha_k k_1}{h_k} \\ &= \frac{(1 - \alpha_k)(1 - \alpha_1)(1 - \alpha_2)k\alpha_k k_1}{h_k} \frac{j_k - h_k j}{h_k} \\ &= \frac{(1 - \alpha_k)(1 - \alpha_1)(1 - \alpha_2)k\alpha_k k_1}{(1 - \alpha_k^3)h_k} \frac{j_k - h_k j}{1 - \alpha_k^3} \leq \alpha_k \end{aligned}$$

Combining the four cases above and Lemma 4.19 yields the result. \square

The next lemma bounds the stochastic model of the reduction of the merit function.

Lemma 4.21. For all $k \geq N$ and $\alpha_k \in (0, 1 - 2\sqrt{\frac{1}{\mu} + \frac{1}{\mu} + \frac{1}{\mu}}]$,

$$\begin{aligned} I(x_k; \alpha_k; g_k; d_k) &\leq \frac{1}{1 - \alpha_k} I(x_k; \alpha_k; g_k; d_k) \\ &+ 2\sqrt{\frac{1}{\mu} + \frac{1}{\mu} + \frac{1}{\mu}} k g_k \leq \frac{1}{1 - \alpha_k} I(x_k; \alpha_k; g_k; d_k) + 2\sqrt{\frac{1}{\mu} + \frac{1}{\mu} + \frac{1}{\mu}} k g_k^2; \end{aligned}$$

where $\frac{1}{1 - \alpha_k} = 2\sqrt{\frac{1}{\mu} + \frac{1}{\mu} + \frac{1}{\mu}} \in \mathbb{R}_{>0}$. Additionally, under Assumption 4.1, for all $k \geq N$,

$$E_k I(x_k; \alpha_k; g_k; d_k) \leq \frac{1}{1 - \alpha_k} I(x_k; \alpha_k; g_k; d_k);$$

where $\frac{1}{1 - \alpha_k} = 2\sqrt{\frac{1}{\mu} + \frac{1}{\mu} + \frac{1}{\mu}} \in \mathbb{R}_{>0}$.

Proof. By (3.5), and Lemmas 4.18 and 4.20, it follows that

$$\begin{aligned} I(x_k; \alpha_k; g_k; d_k) &= \alpha_k g_k^T d_k + k\alpha_k k_1 \alpha_k + J_k d_k k_1 \\ &= \alpha_k g_k^T d_k + k\alpha_k k_1 \\ &+ (\alpha_k - \alpha_k) g_k^T d_k + \alpha_k (g_k^T d_k - g_k^T d_k) \alpha_k + J_k d_k k_1 \\ &= I(x_k; \alpha_k; g_k; d_k) + \alpha_k (g_k^T d_k - g_k^T d_k) + \alpha_k J_k d_k k_1 \\ &= I(x_k; \alpha_k; g_k; d_k) + \frac{1}{1 - \alpha_k} I(x_k; \alpha_k; g_k; d_k) \\ &+ 2\sqrt{\frac{1}{\mu} + \frac{1}{\mu} + \frac{1}{\mu}} k g_k \leq \frac{1}{1 - \alpha_k} I(x_k; \alpha_k; g_k; d_k) + 2\sqrt{\frac{1}{\mu} + \frac{1}{\mu} + \frac{1}{\mu}} k g_k^2 \\ &+ \frac{1}{1 - \alpha_k} I(x_k; \alpha_k; g_k; d_k) + \frac{1}{1 - \alpha_k} I(x_k; \alpha_k; g_k; d_k) \end{aligned}$$

Thus, by choosing $\delta \in (0, \frac{1}{2} - \frac{1}{2} \frac{1}{\gamma})$,

$$\begin{aligned}
 I(x_k; k; g_k; d_k) &= 1 + \frac{(\frac{1}{2} + \frac{1}{2} \frac{1}{\gamma} + \frac{1}{2} \frac{1}{\gamma})}{1 - \frac{1}{2} \frac{1}{\gamma}} \frac{= 2}{= 2} I(x_k; k; g_k; d_k) \\
 &+ \frac{\frac{1}{2} \frac{1}{\gamma} k g_k \quad g_k k_2}{1 - \frac{1}{2} \frac{1}{\gamma}} \frac{p}{I(x_k; k; g_k; d_k) + 2} \frac{L k g_k \quad g_k k_2^2}{= 2} \\
 &= 1 + 2 \left(\frac{1}{2} + \frac{1}{2} \frac{1}{\gamma} + \frac{1}{2} \frac{1}{\gamma} \right) \frac{= 2}{= 2} I(x_k; k; g_k; d_k) \\
 &+ 2 \frac{1}{2} \frac{1}{\gamma} k g_k \quad g_k k_2 \frac{p}{I(x_k; k; g_k; d_k) + 2} \frac{L k g_k \quad g_k k_2^2}{= 2} :
 \end{aligned}$$

The first result follows using the definition of γ . Furthermore, by Assumption 4.1,

$$\begin{aligned}
 E_k \left[I(x_k; k; g_k; d_k) \right] &= 1 + 2 \left(\frac{1}{2} + \frac{1}{2} \frac{1}{\gamma} + \frac{1}{2} \frac{1}{\gamma} \right) \frac{= 2}{= 2} I(x_k; k; g_k; d_k) \\
 &+ 2 \frac{1}{2} \frac{1}{\gamma} k g_k \quad g_k k_2 \frac{p}{I(x_k; k; g_k; d_k) + 2} \frac{L k g_k \quad g_k k_2^2}{= 2} \\
 &= 1 + 2 \left(\frac{1}{2} + \frac{1}{2} \frac{1}{\gamma} + \frac{1}{2} \frac{1}{\gamma} \right) \frac{= 2}{= 2} I(x_k; k; g_k; d_k) \\
 &+ 2 \frac{1}{2} \frac{1}{\gamma} \frac{p}{= 2} \frac{= 2}{= 2} I(x_k; k; g_k; d_k) + 2 \frac{L}{= 2} \frac{= 2}{= 2} I(x_k; k; g_k; d_k) \\
 &= 1 + 2 \left(\frac{1}{2} + \frac{1}{2} \left(\frac{1}{\gamma} + \frac{1}{\gamma} + \frac{1}{\gamma} \right) \frac{p}{= 2} \frac{= 2}{= 2} \right) \frac{= 2}{= 2} I(x_k; k; g_k; d_k);
 \end{aligned}$$

where the second result follows using the definition of γ . □

The next lemma bounds the difference in the merit function after a step.

Lemma 4.22. For all $k \geq N$, there exist $\delta \in R_{>0}$ such that

$$E_k [I(x_{k+1}; k+1) - I(x_k; k)] \leq E_k [I(x_k; k)] \inf_{\delta \in R_{>0}} \left(\frac{1}{2} - \delta \right) I(x_k; k; g_k; d_k);$$

Proof. By Lemmas 4.14, 4.15, and 4.21, and (3.4), it follows that for all $k \geq 2$

$$\begin{aligned}
 & E_k [(x_{k+1}; k+1) - (x_k; k)] \\
 = & E_k [(x_{k+1}; k+1) - (x_{k+1}; k) + (x_{k+1}; k) - (x_k; k)] \\
 = & E_k [(x_{k+1}; k) f_{k+1}] + E_k [(x_{k+1}; k) - (x_k; k)] \\
 & E_k [(x_{k+1}; k)] f_{\text{inf}} \\
 & E_k [(x_k; k; g_k; d_k) - (1 - \frac{1}{L}) (x_k; k; g_k; d_k)] \\
 & + E_k [(x_k; k; g_k^T (d_k - \bar{d}_k)) + E_k [(x_k; k; g_k) g_k^T d_k] + E_k [(x_k; k; J_k (d_k - \bar{d}_k)) k_1] \\
 & E_k [(x_{k+1}; k)] f_{\text{inf}} - E_k [(x_k; k; g_k; d_k)] \\
 & + E_k [(1 - \frac{1}{L}) (x_k; k; g_k; d_k)] \\
 & + E_k [(1 - \frac{1}{L})^2 u^{-1} \frac{1}{g; d; d} k g_k - g_k k_2^p \frac{1}{L(x_k; k; g_k; d_k)} + 2 L k g_k - g_k k_2^2] \\
 & + u^{-2} E_k [j g_k^T (d_k - \bar{d}_k)] \\
 & + u^{-2} E_k [j (x_k; k; g_k) g_k^T d_k] + u^{-2} E_k [J_k (d_k - \bar{d}_k) k_1] :
 \end{aligned}$$

Continuing from the above, by Lemmas 4.14, 4.18, 4.20 and 4.21,

$$\begin{aligned}
 & E_k [(x_{k+1}; k+1) - (x_k; k)] \\
 & E_k [(x_{k+1}; k)] f_{\text{inf}} \\
 & E_k [(x_k; k; g_k; d_k) - (1 - \frac{1}{L}) (x_k; k; g_k; d_k)] \\
 & + 2(1 - \frac{1}{L}) u^{-1} \frac{1}{g; d; d} \frac{1}{L} = 2 + 2 L^{-1} (x_k; k; g_k; d_k) \\
 & + u^{-2} \frac{1}{g; d; d} \frac{1}{1 + g; d; d} (x_k; k; g_k; d_k) + \frac{1}{g; d; d} E_k [(x_k; k; g_k; d_k)] \\
 & + u^{-2} (x_k; k; g_k; d_k) + u^{-2} J; d; d E_k [(x_k; k; g_k; d_k)] \\
 & E_k [(x_{k+1}; k)] f_{\text{inf}} \\
 & - (x_k; k; g_k; d_k) + (1 - \frac{1}{L}) u^{-1} (x_k; k; g_k; d_k) \\
 & + 2(1 - \frac{1}{L}) u^{-1} \frac{1}{g; d; d} \frac{1}{1 + 2 L^{-1}} = 2 (x_k; k; g_k; d_k) \\
 & + u^{-2} \frac{1}{g; d; d} \frac{1}{1 + g; d; d} + \frac{1}{g; d; d} (1 - \frac{1}{L}) = 2 (x_k; k; g_k; d_k) \\
 & + u^{-2} (x_k; k; g_k; d_k) + u^{-2} J; d; d (1 - \frac{1}{L}) = 2 (x_k; k; g_k; d_k) \\
 & E_k [(x_{k+1}; k)] f_{\text{inf}} - (x_k; k; g_k; d_k) + 2 (x_k; k; g_k; d_k) :
 \end{aligned}$$

The result follows by setting

$$\begin{aligned} & := u(1 - \frac{1}{L}) \frac{1}{1 + 2^{-1}} \frac{p}{g;dd} \frac{1}{1 + 2^{-L}} + \\ & + \frac{1}{g;dd} \frac{p}{1 + g;dd} + \left(\frac{1}{g;dd} + \frac{J;dd}{1 + J;dd} \right) \frac{1}{1 + \frac{1}{L}} \quad 2R > 0: \end{aligned}$$

□

Lemma 4.22 provides an upper bound on the change in the merit function across a step. The bound has two terms: the first term is related to the difference in the merit parameter across iterations and the second term is negative, conditioned on being sufficiently small, and proportional to the model of the reduction of the merit function. We are now ready to prove the main theorem of this section.

Theorem 4.23. By choosing $\delta \in (0; 1)$; $\min \left\{ \frac{1}{(2-\delta)^2}, \frac{(1-\delta)}{2} \right\}$; $\frac{1}{(2-\delta g;dd)^2}$ for any $\delta \in (0; 1)$,

$$\lim_{k \rightarrow \infty} E \left[\sum_{j=0}^{k-1} l(x_j; j; g_j; d_j) \right] < 1;$$

from which it follows that, $\lim_{k \rightarrow \infty} E[l(x_k; k; g_k; d_k)] = 0$.

Proof. By Lemma 4.22 and $\delta \in (0; 1)$, it follows that

$$\begin{aligned} E_k [l(x_{k+1}; k+1) - l(x_k; k)] & \leq E_k [l(x_{k+1}; k+1)] - f_{\text{inf}} \frac{1}{(2-\delta)} l(x_k; k; g_k; d_k) \\ & \leq E_k [l(x_{k+1}; k+1)] - f_{\text{inf}} \frac{1}{(2-\delta)} l(x_k; k; g_k; d_k); \end{aligned} \tag{4.12}$$

Applying a telescopic sum to (4.12) and taking the total expectation, it follows that

$$\begin{aligned} 1 & < \inf_{\delta} E[l(x_0; 0) - l(x_k; k)] \\ & = E \left[\sum_{j=0}^{k-1} (l(x_{j+1}; j+1) - l(x_j; j)) \right] \\ & \leq E \left[\sum_{j=0}^{k-1} (l(x_{j+1}; j+1) - f_{\text{inf}} \frac{1}{(2-\delta)} l(x_j; j; g_j; d_j)) \right] \\ & \leq \sum_{j=0}^{k-1} E[l(x_{j+1}; j+1)] - f_{\text{inf}} \frac{1}{(2-\delta)} \sum_{j=0}^{k-1} E[l(x_j; j; g_j; d_j)]; \end{aligned}$$

which completes the proof. □

Theorem 4.23 describes the behavior of the model of the reduction of the merit function evaluated at the iterates generated by the Algorithm 1 in expectation. We connect the result of Theorem 4.23 to feasibility and stationarity measures below.

Corollary 4.24. Under the conditions of Theorem 4.23, Algorithm 1 yields a sequence of iterates $\{x_k; y_k\}$ for which

$$\lim_{k \rightarrow \infty} E \|d_k\|_2^2 = 0; \quad \lim_{k \rightarrow \infty} E \|c_k\|_2 = 0; \quad \text{and} \quad \lim_{k \rightarrow \infty} E \|g_k + J_k^T(y_k - x_k)\|_2 = 0:$$

Proof. By the deterministic analogue of Lemma 4.8, $\|x_k - x^*\|_2 \leq \min_k \|x_k - x^*\|_2$. By Theorem 4.23 and the definitions of $\|c_k\|_2$ and $\|d_k\|_2$, we have $\lim_{k \rightarrow \infty} E \|c_k\|_2 = 0$, and thus, our first two results follow from the deterministic variant of Lemma 4.7. The final result follows from (3.11), i.e., $\|g_k + J_k^T(y_k - x_k)\|_2 = \|H_k d_k\|_2 \leq \|H_k\|_2 \|d_k\|_2 \leq L \|d_k\|_2$, which completes the proof. \square

Corollary 4.24 shows that in the limit and in expectation the search direction, the constraint violation and a first-order stationarity measure converge to zero.

Remark 4.25. We make a few remarks about the main theoretical results (Theorem 4.23 and Corollary 4.24).

- ^ Comparison to deterministic results: The result in Corollary 4.24 is similar, albeit in expectation, to what can be proven for an exact deterministic SQP method, i.e. $\bar{g}_k = g_k$ and $\bar{d}_k = d_k$, under the same assumptions; see e.g., [5, 13].
- ^ Comparison to [5]: The main difference in the result of Corollary 4.24 and similar results for the stochastic SQP algorithm proposed in [5] pertain to the requirements on the $\|f_k\|_2$ sequence. In [5] (and other works, e.g., [4, 21, 22]) a diminishing $\|f_k\|_2$ sequence is required to guarantee convergence, whereas in this work convergence with a constant $\|f_k\|_2$ sequence is derived due to the variance reduction achieved.
- ^ Comparison to [28]: In [28], a stochastic line search SQP method for equality constrained problems that utilizes an exact differentiable merit function is proposed. Under deterministic conditions on the function and derivative approximations and exact solutions to the linear systems, the authors show convergence analogous to that of a deterministic algorithm. We note that under the same deterministic conditions, similar results can be established for our proposed adaptive sampling algorithm.
- ^ Comparison to [6]: A result analogous to Theorem 4.23 is proven in [6]. In both works this is possible due to variance reduction in the approximations employed; the algorithm proposed in [6] makes use of predictive variance reduction via SVRG gradients, whereas in this work achieves variance reduction via adaptive sampling.

The final result we show in this section is a complexity result for our proposed algorithm, i.e., the number of iterations required to achieved an ϵ -accurate solution in expectation. Specifically, we consider the following complexity metric,

$$E[kg_k + J_k^T(y_k + \alpha_k)k_2] \leq L; \text{ and } E[k\alpha_k k_1] \leq c; \quad (4.13)$$

for $L \geq 0$ and $c \geq 0$.

Corollary 4.26. Under the conditions of Theorem 4.23, Algorithm 1 generates an iterate $f(x_k; y_k)$ that satisfies (4.13) in at most

$$K = \frac{1}{\epsilon} \frac{(f(x_0) - f_{\text{inf}}) + kc_0k_1}{x} \max\{L^2, c\} \quad (4.14)$$

iterations. Moreover, if $L = \epsilon$ and $c = \epsilon^2$, then $K = O(\frac{1}{\epsilon^2})$.

Proof. First we show that if $E[kg_k + J_k^T(y_k + \alpha_k)k_2] > L$ or $E[k\alpha_k k_1] > c$, then

$$E[\ell(x_k; y_k; g_k; d_k)] \leq \min\{\frac{L}{2}, c\}; \quad (4.15)$$

where $x = \min\{\frac{1}{2}, \frac{\min\{1, d\}}{H}\} \geq R > 0$. Consider arbitrary $(k; L; c) \in \mathbb{N} \times (0; 1) \times (0; 1)$ for which $E[kg_k + J_k^T(y_k + \alpha_k)k_2] > L$ and/or $E[k\alpha_k k_1] > c$. First, suppose that $E[k\alpha_k k_1] > c$. By the deterministic variant of (3.8),

$$E[\ell(x_k; y_k; g_k; d_k)] \leq E[k \max\{d_k^T H_k d_k; d_k d_k k_2^2 g + \alpha_k k_1\}] \leq E[\alpha_k k_1] > c;$$

Next, suppose that $E[kg_k + J_k^T(y_k + \alpha_k)k_2] > L$. By Assumption 3.1 and (3.11),

$$L < E[kg_k + J_k^T(y_k + \alpha_k)k_2] = E[kH_k d_k k_2] \leq E[kH_k d_k k_2];$$

and thus, by the deterministic variant of (3.8), the fact that α_k is bounded below and the definition of d_k , it follows that

$$E[\ell(x_k; y_k; g_k; d_k)] \leq E[k \max\{d_k^T H_k d_k; d_k d_k k_2^2 g + \alpha_k k_1\}] \leq E[\min\{\frac{1}{2}, \frac{d}{H}\} L];$$

Combining the results of the two cases and using the definition of x yields (4.15).

If (4.13) is violated, then by Lemma 4.14, (3.4), (4.12) and (4.15), we have

$$\begin{aligned} & E[\frac{1}{x}(f(x_0) - f_{\text{inf}}) + kc_0k_1] \leq E[\frac{1}{x}(f(x_0) - f_{\text{inf}}) + kc_0k_1] \\ & E[\frac{1}{x}(f(x_0) - f_{\text{inf}}) + kc_0k_1 - \sum_{j=0}^{k-1} \ell(x_j; y_j; g_j; d_j) + \sum_{j=0}^{k-1} \alpha_j k_1 + (x_k - x_0)f_{\text{inf}}] \\ & = E[\sum_{j=0}^{k-1} (\ell(x_j; y_j; g_j; d_j) - \alpha_j k_1) + (x_k - x_0)f_{\text{inf}}] \\ & \leq E[\sum_{j=0}^{k-1} \ell(x_j; y_j; g_j; d_j)] + E[\sum_{j=0}^{k-1} \alpha_j k_1] \\ & \leq k \frac{L}{x} \min\{\frac{L}{2}, c\}; \end{aligned}$$

which implies that k is bounded above by (4.14). \square

The result of Corollary 4.26 is similar to that of deterministic SQP methods, albeit in expectation, under the same assumptions; see e.g., [21, Theorem 1]. To the best of our knowledge, this is the first time that such complexity results have been derived.

4.2 Predetermined Sublinear Errors

In this subsection, we consider the case where the sample sizes used in the gradient estimation are monotonically increased at a sublinear rate and the accuracy in the linear system solves is increased at a sublinear rate to achieve convergence to the solution in expectation. The reason we include this result is to emphasize that predetermined sampling strategies and error sequences suffice to provide convergence guarantees without the need for stronger assumptions. Also, this result provides guidance on the total sampling complexity. For brevity, we only introduce the assumptions and present the main theoretical results. All technical lemmas and proofs are deferred to Appendix A.

The two assumptions below are analogues of Assumptions 4.1 and 4.2.

Assumption 4.27. For all $k \geq N$, the stochastic gradient estimate $g_k \in \mathbb{R}^n$ satisfies, $E_k \|g_k - \bar{g}_k\|_2^2 \leq \frac{1}{(k+1)^\alpha}$, where $\alpha \in \mathbb{R}_{>0}$, $\beta \in (0; 1)$, $\gamma \in \mathbb{R}_{>1}$ and $\delta \in [2; 4]$. Additionally, for all $k \geq N$, the stochastic gradient estimate $\bar{g}_k \in \mathbb{R}^n$ is an unbiased estimator of the gradient of f at x_k , i.e., $E_k [\bar{g}_k] = \bar{g}_k$.

Assumption 4.28. For all $k \geq N$, the search directions $(d_k; \tilde{d}_k) \in \mathbb{R}^n \times \mathbb{R}^m$ in (3.3) (inexact solutions to (3.2)) satisfy, $\|d_k - \bar{d}_k\|_2^2 \leq \frac{\beta}{(k+1)^\alpha}$, where $\beta \in \mathbb{R}_{>0}$, $\gamma \in (0; 1)$, $\delta \in \mathbb{R}_{>1}$ and $\delta \in [2; 4]$. Note, $(\bar{d}_k; \bar{e}_k)$ and $(\bar{d}_k; \bar{e}_k)$ are the exact and inexact solutions of (3.2), respectively.

Next, we state the main results of this subsection. Lemma 4.29 (analogue of Lemma 4.22) provides a bound on the difference of the merit function across iterations.

Lemma 4.29. For all $k \geq N$, there exist constants $(\eta; \zeta) \in \mathbb{R}_{>0} \times \mathbb{R}_{>0}$ such

$$E_k [f(x_{k+1}; \tilde{d}_{k+1}) - f(x_k; \tilde{d}_k)] \leq E_k [f(x_{k+1}; \tilde{d}_{k+1}) - f(x_k; \tilde{d}_k)] + \frac{\eta}{(k+1)^\alpha};$$

(Note, the definitions of the constants are given in Appendix A)

Lemma 4.29 has an additional as compared to Lemma 4.22. This is due to the fact that we control the gradient error and linear system accuracy at a sublinear rate instead of controlling it relative to the algorithmic progress, as measured in terms of $\|l(x_k; \tilde{d}_k; g_k; d_k)\|$.

That being said, the additional term is proportional to $\frac{1}{\epsilon^2}$ and whose accumulation is finite in the limit.

Theorem 4.30 and Corollary 4.31 (analogues of Theorem 4.23 and Corollary 4.26, respectively) provide convergence and iteration complexity results, respectively. Corollary 4.31 also provides sample complexity results.

Theorem 4.30. By choosing $\epsilon \in (0, 1)$ and $\eta \in (0, \frac{1}{2L})$, $\lim_{k \rightarrow \infty} \mathbb{E} \left[\sum_{j=0}^{k-1} \mathbb{E} \left[\ell(x_j; y_j; g_j; d_j) \right] \right] < \epsilon$, from which it follows that, $\lim_{k \rightarrow \infty} \mathbb{E} \left[\ell(x_k; y_k; g_k; d_k) \right] = 0$.

The conclusion of Theorem 4.30 is the same as that of Theorem 4.23, up to constants, and, as a result, Corollary 4.24 holds for Theorem 4.30.

Corollary 4.31. Under the conditions of Theorem 4.30, Algorithm 1 generates an iterate $(x_k; y_k; g_k; d_k)$ that satisfies (4.13) in at most $K = O\left(\frac{1}{\epsilon^2} \max\{L^2, c^{-1}\}\right)$ iterations and $W = O\left(\frac{1}{\epsilon^2} \max\{L^2, c^{-1}\}^{(1+\alpha)}\right)$ stochastic gradient evaluations. Moreover, if $L = \frac{1}{\epsilon}$ and $c = \frac{1}{\epsilon^2}$, then $K = O\left(\frac{1}{\epsilon^2}\right)$ and $W = O\left(\frac{1}{\epsilon^{2(1+\alpha)}}\right)$.

Corollary 4.31 shows that one can achieve the same iteration complexity as the deterministic variant of the algorithm under Assumptions 4.27 and 4.28 (and other assumptions stated earlier) at an increased overall sample complexity.

5 Practical AIS-SQP

In this section, we present our proposed practical adaptive inexact stochastic SQP method (PAIS-SQP) and describe the sample size selection mechanism, iterative linear system solver and early termination conditions employed.

5.1 Sample Size Selection

We describe the mechanism by which the sample size is selected at every iteration. Condition (4.1) involves computing population variances and deterministic quantities which are not available in our setting, and possibly requires solving multiple linear systems. That being said, one can approximate these quantities with sample variances and sampled stochastic counterparts of the deterministic quantities required following the ideas proposed in [12, 9].

Condition (4.1) is approximated as follows. Let $g_k \in \mathbb{R}^n$ be defined as

$$g_k := \frac{1}{\sqrt{S_k}} \sum_{i \in S_k} r F(x_k; y_i); \quad (5.1)$$

where S_k is a set consisting of indices drawn at random from the distribution of \mathcal{I} , and i_k is a realization of \mathcal{I} . The left-hand-side of (4.1) can be expressed as

$$E_k \|g_k - r - f(x_k)\|_2^2 = \frac{E_k [\|r - F(x_k; i_k)\|_2^2]}{|S_k|}. \quad (5.2)$$

Computing the population variance on the right-hand-side of (5.2) is prohibitively expensive in our setting, and thus we approximate it with the sample variance (see left-hand-side of (5.3)). Moreover, the right-hand-side of (4.1) is approximated with its stochastic counter-part. This results in the following approximation to condition (4.1),

$$\frac{\text{Var}_{i \in S_k} [r - F(x_k; i)]}{|S_k|} \leq \frac{1}{|S_k|} \sum_{i \in S_k} \|r - F(x_k; i; g_k; d_k)\|_2^2, \quad (5.3)$$

where $\text{Var}_{i \in S_k} [r - F(x_k; i)] = \frac{1}{|S_k|} \sum_{i \in S_k} \|r - F(x_k; i; g_k; d_k)\|_2^2$, used in PAIS-SQP

In our practical algorithm (Algorithm 2), if inequality (5.3) is not satisfied (given the set S_k), we increase the sample size to a size that (at least with high probability) will satisfy (4.1). The heuristic we propose to do this is as follows. Suppose we wish to find a larger sample \hat{S}_k ($|\hat{S}_k| > |S_k|$), and let us assume that the change in sample size is gradual enough that for any x_k , such that $\text{Var}_{i \in S_k} [r - F(x_k; i)] \approx \text{Var}_{i \in \hat{S}_k} [r - F(x_k; i)]$ and $\|r - F(x_k; i; g_k; d_k)\|_2 \approx \|r - F(x_k; i; \hat{g}_k; \hat{d}_k)\|_2$, where $\hat{g}_k; \hat{d}_k$ are stochastic realizations of these quantities computed based on the sample \hat{S}_k . Under this assumption, it is clear that (5.3) is satisfied if

$$|\hat{S}_k| \geq \frac{1}{\alpha} \frac{\text{Var}_{i \in S_k} [r - F(x_k; i)]}{\frac{1}{|S_k|} \sum_{i \in S_k} \|r - F(x_k; i; g_k; d_k)\|_2^2} \cdot m. \quad (5.4)$$

This is the final level of approximation where we check (5.4) to set the sample size for the next iteration, i.e., $S_{k+1} = \hat{S}_k$. The ideas above are used in Algorithm 2.

5.2 Inexact Linear System Solutions

Our proposed practical algorithm makes use of an iterative solver with early termination tests to solve the Newton-SQP linear system (3.3). We use the minimum residual (MINRES) method [18, 30], with early termination conditions, for solving the system inexactly, but note that other iterative algorithms could also be used. For all $k \geq 2$, let $(d_{k,t}; r_{k,t}; g_{k,t}; \rho_{k,t})_{t \in \mathcal{I}_k}$ denote the steps (and residuals) generated in iteration $t \in \mathcal{I}_k$ of MINRES, and $(d_k; r_k; g_k; \rho_k) = (d_{k,t^0}; r_{k,t^0}; g_{k,t^0}; \rho_{k,t^0})$ where t^0 is the last MINRES iteration. The MINRES method was terminated for the minimum t such that either condition (a)

$$\|r_k\|_2 < \alpha \|g_k\|_2 \quad \text{and} \quad \|g_k\|_2 < \beta \|r_k\|_2 \quad \text{with} \quad \alpha = \frac{1}{\alpha} \quad \text{and} \quad \beta = \frac{1}{\beta} \quad (5.5)$$

or, condition (b)

$$\|r_k\|_2 < \min(\alpha \|g_k\|_2, \beta \|r_k\|_2) \quad \text{and} \quad \|g_k\|_2 < \min(\alpha \|g_k\|_2, \beta \|r_k\|_2) \quad (5.6)$$

hold. These conditions are inspired by the theory, but relaxed for practicality. Specifically, the additional condition in case (a) is not checked and neither is Assumption 4.9.

5.3 PAIS-SQP

In this section, we present our practical algorithm PAIS-SQP

Algorithm 2 (PAIS-SQP) Practical, Adaptive, Inexact, Stochastic SQP Algorithm

Require: $x_0 \in \mathbb{R}^n$; $y_0 \in \mathbb{R}^m$; $f, H_k \in \mathbb{S}^n$; $\alpha \in \mathbb{R}_{>0}$; $\mathcal{S}_1 \subseteq \mathbb{N}$; $f_1, f_2, \dots, g \in (0, 1)$;
 $\beta_a \in \mathbb{R}_{>0}$; $\beta_b \in \mathbb{R}_{>0}$; $\gamma \in (0, 1]$; $\mu \in \mathbb{R}_{>0}$; $\delta \in (0, \infty)$; $\tau \in [2, 4]$

- 1: for all $k \in \mathbb{N}$ do
- 2: Compute g_k via (5.1) with $S_k = \mathcal{S}_{k-1}$
- 3: Solve (3.2) iteratively using MINRES;
 Compute a step $(d_k; \kappa)$ that satisfies either (5.5) or (5.6)
- 4: Update κ via (3.6) or (3.7)
- 5: Compute a step size α_k via (3.9) or (3.10)
- 6: Update $x_{k+1} = x_k + \alpha_k d_k$, and $y_{k+1} = y_k + \alpha_k \kappa$
- 7: Choose a sample \mathcal{S}_k such that $j_{\mathcal{S}_k} = j_{S_k}$
- 8: If condition (5.3) is not satisfied, augment \mathcal{S}_k using formula (5.4)
- 9: end for

6 Numerical Results

The main goal of this section is to illustrate the efficiency of our proposed practical algorithm (Algorithm 2). We characterize efficiency in terms of two metrics that capture the major costs in solving (2.1). The first metric is the number of objective function gradient evaluations (or accessed data points in the context of the machine learning problems presented below). The second metric is the number of iterations of an iterative solver used to solve the linear system (3.2). To illustrate the efficiency of Algorithm 2, we present results on two classes of problems, constrained classification problems that arise in machine learning (Section 6.2) and standard CUTE collection of nonlinear optimization problems (Section 6.3), and compare exact and inexact variants (all implementation details are given in Section 6.1).

In both Sections 6.2 and 6.3, results are given in terms of feasibility and stationarity errors defined as, $\|q_k\|_1$ and $\|g_k + J_k^T y_k\|_1$, respectively, where the vector $y_k \in \mathbb{R}^m$ is computed as a least-squares multiplier of the above using the true gradient g_k . Moreover, in both sections we terminate all algorithms solely due to iteration, sampled gradient evaluation or linear system iteration budgets.

6.1 Implementation Details

We compare different variants of Algorithm 2. Specifically, we compare variants with different levels of accuracy in the gradient approximations employed, as well as variants with and without the early termination conditions. Precise characterization of the accuracy

levels are given in Sections 6.2 and 6.3. For all variants MINRES was used to solve the linear systems, and all variants employed the same adaptive step size selection strategy described by (3.9)-(3.10). For all problems estimates of L and μ were computed using gradient differences around the initial point, and kept constant throughout the course of optimization. This procedure was performed in such a way that, for each problem instance, all algorithms used the same values for these estimates. For all algorithms $\beta_k = 1$ for all $k \geq 2$, $\alpha_1 = 1$, $\alpha_2 = 1$, $\alpha_3 = 10^2$, $\alpha_4 = 10^2$, $\alpha_5 = 10^2$, $\alpha_6 = 10^2$, $\alpha_7 = 10^2$, $\alpha_8 = 10^2$, $\alpha_9 = 10^2$, $\alpha_{10} = 10^2$, $\alpha_{11} = 10^2$, $\alpha_{12} = 10^2$, $\alpha_{13} = 10^2$, $\alpha_{14} = 10^2$, $\alpha_{15} = 10^2$, $\alpha_{16} = 10^2$, $\alpha_{17} = 10^2$, $\alpha_{18} = 10^2$, $\alpha_{19} = 10^2$, $\alpha_{20} = 10^2$, $\alpha_{21} = 10^2$, $\alpha_{22} = 10^2$, $\alpha_{23} = 10^2$, $\alpha_{24} = 10^2$, $\alpha_{25} = 10^2$, $\alpha_{26} = 10^2$, $\alpha_{27} = 10^2$, $\alpha_{28} = 10^2$, $\alpha_{29} = 10^2$, $\alpha_{30} = 10^2$, $\alpha_{31} = 10^2$, $\alpha_{32} = 10^2$, $\alpha_{33} = 10^2$, $\alpha_{34} = 10^2$, $\alpha_{35} = 10^2$, $\alpha_{36} = 10^2$, $\alpha_{37} = 10^2$, $\alpha_{38} = 10^2$, $\alpha_{39} = 10^2$, $\alpha_{40} = 10^2$, $\alpha_{41} = 10^2$, $\alpha_{42} = 10^2$, $\alpha_{43} = 10^2$, $\alpha_{44} = 10^2$, $\alpha_{45} = 10^2$, $\alpha_{46} = 10^2$, $\alpha_{47} = 10^2$, $\alpha_{48} = 10^2$, $\alpha_{49} = 10^2$, $\alpha_{50} = 10^2$. For PAIS-SQP we additionally set $\alpha_1 = 0.99$. (Note, α_2 and α_3 are not required by the PAIS-SQP algorithm.) The initial sample size j_k is defined in Sections 6.2 and 6.3. For variants that solve the linear system without the early termination conditions, the termination tolerance used in MINRES was $\epsilon_{\text{MINRES}} = 10^{-8}$, in order to obtain accurate solutions.

6.2 Constrained Logistic Regression

In our first set of experiments, we consider constrained logistic regression problems,

$$\min_{x \in \mathbb{R}^n} f(x) = \frac{1}{N} \sum_{i=1}^N \log(1 + e^{y_i(X_i^T x)}) \quad \text{st: } Ax = b_1; \|x\|_2 = b_2; \quad (6.1)$$

where $X \in \mathbb{R}^{n \times N}$ is the data matrix, X_i is the i th column of matrix X , $y \in \{-1, 1\}^N$ contains corresponding label data, $A \in \mathbb{R}^{m \times n}$, $b_1 \in \mathbb{R}^m$ and $b_2 \in \mathbb{R}$. We present results on two data sets from the LIBSVM collection [17]; `australian` and `mushroom` (Results on more data sets from the LIBSVM collection are presented in Appendix B). For the linear constraints, the data was generated as follows: the entries of the matrix A and the vector b_1 were drawn from a standard normal distribution (for each data set the same A and b_1 were used for all methods). For the ℓ_2 -norm constraint, $b_2 = 1$. For all problems and algorithms, the initial primal iterate was set to the vector of all ones of appropriate dimension, and the initial dual variables y_0 were set as the least squares multipliers.

For each data set, we consider exact and inexact (linear system solutions) variants, i.e., variants with and without the early termination conditions, and three different sample sizes $\{S_k\} = \{j_k\} \in \{2, 12, N\}$ for all $k \geq 2$ for a total of 6 variants, and compare against PAIS-SQP. A budget of 50 epochs was given to every method. The results for the two data sets are presented in Figures 1 and 2. For every method, we report the feasibility and stationarity errors in terms of iterations, epochs (gradient evaluations) and linear system iterations. The results indicate that our proposed practical inexact SQP method PAIS-SQP strikes a good balance between reducing constraint violation while attempting to find a point that satisfies approximate first-order stationarity across all three evaluation metrics. In Figure 3 we show the sample size and step size selected by the different variants. While the sample sizes increase relatively quickly, there are significant savings that can be achieved by employing inexact information. The step size figures illustrate that the adaptive step size mechanism (3.9)-(3.10) is stable.

(a) Feasibility vs. Iterations (b) Feasibility vs. Epochs (c) Feasibility vs. LS Iters

(d) Stationarity vs. Iterations (e) Stationarity vs. Epochs (f) Stationarity vs. LS Iters

Figure 1: `australian` : Feasibility & stationarity errors versus iterations/epochs/linear system iterations for exact and inexact variants of Algorithm 2 on (6.1).

(a) Feasibility vs. Iterations (b) Feasibility vs. Epochs (c) Feasibility vs. LS Iters

(d) Stationarity vs. Iterations (e) Stationarity vs. Epochs (f) Stationarity vs. LS Iters

Figure 2: mushroom Feasibility/stationarity errors versus iterations/epochs/linear system iterations for exact and inexact variants of Algorithm 2 on (6.1).

(a) Step Size vs. Iterations (b) Batch Size vs. Iterations (c) Step Size vs. Iterations (d) Batch Size vs. Iterations

Figure 3: australian ((a), (b)); mushroom((c), (d)): Step sizes and batch sizes versus iterations.

6.3 CUTE Problems

Next, we consider equality constrained problems from the CUTE collection of nonlinear optimization problems [11]. Specifically, of the 123 such problems in the collection we considered 49 problems. (We only used those for which: (i) LICQ held for all iterations of all algorithms, (ii) f is not a constant function, and (iii) $n + m = 1000$.) We used the prescribed starting point for all problems and all algorithms. The CUTE problems are deterministic, so we added noise to the gradient computations to make the problems stochastic. Specifically, we consider additive noise where the gradient was computed as

$$g_k = \frac{1}{\sqrt{S_k}} \sum_{i \in S_k} (r f(x_k) + N(0; N_{i,i} I));$$

where $N_{i,i} = 10^{-1}$ for all i , and S_k is prescribed by the variant and determines the level of noise.

For each problem, we again consider variants that compute exact and inexact (early termination conditions) linear system solutions. We compare PAIS-SQP to non-adaptive sampling variants with $|S_k| = |S| \in \{2, 128, 1024\}$ for all $k \in N$, and limit the maximum sample size employed by PAIS-SQP to 1024. For each problem, we ran 10 instances with different random seeds. This led to a total of 490 runs of each algorithm for each noise level. We terminated the methods on the following budget: $1024 \cdot 10^3$ gradient evaluations or $1024 \cdot 10^2$ linear system iterations (whichever comes first).

The results of these experiments are reported in Figure 4 in the form of performance profiles [27]. We present results in terms of feasibility and stationarity with respect to gradient evaluations and linear system iterations. The performance profiles were constructed as follows. For each problem, method and seed, the iterate used in the performance profile x_{pp} was chosen as: either the point with minimum $\|g_k + J_k^T y_k\|_1$ among all points with $\|c_k\|_1 \leq 10^{-6}$, or if no such point exists, then the point with minimum $\|c_k\|_1$. Following [27], for the two metrics an algorithm was deemed to have solved a given problem for a given seed if $m(x_0) - m(x_{pp}) \leq (1 - \text{pp})(m(x_0) - m(x_b))$, where $m(x_1)$ is $\|g_1 + J_1^T y_1\|_1$ (for stationarity) and $\|c_1\|_1$ (for feasibility), respectively, $m(x_b)$ denotes the best possible value of either metric for each problem and seed, and $\text{tolerance}_{pp} \in (0; 1)$. Overall, across all tolerances and metrics, the PAIS-SQP method appears to be the most robust (as seen by the right-most points on the figures). The ability of PAIS-SQP to make sufficient progress with inexact information in the initial stages of the optimization, combined with its ability to increase accuracy of the approximations employed, as needed, as the optimization progresses allows the algorithm to balance convergence and cost. As a result, PAIS-SQP is efficient and robust in terms of all metrics.

(a) Feas. vs. Grad. (b) Feas. vs. LS Iters (c) Stat. vs. Grad. (d) Stat. vs. LS Iters

(e) Feas. vs. Grad. (f) Feas. vs. LS Iters (g) Stat. vs. Grad. (h) Stat. vs. LS Iters

(i) Feas. vs. Grad. (j) Feas. vs. LS Iters (k) Stat. vs. Grad. (l) Stat. vs. LS Iters

Figure 4: CUTE Performance profiles for exact and inexact variants of Algorithm 2 on CUTE collection. First row accuracy $\rho_{pp} = 10^{-1}$; Second row accuracy $\rho_{pp} = 10^{-3}$; Third row accuracy $\rho_{pp} = 10^{-5}$.

7 Final Remarks

In this paper, we have designed and analyzed a stochastic SQP algorithm (AIS-SQP) for solving optimization problems involving deterministic nonlinear equality constraints and a stochastic objective function. At each iteration, the AIS-SQP method computes a stochastic approximation of the gradient of the objective function, computes a step by solving a stochastic Newton-SQP linear system inexactly, potentially updates the merit parameter, and adaptively selects a step size and updates the iterate. Our algorithm is adaptive in several ways. We have proposed accuracy conditions for the stochastic gradient approximation and the quality of the linear system solutions. Moreover, we have proposed adaptive update rules for the merit and step size parameters. Our algorithmic development and analysis has revealed an intrinsic relationship between the accuracy of stochastic objective gradient realizations, the quality of inexact solutions to the Newton-SQP linear systems, and the adaptive step sizes selected. That is, higher accuracy in the gradient estimation

and linear system solution can potentially lead to the acceptance of larger step sizes as in the deterministic counterparts.

We have proved that our algorithm generates a sequence of iterates whose first-order stationarity measure converges to zero in expectation. While similar results have been established in the literature, e.g., [5, 21], these works only consider asymptotic regimes after which the penalty parameter is sufficiently small and has stabilized. In this work, we have analyzed the complete behavior of the algorithm across all potential merit parameter changes, under an additional assumption, and have provided iteration complexity analysis for AIS-SQP, which matches that of deterministic SQP methods [21], in expectation. We have also established sublinear (gradient) sample complexity results of the proposed algorithm when the gradient and linear system accuracies are controlled at predetermined sublinear rates.

Inspired by AIS-SQP, we have developed, implemented and tested a practical variant PAIS-SQP of the adaptive stochastic SQP method. Our results on two different sets of experiments, constrained logistic regression and standard nonlinear optimization test problems, suggest that our practical algorithm strikes a good balance between minimizing constraint violation while also minimizing the objective function in terms of importance evaluation metrics such as iterations, gradient evaluations and linear system iterations.

Acknowledgements

We would like to thank Professor Frank E. Curtis for all his useful suggestions and feedback. Moreover, we would like to thank the Office of Naval Research for their support of this project (award number: N00014-21-1-2532).

References

- [1] Joshua Achiam, David Held, Aviv Tamar, and Pieter Abbeel. Constrained policy optimization. In International Conference on Machine Learning, pages 22{31. PMLR, 2017.
- [2] Florian Beiser, Brendan Keith, Simon Urbainczyk, and Barbara Wohlmuth. Adaptive sampling strategies for risk-averse stochastic optimization with constraints. arXiv preprint arXiv:2012.03844, 2020.
- [3] Albert S Berahas, Liyuan Cao, and Katya Scheinberg. Global convergence rate analysis of a generic line search algorithm with noise. SIAM Journal on Optimization , 31(2):1489{1518, 2021.
- [4] Albert S Berahas, Frank E Curtis, Michael J O'Neill, and Daniel P Robinson. A stochastic sequential quadratic optimization algorithm for nonlinear equality constrained optimization with rank-de cient jacobians. arXiv preprint arXiv:2106.13015, 2021.
- [5] Albert S Berahas, Frank E Curtis, Daniel Robinson, and Baoyu Zhou. Sequential quadratic optimization for nonlinear equality constrained stochastic optimization. SIAM Journal on Optimization , 31(2):1352{1379, 2021.
- [6] Albert S Berahas, Jiahao Shi, Zihong Yi, and Baoyu Zhou. Accelerating Stochastic Sequential Quadratic Programming for Equality Constrained Optimization using Predictive Variance Reduction. arXiv preprint arXiv:2204.04161 , 2022.
- [7] Dimitri Bertsekas. Network optimization: continuous and discrete models Athena Scientific, 1998.
- [8] John T Betts. Practical methods for optimal control and estimation using nonlinear programming. SIAM, 2010.
- [9] Raghu Bollapragada, Richard Byrd, and Jorge Nocedal. Adaptive sampling strategies for stochastic optimization. SIAM Journal on Optimization , 28(4):3312{3343, 2018.
- [10] Raghu Bollapragada, Jorge Nocedal, Dheevatsa Mudigere, Hao-Jun Shi, and Ping Tak Peter Tang. A progressive batching l-bfgs method for machine learning. In International Conference on Machine Learning, pages 620{629. PMLR, 2018.
- [11] Ingrid Bongartz, Andrew R Conn, Nick Gould, and Ph L Toint. Cute: Constrained and unconstrained testing environment. ACM Transactions on Mathematical Software (TOMS) , 21(1):123{160, 1995.

- [12] Richard H Byrd, Gillian M Chin, Jorge Nocedal, and Yuchen Wu. Sample size selection in optimization methods for machine learning. *Mathematical programming*, 134(1):127{155, 2012.
- [13] Richard H Byrd, Frank E Curtis, and Jorge Nocedal. An inexact sqp method for equality constrained optimization. *SIAM Journal on Optimization* , 19(1):351{369, 2008.
- [14] Richard H Byrd, Frank E Curtis, and Jorge Nocedal. An inexact newton method for nonconvex equality constrained optimization. *Mathematical programming*, 122(2):273{299, 2010.
- [15] Richard G Carter. On the global convergence of trust region algorithms using inexact gradient information. *SIAM Journal on Numerical Analysis* , 28(1):251{265, 1991.
- [16] Coralia Cartis and Katya Scheinberg. Global convergence rate analysis of unconstrained optimization methods based on probabilistic models. *Mathematical Programming*, 169(2):337{375, 2018.
- [17] Chih-Chung Chang and Chih-Jen Lin. LIBSVM: a library for support vector machines. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 2(3):1{27, 2011.
- [18] Sou-Cheng T Choi, Christopher C Paige, and Michael A Saunders. Minres-qlp: A krylov subspace method for inde nite or singular symmetric systems. *SIAM Journal on Scienti c Computing* , 33(4):1810{1836, 2011.
- [19] Andrew Cotter, Maya Gupta, and Jan Pfeifer. A light touch for heavily constrained sgd. In *Conference on Learning Theory*, pages 729{771. PMLR, 2016.
- [20] Frank E Curtis, Jorge Nocedal, and Andreas Wachter. A matrix-free algorithm for equality constrained optimization problems with rank-de cient jacobians. *SIAM Journal on Optimization*, 20(3):1224{1249, 2010.
- [21] Frank E Curtis, Michael J O'Neill, and Daniel P Robinson. Worst-case complexity of an sqp method for nonlinear equality constrained stochastic optimization. *arXiv preprint arXiv:2112.14799*, 2021.
- [22] Frank E Curtis, Daniel P Robinson, and Baoyu Zhou. Inexact sequential quadratic optimization for minimizing a stochastic objective function subject to deterministic nonlinear equality constraints. *arXiv preprint arXiv:2107.03512*, 2021.
- [23] Michael P Friedlander and Mark Schmidt. Hybrid deterministic-stochastic methods for data tting. *SIAM Journal on Scienti c Computing* , 34(3):A1380{A1405, 2012.

- [24] Fatemeh S Hashemi, Soumyadip Ghosh, and Raghu Pasupathy. On adaptive sampling rules for stochastic recursions. In *Proceedings of the Winter Simulation Conference 2014*, pages 3959{3970. IEEE, 2014.
- [25] Billy Jin, Katya Scheinberg, and Miaolan Xie. High probability complexity bounds for line search based on stochastic oracles. *arXiv preprint arXiv:2106.06454*, 2021.
- [26] Mehrdad Mahdavi, Tianbao Yang, Rong Jin, Shenghuo Zhu, and Jinfeng Yi. Stochastic gradient descent with only one projection. *Advances in neural information processing systems* 25:494{502, 2012.
- [27] Jorge J Moré and Stefan M Wild. Benchmarking derivative-free optimization algorithms. *SIAM Journal on Optimization* , 20(1):172{191, 2009.
- [28] Sen Na, Mihai Anitescu, and Mladen Kolar. An adaptive stochastic sequential quadratic programming with differentiable exact augmented lagrangians. *arXiv preprint arXiv:2102.05320*, 2021.
- [29] Yatin Nandwani, Abhishek Pathak, Mausam Singla, and Parag Singla. A primal dual formulation for deep learning with constraints. In *Advances in Neural Information Processing Systems* pages 12157{12168, 2019.
- [30] Christopher C Paige and Michael A Saunders. Solution of sparse indefinite systems of linear equations. *SIAM journal on numerical analysis*, 12(4):617{629, 1975.
- [31] Raghu Pasupathy, Peter Glynn, Soumyadip Ghosh, and Fatemeh S Hashemi. On sampling rates in simulation-based recursions. *SIAM Journal on Optimization* , 28(1):45{73, 2018.
- [32] Sathya N Ravi, Tuan Dinh, Vishnu Suresh Lokhande, and Vikas Singh. Explicitly imposing constraints in deep networks via conditional gradients gives improved generalization and faster convergence. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 4772{4779, 2019.
- [33] Tyrone Rees, H Sue Dollar, and Andrew J Wathen. Optimal solvers for pde-constrained optimization. *SIAM Journal on Scientific Computing* , 32(1):271{298, 2010.
- [34] Lloyd N Trefethen and David Bau III. *Numerical linear algebra*, volume 50. Siam, 1997.
- [35] Yuchen Xie, Raghu Bollapragada, Richard Byrd, and Jorge Nocedal. Constrained and composite optimization via adaptive sampling methods. *arXiv preprint arXiv:2012.15411*, 2020.

- [36] Yinhao Zhu, Nicholas Zabaras, Phaedon-Stelios Koutsourelakis, and Paris Perdikaris. Physics-constrained deep learning for high-dimensional surrogate modeling and uncertainty quantification without labeled data. *Journal of Computational Physics*, 394:56{81, 2019.

A Technical Results (Section 4.2)

Here, we present the complete theoretical results for Algorithm 1 under Assumptions 4.27 and 4.28. We present all the technical lemmas required to establish the main theoretical results presented in Section 4.2. We unify the assumptions in the following assumption.

Assumption A.1. There exist universal quantities such that Assumptions 2.1, 3.1, 3.3, 4.27 and 4.28 hold for any realization of Algorithm 1.

Lemmas 4.5-4.8 from Section 4 hold without change. We unify them in the lemma below for completeness of the analysis under Assumption A.1.

Lemma A.2. Under Assumption A.1, for all $k \geq N$, (3.2) has a unique solution. In addition, for the same constants $(\underline{v}; \underline{uv}; \underline{\gamma}; \underline{\mu}) \in \mathbb{R}_{>0} \times \mathbb{R}_{>0} \times \mathbb{R}_{>0} \times \mathbb{R}_{>0}$ that appear in the respective Lemmas (listed in parenthesis below), the following statements hold true for all $k \geq N$:

- (Lemma 4.5) The normal component v_k is bounded as $\max\{k v_k k_2; k v_k k_2^2\} \leq \underline{v} \max\{k \alpha_k k_2; k r_k k_2\} g$.
- (Lemma 4.6) If $k u_k k_2^2 \leq \underline{uv} k v_k k_2^2$, then $d_k^T H_k d_k \leq \frac{1}{2} k u_k k_2^2$ and $d_k^T H_k d_k \geq \underline{\mu} k d_k k_2^2$ where $\underline{\mu} \in (0; \frac{1}{2})$ is an user-defined parameter in Algorithm 1.
- (Lemma 4.7) The search direction satisfies $k d_k k_2^2 \leq \underline{\gamma} k$ and $k d_k k_2^2 + k \alpha_k k_2 \leq (\underline{\mu} + 1) k$.
- (Lemma 4.8) The model reduction satisfies $|l(x_k; k; g_k; d_k) - l_k| \leq \underline{\mu} k$.

Additionally, if Assumption 4.9 holds, then the following results also hold, whose proofs are from the corresponding lemmas in parentheses.

Lemma A.3. Under Assumptions A.1 and 4.9, the following statements hold true for all $k \geq N$ with the same constants appearing in the respective Lemmas (listed in parenthesis below).

- (Lemma 4.11) There exists a constant $\underline{\mu} \in \mathbb{R}_{>0}$ such that Algorithm 1 generates a sequence of $\{k, g_k\}$, where $k \geq \underline{\mu}$.
- (Lemma 4.12) There exist constants $\underline{\mu}; \underline{\mu}_d; \underline{\mu}_g \in \mathbb{R}_{>0}$ such that for all $k \geq N$, $|l(x_k; k; g_k; d_k) - (k L_k + k) k d_k k_2^2| \leq \underline{\mu}_d k d_k k_2^2$.
- (Lemma 4.14) For all $k \geq N$, there exists a constant $\underline{\mu} \in \mathbb{R}_{>0}$ such that, $\underline{\mu} \leq k$ ($\underline{\mu} = 2$).
- (Lemma 4.15) Finally, for all $k \geq N$, it follows that

$$\begin{aligned} & (x_k + k d_k; k) - (x_k; k) \\ & \leq k |l(x_k; k; g_k; d_k) - (k L_k + k) k d_k k_2^2| + \underline{\mu}_d k d_k k_2^2 \\ & + k g_k^T (d_k - \underline{\mu} d_k) + k (k - \underline{\mu}) g_k^T d_k + k J_k (d_k - \underline{\mu} d_k) k_1 \end{aligned} \quad (\text{A.1})$$

We now proceed to state and prove a series of Lemmas that are analogues to those proven in Section 4.1. For brevity, in the rest of this appendix, we always assume that Assumptions A.1 and 4.9 hold without explicitly stating this fact. We begin by bounding the variance in the stochastic gradient approximations and the variance in the search directions computed in Lemmas A.4 and A.5.

Lemma A.4. For all $k \geq N$, $E_k [k g_k - g_k k_2] \leq \frac{p-1}{(k+1)^2}$.

Proof. By Assumption 4.27 and Jensen's inequality, we have

$$E_k [k g_k - g_k k_2] \leq \sqrt{E_k [k g_k - g_k k_2^2]} \leq \frac{p-1}{(k+1)^2};$$

which completes the proof. \square

Lemma A.5. For all k , $E_k[d_k] = d_k$, $E_k[u_k] = u_k$, and $E_k[\tilde{v}_k] = v_k$. Moreover, there exists $L \geq R_{>0}$, independent of k , such that $k d_k - d_k k_2 \leq L k g_k - g_k k_2$ and $E_k [k d_k - d_k k_2] \leq \frac{p-1}{(k+1)^2}$, where $\sigma = L \geq R_{>0}$.

Proof. The first two statements follow the same arguments as in the proof of Lemma 4.17. By (3.2) and Lemma A.4,

$$E_k [k d_k - d_k k_2] \leq L E_k [k g_k - g_k k_2] \leq \frac{p-1}{(k+1)^2};$$

which proves the last statement. \square

Similar to Lemma 4.18, the next lemma provides bounds on the differences between stochastic and deterministic gradient approximations and exact and inexact search directions.

Lemma A.6. For all $k \geq N$,

$$\begin{aligned}
jg_k^T(d_k - \hat{d}_k)j & \leq g_{;dd}^{p-} kg_k \frac{g_k k_2^p}{l(x_k; k; g_k; d_k)}; \\
jg_k^T \hat{d}_k - g_k^T d_k j & \leq gg_{;dd}^{p-} kg_k \frac{g_k k_2^p}{l(x_k; k; g_k; d_k)} + L kg_k g_k k_2^2; \\
jg_k^T(d_k - \hat{d}_k)j & \leq g_{;dd}^{p-} = 2 \frac{p}{(k+1)} \frac{l(x_k; k; g_k; d_k)}{l(d)} + p \frac{2}{(k+1)} kg_k g_k k_2 \\
& \quad + g_{;dd}^{p-} = 2 l(x_k; k; g_k; d_k); \\
jg_k^T(d_k - d_k)j & \leq g_{;dd}^{p-} kg_k \frac{g_k k_2^p}{l(x_k; k; g_k; d_k)} \\
& \quad + g_{;dd}^{p-} = 2 \frac{p}{(k+1)} \frac{l(x_k; k; g_k; d_k)}{l(d)} + g_{;dd}^{p-} = 2 l(x_k; k; g_k; d_k); \\
jg_k^T d_k - g_k^T d_k j & \leq gg_{;dd}^{p-} kg_k \frac{g_k k_2^p}{l(x_k; k; g_k; d_k)} + L kg_k g_k k_2^2 \\
& \quad + g_{;dd}^{p-} = 2 \frac{p}{(k+1)} \frac{l(x_k; k; g_k; d_k)}{l(d)} + p \frac{2}{(k+1)} kg_k g_k k_2 \\
& \quad + g_{;dd}^{p-} = 2 l(x_k; k; g_k; d_k); \\
\text{and } kJ_k(d_k - \hat{d}_k)k_1 & \leq J_{;dd}^{p-} = 2 l(x_k; k; g_k; d_k);
\end{aligned}$$

where $g_{;dd}^{p-} = p \frac{H-1}{l,d} \geq 2 R_{>0}$, $gg_{;dd}^{p-} = g_{;dd}^{p-} + p \frac{1}{l,d} \geq 2 R_{>0}$, $g_{;dd}^{p-} = H \frac{q}{\frac{2}{l,d}} \geq 2 R_{>0}$, $g_{;dd}^{p-} = \gamma! a \geq 2 R_{>0}$, $g_{;dd}^{p-} = H \frac{2}{l,d} \geq 2 R_{>0}$, $g_{;dd} = \gamma! a \geq 2 R_{>0}$, and $J_{;dd} = ! a \geq 2 R_{>0}$.

Proof. (First inequality) By Assumption 3.1, Lemma A.5 and (3.2), (3.11), (4.3),

$$\begin{aligned}
jg_k^T(d_k - \hat{d}_k)j & = j(g_k + J_k^T(y_k + \kappa))^T(d_k - \hat{d}_k)j \\
& = j(H_k d_k)^T(d_k - \hat{d}_k)j \leq H kd_k k_2 kd_k - \hat{d}_k k_2 \\
& \leq H \frac{p}{l(d)} \frac{l(x_k; k; g_k; d_k)}{l(d)} L kg_k g_k k_2;
\end{aligned} \tag{A.2}$$

where the result follows using the definition of $g_{;dd}^{p-}$.

(Second inequality) By the Cauchy-Schwarz inequality, the triangle inequality, Assumption 3.1, (3.2), (3.11), (4.3) and (A.2), and Lemma A.5, it follows that

$$\begin{aligned}
& jg_k^T \hat{d}_k - g_k^T d_k j \\
& \leq j(g_k - g_k)^T d_k j + jg_k^T(\hat{d}_k - d_k)j + j(g_k - g_k)^T(\hat{d}_k - d_k)j \\
& \leq k g_k - g_k k_2 kd_k k_2 + jg_k^T(\hat{d}_k - d_k)j + kg_k - g_k k_2 kd_k - \hat{d}_k k_2 \\
& \leq g_{;dd}^{p-} + p \frac{1}{l,d} kg_k - g_k k_2 \frac{p}{l(x_k; k; g_k; d_k)} + L kg_k - g_k k_2^2;
\end{aligned} \tag{A.3}$$

where the result follows using the definitions of $gg_{;dd}^{p-}$.

(Third inequality) By Assumption 3.1, Lemmas 3.7, 3.8 and A.5, (3.2), (3.3), (3.11), and (4.3), we have for all $k \geq 2$ that

$$\begin{aligned}
jg_k^T(d_k - \hat{d}_k)j &= j(g_k + J_k^T(y_k + \kappa) - J_k^T(y_k + \kappa))^T(d_k - \hat{d}_k)j \\
&\quad j(g_k + J_k^T(y_k + \kappa))^T(d_k - \hat{d}_k)j + j(J_k^T(y_k + \kappa))^T(d_k - \hat{d}_k)j \\
&\quad j(g_k + J_k^T(y_k + \kappa))^T(d_k - \hat{d}_k)j + j(g_k - g_k)^T(d_k - \hat{d}_k)j \\
&\quad + j(J_k^T(y_k + \kappa))^T(d_k - \hat{d}_k)j \\
&= j(H_k d_k)^T(d_k - \hat{d}_k)j + j(g_k - g_k)^T(d_k - \hat{d}_k)j + j(y_k + \kappa)^T r_k j \quad (\text{A.4}) \\
&\quad j(H_k d_k)^T(d_k - \hat{d}_k)j + kg_k - g_k k_2 k d_k - \hat{d}_k k_2 + ky_k + \kappa k_1 k r_k k_1 \\
&\quad \leq \frac{L(x_k; \kappa; g_k; d_k)}{l; d} p^{\frac{2P-2}{(k+1)}} + p^{\frac{2P-2}{(k+1)}} kg_k - g_k k_2 \\
&\quad + \gamma! a^{-2} l(x_k; \kappa; g_k; d_k);
\end{aligned}$$

where the result follows using the definitions of $g; d; p$ and $g; d$.

(Fourth inequality) By Assumption 3.1, Lemmas 3.7 and 3.8, (3.2), (3.3), (3.11) and (4.3), it follows that

$$\begin{aligned}
jg_k^T(d_k - \hat{d}_k)j &= j(g_k + J_k^T(y_k + \kappa) - J_k^T(y_k + \kappa))^T(d_k - \hat{d}_k)j \\
&\quad j(g_k + J_k^T(y_k + \kappa))^T(d_k - \hat{d}_k)j + j(J_k^T(y_k + \kappa))^T(d_k - \hat{d}_k)j \\
&\quad j(H_k d_k)^T(d_k - \hat{d}_k)j + ky_k + \kappa k_1 k r_k k_1 \\
&\quad \leq \frac{L(x_k; \kappa; g_k; d_k)}{l; d} = 2 p^{\frac{P-2}{(k+1)}} + \gamma! a^{-2} l(x_k; \kappa; g_k; d_k);
\end{aligned}$$

By the triangle inequality and (A.2),

$$\begin{aligned}
jg_k^T(d_k - d_k)j &\leq jg_k^T(\hat{d}_k - d_k)j + jg_k^T(d_k - \hat{d}_k)j \\
&\quad \leq \frac{g; d; p - kg_k - g_k k_2}{q} p^{\frac{P-2}{(k+1)}} l(x_k; \kappa; g_k; d_k) \\
&\quad + \frac{L(x_k; \kappa; g_k; d_k)}{l; d} = 2 p^{\frac{P-2}{(k+1)}} l(x_k; \kappa; g_k; d_k);
\end{aligned}$$

where the result follows using the definition of $g; d; p$ and $g; d$.

(Fifth inequality) By (A.3), (A.4),

$$\begin{aligned}
jg_k^T(d_k - g_k^T d_k)j &\leq jg_k^T(d_k - g_k^T \hat{d}_k)j + jg_k^T(\hat{d}_k - d_k)j \\
&\quad \leq \frac{g; d; p - kg_k - g_k k_2}{q} p^{\frac{P-2}{(k+1)}} l(x_k; \kappa; g_k; d_k) + L kg_k - g_k k_2^2 \\
&\quad + \frac{g; d; p - kg_k - g_k k_2}{q} p^{\frac{P-2}{(k+1)}} l(x_k; \kappa; g_k; d_k) + p^{\frac{P-2}{(k+1)}} kg_k - g_k k_2 \\
&\quad + \gamma! a^{-2} l(x_k; \kappa; g_k; d_k);
\end{aligned}$$

(Sixth inequality) By Lemma 3.7, and (3.2), (3.3),

$$k J_k(d_k - \bar{d}_k) k_1 = k r_k k_1 - \alpha = 2 l(x_k; k; g_k; d_k);$$

where the result follows using the definition of J_k . \square

The next lemma provides a useful upper bound on the difference between the deterministic and stochastic merit parameters.

Lemma A.7. For all $k \in \mathbb{N}$, $j \in \{1, \dots, k\}$, $g_k^T d_k = 2 l(x_k; k; g_k; d_k)$, where $\alpha = 2 \beta - 1$ and $\beta \geq 0$; $\frac{1}{(2 - \beta)^2} > 0$.

Proof. The proof of this statement is identical to the proof of Lemma 4.20. \square

Similar to Lemma 4.21, the next lemma bounds the stochastic model reduction function with respect to its deterministic counterpart (with additional terms).

Lemma A.8. For all $k \in \mathbb{N}$, by choosing $\beta \geq 0$; $\frac{1}{(2 - \beta)^2} > 0$, it follows that

$$\begin{aligned} l(x_k; k; g_k; d_k) &= (1 + \frac{\beta}{2 - \beta}) l(x_k; k; g_k; d_k) \\ &+ 2 \beta - 1 \frac{1}{2} k g_k^T g_k + \frac{\beta}{2 - \beta} \frac{1}{(k+1)} \frac{l(x_k; k; g_k; d_k)}{p} \\ &+ \frac{\beta}{2 - \beta} \frac{1}{(k+1)} \frac{l(x_k; k; g_k; d_k)}{p} + \frac{\beta}{2 - \beta} \frac{1}{(k+1)} \frac{2 k g_k^T g_k}{p} ; \end{aligned}$$

where $\frac{\beta}{2 - \beta} = 2(1 + \frac{\beta}{2 - \beta})^2 R > 0$. Additionally, under Assumption 4.1, for all $k \in \mathbb{N}$

$$\begin{aligned} E_k l(x_k; k; g_k; d_k) &= (1 + \frac{\beta}{2 - \beta}) l(x_k; k; g_k; d_k) \\ &+ \frac{\beta}{2 - \beta} \frac{1}{(k+1)} \frac{l(x_k; k; g_k; d_k)}{p} + \frac{\beta}{2 - \beta} \frac{1}{(k+1)} ; \end{aligned}$$

where $\frac{\beta}{2 - \beta} = \frac{\beta}{2 - \beta} \frac{2}{p} R > 0$, $\frac{\beta}{2 - \beta} \frac{1}{(k+1)} = 2(1 + \frac{\beta}{2 - \beta})^2 R > 0$, and $\frac{\beta}{2 - \beta} = 2(1 + \frac{\beta}{2 - \beta})^2 R > 0$.

Proof. By (3.5), and Lemmas A.6 and A.7, it follows that

$$\begin{aligned} l(x_k; k; g_k; d_k) &= k g_k^T d_k + k \alpha_k k_1 - k \alpha_k + J_k d_k k_1 \\ &+ k g_k^T d_k + k \alpha_k k_1 + (k - k) g_k^T d_k + k (g_k^T d_k - g_k^T d_k) \\ l(x_k; k; g_k; d_k) &+ j \in \{1, \dots, k\} g_k^T d_{kj} + \frac{1}{j} g_k^T d_k - g_k^T d_{kj} \\ l(x_k; k; g_k; d_k) &+ \frac{\beta}{2 - \beta} l(x_k; k; g_k; d_k) \\ &+ 2 \beta - 1 \frac{1}{2} k g_k^T g_k + \frac{\beta}{2 - \beta} \frac{1}{(k+1)} \frac{l(x_k; k; g_k; d_k)}{p} + \frac{\beta}{2 - \beta} \frac{1}{(k+1)} \frac{2 k g_k^T g_k}{p} \\ &+ \frac{\beta}{2 - \beta} \frac{1}{(k+1)} \frac{l(x_k; k; g_k; d_k)}{p} + \frac{\beta}{2 - \beta} \frac{1}{(k+1)} \frac{2 k g_k^T g_k}{p} + \frac{\beta}{2 - \beta} \frac{1}{(k+1)} \frac{l(x_k; k; g_k; d_k)}{p} ; \end{aligned}$$

The first result follows by re-arranging the above, invoking the restriction on β , and the definition of β_{k+1} .

Taking the conditional expectation, by Assumption 4.27 and Lemma A.4,

$$\begin{aligned} E_k [l(x_k; \beta_k; g_k; d_k)] &= 1 + 2(\beta_k + \beta_{k+1} g_k; d_k)^2 = 2 l(x_k; \beta_k; g_k; d_k) \\ &+ 2 \beta_{k+1} \frac{\beta_{k+1}}{(k+1)} + \beta_{k+1} g_k; d_k; \beta_{k+1} = 2 \frac{\beta_{k+1} l(x_k; \beta_k; g_k; d_k)}{(k+1)} \\ &+ \beta_{k+1} g_k; d_k; \beta_{k+1} = 2 \frac{\beta_{k+1} l(x_k; \beta_k; g_k; d_k)}{(k+1)} + \frac{\beta_{k+1}}{(k+1)} : \end{aligned}$$

Using the definitions of β_{k+1} , β_{k+1} , β_{k+1} , and β_{k+1} , completes the proof. \square

Finally, we restate and prove the theoretical results stated in Section 4.2. Specifically, we state and prove Lemma 4.29, Theorem 4.30 and Corollary 4.31.

Lemma A.9. (Lemma 4.29) For all $k \geq 2$, N ,

$$E_k [l(x_{k+1}; \beta_{k+1}) - l(x_k; \beta_k)] \leq E_k [l(x_{k+1}; \beta_{k+1}) - l(x_k; \beta_k)] + \frac{\beta_{k+1}}{(k+1)};$$

where $\beta_{k+1} = \beta_k + \frac{\beta_{k+1}}{2} \geq 0$ and $\beta_{k+1} = \beta_k + \frac{\beta_{k+1}}{2} \geq 0$, and

$$\begin{aligned} \beta_{k+1} &= \beta_k (1 - \beta_{k+1}) + (\beta_{k+1} g_k; d_k + \beta_{k+1} g_k; d_k)(1 + \beta_{k+1}) \geq 0; \\ \beta_{k+1} &= \beta_k (1 - \beta_{k+1}) (1 + \beta_{k+1}) + (\beta_{k+1} g_k; d_k + \beta_{k+1} g_k; d_k) \beta_{k+1} \geq 0; \\ \beta_{k+1} &= \beta_k (1 - \beta_{k+1}) \beta_{k+1} + \beta_{k+1} g_k; d_k; \beta_{k+1} \\ &+ \beta_{k+1} g_k; d_k; \beta_{k+1} + (\beta_{k+1} g_k; d_k + \beta_{k+1} g_k; d_k) \beta_{k+1} \geq 0; \end{aligned}$$

Proof. By Lemmas A.3, A.7, A.6, A.8 and 4.14, and (3.4), it follows that for all $k \geq 2$ N

$$\begin{aligned}
 & E_k [(x_{k+1}; k+1) - (x_k; k)] \\
 = & E_k [(x_{k+1}; k) f_{k+1}] + E_k [(x_{k+1}; k) - (x_k; k)] \\
 & E_k [(x_{k+1}; k)] f_{inf} \\
 & E_k [(x_k; k; g_k; d_k) - (1 - \epsilon) (x_k; k; g_k; d_k)] \\
 & + E_k [(x_k; k; g_k^T (d_k - d_k) + E_k [(x_k; k; g_k^T d_k) + E_k [(x_k; k; g_k^T d_k)]]] \\
 & E_k [(x_{k+1}; k)] f_{inf} - E_k [(x_k; k; g_k; d_k)] \\
 & + E_k [(1 - \epsilon) (x_k; k; g_k; d_k) - (1 + \epsilon) (x_k; k; g_k; d_k)] \\
 & + 2 \epsilon (1 - \epsilon) L k g_k - g_k k_2^2 + \frac{p}{g;dd;^p} - \frac{= 2^p}{p} \frac{l(x_k; k; g_k; d_k)}{(k+1)} \\
 & + \frac{p}{g;dd;^p} - k g_k - g_k k_2^2 \frac{p}{l(x_k; k; g_k; d_k)} + \frac{p}{2} \frac{= 2^p k g_k - g_k k_2}{(k+1)} \\
 & + \frac{u}{h} (2 - \epsilon) E_k [j g_k^T (d_k - d_k)] + \frac{u}{h} (2 - \epsilon) E_k [j (x_k; k) g_k^T d_k] \\
 & + \frac{u}{h} (2 - \epsilon) E_k [k J_k (d_k - d_k) k_1] \\
 & E_k [(x_{k+1}; k)] f_{inf} - E_k [(x_k; k; g_k; d_k)] \\
 & + (1 - \epsilon) (1 - \epsilon) u (1 + \epsilon) l(x_k; k; g_k; d_k) \\
 & + 2 \epsilon (1 - \epsilon) u L \frac{1}{(k+1)} + \frac{p}{g;dd;^p} - \frac{= 2^p}{p} \frac{l(x_k; k; g_k; d_k)}{(k+1)} \\
 & + \frac{p}{g;dd;^p} - \frac{p}{(k+1)} \frac{= 2^p}{p} \frac{l(x_k; k; g_k; d_k)}{(k+1)} + \frac{p}{(k+1)} \frac{= 2^p}{p} \frac{l(x_k; k; g_k; d_k)}{(k+1)} \\
 & + \frac{u}{h} (2 - \epsilon) \frac{1}{g;dd;^p} - \frac{p}{(k+1)} \frac{= 2^p}{p} \frac{l(x_k; k; g_k; d_k)}{(k+1)} + \frac{p}{g;dd;^p} - \frac{= 2^p}{p} \frac{l(x_k; k; g_k; d_k)}{(k+1)} \\
 & + \frac{p}{g;dd;^p} - \frac{= 2^p}{p} E_k [l(x_k; k; g_k; d_k)]
 \end{aligned}$$

Continuing from the above, by Lemmas A.8 and 4.14, it follows that for all $k \geq N$

$$\begin{aligned}
 & E_k [I(x_{k+1}; k+1) - I(x_k; k)] \\
 & E_k [(I(x_{k+1}; k+1) - I(x_k; k))] f_{\text{inf}} \\
 & \quad - I(x_k; k; g_k; d_k) + (1 - \bar{l}_i) u^{(1+)=2} I(x_k; k; g_k; d_k) \\
 & + 2^{-1} (1 - \bar{l}_i) u^{(1+)=2} \frac{p}{g; d; p - 1} \frac{p}{g; d; p - 1} \frac{p}{\beta^{(k+1)}(x_k; k; g_k; d_k)} \\
 & \quad + \frac{p}{L 1 + 1} \frac{p}{1 2} \frac{= 2}{(k+1)} \\
 & + u^{-1} 2^{-1} \frac{p}{g; d; p - 1} \frac{p}{g; d; p - 1} \frac{p}{\beta^{(k+1)}(x_k; k; g_k; d_k)} + u^{-2} I(x_k; k; g_k; d_k) \\
 & + u^{-2} (1 - g; d; d + J; d; d) (1 + \bar{l}_i) u^{(1+)=2} I(x_k; k; g_k; d_k) \\
 & \quad + \bar{l}_i; l_i; p - = 2 \frac{p}{\beta^{(k+1)}(x_k; k; g_k; d_k)} + \bar{l}_i; l_i; \frac{1}{(k+1)} \\
 & E_k [(I(x_{k+1}; k+1) - I(x_k; k))] f_{\text{inf}} - I(x_k; k; g_k; d_k) \\
 & + u (1 - \bar{l}_i) u^{(1+)=2} + (1 - g; d; d + J; d; d) (1 + \bar{l}_i) u^{(1+)=2} I(x_k; k; g_k; d_k) \\
 & + u^{-2} (1 - \bar{l}_i) \frac{p}{g; d; p - 1} \frac{p}{g; d; p - 1} \frac{p}{g; d; p - 1} \frac{p}{g; d; p - 1} \frac{p}{\beta^{(k+1)}(x_k; k; g_k; d_k)} \\
 & + (1 - g; d; d + J; d; d) \bar{l}_i; l_i; p - 2 \frac{p}{\beta^{(k+1)}(x_k; k; g_k; d_k)} \\
 & + u^{-2} (1 - \bar{l}_i) (\frac{p}{L 1 + 1} \frac{p}{1 2}) + (1 - g; d; d + J; d; d) \bar{l}_i; l_i; \frac{2}{(k+1)} \frac{1}{(k+1)} \\
 & = E_k [(I(x_{k+1}; k+1) - I(x_k; k))] f_{\text{inf}} - I(x_k; k; g_k; d_k) + u^{(1+)=2} I(x_k; k; g_k; d_k) \\
 & \quad + \bar{l}_i; l_i; p - 2 \frac{p}{\beta^{(k+1)}(x_k; k; g_k; d_k)} + \bar{l}_i; l_i; \frac{2}{(k+1)} \frac{1}{(k+1)} \\
 & E_k [(I(x_{k+1}; k+1) - I(x_k; k))] f_{\text{inf}} - I(x_k; k; g_k; d_k) + \frac{p - 2}{2} u^{(1+)=2} I(x_k; k; g_k; d_k) \\
 & \quad + \bar{l}_i; l_i; \frac{p - 2}{2} \frac{1}{(k+1)};
 \end{aligned}$$

where the result follows using the definitions of \bar{l}_i , l_i , $p - 2$, $u^{(1+)=2}$, and u^{-2} . □

We are now ready to prove the main theorem of this section.

Theorem A.10. (Theorem 4.30) For some $\alpha \in (0; 1)$ and $\beta > 1$, and choosing $\epsilon \in (0; \min\{\frac{1}{(2 - \alpha)^2}; \frac{(1 - \alpha)}{2}\})$, it follows that

$$\lim_{k \rightarrow \infty} E \sum_{j=0}^{k-1} I(x_j; j; g_j; d_j)^5 < 1;$$

from which it follows that $\lim_{k \rightarrow \infty} E[\|x_k - x^*\|] = 0$.

Proof. By Lemma A.9 and $\alpha > 0$, it follows that

$$E_k[\|x_{k+1} - x_k\|] \leq E_k[\|x_k - x^*\|] + \frac{\alpha}{k+1} \quad (\text{A.5})$$

Applying a telescopic sum to (A.5) and taking the total expectation, it follows that

$$\begin{aligned} 1 &< \sum_{j=0}^k E[\|x_{j+1} - x_j\|] \\ &= E\left[\sum_{j=0}^k (\|x_{j+1} - x_j\|)\right] \\ &= E\left[\sum_{j=0}^k \left(\|x_{j+1} - x_j\| + \frac{\alpha}{j+1}\right)\right] \\ &= E\left[\sum_{j=0}^k \|x_{j+1} - x_j\| + \alpha \sum_{j=0}^k \frac{1}{j+1}\right] \end{aligned}$$

Finally, using the fact that $\sum_{j=0}^k \frac{1}{j+1} < 1$ for any $k > 0$, we may complete the proof. \square

Corollary A.11. Under Assumptions A.1 and 4.9, and the conditions of Theorem A.10, Algorithm 1 yields a sequence of iterates (x_k, y_k) for which

$$\lim_{k \rightarrow \infty} E\|kx_k\|_2^2 = 0; \quad \lim_{k \rightarrow \infty} E\|ky_k\|_2 = 0; \quad \text{and} \quad \lim_{k \rightarrow \infty} E\|kg_k + J_k^T(y_k - x_k)\|_2 = 0:$$

Proof. The proof of this corollary follows the exact same arguments as the proof of Corollary 4.24 (Section 4). \square

The final result we show in this section is a complexity result for our proposed algorithm, i.e., the number of iterations and the total number of stochastic gradient evaluations required to achieved an ϵ -accurate solution in expectation. Specifically, to measure the complexity of our algorithm, we consider the minimum number of iterations, K , and the minimum total number of stochastic gradient evaluations, W , required to achieve the following approximate stationary measure

$$E\|kg_k + J_k^T(y_k - x_k)\|_2 \leq L \quad \text{and} \quad E\|kx_k\|_2 \leq c; \quad (\text{A.6})$$

for $L \in (0, 1)$ and $c \in (0, 1)$.

Corollary A.12. (Corollary 4.31) Under the conditions of Theorem 4.30, Algorithm 1 generates an iterate $(x_k; y_k)$ that satisfies (A.6) in at most $K = O(\max\{L^2; c^{-1}\})$ iterations and $W = O(\max\{L^2; c^{-1}\}^{(+)})$, $2R > 1$, stochastic gradient evaluations. Moreover, if $L = 1$ and $c = 2$, then $K = O(2)$ and $W = O(2^{(+)})$.

Proof. Using the same logic as the proof of Corollary 4.26, we know that if (A.6) is violated, then for all $k \geq 0; \dots; K-1$

$$E[\sum_{k=0}^{K-1} l(x_k; y_k; g_k; d_k)] \geq \sum_{k=0}^{K-1} \min_{x,y} \{ \frac{1}{2} L \|g_k\|^2; c \|g_k\| \} \tag{A.7}$$

where $x = \min\{1; \frac{\min\{1; d\}}{2R}\}$, $2R > 0$. Re-arranging terms in (A.5), and using (3.4) and (A.7) to do a telescoping, it follows that

$$\begin{aligned} & E[\sum_{k=0}^{K-1} (f(x_k) - f_{\text{inf}}) + kc_0k_1] = E[\sum_{k=0}^{K-1} (f(x_k) - f_{\text{inf}}) + kc_0k_1] \\ & E[\sum_{k=0}^{K-1} (f(x_k) - f_{\text{inf}}) + kc_0k_1 - \sum_{k=0}^{K-1} (f(x_{k+1}) - f_{\text{inf}}) + (k+1)c_0k_1 + (k+1)f_{\text{inf}})] \\ & = E[\sum_{k=0}^{K-1} (f(x_k) - f(x_{k+1})) + kc_0k_1 - (k+1)c_0k_1 - (k+1)f_{\text{inf}} + (k+1)f_{\text{inf}})] \\ & = E[\sum_{k=0}^{K-1} (f(x_k) - f(x_{k+1})) + (k+1)c_0k_1 - (k+1)f_{\text{inf}})] \\ & = E[\sum_{k=0}^{K-1} l(x_k; y_k; g_k; d_k) - \sum_{k=0}^{K-1} \min_{x,y} \{ \frac{1}{2} L \|g_k\|^2; c \|g_k\| \}] \\ & = \sum_{k=0}^{K-1} \min_{x,y} \{ \frac{1}{2} L \|g_k\|^2; c \|g_k\| \} - \sum_{k=0}^{K-1} \min_{x,y} \{ \frac{1}{2} L \|g_k\|^2; c \|g_k\| \} \\ & = \sum_{k=0}^{K-1} \min_{x,y} \{ \frac{1}{2} L \|g_k\|^2; c \|g_k\| \} - \sum_{k=0}^{K-1} \min_{x,y} \{ \frac{1}{2} L \|g_k\|^2; c \|g_k\| \} \end{aligned}$$

It further implies that K is bounded as

$$K \leq \frac{\sum_{k=0}^{K-1} (f(x_k) - f_{\text{inf}}) + kc_0k_1 + \sum_{k=0}^{K-1} \min_{x,y} \{ \frac{1}{2} L \|g_k\|^2; c \|g_k\| \}}{\sum_{k=0}^{K-1} \min_{x,y} \{ \frac{1}{2} L \|g_k\|^2; c \|g_k\| \}} = O(\min\{L^2; c\}) \tag{A.8}$$

under the condition that $2R > 1$. Next, we analyze the sampling complexity. Suppose, we use $|S_k|$ samples to estimate the stochastic gradient g_k , then,

$$E_k \|kg_k - g_k\|_2^2 = \frac{\sigma^2}{|S_k|}$$

where σ^2 is the population variance. Therefore, the minimum number of samples required to satisfy Assumption 4.27 is

$$|S_k| = \frac{2(k+1)}{1}$$

By (A.8), it follows that the total number of stochastic gradient evaluations required to satisfy (A.6) can be expressed as

$$W = \sum_{k=0}^{K-1} \sum_j S_{kj} = \sum_{k=0}^{K-1} \frac{2^{k+1}}{1} = \sum_{k=0}^{K-1} 2^{k+1} :$$

Using Faulhaber's formula, sum of k -th power of first K positive integers is a function of polynomial $K+1$. Therefore, there exists some constant $c_k \in (0; 1)$ such that $\sum_{k=0}^{K-1} 2^{k+1} = c_k K^{k+1}$. Therefore,

$$W = \frac{2^{k+1}}{1} K^{k+1}$$

Substituting $K = O(\epsilon^{-2})$ yields the desired result. □

B Additional Numerical Results: Constrained Logistic Regression

In this section, we provide numerical results corresponding to the binary classification data sets given in Table 1 from the [17] collection.

Table 1: Binary classification data set details. For more information see [17].

data set	dimension (n)	datapoints (N)
australian	14	690
ionosphere	34	351
mushrooms	112	8,124
sonar	60	208
splice	60	1,000

(a) Feasibility vs. Iterations (b) Feasibility vs. Epochs (c) Feasibility vs. LS Iters

(d) Stationarity vs. Iterations (e) Stationarity vs. Epochs (f) Stationarity vs. LS Iters

(g) Step Size vs. Iterations (h) Batch Size vs. Iterations

Figure 5: *australian* : First & Second Row: Feasibility & stationarity errors versus iterations/epochs/linear system iterations for exact and inexact variants of Algorithm 2 on (6.1). Last Row: Step sizes and Batch sizes versus iterations.

(a) Feasibility vs. Iterations (b) Feasibility vs. Epochs (c) Feasibility vs. LS Iters

(d) Stationarity vs. Iterations (e) Stationarity vs. Epochs (f) Stationarity vs. LS Iters

(g) Step Size vs. Iterations (h) Batch Size vs. Iterations

Figure 6: ionosphere : First & Second Row: Feasibility & stationarity errors versus iterations/epochs/linear system iterations for exact and inexact variants of Algorithm 2 on (6.1). Last Row: Step sizes and Batch sizes versus iterations.

(a) Feasibility vs. Iterations (b) Feasibility vs. Epochs (c) Feasibility vs. LS Iters

(d) Stationarity vs. Iterations (e) Stationarity vs. Epochs (f) Stationarity vs. LS Iters

(g) Step Size vs. Iterations (h) Batch Size vs. Iterations

Figure 7: mushroom First & Second Row: Feasibility & stationarity errors versus iterations/epochs/linear system iterations for exact and inexact variants of Algorithm 2 on (6.1). Last Row: Step sizes and Batch sizes versus iterations.

(a) Feasibility vs. Iterations (b) Feasibility vs. Epochs (c) Feasibility vs. LS Iters

(d) Stationarity vs. Iterations (e) Stationarity vs. Epochs (f) Stationarity vs. LS Iters

(g) Step Size vs. Iterations (h) Batch Size vs. Iterations

Figure 8: sonar: First & Second Row: Feasibility & stationarity errors versus iterations/epochs/linear system iterations for exact and inexact variants of Algorithm 2 on (6.1). Last Row: Step sizes and Batch sizes versus iterations.

