

Accelerating Stochastic Sequential Quadratic Programming for Equality Constrained Optimization using Predictive Variance Reduction

Albert S. Berahas^{*†} Jiahao Shi[†] Zihong Yi[‡] Baoyu Zhou[§]

March 28, 2023

Abstract

In this paper, we propose a stochastic method for solving equality constrained optimization problems that utilizes predictive variance reduction. Specifically, we develop a method based on the sequential quadratic programming paradigm that employs variance reduction in the gradient approximations. Under reasonable assumptions, 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, for both constant and adaptive step size strategies. Finally, we demonstrate the practical performance of our proposed algorithm on constrained binary classification problems that arise in machine learning.

1 Introduction

We consider the design of algorithms for solving equality constrained finite-sum problems of the form

$$\min_{x \in \mathbb{R}^n} f(x) \quad \text{s.t.} \quad c(x) = 0, \quad \text{with} \quad f(x) = \frac{1}{N} \sum_{i=1}^N f_i(x), \quad (1.1)$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$ for all $i \in \{1, \dots, N\}$, and $c : \mathbb{R}^n \rightarrow \mathbb{R}^m$ are smooth nonlinear (possibly nonconvex) functions. Such problems arise in a plethora of areas such as machine/deep learning [1, 21, 24, 30, 35, 45], statistics [9, 14], and stochastic optimal

^{*}Corresponding author.

[†]Department of Industrial and Operations, University of Michigan. (albertberahas@gmail.com, jiahaos@umich.edu)

[‡]Department of Computer Science and Engineering, University of Michigan. (zihongy@umich.edu)

[§]Booth School of Business, University of Chicago. (baoyu.zhou@chicagobooth.edu)

control [19, 20], as well as other science and engineering applications such as optimal power flow [40, 42, 43], multi-stage modeling [38], and portfolio optimization [41, 46].

Numerous algorithms have been developed over the last half century for solving deterministic equality constrained optimization problems, such as that in (1.1). A few classical examples are penalty methods, projection methods and sequential quadratic programming (SQP), each of which have their associated merits and limitations [28]. Penalty methods are intuitive and simple to implement, however, their performance critically relies on the choice of the penalty function and penalty parameter, and, in practice, often suffers from ill-conditioning issues and subproblems' nonsmoothness. On the other hand, projection methods are powerful *feasible* methods, however, they assume that projections can be efficiently computed at every iteration, something that is often not the case with general nonlinear constraints. SQP methods attempt to alleviate these issues by solving a sequence of subproblems that minimize a quadratic model of the objective function subject to a linearization of the constraints, and, as such can handle general nonlinear constraints, however, this comes at the cost of more expensive iterations (SQP methods require solving a linear system at every iteration). That being said, all aforementioned deterministic methods require the computation of the true gradient of the objective function (and constraints) at every iteration, which can be prohibitively expensive in settings in which n and/or N are large.

Rather than minimizing the finite-sum optimization problem (1.1) with a deterministic method, one can employ stochastic methods that utilize a stochastic approximation of the gradient in lieu of the true gradient in order to reduce the per iteration computational cost. In this direction, several stochastic penalty and projection methods have been proposed [10, 15, 17, 18, 24, 25, 32, 39]. Another line of research considers stochastic alternating direction method of multipliers (ADMM) algorithms [29, 44], and variants of ADMM that utilize variance reduction [2, 5]. Following the SQP paradigm, recent work [4] proposed a stochastic SQP method with an adaptive step size selection rule for solving equality constrained stochastic optimization problems endowed with theoretical guarantees (convergence in expectation) analogous of those of the stochastic gradient (SG) method for unconstrained problems, and empirical performance superior to that of the stochastic sub-gradient method. Several extensions of this work have been developed; namely, in [3] the authors relax requirements on the constraints (relax constraint qualifications), in [12] the authors develop an inexact stochastic SQP method (linear system solved inexactly at every iteration), and in [11] the authors analyze the complexity of the stochastic SQP algorithm proposed in [4]. Along a slightly different direction, under the assumption that the error in the stochastic gradient approximations employed can be diminished as needed, the authors in [22, 23] proposed stochastic line search SQP methods for equality and inequality constrained stochastic optimization problems, respectively, that utilize a differentiable exact augmented Lagrangian function as its merit function.

In the last decade, several stochastic first-order algorithms have been proposed for solving unconstrained finite-sum optimization problems. One such class of algorithms are

variance-reduction methods, that attempt to reduce the the variance in the stochastic gradient approximation employed as the optimization progresses. Examples of popular variance reduction methods include, but are not limited to, the Stochastic Average Gradient (SAG/SAGA) method [13, 36], the Stochastic Variance Reduced Gradient (SVRG) method [16], the Stochastic Recursive Gradient Algorithm (SARAH) method [27], and the Stochastic Dual Coordinate Ascent (SDCA) method [37]. As a result of the variance reduction, these methods enjoy improved convergence results as compared to their classical counter-parts (e.g., SG method [6, 33]), and these benefits are very often also observed in practice. Motivated by this fact, we design and analyze a stochastic SQP method that employs variance reduced gradients for solving (1.1).

1.1 Contributions

The contributions of our work can be categorized as follows:

- **Algorithmic.** We present a stochastic sequential quadratic optimization algorithm that uses variance reduced gradients. Specifically, inspired by the theoretical and empirical advantages of *variance reduced methods* (unconstrained finite-sum problems) and SQP methods (deterministic equality constrained problems), we propose a stochastic SQP method that uses variance reduced gradients (SVRG-type, [16]) in lieu of the true gradient (SVR-SQP). We propose one algorithm with two possible step size selection strategies; a constant step size scheme (similar to that in SVRG [16, 31]), and an adaptive step size scheme (similar to that in [4]).

Our proposed algorithm is based on a stochastic SQP framework, similar to the stochastic algorithm proposed in [4], but with several distinguishing algorithmic and theoretical differences. In particular, our proposed algorithm is of nested form, due to the nature of the construction of the SVRG gradients, and operates with two types of iterations (inner and outer). As a direct consequence of the use of variance reduced gradients, the proposed step size selection strategy only requires minimal safe-guarding, as compared to the safe-guards imposed in [4], e.g., safe-guarding parameters or sequences and projections.

We should note that while in this work we chose to employ SVRG-type gradient approximations, others, e.g., [13, 27, 36, 37], may also be employed. We chose SVRG because it is based on an intuitive idea (control variates [34]), has no additional storage requirements, has proven robust and efficient in practice, and, perhaps most importantly because the unbiasedness of the SVRG gradients allows for simple convergence analysis.

- **Theoretical.** We provide convergence guarantees for the SVR-SQP method with the two different step size strategies (constant and adaptive). For both strategies, we present strong theoretical guarantees in the sense that a measure of first-order

Table 1: Summary of asymptotic results for different problem settings and different methods (and step size choices) for nonconvex functions. For unconstrained finite sum problems the stationarity measure is $\|\nabla f(x)\|_2^2$, whereas for equality constrained finite sum problems the stationarity measure is $\|\nabla f(x) + \nabla c(x)y\|_2^2 + \|c(x)\|_2$ (where y are least-squares Lagrange multipliers). In the table, “exact” and “neighborhood” denote convergence to the stationarity measure in expectation and convergence to a neighborhood of the stationarity measure in expectation, respectively, and “-” denotes that no result exists. For brevity and ease of exposition, we do not state the exact constants in the results below.

Setting	Method	Step size		
		Diminishing	Constant	Adaptive
Unconstrained	SG [6]	exact	neighborhood	-
Finite Sum	SVRG [31]	-	exact	-
Equality	Stoch. SQP [4]	exact	neighborhood	neighborhood
Constrained	SVR-SQP	-	exact	exact
Finite Sum	(this paper)	-	exact	exact

stationarity evaluated at the iterates generated by SVR-SQP vanishes in expectation with *non-diminishing* step size sequences. This is in contrast with the results in [4] where a *diminishing* step size sequence is required to ensure exact convergence in expectation. Our result (equality constrained finite sum setting) for the SVR-SQP method with a constant step size can be viewed as analogues of the results that can be proven for the SVRG method [16, 31] in the unconstrained finite sum setting. Similar convergence guarantees are established for the more flexible variant with adaptive step sizes. Table 1 summarizes our results.

- **Empirical.** We illustrate the performance of our proposed method on constrained binary classification problems, and we compare our proposed algorithm against other popular methods, such as the adaptive stochastic SQP method proposed in [4] and a stochastic subgradient method that utilizes variance reduction. We provide evidence illustrating the benefits of using variance reduced gradients within the stochastic SQP framework for solving equality constrained finite sum optimization problems.

1.2 Organization

The paper is organized as follows. We conclude this section by setting the notation that will be used throughout the paper. In Section 2 we introduce the assumptions and main algorithmic components of our proposed method. The stochastic variance reduced sequential quadratic optimization method is presented in Section 3, and its associated convergence

guarantees are presented in Section 4. In Section 5, we demonstrate the empirical performance of the proposed algorithm. Finally, in Section 6 we make some concluding remarks.

1.3 Notation

Let \mathbb{N} denote the set of natural numbers, \mathbb{R} denote the set of real numbers and $\mathbb{R}_{>0}$ denote the set of positive real numbers. For any $m \in \mathbb{N}$, let $[m]$ denote the set of integers $\{1, \dots, m\}$, and $[\bar{m}]$ denote the set of integers $\{0, 1, \dots, m-1\}$. Let \mathbb{R}^n denote the set of n -dimensional real vectors, $\mathbb{R}^{m \times n}$ denote the set of m -by- n -dimensional real matrices, and \mathbb{S}^n denote the set of n -by- n -dimensional symmetric matrices.

The algorithms described in this paper will either operate with a single type of iteration and produce sequences of iterates $\{x_k\}$ where $k \in \mathbb{N}$ is the index of iterations, or will operate with two types of iterations (e.g., inner and outer) and produce sequences of iterates $\{x_{k,s}\}$, where $k \in \mathbb{N}$ is the index of outer iterations and $s \in [\bar{S}]$ is the index of inner iterations. The index of iteration number is also appended as a subscript to other quantities corresponding to each iteration; e.g., $f_{k,s} = f(x_{k,s})$, respectively for the single iteration algorithms. Throughout the paper, we use the overline to denote stochastic quantities; e.g., $\bar{g}_{k,s}$ is an estimate of $g_{k,s} := \nabla f(x_{k,s})$.

2 Assumptions and Algorithm Preliminaries

Throughout the paper, we assume that the constraint function and its associated first-order derivatives can be computed exactly. With regards to the objective function and its associated derivatives, we assume that those quantities are prohibitively expensive to compute at every iteration, however, exact evaluations can be accessed as required by the algorithm. We formalize our assumptions with regards to (1.1) and the iterates generated by our algorithm $\{x_{k,s}\}$ below.

Assumption 2.1. *Let $\mathcal{X} \subseteq \mathbb{R}^n$ be an open convex set containing the iterates $\{x_{k,s}\}$ generated by any run of the algorithm. The objective function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and its gradient $g := \nabla f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ are bounded over \mathcal{X} . For each $i \in [N]$, the component objective function $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$ is continuously differentiable, and each component gradient $\nabla f_i : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is Lipschitz continuous with constant L . For each $i \in [m]$, the constraint function $c_i : \mathbb{R}^n \rightarrow \mathbb{R}$ is continuously differentiable and bounded over \mathcal{X} , and its gradient $\nabla c_i : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is Lipschitz continuous with constant $\gamma_i \in \mathbb{R}_{>0}$. We define $\Gamma := \sum_{i=1}^m \gamma_i$. The Jacobian function $J := \nabla c^T : \mathbb{R}^n \rightarrow \mathbb{R}^{m \times n}$ is bounded over \mathcal{X} , and has singular values bounded away from zero over \mathcal{X} .*

Remark 2.2. *Assumption 2.1 implies that the objective function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is continuously differentiable, and that its gradient $g := \nabla f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is Lipschitz continuous with constant L . Under Assumption 2.1, it follows that, for all $(k, s) \in \mathbb{N} \times [\bar{S}]$, there exist*

positive real numbers $(f_{\inf}, f_{\sup}, \kappa_g, \kappa_c, \kappa_J, \kappa_\sigma) \in \mathbb{R} \times \mathbb{R} \times \mathbb{R}_{>0} \times \mathbb{R}_{>0} \times \mathbb{R}_{>0} \times \mathbb{R}_{>0}$ such that

$$\begin{aligned} f_{\inf} &\leq f_{k,s} \leq f_{\sup}, & \|g_{k,s}\|_2 &\leq \kappa_g, \\ \text{and } \|c_{k,s}\|_1 &\leq \kappa_c, & \|J_{k,s}\|_2 &\leq \kappa_J, & \|(J_{k,s}J_{k,s}^T)^{-1}\|_2 &\leq \kappa_\sigma^{-2}. \end{aligned} \quad (2.1)$$

Assumption 2.1 ensures the smoothness of the objective function and constraint functions. Unlike many projection methods aimed to solve stochastic optimization problems [26, 39], we do not assume that \mathcal{X} is bounded. We remark that the boundedness assumption of the singular values of ∇c^T guarantees the linear independence constraint qualification (LICQ). Note that it is generally not ideal to assume that the objective and constraint function and derivative values are bounded over \mathcal{X} containing stochastic iterates $\{x_{k,s}\}$. However, this assumption is reasonable in our problem setting if we assume the component functions $\{f_i\}$ have bounded derivatives over \mathcal{X} . In addition, this assumption can be loosened if one chooses to use constant step sizes. This assumption is similar to those in [4, 22].

We define the Lagrangian, $\mathcal{L} : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$, of (1.1) as $\mathcal{L}(x, y) := f(x) + y^T c(x)$, where $y \in \mathbb{R}^m$ represents a vector of Lagrange multipliers. Under Assumption 2.1 (and as a result of LICQ), necessary conditions for first-order stationarity with respect to (1.1) are given by

$$0 = \begin{bmatrix} \nabla_x \mathcal{L}(x, y) \\ \nabla_y \mathcal{L}(x, y) \end{bmatrix} = \begin{bmatrix} \nabla f(x) + \nabla c(x)y \\ c(x) \end{bmatrix}.$$

Next, we formalize our assumption on the gradient approximation employed by the SVR-SQP method. Given an iterate $x_{k,s} \in \mathbb{R}^n$ (where $(k, s) \in \mathbb{N} \times [\bar{S}]$), let $\tilde{g}_{k,s} \in \mathbb{R}^n$ be defined as

$$\tilde{g}_{k,s} := \frac{1}{b} \sum_{i \in I_{k,s}} \nabla f_i(x_{k,s}), \quad (2.2)$$

where $I_{k,s} \subset [N]$ of size b is a mini-batch (subset) of all the data. Throughout the paper, we refer to the gradient approximation in (2.2) as the *stochastic gradient*.

Assumption 2.3. *The gradient approximation (2.2) is an unbiased estimator of the true gradient of the objective function, i.e., we have that $\mathbb{E}_{k,s}[\tilde{g}_{k,s}] = g_{k,s}$, where $\mathbb{E}_{k,s}$ denotes the expectation taken conditioned on the event that the algorithm has reached $x_{k,s} \in \mathbb{R}^n$ in iteration $(k, s) \in \mathbb{N} \times [\bar{S}]$. (We impose an additional condition on this expectation in subsequent sections of the paper; see Lemma 4.2.) The unbiasedness assumption of $\tilde{g}_{k,s}$ can be easily satisfied, e.g., when each sample in the mini-batch $I_{k,s} \subset [N]$ is selected uniformly at random.*

Finally, the variance reduced gradient approximation employed by the SVR-SQP method $\bar{g}_{k,s} \in \mathbb{R}^n$ is defined as

$$\begin{aligned}\bar{g}_{k,s} &:= \frac{1}{b} \sum_{i \in I_{k,s}} (\nabla f_i(x_{k,s}) - (\nabla f_i(x_{k,0}) - \nabla f(x_{k,0}))) \\ &= \tilde{g}_{k,s} - \tilde{g}_{k,0} + g_{k,0},\end{aligned}\tag{2.3}$$

where $I_{k,s} \subset [N]$ of size b , and $x_{k,0}$ is known as the reference point (the initial point for the inner iterations of the k th outer iteration). Throughout the paper, we refer to the gradient approximation in (2.3) as the *SVRG gradient*. Under Assumption 2.3, it follows that the SVRG gradient is an unbiased estimate of the true gradient, i.e., $\mathbb{E}_{k,s}[\bar{g}_{k,s}] = g_{k,s}$.

3 Stochastic Variance Reduced Sequential Quadratic Programming

Our proposed algorithm (SVR-SQP) is based on the Sequential Quadratic Programming (SQP) paradigm. A high level description of the SVR-SQP method is as follows: SVR-SQP operates with two types of iterations (inner and outer), employs variance reduced approximations of the gradient of the objective function following (2.3) in lieu of the true gradient, and updates the iterates in SQP fashion.

Given an iterate $x_{k,s}$ for all $(k, s) \in \mathbb{N} \times [\bar{S}]$, the SVR-SQP methods proceeds to compute a search direction $\bar{d}_{k,s} \in \mathbb{R}^n$ by solving the following subproblem

$$\min_{\bar{d} \in \mathbb{R}^n} \bar{g}_{k,s}^T \bar{d} + \frac{1}{2} \bar{d}^T H_{k,s} \bar{d} \quad \text{s.t.} \quad c_{k,s} + J_{k,s} \bar{d} = 0,\tag{3.1}$$

where $\bar{g}_{k,s}$ is defined in (2.3) and $H_{k,s} \in \mathbb{S}^n$ satisfies Assumption 3.1 below.

Assumption 3.1. *The sequence $\{H_{k,s}\}$ is independent of $\{\bar{g}_{k,s}\}$ and is bounded in norm by $\kappa_H \in \mathbb{R}_{>0}$. In addition, there exists a constant $\zeta \in \mathbb{R}_{>0}$ such that, for all $(k, s) \in \mathbb{N} \times [\bar{S}]$, the matrix $H_{k,s}$ has the property that $u^T H_{k,s} u \geq \zeta \|u\|_2^2$ for all $u \in \mathbb{R}^n$ such that $J_{k,s} u = 0$.*

Under Assumptions 2.1 and 3.1, the solution of (3.1), denoted by $\bar{d}_{k,s} \in \mathbb{R}^n$, can be equivalently computed by solving the following linear system of equations

$$\begin{bmatrix} H_{k,s} & J_{k,s}^T \\ J_{k,s} & 0 \end{bmatrix} \begin{bmatrix} \bar{d}_{k,s} \\ \bar{y}_{k,s} \end{bmatrix} = - \begin{bmatrix} \bar{g}_{k,s} \\ c_{k,s} \end{bmatrix},\tag{3.2}$$

where $\bar{y}_{k,s} \in \mathbb{R}^m$ is the vector of associated Lagrange multipliers of (3.1). The linear system in (3.2) has a unique solution under Assumptions 2.1 and 3.1; see [28].

With a search direction $\bar{d}_{k,s} \in \mathbb{R}^n$ in hand, SVR-SQP proceeds to utilize a merit function, $\phi : \mathbb{R}^n \times \mathbb{R}_{>0} \rightarrow \mathbb{R}$, to judge the quality of the computed step (in terms of stationarity

and feasibility), and then compute a positive step size $\bar{\alpha}_{k,s} \in \mathbb{R}_{>0}$ in order to update the current iterate $x_{k,s} \in \mathbb{R}^n$ via

$$x_{k,s+1} = x_{k,s} + \bar{\alpha}_{k,s} \bar{d}_{k,s}. \quad (3.3)$$

Similar to [4], our algorithm makes use of, possibly the most common merit (penalty) function, the l_1 -norm merit function, defined as

$$\phi(x, \tau) := \tau f(x) + \|c(x)\|_1, \quad (3.4)$$

where $\tau \in \mathbb{R}_{>0}$ is known as the merit (penalty) parameter and whose value is set adaptively as the optimization progresses. Before we proceed, we introduce two quantities that are used in our proposed algorithm, and that are vital to the analysis. First, a local linear model of the merit function $l : \mathbb{R}^n \times \mathbb{R}_{>0} \times \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ is defined by

$$l(x, \tau, g, d) := \tau(f(x) + g^T d) + \|c(x) + \nabla c(x)^T d\|_1. \quad (3.5)$$

Second, the reduction function of the local linear model of the merit function $\Delta l : \mathbb{R}^n \times \mathbb{R}_{>0} \times \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$, given $d \in \mathbb{R}^n$ with $c(x) + \nabla c(x)^T d = 0$, is defined by

$$\Delta l(x, \tau, g, d) := l(x, \tau, g, 0) - l(x, \tau, g, d) = -\tau g^T d + \|c(x)\|_1. \quad (3.6)$$

Given a search direction $\bar{d}_{k,s} \in \mathbb{R}^n$, the merit parameter update strategy goes as follows. To begin with, a trial merit parameter is defined as

$$\bar{\tau}_{k,s}^{trial} \leftarrow \begin{cases} \infty & \text{if } \bar{g}_{k,s}^T \bar{d}_{k,s} + \max\{\bar{d}_{k,s}^T H_{k,s} \bar{d}_{k,s}, 0\} \leq 0; \\ \frac{(1-\sigma)\|c_{k,s}\|_1}{\bar{g}_{k,s}^T \bar{d}_{k,s} + \max\{\bar{d}_{k,s}^T H_{k,s} \bar{d}_{k,s}, 0\}} & \text{otherwise,} \end{cases} \quad (3.7)$$

where the parameter $\sigma \in (0, 1)$ is user-defined. It follows that $\bar{\tau}_{k,s}^{trial} > 0$ since if $\|c_{k,s}\|_1 = 0$, by Assumption 3.1 and (3.2), $\bar{g}_{k,s}^T \bar{d}_{k,s} + \max\{\bar{d}_{k,s}^T H_{k,s} \bar{d}_{k,s}, 0\} = \bar{g}_{k,s}^T \bar{d}_{k,s} + \bar{d}_{k,s}^T H_{k,s} \bar{d}_{k,s} = -\bar{d}_{k,s}^T J_{k,s}^T \bar{y}_{k,s} = c_{k,s}^T \bar{y}_{k,s} = 0$. Next, for some user-defined parameter $\epsilon_\tau \in (0, 1)$, $\bar{\tau}_{k,s}$ is computed via

$$\bar{\tau}_{k,s} \leftarrow \begin{cases} \bar{\tau}_{k,s-1} & \text{if } \bar{\tau}_{k,s-1} \leq \bar{\tau}_{k,s}^{trial} \\ (1 - \epsilon_\tau) \bar{\tau}_{k,s}^{trial} & \text{otherwise.} \end{cases} \quad (3.8)$$

Note, the above rule ensures that $\bar{\tau}_{k,s} \leq \bar{\tau}_{k,s}^{trial}$. Moreover, and more importantly, the updates (3.7)–(3.8) ensure that

$$\Delta l(x_{k,s}, \bar{\tau}_{k,s}, \bar{g}_{k,s}, \bar{d}_{k,s}) \geq \bar{\tau}_{k,s} \max\{\bar{d}_{k,s}^T H_{k,s} \bar{d}_{k,s}, 0\} + \sigma \|c_{k,s}\|_1. \quad (3.9)$$

The above inequality plays a critical role in our algorithm and analysis.

Finally, the SVR-SQP method computes a positive step size. We propose two different step size selection strategies; a constant step size strategy and an adaptive step size strategy. The constant step size strategy, similar to that in [31], specifies an upper bound on acceptable step sizes (see Theorem 4.12 for the exact specification).

The adaptive step size strategy, inspired by [4], is motivated by the desire to select a step size that minimizes an upper bound on the change in the merit function. By the definition of the merit function (3.4) and under Assumption 2.1, the upper bound on the change in the merit function is a convex (strongly-convex when $\|\bar{d}_{k,s}\| \neq 0$), piece-wise quadratic function in $\bar{\alpha}_{k,s} \in \mathbb{R}_{>0}$,

$$\begin{aligned} & \phi(x_{k,s+1}, \bar{\tau}_{k,s}) - \phi(x_{k,s}, \bar{\tau}_{k,s}) \\ & \leq \bar{\alpha}_{k,s} \bar{\tau}_{k,s} g_{k,s}^T \bar{d}_{k,s} + (|1 - \bar{\alpha}_{k,s}| - 1) \|c_{k,s}\|_1 + \frac{1}{2} (\bar{\tau}_{k,s} L + \Gamma) \bar{\alpha}_{k,s}^2 \|\bar{d}_{k,s}\|_2^2. \end{aligned} \quad (3.10)$$

(See [7, Lemma 3.1] for derivation of above inequality.) Our adaptive strategy attempts to select a step size that approximately minimizes this upper bound. To this end, at iteration $(k, s) \in \mathbb{N} \times [\bar{S}]$, two trial step sizes are computed, specifically,

$$\bar{\tilde{\alpha}}_{k,s} \leftarrow \min \left\{ \frac{\Delta l(x_{k,s}, \bar{\tau}_{k,s}, \bar{g}_{k,s}, \bar{d}_{k,s})}{(\bar{\tau}_{k,s} L_{k,s} + \Gamma_{k,s}) \|\bar{d}_{k,s}\|_2^2}, \alpha_u \right\} \beta \quad (3.11)$$

$$\text{and } \tilde{\alpha}_{k,s} \leftarrow \bar{\tilde{\alpha}}_{k,s} - \frac{4 \|c_{k,s}\|_1}{(\bar{\tau}_{k,s} L_{k,s} + \Gamma_{k,s}) \|\bar{d}_{k,s}\|_2^2}, \quad (3.12)$$

where $\alpha_u \in \mathbb{R}_{>0}$ is a user-defined parameter that is introduced here to avoid the step size being arbitrarily large (the precise specification of α_u is given in Section 4.4). Due to the nonsmoothness of the upper bound (notice, nonsmooth point at $\bar{\alpha}_{k,s} = 1$), the approximate minimizer, and the step size used by the SVR-SQP, is set as

$$\bar{\alpha}_{k,s} \leftarrow \begin{cases} \bar{\tilde{\alpha}}_{k,s} & \text{if } \bar{\tilde{\alpha}}_{k,s} < 1 \\ 1 & \text{if } \tilde{\alpha}_{k,s} \leq 1 \leq \bar{\tilde{\alpha}}_{k,s} \\ \tilde{\alpha}_{k,s} & \text{if } \tilde{\alpha}_{k,s} > 1 \end{cases} \quad (3.13)$$

Our proposed algorithm SVR-SQP is fully described in Algorithm 1. Similar to the SVRG method [16, 31], SVR-SQP operates with inner and outer iterations. Each outer iteration commences with the computation of the full gradient of the objective function at the *reference point* $x_{k,0}$, i.e., $g_{k,0}$. Given the $g_{k,0}$ at every inner iteration, a stochastic variance reduced gradient is computed via (2.3), and then paralleling the SQP paradigm, the search direction is computed by solving the linear system given in (3.2). Finally, similar to the stochastic SQP algorithm proposed in [4], the merit parameter is updated following (3.7)–(3.8), a step size is computed, and the current iterate is updated. The algorithm allows for two different step size choices: constant step size (**Option I**) and adaptive step size (**Option II**) via the equations (3.11)–(3.13).

Remark 3.2. *Due to the nature of the SVRG gradient estimate, our proposed algorithm is of nested nature (inner and outer iterations), and the full batch gradient is computed once*

Algorithm 1 Stochastic Variance Reduced SQP (SVR-SQP)

Require: $x_{-1,S} \in \mathbb{R}^n$ (initial iterate); $\bar{\tau}_{-1,S-1} \in \mathbb{R}_{>0}$ (initial merit parameter value); $\epsilon_\tau \in (0, 1)$ (merit decrease parameter); $\sigma \in (0, 1)$ (model reduction parameter), $b \in [N - 1]$ (batch size)

Require (Option I: Constant step size algorithm): $\alpha \in (0, 1]$ (constant step size parameter)

Require (Option II: Adaptive step size algorithm): $\beta \in (0, 1]$ (adaptive step size parameter), $\alpha_u \in \mathbb{R}_{>0}$ (adaptive step size bound)

```
1: for  $k = 0, 1, \dots$ , do
2:   Set  $x_{k,0} = x_{k-1,S}$ ;  $\bar{\tau}_{k,-1} = \bar{\tau}_{k-1,S-1}$ 
3:   Compute gradient  $g_{k,0} = \nabla f(x_{k,0})$ 
4:   for  $s = 0, 1, \dots, S - 1$  do
5:     Choose a mini-batch  $I_{k,s} \subset [N]$  of size  $b$ , and compute  $\bar{g}_{k,s}$  via (2.3)
6:     Compute  $(\bar{d}_{k,s}, \bar{y}_{k,s})$  as the solution of (3.2)
7:     if  $\bar{d}_{k,s} = 0$  then set  $x_{k,s+1} \leftarrow x_{k,s}$ ,  $\bar{\tau}_{k,s} \leftarrow \bar{\tau}_{k,s-1}$ ; go to Line 5
8:     end if
9:     Set  $\bar{\tau}_{k,s}^{trial}$  via (3.7) and  $\bar{\tau}_{k,s}$  via (3.8)
10:    Set step size parameter  $\bar{\alpha}_{k,s}$ 
11:     Option I: Set  $\bar{\alpha}_{k,s} = \alpha$ 
12:     Option II: Compute  $\tilde{\alpha}_{k,s}$  and  $\tilde{\tilde{\alpha}}_{k,s}$  via (3.11)–(3.12)
13:     Set  $\bar{\alpha}_{k,s}$  via (3.13)
14:    Set  $x_{k,s+1} \leftarrow x_{k,s} + \bar{\alpha}_{k,s} \bar{d}_{k,s}$ 
15:   end for
16: end for
```

every outer iteration in order to reduce the variance of the gradient estimate. Our proposed algorithm has two options for selecting the step size. Algorithm 1 with **Option I** (constant step size) can be considered a natural extension of [31] to the equality constrained setting. Algorithm 1 with **Option II** (adaptive step size) can be considered a natural extension of [4] where the stochastic gradient estimate is replaced by the SVRG gradient estimate.

4 Convergence Analysis

In this section, we present convergence guarantees for SVR-SQP (Algorithm 1) under the two step size regimes. We begin with some general technical lemmas (Section 4.1), then discuss the behavior of the merit parameter (Section 4.2), and finally present our main theoretical results for constant and adaptive step size choices (Sections 4.3 and 4.4, respectively). Throughout this section we assume that Assumptions 2.1, 2.3 and 3.1 hold; for brevity, we do not remind the reader of this fact within the statement of each result.

For the purposes of the analysis, we define several deterministic quantities that are never explicitly computed in Algorithm 1. First, $d_{k,s} \in \mathbb{R}^n$ and $y_{k,s} \in \mathbb{R}^m$ for all $(k, s) \in \mathbb{N} \times [\bar{S}]$ are defined as

$$\begin{bmatrix} H_{k,s} & J_{k,s}^T \\ J_{k,s} & 0 \end{bmatrix} \begin{bmatrix} d_{k,s} \\ y_{k,s} \end{bmatrix} = - \begin{bmatrix} g_{k,s} \\ c_{k,s} \end{bmatrix}. \quad (4.1)$$

We note that the only difference between (3.2) and (4.1) is the right-hand-side, where the gradient approximation $\bar{g}_{k,s}$ is replaced by the true gradient $g_{k,s}$. Moreover, for all $(k, s) \in \mathbb{N} \times [\bar{S}]$, $\tau_{k,s}^{trial}$ and $\tau_{k,s}$ are the deterministic analogues of stochastic merit parameters values $\bar{\tau}_{k,s}^{trial}$ and $\bar{\tau}_{k,s}$, respectively, where $\bar{g}_{k,s}$ and $\bar{d}_{k,s}$ are replaced by $g_{k,s}$ and $d_{k,s}$ in (3.7) and (3.8).

4.1 General results

The first lemma of this section consists of several technical conditions that are used for the convergence analysis of the SVR-SQP method. These conditions are analogues of those in [4, Lemma 3.4]¹.

Lemma 4.1. *There exists a constant $\kappa_l \in \mathbb{R}_{>0}$ such that the following statements hold true for all $(k, s) \in \mathbb{N} \times [\bar{S}]$:*

- (a) $\Delta l(x_{k,s}, \bar{\tau}_{k,s}, g_{k,s}, d_{k,s}) \geq \kappa_l \bar{\tau}_{k,s} \|d_{k,s}\|_2^2$;
- (b) $\Delta l(x_{k,s}, \bar{\tau}_{k,s}, \bar{g}_{k,s}, \bar{d}_{k,s}) \geq \kappa_l \bar{\tau}_{k,s} \|\bar{d}_{k,s}\|_2^2$;
- (c) and,

$$\begin{aligned} & \phi(x_{k,s+1}, \bar{\tau}_{k,s}) - \phi(x_{k,s}, \bar{\tau}_{k,s}) \\ & \leq \bar{\alpha}_{k,s} \bar{\tau}_{k,s} g_{k,s}^T \bar{d}_{k,s} + |1 - \bar{\alpha}_{k,s}| \|c_{k,s}\|_1 - \|c_{k,s}\|_1 + \frac{1}{2} (\bar{\tau}_{k,s} L + \Gamma) \bar{\alpha}_{k,s}^2 \|\bar{d}_{k,s}\|_2^2. \end{aligned} \quad (4.2)$$

Proof. First note that condition (b) is the stochastic analogue of condition (a). Conditions (a) and (b) can be derived directly from [4, Lemmas 2.11, 2.12 & 3.4(c),(d)]. Inequality (c) is identical to that in [4, Lemma 3.4], but accounts for the double indices of the SVR-SQP algorithm. \square

In the next lemma, we bound the error in the gradient approximation $\bar{g}_{k,s}$ employed by the SVR-SQP algorithm.

Lemma 4.2. *Let $\bar{g}_{k,s} \in \mathbb{R}^n$ be the gradient approximation computed by Algorithm 1 via (2.3). Then, for all $(k, s) \in \mathbb{N} \times [\bar{S}]$,*

$$\mathbb{E}_{k,s} [\|\bar{g}_{k,s} - \nabla f(x_{k,s})\|_2^2] \leq M_{k,s}, \quad (4.3)$$

¹For lemmas with proofs equivalent to those in [4], we refer interested reader to the appropriate sections.

where $M_{k,s} = \frac{L^2}{b} \|x_{k,s} - x_{k,0}\|_2^2$, and $\mathbb{E}_{k,s}$ denotes the expectation taken conditioned on the event that the algorithm has reached $x_{k,0} \in \mathbb{R}^n$ in (outer) iteration $k \in \mathbb{N}$ and $x_{k,s}$ in (outer-inner) iteration $(k, s) \in \mathbb{N} \times [\bar{S}]$.

Proof. For the ease of exposition, we introduce the following notation,

$$\zeta_{k,s} = \frac{1}{b} \sum_{i \in I_{k,s}} (\nabla f_i(x_{k,s}) - \nabla f_i(x_{k,0})). \quad (4.4)$$

It follows by Assumption 2.3 and (4.4) that $\mathbb{E}_{k,s}[\zeta_{k,s}] = \nabla f(x_{k,s}) - \nabla f(x_{k,0})$. By the definition of $\bar{g}_{k,s}$ (2.3) and Assumption 2.1, and the facts that $\mathbb{E}[\|z - \mathbb{E}[z]\|^2] \leq \mathbb{E}[\|z\|^2]$ for random variable z and $\mathbb{E}[\|z_1 + \dots + z_r\|^2] = \mathbb{E}[\|z_1\|^2 + \dots + \|z_r\|^2]$ for independent mean zero random variables z_1, \dots, z_r [31], it follows that

$$\begin{aligned} \mathbb{E}_{k,s} [\|\bar{g}_{k,s} - \nabla f(x_{k,s})\|_2^2] &= \mathbb{E}_{k,s} [\|\zeta_{k,s} + \nabla f(x_{k,0}) - \nabla f(x_{k,s})\|_2^2] \\ &= \mathbb{E}_{k,s} [\|\zeta_{k,s} - \mathbb{E}_{k,s}[\zeta_{k,s}]\|_2^2] \\ &= \frac{1}{b^2} \mathbb{E}_{k,s} \left[\left\| \sum_{i \in I_{k,s}} (\nabla f_i(x_{k,s}) - \nabla f_i(x_{k,0}) - \mathbb{E}_{k,s}[\zeta_{k,s}]) \right\|_2^2 \right] \\ &= \frac{1}{b^2} \mathbb{E}_{k,s} \left[\sum_{i \in I_{k,s}} \|\nabla f_i(x_{k,s}) - \nabla f_i(x_{k,0}) - \mathbb{E}_{k,s}[\zeta_{k,s}]\|_2^2 \right] \\ &\leq \frac{1}{b^2} \mathbb{E}_{k,s} \left[\sum_{i \in I_{k,s}} \|\nabla f_i(x_{k,s}) - \nabla f_i(x_{k,0})\|_2^2 \right] \\ &\leq \frac{L^2}{b} \mathbb{E}_{k,s} [\|x_{k,s} - x_{k,0}\|_2^2] = \frac{L^2}{b} \|x_{k,s} - x_{k,0}\|_2^2. \end{aligned}$$

□

Lemma 4.2 is one of the major differences between this work and that in [4], and in [3, 11, 12]. Specifically, in [4] (and in [3, 11, 12]) it is assumed that the variance in the stochastic gradients employed is bounded uniformly by a constant $M \in \mathbb{R}_{>0}$ (i.e., this would be equivalent to having $M_{k,s} = M$ for all $(k, s) \in \mathbb{N} \times [\bar{S}]$). This is a classical assumption for the convergence analysis of the SG method [6, 33], which leads to the fact that the algorithm can only converge to a neighborhood depending on M in expectation when a constant step size is employed. By employing variance reduced gradients, this allows us to control the variance, and diminish it as needed, in order to prove exact convergence of first-order stationary measure in expectation.

In the next two lemmas, we present some useful bounds pertaining to the solutions of the linear system (3.2).

Lemma 4.3. For all $(k, s) \in \mathbb{N} \times [\bar{S}]$, we always have $\mathbb{E}_{k,s}[\bar{d}_{k,s}] = d_{k,s}$ and $\mathbb{E}_{k,s}[\bar{y}_{k,s}] = y_{k,s}$. In addition, there exists some constant $\kappa_d \in \mathbb{R}_{>0}$, independent of (k, s) and any run of the algorithm, with $\mathbb{E}_{k,s}[\|\bar{d}_{k,s} - d_{k,s}\|_2] \leq \kappa_d \sqrt{M_{k,s}}$.

Proof. The proof of this lemma is similar to that in [4, Lemma 3.8]. The first statement follows from the facts that (i) conditioned on $x_{k,s}$, the matrix on the left-hand-side of (3.2) is deterministic; (ii) under Assumption 2.1, the matrix is invertible; (iii) under Assumption 2.3, $\mathbb{E}_{k,s}[\bar{g}_{k,s}] = \nabla f(x_{k,s})$; and (iv) expectation is a linear operator. By (3.2) for any realization $\bar{g}_{k,s}$, it follows that

$$\begin{bmatrix} \bar{d}_{k,s} - d_{k,s} \\ \bar{y}_{k,s} - y_{k,s} \end{bmatrix} = - \begin{bmatrix} H_{k,s} & J_{k,s}^T \\ J_{k,s} & 0 \end{bmatrix}^{-1} \begin{bmatrix} \bar{g}_{k,s} - \nabla f(x_{k,s}) \\ 0 \end{bmatrix}. \quad (4.5)$$

The second result follows by Jensen's inequality, the concavity of the square root, and Lemma 4.2, and where $\kappa_d \in \mathbb{R}_{>0}$ is an upper bound on the norm of the matrix in (4.5). \square

Lemma 4.4. For all $(k, s) \in \mathbb{N} \times [\bar{S}]$, it follows that

$$g_{k,s}^T d_{k,s} \geq \mathbb{E}_{k,s}[\bar{g}_{k,s}^T \bar{d}_{k,s}] \geq g_{k,s}^T d_{k,s} - \zeta^{-1} M_{k,s}, \quad (4.6)$$

Proof. The proof is identical to [4, Lemma 3.9 (proof)] with M replaced by $M_{k,s}$ ($M_{k,s}$ defined in Lemma 4.2). \square

We conclude this subsection by defining a Lyapunov function $R : \mathbb{R}^n \times \mathbb{R}^n \times \mathbb{R}_{>0} \times \mathbb{R}_{>0} \rightarrow \mathbb{R}$ that will be used in the analysis. Specifically,

$$R_{k,s} := R(x_{k,s}, x_{k,0}, \bar{\tau}_{k,s}, \lambda_s) = \mathbb{E}_{k,s} [\phi(x_{k,s}, \bar{\tau}_{k,s}) + \lambda_s \|x_{k,s} - x_{k,0}\|_2^2], \quad (4.7)$$

where $x_{k,s} \in \mathbb{R}^n$ and $\bar{\tau}_{k,s} \in \mathbb{R}_{>0}$ are the iterate and merit parameter at outer-inner iteration $(k, s) \in \mathbb{N} \times [\bar{S}]$, respectively, $x_{k,0} \in \mathbb{R}^n$ is the reference point at the k th outer iteration and $\lambda_s \in \mathbb{R}_{>0}$ is a parameter (defined explicitly later in the analysis). The Lyapunov function is defined as the expected value of the merit function plus the distance squared between any inner iterate and the reference iterate parameterized by a constant. When $s = 0$ the Lyapunov function only involves the merit function. Moreover, the last term in the Lyapunov function is similar to that of the upper bound in the variance of the SVRG gradient (see Lemma 4.2), and, if the iterates converge, the Lyapunov function reduces to the expected value of the merit function. This is by construction, and will allow us to prove strong theoretical guarantees for SVR-SQP.

4.2 Merit Parameter behavior

The behavior of the merit parameter $\bar{\tau}_{k,s}$ requires careful considerations as it is a crucial component of the SVR-SQP method and the analysis. Specifically, what is important is the behavior of $\bar{\tau}_{k,s}$ for large $k \in \mathbb{N}$. As described in [4], there are three possible outcomes for $\bar{\tau}_{k,s}$: (i) converges to zero (vanishes); (ii) remains constant at a large positive value; (iii) remains constant at a sufficiently small positive value. We argue that in the finite-sum setting (1.1) and under reasonable assumptions, outcome (i) is not possible, and outcome (ii) occurs with probability zero. To show the former, i.e., outcome (i) is not possible, we make the following assumption.

Assumption 4.5. *Each component $\{f_i\}$ of the objective function f in (1.1) has bounded derivatives over \mathcal{X} (defined in Assumption 2.1).*

Under Assumption 4.5, the merit parameter cannot vanish.

Lemma 4.6. *Suppose Assumption 4.5 holds, then there exists $\bar{k}_\tau \in \mathbb{N}$ and $\bar{\tau}_{const} \in \mathbb{R}_{>0}$ such that $\bar{\tau}_{k,s} = \bar{\tau}_{const}$ for all $k \geq \bar{k}_\tau$ and $s \in [\bar{S}]$.*

Proof. Under Assumption 4.5, there exists $g_{\max} \in \mathbb{R}_{>0}$ such that $\|\bar{g}_{k,s} - g_{k,s}\| \leq g_{\max}$ for all $k \geq \bar{k}_\tau$ and $s \in [\bar{S}]$. The desired conclusion follows using similar arguments as in [4, Proposition 3.18]. \square

Following a similar argument as that in [4], we argue the latter, i.e., (ii) occurs with probability zero.

Lemma 4.7. *Suppose event $E_{\tau\uparrow}$ occurs in the sense that there exists infinite $(\bar{\mathcal{K}}_\tau, \bar{\mathcal{S}}_\tau) \subseteq \mathbb{N} \times [\bar{S}]$ and $\bar{\tau}_{big} \in \mathbb{R}_{>0}$ such that*

$$\bar{\tau}_{k,s} = \bar{\tau}_{big} > \tau_{k,s}^{trial} \text{ for all } (k,s) \in (\bar{\mathcal{K}}_\tau, \bar{\mathcal{S}}_\tau). \quad (4.8)$$

Moreover, suppose that $\bar{d}_{k,s}^T H_{k,s} \bar{d}_{k,s} \geq 0$ for all $(k,s) \in \mathbb{N} \times [\bar{S}]$. Then, $E_{\tau\uparrow}$ occurs with probability zero.

Proof. By (3.2) and $\bar{d}_{k,s}^T H_{k,s} \bar{d}_{k,s} \geq 0$ it follows that $\bar{g}_{k,s}^T \bar{d}_{k,s} + \max\{\bar{d}_{k,s}^T H_{k,s} \bar{d}_{k,s}, 0\} = \bar{g}_{k,s}^T \bar{d}_{k,s} + \bar{d}_{k,s}^T H_{k,s} \bar{d}_{k,s} = c_{k,s}^T \bar{y}_{k,s}$. Similarly, by (4.1) it follows that $g_{k,s}^T d_{k,s} + \max\{d_{k,s}^T H_{k,s} d_{k,s}, 0\} = c_{k,s}^T y_{k,s}$. Since we are considering an objective function composed of a finite number of components, there are a finite number of realizations for $\bar{g}_{k,s}$. Among $\binom{N}{b}$ possible realizations of $\bar{g}_{k,s}$, there should at least be one realization such that $c_{k,s}^T \bar{y}_{k,s}$ is no smaller than $c_{k,s}^T y_{k,s}$, since $\mathbb{E}_{k,s}[c_{k,s}^T \bar{y}_{k,s}] = c_{k,s}^T y_{k,s}$ by Lemma 4.3. Hence, it follows that

$$\begin{aligned} & \mathbb{P}[\bar{g}_{k,s}^T \bar{d}_{k,s} + \max\{\bar{d}_{k,s}^T H_{k,s} \bar{d}_{k,s}, 0\} \geq g_{k,s}^T d_{k,s} + \max\{d_{k,s}^T H_{k,s} d_{k,s}, 0\}] \\ &= \mathbb{P}[c_{k,s}^T \bar{y}_{k,s} \geq c_{k,s}^T y_{k,s}] \\ &\geq \frac{1}{\binom{N}{b}} \end{aligned}$$

The desired conclusion then follows from [4, Proposition 3.16]. \square

If Assumption 4.5 holds and $\bar{d}_{k,s}^T H_{k,s} \bar{d}_{k,s} \geq 0$ for all $(k, s) \in \mathbb{N} \times [\bar{S}]$, then $\bar{\tau}_{k,s}$ is guaranteed to remain constant at a sufficiently small positive value eventually with probability 1. While one could prove such a corollary, we instead assume that the merit parameter remains constant at a sufficiently small positive value because the above only provides sufficient conditions, and this merit parameter behavior can potentially be exhibited on a wider class of problems. For the remainder of the paper, we will assume that the merit parameter remains constant at a sufficiently small positive value, and formalize this assumption below.

Assumption 4.8. *Suppose event $E_{\tau_{\min}}$ occurs in the sense that there exists an iteration number $\bar{k}_\tau \in \mathbb{N}$ and a merit parameter value $\bar{\tau}_{\min} \in \mathbb{R}_{>0}$ such that,*

$$\bar{\tau}_{k,s} = \bar{\tau}_{\min} \leq \tau_{k,s}^{\text{trial}} \text{ for all } k \geq \bar{k}_\tau \text{ and } s \in [\bar{S}]. \quad (4.9)$$

In addition, we further assume that the stochastic gradient sequence $\{\bar{g}_{k,s}\}_{k \geq \bar{k}_\tau, s \in [\bar{S}]}$ satisfies $\mathbb{E}_{k,s,\tau_{\min}}[\bar{g}_{k,s}] = g_{k,s}$, where $\mathbb{E}_{k,s,\tau_{\min}}$ denotes the expectation taken conditioned on the event that $E_{\tau_{\min}}$ occurs and that the algorithm has reached $x_{k,0} \in \mathbb{R}^n$ in (outer) iteration $k \in \mathbb{N}$ and $x_{k,s}$ in (outer-inner) iteration $(k, s) \in \mathbb{N} \times [\bar{S}]$.

Assumption 4.8 is a critical assumption in proving the convergence of the SVR-SQP method, and will be assumed to hold throughout the remainder of this section. For ease of exposition, we use $\mathbb{E}_{k,s}$ to denote $\mathbb{E}_{k,s,\tau_{\min}}$, and we define the following quantity

$$\mathbb{E}_{\tau_{\min}}[\cdot] := \mathbb{E}[\cdot | \text{Assumption 4.8}],$$

i.e., the total expectation conditioned on the event $E_{\tau_{\min}}$. Moreover, we define a constant $\phi_{\inf} > -\infty$ as

$$\phi_{\inf} := \inf_{x \in \mathcal{X}} \phi(x, \bar{\tau}_{\min}),$$

and whose existence is guaranteed under Assumptions 2.1 and 4.8.

Before we proceed, we state and prove one more technical lemma that will be used in the analysis in Sections 4.3 and 4.4.

Lemma 4.9. *Suppose that Assumption 4.8 holds. For all $k \geq \bar{k}_\tau$ and $s \in [\bar{S}]$, it follows that*

$$\mathbb{E}_{k,s}[\Delta l(x_{k,s}, \bar{\tau}_{k,s}, \bar{g}_{k,s}, \bar{d}_{k,s})] \leq \Delta l(x_{k,s}, \bar{\tau}_{\min}, g_{k,s}, d_{k,s}) + \bar{\tau}_{\min} \zeta^{-1} M_{k,s}.$$

Proof. For $k \geq \bar{k}_\tau$ and $s \in [\bar{S}]$, it follows by (3.6) and Lemma 4.4 that

$$\begin{aligned} \mathbb{E}_{k,s}[\Delta l(x_{k,s}, \bar{\tau}_{k,s}, \bar{g}_{k,s}, \bar{d}_{k,s})] &= \mathbb{E}_{k,s} \left[-\bar{\tau}_{\min} \bar{g}_{k,s}^T \bar{d}_{k,s} + \|c_{k,s}\|_1 \right] \\ &= \mathbb{E}_{k,s} \left[-\bar{\tau}_{\min} \bar{g}_{k,s}^T \bar{d}_{k,s} + \bar{\tau}_{\min} g_{k,s}^T d_{k,s} \right. \\ &\quad \left. - \bar{\tau}_{\min} g_{k,s}^T d_{k,s} + \|c_{k,s}\|_1 \right] \\ &\leq \Delta l(x_{k,s}, \bar{\tau}_{\min}, g_{k,s}, d_{k,s}) + \bar{\tau}_{\min} \zeta^{-1} M_{k,s}. \end{aligned}$$

□

4.3 Constant step size analysis

In this subsection, we present convergence results for Algorithm 1 with the constant step size strategy (**Option I**) under Assumption 4.8. The first lemma provides a useful upper bound for the difference in merit function after a step.

Lemma 4.10. *Suppose that Assumption 4.8 holds and $\alpha \in (0, 1]$. For all $k \geq \bar{k}_\tau$ and $s \in [\bar{S}]$, it follows that*

$$\begin{aligned} & \phi(x_{k,s+1}, \bar{\tau}_{k,s}) - \phi(x_{k,s}, \bar{\tau}_{k,s}) \\ & \leq -\alpha \Delta l(x_{k,s}, \bar{\tau}_{\min}, g_{k,s}, d_{k,s}) + \frac{\alpha^2(\bar{\tau}_{\min}L + \Gamma)}{2\bar{\tau}_{\min}\kappa_l} \Delta l(x_{k,s}, \bar{\tau}_{\min}, \bar{g}_{k,s}, \bar{d}_{k,s}) \\ & \quad + \alpha \bar{\tau}_{\min} g_{k,s}^T (\bar{d}_{k,s} - d_{k,s}). \end{aligned}$$

Proof. For $\alpha \in (0, 1]$, $k \geq \bar{k}_\tau$ and $s \in [\bar{S}]$, by (3.6) and Lemma 4.1(b) it follows that

$$\begin{aligned} & \phi(x_{k,s+1}, \bar{\tau}_{k,s}) - \phi(x_{k,s}, \bar{\tau}_{k,s}) \\ & \leq \alpha(\bar{\tau}_{\min} g_{k,s}^T \bar{d}_{k,s} - \|c_{k,s}\|_1) + \frac{1}{2}(\bar{\tau}_{\min}L + \Gamma)\alpha^2 \|\bar{d}_{k,s}\|_2^2 \\ & = \alpha(\bar{\tau}_{\min} g_{k,s}^T d_{k,s} - \|c_{k,s}\|_1) + \alpha \bar{\tau}_{\min} g_{k,s}^T (\bar{d}_{k,s} - d_{k,s}) + \frac{1}{2}(\bar{\tau}_{\min}L + \Gamma)\alpha^2 \|\bar{d}_{k,s}\|_2^2 \\ & \leq -\alpha \Delta l(x_{k,s}, \bar{\tau}_{\min}, g_{k,s}, d_{k,s}) + \alpha \bar{\tau}_{\min} g_{k,s}^T (\bar{d}_{k,s} - d_{k,s}) \\ & \quad + \frac{\alpha^2(\bar{\tau}_{\min}L + \Gamma)}{2\bar{\tau}_{\min}\kappa_l} \Delta l(x_{k,s}, \bar{\tau}_{\min}, \bar{g}_{k,s}, \bar{d}_{k,s}). \end{aligned}$$

□

The next lemma is the central lemma of this subsection; it provides a useful upper bound on the expected value of the sum (over all inner and outer iterations) of the model reduction function of the merit function.

Lemma 4.11. *Suppose that Assumption 4.8 holds. Let $\lambda_S = 0$, and*

$$\begin{aligned} \lambda_s &= \lambda_{s+1} \left(\frac{\alpha^2 L^2}{\kappa_l b \zeta} + \alpha z + 1 \right) + \frac{\alpha^2 (\bar{\tau}_{\min} L + \Gamma) L^2}{2 \kappa_l b \zeta}, \\ \Lambda_s &= \alpha - \frac{\alpha^2 (\bar{\tau}_{\min} L + \Gamma)}{2 \bar{\tau}_{\min} \kappa_l} - \lambda_{s+1} \frac{\alpha}{\bar{\tau}_{\min} \kappa_l} \left(\alpha + \frac{1}{z} \right), \end{aligned} \tag{4.10}$$

for $s \in [\bar{S}]$, where $\alpha \in (0, 1]$, $z \in \mathbb{R}_{>0}$, $\lambda_s \in \mathbb{R}_{>0}$ are chosen such that $\Lambda_s \in \mathbb{R}_{>0}$, and $\Lambda_{\min} = \min_{s \in [\bar{S}]} \Lambda_s$. Then, for all $k \geq \bar{k}_\tau$ and $s \in [\bar{S}]$, the sequence of iterates $\{x_{k,s}\}$ generated by Algorithm 1 (**Option I**) satisfy

$$\mathbb{E}_{\tau_{\min}} \left[\frac{1}{(k - \bar{k}_\tau + 1) \bar{S}} \sum_{j=\bar{k}_\tau}^k \sum_{s=0}^{S-1} \Delta l(x_{j,s}, \bar{\tau}_{\min}, g_{j,s}, d_{j,s}) \right] \leq \frac{\mathbb{E}_{\tau_{\min}} [\phi(x_{\bar{k}_\tau, 0}, \bar{\tau}_{\min})] - \phi_{\inf}}{(k - \bar{k}_\tau + 1) S \Lambda_{\min}}. \tag{4.11}$$

Proof. Consider arbitrary $k \geq \bar{k}_\tau$ and $s \in [\bar{S}]$. By Lemmas 4.3, 4.9 and 4.10, we have

$$\begin{aligned} \mathbb{E}_{k,s}[\phi(x_{k,s+1}, \bar{\tau}_{k,s})] &\leq \mathbb{E}_{k,s}[\phi(x_{k,s}, \bar{\tau}_{k,s})] - \alpha \Delta l(x_{k,s}, \bar{\tau}_{\min}, g_{k,s}, d_{k,s}) \\ &\quad + \frac{(\bar{\tau}_{\min} L + \Gamma) \alpha^2}{2 \bar{\tau}_{\min} \kappa_l} (\Delta l(x_{k,s}, \bar{\tau}_{\min}, g_{k,s}, d_{k,s}) + \bar{\tau}_{\min} \zeta^{-1} M_{k,s}) \\ &= \mathbb{E}_{k,s}[\phi(x_{k,s}, \bar{\tau}_{k,s})] - \left(\alpha - \frac{\alpha^2 (\bar{\tau}_{\min} L + \Gamma)}{2 \bar{\tau}_{\min} \kappa_l} \right) \Delta l(x_{k,s}, \bar{\tau}_{\min}, g_{k,s}, d_{k,s}) \\ &\quad + \frac{\alpha^2 (\bar{\tau}_{\min} L + \Gamma) L^2}{2 \kappa_l b \zeta} \|x_{k,s} - x_{k,0}\|_2^2. \end{aligned}$$

Moreover, by Lemmas 4.1, 4.3 and 4.9, and the fact that $2XY = 2(\sqrt{z}X)(Y/\sqrt{z}) \leq zX^2 + Y^2/z$ for $\{X, Y\} \subset \mathbb{R}$ and $z \in \mathbb{R}_{>0}$, it follows that

$$\begin{aligned} \mathbb{E}_{k,s} [\|x_{k,s+1} - x_{k,0}\|_2^2] &= \mathbb{E}_{k,s} [\|x_{k,s+1} - x_{k,s} + x_{k,s} - x_{k,0}\|_2^2] \\ &= \mathbb{E}_{k,s} [\alpha^2 \|\bar{d}_{k,s}\|_2^2] + \|x_{k,s} - x_{k,0}\|_2^2 + 2\alpha d_{k,s}^T (x_{k,s} - x_{k,0}) \\ &\leq \mathbb{E}_{k,s} [\alpha^2 \|\bar{d}_{k,s}\|_2^2] + \|x_{k,s} - x_{k,0}\|_2^2 \\ &\quad + 2\alpha \left(\frac{1}{2z} \|d_{k,s}\|_2^2 + \frac{z}{2} \|x_{k,s} - x_{k,0}\|_2^2 \right) \\ &\leq \mathbb{E}_{k,s} \left[\frac{\alpha^2}{\bar{\tau}_{\min} \kappa_l} \Delta l(x_{k,s}, \bar{\tau}_{\min}, \bar{g}_{k,s}, \bar{d}_{k,s}) \right] + \|x_{k,s} - x_{k,0}\|_2^2 \\ &\quad + \frac{\alpha}{z \bar{\tau}_{\min} \kappa_l} \Delta l(x_{k,s}, \bar{\tau}_{\min}, g_{k,s}, d_{k,s}) + \alpha z \|x_{k,s} - x_{k,0}\|_2^2 \\ &\leq \frac{\alpha}{\bar{\tau}_{\min} \kappa_l} \left(\alpha + \frac{1}{z} \right) \Delta l(x_{k,s}, \bar{\tau}_{\min}, g_{k,s}, d_{k,s}) \\ &\quad + \left(\frac{\alpha^2 L^2}{\kappa_l b \zeta} + \alpha z + 1 \right) \|x_{k,s} - x_{k,0}\|_2^2. \end{aligned}$$

Taking total expectation conditioned on the event $E_{\tau_{\min}}$, for all $k \geq \bar{k}_\tau$ and $s \in [\bar{S}]$, combining the results above and by the definitions of λ_s , $R_{k,s}$ and Λ_s and the fact that $\bar{\tau}_{k,s+1} = \bar{\tau}_{k,s} = \bar{\tau}_{\min}$, it follows that

$$\begin{aligned} \mathbb{E}_{\tau_{\min}} [R_{k,s+1}] &= \mathbb{E}_{\tau_{\min}} [\phi(x_{k,s+1}, \bar{\tau}_{k,s+1}) + \lambda_{s+1} \|x_{k,s+1} - x_{k,0}\|_2^2] \\ &\leq \mathbb{E}_{\tau_{\min}} [R_{k,s}] - \Lambda_s \mathbb{E}_{\tau_{\min}} [\Delta l(x_{k,s}, \bar{\tau}_{\min}, g_{k,s}, d_{k,s})] \\ &\leq \mathbb{E}_{\tau_{\min}} [R_{k,s}] - \Lambda_{\min} \mathbb{E}_{\tau_{\min}} [\Delta l(x_{k,s}, \bar{\tau}_{\min}, g_{k,s}, d_{k,s})]. \end{aligned}$$

Summing over all inner iterations ($s \in [\bar{S}]$), we have

$$\begin{aligned} \sum_{s=0}^{S-1} \mathbb{E}_{\tau_{\min}} [\Delta l(x_{k,s}, \bar{\tau}_{\min}, g_{k,s}, d_{k,s})] &\leq \frac{\mathbb{E}_{\tau_{\min}} [R_{k,0} - R_{k,S}]}{\Lambda_{\min}} \\ &= \frac{\mathbb{E}_{\tau_{\min}} [\phi(x_{k,0}, \bar{\tau}_{\min}) - \phi(x_{k+1,0}, \bar{\tau}_{\min})]}{\Lambda_{\min}}. \end{aligned}$$

The equality follows from the fact that $\lambda_S = 0$ and $x_{k,S} = x_{k+1,0}$. Summing this inequality for $j \in \{\bar{k}_\tau, \bar{k}_\tau + 1, \dots, k\}$, we have

$$\sum_{j=\bar{k}_\tau}^k \sum_{s=0}^{S-1} \mathbb{E}_{\tau_{\min}} [\Delta l(x_{j,s}, \bar{\tau}_{\min}, g_{j,s}, d_{j,s})] \leq \frac{\mathbb{E}_{\tau_{\min}} [\phi(x_{\bar{k}_\tau,0}, \bar{\tau}_{\min})] - \phi_{\inf}}{\Lambda_{\min}},$$

for which the desired conclusion (4.11) follows. \square

As a consequence of Lemma 4.11, in Theorem 4.12 we present the main convergence result of this subsection, along with a specification of the controlled parameters (e.g., step size, inner iteration length, etc).

Theorem 4.12. *Suppose Assumption 4.8 holds. Let λ_s , Λ_s and Λ_{\min} be defined as in Lemma 4.11. Suppose $\alpha = \frac{\mu_0 b}{(\bar{\tau}_{\min} L + \Gamma) N^\gamma} \in (0, 1]$ with $\mu_0 \in (0, 1]$, $z = \frac{\bar{\tau}_{\min} L + \Gamma}{N^{\gamma/2}}$, $\gamma \in (0, 1]$,*

$b < N^\gamma$, and $S \leq \left\lfloor \frac{N^{3\gamma/2}}{\mu_0 \left(b + \frac{L^2}{(\bar{\tau}_{\min} L + \Gamma)^2 \kappa_l \zeta} \right)} \right\rfloor$. Then, for all $k \geq \bar{k}_\tau$ and $s \in [\bar{S}]$, there exist universal constants μ_0 and $\nu_0 \in (0, 1)$ such that $\Lambda_{\min} \geq \frac{\nu_0 b}{(\bar{\tau}_{\min} L + \Gamma) N^\gamma}$ and,

$$\begin{aligned} \mathbb{E}_{\tau_{\min}} \left[\frac{1}{(k - k_\tau + 1)S} \sum_{j=\bar{k}_\tau}^k \sum_{s=0}^{S-1} \Delta l(x_{j,s}, \bar{\tau}_{\min}, g_{j,s}, d_{j,s}) \right] \\ \leq \frac{(\bar{\tau}_{\min} L + \Gamma) N^\gamma (\mathbb{E}_{\tau_{\min}} [\phi(x_{\bar{k}_\tau, 0}, \bar{\tau}_{\min})] - \phi_{\inf})}{(k - k_\tau + 1) S \nu_0 b}. \end{aligned}$$

Proof. By the recursive definition of λ_s (4.10) and the fact that $\lambda_S = 0$, we have that

$$\lambda_0 = \frac{1}{2} (\bar{\tau}_{\min} L + \Gamma) \frac{\alpha^2 L^2}{\kappa_l b \zeta} \frac{((1+\rho)^S - 1)}{\rho}, \quad (4.12)$$

with

$$\begin{aligned} \rho &= \alpha z + \frac{\alpha^2 L^2}{\kappa_l b \zeta} \\ &= \frac{\mu_0 b}{N^{3\gamma/2}} + \frac{\mu_0^2 b L^2}{(\bar{\tau}_{\min} L + \Gamma)^2 N^{2\gamma} \kappa_l \zeta} \\ &\leq \mu_0 N^{-3\gamma/2} \left(b + \frac{b L^2}{(\bar{\tau}_{\min} L + \Gamma)^2 \kappa_l \zeta} \right), \end{aligned}$$

where $\mu_0 \in (0, 1]$ and $N \geq 1$. (Note, without loss of generality, we assume that the user defined constants are chosen such that $\rho \in (0, 1)$.) Plugging α and ρ into equation (4.12), it follows that

$$\begin{aligned} \lambda_0 &= \frac{1}{2} \frac{L^2 \mu_0^2 b}{\kappa_l (\bar{\tau}_{\min} L + \Gamma) N^{2\gamma} \zeta} \frac{(1+\rho)^S - 1}{\frac{\mu_0 b}{N^{3\gamma/2}} + \frac{\mu_0^2 b L^2}{(\bar{\tau}_{\min} L + \Gamma)^2 N^{2\gamma} \zeta}} \\ &\leq \frac{L^2 \mu_0 (e-1)}{2 \kappa_l (\bar{\tau}_{\min} L + \Gamma) \zeta} N^{-\gamma/2}, \end{aligned}$$

where the inequality is obtained by noticing that for $l > 0$, $(1 + \frac{1}{l})^l$ is an increasing function and $(1 + \frac{1}{l})^l \rightarrow e$ as $l \rightarrow \infty$. Hence, $(1 + \rho)^S \leq e$ by the definition of S . Now, with the

upper bound of λ_0 , the fact that λ_s is decreasing as s increases from 0 to S , and $\mu_0 \in (0, 1]$, $b < N$ and $N \geq 1$, it follows that Λ_{\min} can be lower bounded by

$$\begin{aligned}\Lambda_{\min} &= \min_{0 \leq s \leq S-1} \left\{ -\frac{(\bar{\tau}_{\min}L+\Gamma)\alpha^2}{2\bar{\tau}_{\min}\kappa_l} + \alpha - \lambda_{s+1} \frac{\alpha}{\bar{\tau}_{\min}\kappa_l} \left(\alpha + \frac{1}{z}\right) \right\} \\ &> -\frac{(\bar{\tau}_{\min}L+\Gamma)\alpha^2}{2\bar{\tau}_{\min}\kappa_l} + \alpha - \frac{\lambda_0\alpha}{\bar{\tau}_{\min}\kappa_l} \left(\alpha + \frac{1}{z}\right) \\ &\geq -\frac{\mu_0 b}{2\bar{\tau}_{\min}\kappa_l N^\gamma} \alpha + \alpha - \frac{L^2 \mu_0^2 (e-1) b}{2\kappa_l^2 (\bar{\tau}_{\min}L+\Gamma)^2 \bar{\tau}_{\min}\zeta} N^{-3\gamma/2} \alpha - \frac{L^2 \mu_0 (e-1)}{2\kappa_l^2 (\bar{\tau}_{\min}L+\Gamma)^2 \bar{\tau}_{\min}\zeta} \alpha \\ &\geq \alpha \left[1 - \frac{\mu_0 b}{\bar{\tau}_{\min}\kappa_l} - \frac{L^2 \mu_0 (e-1)}{2\kappa_l^2 (\bar{\tau}_{\min}L+\Gamma)^2 \bar{\tau}_{\min}\zeta} \right].\end{aligned}$$

Let $\nu_0 = 1 - \frac{\mu_0 b}{2\bar{\tau}_{\min}\kappa_l} - \frac{L^2 \mu_0 (e-1)}{\kappa_l^2 (\bar{\tau}_{\min}L+\Gamma)^2 \bar{\tau}_{\min}\zeta}$. By choosing μ_0 (independent of N) such that $\nu_0 > 0$, it follows that $\Lambda_{\min} \geq \frac{b\nu_0}{(\bar{\tau}_{\min}L+\Gamma)N^\gamma}$. Combining this lower bound with Lemma 4.11 yields the desired result. \square

Finally, we conclude this section by presenting a corollary to Theorem 4.12; this result shows that SVR-SQP generates a sequence of iterates whose first order stationary measure (corresponding to (1.1)) converges to zero.

Corollary 4.13. *Under the conditions of Theorem 4.12, there exists universal constants μ_0, ν_0 such that $\Lambda_{\min} \geq \frac{\nu_0 b}{(\bar{\tau}_{\min}L+\Gamma)N^\gamma}$ and*

$$\begin{aligned}\mathbb{E}_{\tau_{\min}} \left[\frac{1}{(k-k_\tau+1)S} \sum_{j=\bar{k}_\tau}^k \sum_{s=0}^{S-1} \frac{\|g_{j,s} + J_{j,s} y_{j,s}\|_2^2}{\kappa_H^2} + \|c_{j,s}\|_2 \right] \\ \leq \frac{(\bar{\tau}_{\min}L+\Gamma)N^\gamma (\mathbb{E}_{\tau_{\min}}[\phi(x_{\bar{k}_\tau,0}, \bar{\tau}_{\min})] - \phi_{\inf})}{(k-k_\tau+1)S\nu_0 b}.\end{aligned}$$

Moreover, if for some $(k, s) \in \mathbb{N} \times [\bar{S}]$, $\|x_{k,s} - x^*\|_2 \leq \delta_x$, $\|x_{k,0} - x^*\|_2 \leq \delta_{x,0}$ and $\|c_{k,s}\|_2 \leq \delta_c$, for $(\delta_x, \delta_{x,0}, \delta_c) \in \mathbb{R}_{>0} \times \mathbb{R}_{>0} \times \mathbb{R}_{>0}$, and some stationary point $(x^*, y^*) \in \mathbb{R}^n \times \mathbb{R}^m$ of (1.1), then, there exists $\kappa_{g^*} \in \mathbb{R}_{>0}$, such that

$$\mathbb{E}_{k,s} \left[\left\| \begin{bmatrix} \bar{g}_{k,s} - \nabla f(x_*) \\ c_{k,s} \end{bmatrix} \right\|_2 \right] \leq \delta_y \quad \text{and} \quad \mathbb{E}_{k,s} [\|\bar{y}_{k,s} - y_{k,s}\|_2] \leq \kappa_d \delta_y + 2\kappa_d^2 \Gamma \delta_x \|g^*\|$$

where $\delta_y = \delta_c + L(\delta_{x,0} + \delta_x)/\sqrt{b} + L\delta_x \in \mathbb{R}_{>0}$ and $\kappa_d \in \mathbb{R}_{>0}$ is an upper bound for $\left\| \begin{bmatrix} H_{k,s} & J_{k,s}^T \\ J_{k,s} & 0 \end{bmatrix} \right\|^{-1}$.

Proof. The first part follows by Lemma 4.1 and Theorem 4.12, and the fact that by the deterministic variant of (3.2) for all $(k, s) \in \mathbb{N} \times [\bar{S}]$

$$\|g_{k,s} + J_{k,s} y_{k,s}\|_2 \leq \|H_{k,s} d_{k,s}\|_2 \leq \|H_{k,s}\|_2 \|d_{k,s}\|_2 \leq \kappa_H \|d_{k,s}\|_2.$$

The second part follows by Assumptions 2.1 and 2.3, the definitions of $(\delta_x, \delta_{x,0}, \delta_c)$, and the triangle inequality that $\|x_{k,s} - x_{k,0}\|_2 \leq \|x_{k,0} - x^*\|_2 + \|x_{k,s} - x^*\|_2 \leq \delta_{x,0} + \delta_x$

$$\begin{aligned} \mathbb{E}_{k,s} \left[\left\| \begin{bmatrix} \bar{g}_{k,s} - \nabla f(x_*) \\ c_{k,s} \end{bmatrix} \right\|_2 \right] &\leq \|c_{k,s}\|_2 + \mathbb{E}_{k,s} [\|\bar{g}_{k,s} - \nabla f(x_{k,s})\|_2] + \|\nabla f(x_{k,s}) - \nabla f(x_*)\|_2 \\ &\leq \|c_{k,s}\|_2 + \sqrt{\mathbb{E}_{k,s} [\|\bar{g}_{k,s} - \nabla f(x_{k,s})\|_2^2]} + \|\nabla f(x_{k,s}) - \nabla f(x_*)\|_2 \\ &\leq \delta_c + L(\delta_{x,0} + \delta_x)/\sqrt{b} + L\delta_x = \delta_y. \end{aligned}$$

Let $g_* := \nabla f(x_*)$ and $c_* := c(x_*) = 0$, for $x_{k,s}$ sufficiently close to x_* there exists $\kappa_{g_*} \in \mathbb{R}_{>0}$ such that

$$\begin{aligned} &\mathbb{E}_{k,s} [\|\bar{y}_{k,s} - y_*\|_2] \\ &\leq \mathbb{E}_{k,s} \left[\left\| \begin{bmatrix} H_{k,s} & J_{k,s}^T \\ J_{k,s} & 0 \end{bmatrix}^{-1} \begin{bmatrix} \bar{g}_{k,s} \\ c_{k,s} \end{bmatrix} - \begin{bmatrix} H_{k,s} & J_*^T \\ J_* & 0 \end{bmatrix}^{-1} \begin{bmatrix} g_* \\ c_* \end{bmatrix} \right\|_2 \right] \\ &= \mathbb{E}_{k,s} \left[\left\| \begin{bmatrix} H_{k,s} & J_{k,s}^T \\ J_{k,s} & 0 \end{bmatrix}^{-1} \begin{bmatrix} \bar{g}_{k,s} - g_* \\ c_{k,s} \end{bmatrix} + \left(\begin{bmatrix} H_{k,s} & J_{k,s}^T \\ J_{k,s} & 0 \end{bmatrix}^{-1} - \begin{bmatrix} H_{k,s} & J_*^T \\ J_* & 0 \end{bmatrix}^{-1} \right) \begin{bmatrix} g_* \\ 0 \end{bmatrix} \right\|_2 \right] \\ &\leq \left\| \begin{bmatrix} H_{k,s} & J_{k,s}^T \\ J_{k,s} & 0 \end{bmatrix}^{-1} \right\|_2 \mathbb{E}_{k,s} \left[\left\| \begin{bmatrix} \bar{g}_{k,s} - g_* \\ c_{k,s} \end{bmatrix} \right\|_2 \right] \\ &\quad + 2 \left\| \begin{bmatrix} H_{k,s} & J_{k,s}^T \\ J_{k,s} & 0 \end{bmatrix}^{-1} \right\|_2 \left\| \begin{bmatrix} 0 & (J_* - J_{k,s})^T \\ J_* - J_{k,s} & 0 \end{bmatrix}^{-1} \right\|_2 \left\| \begin{bmatrix} H_{k,s} & J_{k,s}^T \\ J_{k,s} & 0 \end{bmatrix}^{-1} \right\|_2 \left\| \begin{bmatrix} g_* \\ 0 \end{bmatrix} \right\|_2 \\ &\leq \kappa_d \delta_y + 2\kappa_d^2 \Gamma \delta_x \|g_*\| \end{aligned}$$

where the last inequality is satisfied since $(A + \Delta)^{-1} = A^{-1} - A^{-1}\Delta A^{-1} + \mathcal{O}(\|\Delta\|^2)$, and we assume that $\|\Delta\|_2 = \|J_* - J_{k,s}\|_2 \leq \Gamma \delta_x$ is small enough such that $\mathcal{O}(\|\Delta\|^2) \leq \|A^{-1}\Delta A^{-1}\|_2$. This completes the proof. \square

Corollary 4.13 characterizes the behavior of optimality measure $\|g_{k,s} + J_{k,s}y_{k,s}\|_2^2$ and feasibility measure $\|c_{k,s}\|_2$ for all $k \geq \bar{k}_\tau$ and $s \in [\bar{S}]$. The result of Corollary 4.13 reveals that, under the assumption that merit parameter $\bar{\tau}_k$ has stabilized at a sufficiently small value, both measures converge to zero in expectation, which justifies our summary in Table 1. It is important to note the difference in nature of the results of Corollary 4.13 and the analogues proven in the unconstrained setting for the SVRG method [31]. The first result in Corollary 4.13 is with respect to the expectation of the averaged optimality/feasibility measure across iterations, whereas in [31] the results are with respect to the minimal optimality measure ([31] considers the unconstrained setting, and so the optimality measure is the norm of the gradient) over the iterations. One can easily derive similar convergence results for SVR-SQP. Moreover, if the output of Algorithm 1 is uniformly chosen from (k, s) ,

where $k \geq \bar{k}_\tau$, one can derive a bound for $\mathbb{E}_{\tau_{\min}} \left[\frac{\|g_{k,s} + J_{k,s} y_{k,s}\|_2^2}{\kappa_H^2} + \|c_{k,s}\|_2 \right]$. Finally, we provide an upper bound for the error of the Lagrange multiplier estimates which is dependent on a feasibility measure and the distances from $x_{k,s}$ and $x_{k,0}$ to the optimal solution. As a result, if the primal iterates converge to a feasible point and in expectation the SVRG gradient approximation converges to the true gradient of the objective function at the optimal solution, then the Lagrange multipliers also converge.

We conclude this section with a remark about the iteration complexity of SVR-SQP. In the unconstrained setting, by employing variance reduced gradients SVRG is able to improve upon the iteration complexity of the stochastic gradient (SG) method in term of the dependence on ϵ (the termination tolerance). Specially, in the nonconvex setting the iteration complexity for SVRG is $\mathcal{O}(n + \frac{n^{2/3}}{\epsilon})$ [31] whereas the iteration complexity for SG is $\mathcal{O}(\frac{1}{\epsilon^2})$. In the constrained setting, deriving such results is significantly more difficult due to the fact that one needs to consider two measures of optimality (feasibility and stationarity) and the fact that the merit function (measure of progress) is changing over the course of the optimization (the merit parameter changes adaptively). Under the assumption that the merit parameter has stabilized at a sufficient small positive value, as a result of Corollary 4.13 one can show that the number of iteration to achieve ϵ -optimality (where the optimality measure is a combination of stationarity and feasibility, i.e., $\max\{\|g_{k,s} + J_{k,s} y_{k,s}\|_2^2, \|c_{k,s}\|_2\} \leq \epsilon$) is $\mathcal{O}(n + \frac{n^{2/3}}{\epsilon^2})$. To contrast this result, the algorithm in [4] (the analogue of the SG method in the equality constrained setting), after the merit parameter has stabilized, requires $\mathcal{O}(\frac{1}{\epsilon^4})$. As a result, it is clear that variance reduction does have an effect, albeit the limited setting under which the result has been derived and the fact that this result does not say anything about the iterations before the merit parameter stabilizes. To the best of our knowledge, the only work that has analyzed the iteration complexity with regards to the whole sequence of merit parameter is [11]. One can certainly extend that analysis for our algorithm. We defer such analysis to a different study since it would require extending the paper significantly.

4.4 Adaptive step size

In this subsection, we present convergence results for Algorithm 1 with the adaptive step size strategy (**Option II** in the Algorithm 1) under Assumption 4.8. The analysis in this section is significantly more involved than the analysis in Section 4.3 primarily due to the stochastic nature of the step size rule ((3.11)–(3.13)). Paralleling the analysis of the constant step size strategy, we first provide an upper bound for the difference in merit function after a step.

Lemma 4.14. *Suppose that Assumption 4.8 holds. For all $k \geq \bar{k}_\tau$ and $s \in [\bar{S}]$, it follows*

that

$$\begin{aligned}
& \phi(x_{k,s+1}, \bar{\tau}_{k,s}) - \phi(x_{k,s}, \bar{\tau}_{k,s}) \\
& \leq -\bar{\alpha}_{k,s} \Delta l(x_{k,s}, \bar{\tau}_{\min}, g_{k,s}, d_{k,s}) + \frac{1}{2} \bar{\alpha}_{k,s} \beta \Delta l(x_{k,s}, \bar{\tau}_{\min}, \bar{g}_{k,s}, \bar{d}_{k,s}) \\
& \quad + \bar{\alpha}_{k,s} \bar{\tau}_{\min} g_{k,s}^T (\bar{d}_{k,s} - d_{k,s})
\end{aligned}$$

Proof. For $k \geq \bar{k}_\tau$ and $s \in [\bar{S}]$, we consider three cases depending on how the step size is set in Algorithm 1 (Option II).

Case 1: Suppose in Algorithm 1 (Option II) that $\bar{\alpha}_{k,s} < 1$, meaning that $\bar{\alpha}_{k,s} \leftarrow \bar{\alpha}_{k,s} \leq \frac{\beta \Delta l(x_{k,s}, \bar{\tau}_{\min}, \bar{g}_{k,s}, \bar{d}_{k,s})}{(\bar{\tau}_{\min} L + \Gamma) \|\bar{d}_{k,s}\|_2^2}$. It then follows from (3.9) and Lemma 4.1 that

$$\begin{aligned}
& \phi(x_{k,s} + \bar{\alpha}_{k,s} \bar{d}_{k,s}, \bar{\tau}_{k,s}) - \phi(x_{k,s}, \bar{\tau}_{k,s}) \\
& \leq \bar{\alpha}_{k,s} (\bar{\tau}_{\min} g_{k,s}^T \bar{d}_{k,s} - \|c_{k,s}\|_1) + \frac{1}{2} (\bar{\tau}_{\min} L + \Gamma) \bar{\alpha}_{k,s}^2 \|\bar{d}_{k,s}\|_2^2 \\
& = \bar{\alpha}_{k,s} (\bar{\tau}_{\min} g_{k,s}^T d_{k,s} - \|c_{k,s}\|_1) + \frac{1}{2} (\bar{\tau}_{\min} L + \Gamma) \bar{\alpha}_{k,s}^2 \|\bar{d}_{k,s}\|_2^2 + \bar{\alpha}_{k,s} \bar{\tau}_{\min} g_{k,s}^T (\bar{d}_{k,s} - d_{k,s}) \\
& = -\bar{\alpha}_{k,s} \Delta l(x_{k,s}, \bar{\tau}_{\min}, g_{k,s}, d_{k,s}) + \frac{1}{2} (\bar{\tau}_{\min} L + \Gamma) \bar{\alpha}_{k,s}^2 \|\bar{d}_{k,s}\|_2^2 \\
& \quad + \bar{\alpha}_{k,s} \bar{\tau}_{\min} g_{k,s}^T (\bar{d}_{k,s} - d_{k,s}) \\
& \leq -\bar{\alpha}_{k,s} \Delta l(x_{k,s}, \bar{\tau}_{\min}, g_{k,s}, d_{k,s}) + \frac{1}{2} \bar{\alpha}_{k,s} \beta \Delta l(x_{k,s}, \bar{\tau}_{\min}, \bar{g}_{k,s}, \bar{d}_{k,s}) \\
& \quad + \bar{\alpha}_{k,s} \bar{\tau}_{\min} g_{k,s}^T (\bar{d}_{k,s} - d_{k,s})
\end{aligned}$$

Case 2: Suppose in Algorithm 1 (Option II) that $\bar{\alpha}_{k,s} \leq 1 \leq \bar{\alpha}_{k,s}$, meaning that $\bar{\alpha}_{k,s} \leftarrow 1 \leq \frac{\beta \Delta l(x_{k,s}, \bar{\tau}_{\min}, g_{k,s}, d_{k,s})}{(\bar{\tau}_{\min} L + \Gamma) \|d_{k,s}\|_2^2}$. Similar to Case 1, it follows that

$$\begin{aligned}
& \phi(x_{k,s} + \bar{\alpha}_{k,s} \bar{d}_{k,s}, \bar{\tau}_{k,s}) - \phi(x_{k,s}, \bar{\tau}_{k,s}) \\
& \leq -\bar{\alpha}_{k,s} \Delta l(x_{k,s}, \bar{\tau}_{\min}, g_{k,s}, d_{k,s}) + \frac{1}{2} (\bar{\tau}_{\min} L + \Gamma) \bar{\alpha}_{k,s}^2 \|\bar{d}_{k,s}\|_2^2 \\
& \quad + \bar{\alpha}_{k,s} \bar{\tau}_{\min} g_{k,s}^T (\bar{d}_{k,s} - d_{k,s}) \\
& \leq -\bar{\alpha}_{k,s} \Delta l(x_{k,s}, \bar{\tau}_{\min}, g_{k,s}, d_{k,s}) + \frac{1}{2} \bar{\alpha}_{k,s} \beta \Delta l(x_{k,s}, \bar{\tau}_{\min}, \bar{g}_{k,s}, \bar{d}_{k,s}) \\
& \quad + \bar{\alpha}_{k,s} \bar{\tau}_{\min} g_{k,s}^T (\bar{d}_{k,s} - d_{k,s})
\end{aligned}$$

Case 3: Suppose in Algorithm 1 (Option II) that $\bar{\alpha}_{k,s} > 1$, meaning that $\bar{\alpha}_{k,s} \leftarrow \bar{\alpha}_{k,s} \leq$

$\frac{\beta \Delta l(x_{k,s}, \bar{\tau}_{\min}, g_{k,s}, d_{k,s}) - 4 \|c_{k,s}\|_1}{(\bar{\tau}_{k,s} L + \Gamma) \|d_{k,s}\|_2^2}$. It follows from (3.9) and Lemma 4.1 that

$$\begin{aligned}
& \phi(x_{k,s} + \bar{\alpha}_{k,s} \bar{d}_{k,s}, \bar{\tau}_{k,s}) - \phi(x_{k,s}, \bar{\tau}_{k,s}) \\
& \leq \bar{\alpha}_{k,s} \bar{\tau}_{\min} g_{k,s}^T \bar{d}_{k,s} + (\bar{\alpha}_{k,s} - 1) \|c_{k,s}\|_1 - \|c_{k,s}\|_1 + \frac{1}{2} (\bar{\tau}_{\min} L + \Gamma) \bar{\alpha}_{k,s}^2 \|\bar{d}_{k,s}\|_2^2 \\
& = \bar{\alpha}_{k,s} (\bar{\tau}_{\min} g_{k,s}^T \bar{d}_{k,s} - \|c_{k,s}\|_1) + 2(\bar{\alpha}_{k,s} - 1) \|c_{k,s}\|_1 + \frac{1}{2} (\bar{\tau}_{\min} L + \Gamma) \bar{\alpha}_{k,s}^2 \|\bar{d}_{k,s}\|_2^2 \\
& \leq \bar{\alpha}_{k,s} (\bar{\tau}_{\min} g_{k,s}^T d_{k,s} - \|c_{k,s}\|_1) + 2\bar{\alpha}_{k,s} \|c_{k,s}\|_1 + \frac{1}{2} (\bar{\tau}_{\min} L + \Gamma) \bar{\alpha}_{k,s}^2 \|\bar{d}_{k,s}\|_2^2 \\
& \quad + \bar{\alpha}_{k,s} \bar{\tau}_{\min} g_{k,s}^T (\bar{d}_{k,s} - d_{k,s}) \\
& \leq -\bar{\alpha}_{k,s} \Delta l(x_{k,s}, \bar{\tau}_{\min}, g_{k,s}, d_{k,s}) + \frac{1}{2} \bar{\alpha}_{k,s} \beta \Delta l(x_{k,s}, \bar{\tau}_{\min}, \bar{g}_{k,s}, \bar{d}_{k,s}) \\
& \quad + \bar{\alpha}_{k,s} \bar{\tau}_{\min} g_{k,s}^T (\bar{d}_{k,s} - d_{k,s})
\end{aligned}$$

The result follows by combining the three cases. \square

While the upper bounds for the difference in merit function after a step for the two step size strategies are very similar (Lemmas 4.10 and 4.14, respectively), a key difference pertains to the fact step sizes computed by the adaptive algorithm (**Option II**) are stochastic, and as such the last term in the bound in Lemma 4.14 is nonzero in expectation. Moreover, due to the adaptive and stochastic nature of the step size strategy, an additional user-defined parameter $\alpha_u \in \mathbb{R}_{>0}$ is required. Before we proceed, we make the following remark with regards to the selection of α_u .

Remark 4.15. Under Assumption 4.8 and by Lemma 4.1(b), it follows that for all $k \geq \bar{k}_\tau$ and $s \in [\bar{S}]$

$$\frac{\Delta l(x_{k,s}, \bar{\tau}_{\min}, \bar{g}_{k,s}, \bar{d}_{k,s})}{(\bar{\tau}_{\min} L + \Gamma) \|\bar{d}_{k,s}\|_2^2} \geq \frac{\kappa_l \bar{\tau}_{\min}}{\bar{\tau}_{\min} L + \Gamma} \in \mathbb{R}_{>0}.$$

When the user-defined parameters $\alpha_u \in \mathbb{R}_{>0}$ and $\beta \in (0, 1]$ are chosen such that $\alpha_u \beta \leq \frac{\beta \kappa_l \bar{\tau}_{\min}}{\bar{\tau}_{\min} L + \Gamma} \in (0, 1]$, it follows that **Option II** in Algorithm 1 always selects a constant step size $\alpha_u \beta \in (0, 1]$, whose analysis has already been discussed in Section 4.3. Therefore, under Assumption 4.8, for the rest of this subsection we only consider the case where $\alpha_u > \frac{\kappa_l \bar{\tau}_{\min}}{\bar{\tau}_{\min} L + \Gamma}$ and $\beta \in \mathbb{R}_{>0}$ is chosen such that $\frac{\beta \kappa_l \bar{\tau}_{\min}}{\bar{\tau}_{\min} L + \Gamma} \in (0, 1]$.

Next, we provide upper and lower bounds for the step sizes $\bar{\alpha}_{k,s} \in \mathbb{R}_{>0}$ chosen by SVR-SQP.

Lemma 4.16. Suppose that Assumption 4.8 holds. Let $\bar{\alpha}_{k,s}$ be defined as in (3.11)–(3.13), and consider $\alpha_l := \frac{\kappa_l \bar{\tau}_{\min}}{\bar{\tau}_{\min} L + \Gamma} \in \mathbb{R}_{>0}$ with $\alpha_l < \alpha_u$. Suppose $\beta \in (0, 1]$ is chosen such that $\frac{\beta \kappa_l \bar{\tau}_{\min}}{\bar{\tau}_{\min} L + \Gamma} \in (0, 1]$, then for all $k \geq \bar{k}_\tau$ and $s \in [\bar{S}]$, it follows that $\bar{\alpha}_{k,s} \in [\alpha_l \beta, \alpha_u \beta]$.

Proof. By Lemma 4.1(b), it follows that $\frac{\Delta l(x_{k,s}, \bar{\tau}_{\min}, \bar{g}_{k,s}, \bar{d}_{k,s})}{(\bar{\tau}_{\min} L + \Gamma) \|\bar{d}_{k,s}\|_2^2} \geq \frac{\kappa_l \bar{\tau}_{\min}}{\bar{\tau}_{\min} L + \Gamma}$ for all $k \geq \bar{k}_\tau$ and $s \in [\bar{S}]$. By (3.11)–(3.13), the desired conclusion follows by considering the following three cases.

Case 1: Suppose that $\tilde{\alpha}_{k,s} = \min \left\{ \frac{\Delta l(x_{k,s}, \bar{\tau}_{k,s}, \bar{g}_{k,s}, \bar{d}_{k,s})}{(\bar{\tau}_{k,s} L_{k,s} + \Gamma_{k,s}) \|d_{k,s}\|_2^2}, \alpha_u \right\} \beta < 1$, in which case the algorithm sets $\bar{\alpha}_{k,s} = \tilde{\alpha}_{k,s}$. It follows that

$$\alpha_l \beta = \frac{\kappa_l \bar{\tau}_{\min}}{\bar{\tau}_{\min} L + \Gamma} \beta \leq \min \left\{ \frac{\kappa_l \bar{\tau}_{\min}}{\bar{\tau}_{\min} L + \Gamma}, \alpha_u \right\} \beta \leq \bar{\alpha}_{k,s} \leq \alpha_u \beta.$$

Case 2: Suppose that $\tilde{\alpha}_{k,s} = \bar{\alpha}_{k,s} - \frac{4\|c_{k,s}\|_1}{(\bar{\tau}_{k,s} L_{k,s} + \Gamma_{k,s}) \|d_{k,s}\|_2^2} \leq 1 \leq \bar{\alpha}_{k,s}$, in which case the algorithm sets $\bar{\alpha}_{k,s} = 1$. It follows that

$$\alpha_l \beta \leq 1 = \bar{\alpha}_{k,s} \leq \tilde{\alpha}_{k,s} \leq \alpha_u \beta.$$

Case 3: Suppose that $\tilde{\alpha}_{k,s} > 1$, in which case the algorithm sets $\bar{\alpha}_{k,s} = \tilde{\alpha}_{k,s}$. It follows that

$$\alpha_l \beta \leq 1 < \bar{\alpha}_{k,s} = \tilde{\alpha}_{k,s} - \frac{4\|c_{k,s}\|_1}{(\bar{\tau}_{k,s} L_{k,s} + \Gamma_{k,s}) \|d_{k,s}\|_2^2} \leq \tilde{\alpha}_{k,s} \leq \alpha_u \beta$$

□

As mentioned above, due to the adaptive (and stochastic) nature of the step size strategy, the third term on the right-hand-side of the bound in Lemma 4.14 is nonzero in expectation. We provide an upper bound for this quantity in the next lemma.

Lemma 4.17. *Suppose that Assumption 4.8 holds. For all $k \geq \bar{k}_\tau$ and $s \in [\bar{S}]$, it follows that*

$$\begin{aligned} & \mathbb{E}_{k,s} \left[\bar{\alpha}_{k,s} \bar{\tau}_{k,s} g_{k,s}^T (\bar{d}_{k,s} - d_{k,s}) \right] \\ & \leq \frac{\alpha_u \kappa_H \kappa_d \beta^2}{2\kappa_l} \Delta l(x_{k,s}, \bar{\tau}_{\min}, g_{k,s}, d_{k,s}) + \frac{\alpha_u \bar{\tau}_{\min} \kappa_H \kappa_d L^2}{2b} \|x_{k,s} - x_{k,0}\|_2^2. \end{aligned}$$

Proof. For all $k \geq \bar{k}_\tau$ and $s \in [\bar{S}]$, by Lemmas 4.1, 4.3 and 4.16, Cauchy–Schwarz inequality and the fact of $2XY = 2(\sqrt{\beta}X)(Y/\sqrt{\beta}) \leq \beta X^2 + Y^2/\beta$ for any $\{X, Y\} \subset \mathbb{R}$, it follows that

$$\begin{aligned} & \mathbb{E}_{k,s} \left[\bar{\alpha}_{k,s} \bar{\tau}_{k,s} g_{k,s}^T (\bar{d}_{k,s} - d_{k,s}) \right] \\ & = \mathbb{E}_{k,s} \left[\bar{\alpha}_{k,s} \bar{\tau}_{\min} (g_{k,s} + J_{k,s}^T y_{k,s})^T (d_{k,s} - \bar{d}_{k,s}) \right] \\ & = \mathbb{E}_{k,s} \left[\bar{\alpha}_{k,s} \bar{\tau}_{\min} (-H_{k,s} d_{k,s})^T (d_{k,s} - \bar{d}_{k,s}) \right] \\ & \leq \alpha_u \beta \bar{\tau}_{\min} \kappa_H \|d_{k,s}\|_2 \mathbb{E}_{k,s} [\|\bar{d}_{k,s} - d_{k,s}\|_2] \\ & \leq \alpha_u \beta \bar{\tau}_{\min} \kappa_H \kappa_d \|d_{k,s}\|_2 \sqrt{M_{k,s}} \\ & \leq \alpha_u \beta \bar{\tau}_{\min} \kappa_H \kappa_d \left(\frac{\|d_{k,s}\|_2^2 \beta}{2} + \frac{M_{k,s}}{2\beta} \right) \\ & \leq \frac{\alpha_u \kappa_H \kappa_d \beta^2}{2\kappa_l} \Delta l(x_{k,s}, \bar{\tau}_{\min}, g_{k,s}, d_{k,s}) + \frac{\alpha_u \bar{\tau}_{\min} \kappa_H \kappa_d L^2}{2b} \|x_{k,s} - x_{k,0}\|_2^2. \end{aligned}$$

□

Lemma 4.18 and Theorem 4.19 (below) are the analogues of Lemma 4.11 and Theorem 4.12, respectively, for the adaptive step size case.

Lemma 4.18. *Suppose that Assumption 4.8 holds. Let α_l and α_u be defined as Lemma 4.16, and $\beta \in (0, 1]$ chosen such that $\frac{\beta\kappa_l\bar{\tau}_{\min}}{(\bar{\tau}_{\min}L+\Gamma)} \in (0, 1]$. In addition, let $\lambda_S = 0$, and*

$$\begin{aligned}\lambda_s &= \lambda_{s+1}(1+z)\left(1 + \frac{a_u^2\beta^2L^2}{\kappa_lzb\zeta}\right) + \frac{\alpha_u\bar{\tau}_{\min}L^2}{2b}(\kappa_H\kappa_d + \frac{\beta^2}{\zeta}) \\ \Lambda_s &= \alpha_l\beta - \frac{1}{2}\alpha_u\beta^2 - \frac{\alpha_u\kappa_H\kappa_d\beta^2}{2\kappa_l} - \lambda_{s+1}(1+z)\frac{\alpha_u^2\beta^2}{\bar{\tau}_{\min}\kappa_lz}\end{aligned}\tag{4.13}$$

for $s \in [\bar{S}]$, where $\beta, z \in \mathbb{R}_{>0}$, $\lambda_s \in \mathbb{R}_{>0}$ are chosen such that $\Lambda_s \in \mathbb{R}_{>0}$, and $\Lambda_{\min} = \min_{s \in [\bar{S}]} \Lambda_s$. Then, for all $k \geq \bar{k}_\tau$ and $s \in [\bar{S}]$, the sequence of iterates $\{x_{k,s}\}$ generated by Algorithm 1 (**Option II**) satisfy

$$\mathbb{E}_{\tau_{\min}} \left[\frac{1}{(k-\bar{k}_\tau+1)S} \sum_{j=\bar{k}_\tau}^k \sum_{s=0}^{S-1} \Delta l(x_{j,s}, \bar{\tau}_{\min}, g_{j,s}, d_{j,s}) \right] \leq \frac{\mathbb{E}_{\tau, small}[\phi(x_{\bar{k}_\tau, 0}, \bar{\tau}_{\min})] - \phi_{\inf}}{(k-\bar{k}_\tau+1)S\Lambda_{\min}}.\tag{4.14}$$

Proof. Consider arbitrary $k \geq \bar{k}_\tau$ and $s \in [\bar{S}]$. By Lemmas 4.2, 4.9, 4.17, we have

$$\begin{aligned}& \mathbb{E}_{k,s}[\phi(x_{k,s+1}, \bar{\tau}_{k,s})] \\ & \leq \mathbb{E}_{k,s}[\phi(x_{k,s}, \bar{\tau}_{k,s})] - \alpha_l\beta\Delta l(x_{k,s}, \bar{\tau}_{\min}, g_{k,s}, d_{k,s}) \\ & \quad + \frac{1}{2}\alpha_u\beta^2\mathbb{E}_{k,s}[\Delta l(x_{k,s}, \bar{\tau}_{\min}, \bar{g}_{k,s}, \bar{d}_{k,s})] + \mathbb{E}_k[\bar{\alpha}_{k,s}\bar{\tau}_{k,s}g_{k,s}^T(\bar{d}_{k,s} - d_{k,s})] \\ & \leq \mathbb{E}_{k,s}[\phi(x_{k,s}, \bar{\tau}_{k,s})] - \beta\left(\alpha_l - \frac{1}{2}\alpha_u\beta - \frac{\alpha_u\kappa_H\kappa_d\beta}{2\kappa_l}\right)\Delta l(x_{k,s}, \bar{\tau}_{\min}, g_{k,s}, d_{k,s}) \\ & \quad + \frac{\alpha_u\bar{\tau}_{\min}L^2}{2b}\left(\kappa_H\kappa_d + \frac{\beta^2}{\zeta}\right)\|x_{k,s} - x_{k,0}\|^2.\end{aligned}$$

Moreover, similar to the proof of lemma 4.11, we have

$$\begin{aligned}& \mathbb{E}_{k,s}[\|x_{k,s+1} - x_{k,0}\|_2^2] \\ & = \mathbb{E}_{k,s}[\|x_{k,s+1} - x_{k,s} + x_{k,s} - x_{k,0}\|_2^2] \\ & = \mathbb{E}_{k,s}[\|\bar{\alpha}_{k,s}\bar{d}_{k,s}\|_2^2 + \|x_{k,s} - x_{k,0}\|_2^2 + 2\bar{\alpha}_{k,s}\bar{d}_{k,s}^T(x_{k,s} - x_{k,0})] \\ & \leq \mathbb{E}_{k,s}[\|(1 + \frac{1}{z})\bar{\alpha}_{k,s}\bar{d}_{k,s}\|_2^2 + (1+z)\|x_{k,s} - x_{k,0}\|_2^2] \\ & \leq (1+z)\frac{a_u^2\beta^2}{\bar{\tau}_{\min}\kappa_lz}\mathbb{E}_{k,s}[\Delta l(x_{k,s}, \bar{\tau}_{\min}, \bar{g}_{k,s}, \bar{d}_{k,s})] + (1+z)\|x_{k,s} - x_{k,0}\|_2^2 \\ & \leq (1+z)\frac{a_u^2\beta^2}{\bar{\tau}_{\min}\kappa_lz}\Delta l(x_{k,s}, \bar{\tau}_{\min}, g_{k,s}, d_{k,s}) + \left[1+z+(1+z)\frac{a_u^2\beta^2L^2}{\kappa_lzb\zeta}\right]\|x_{k,s} - x_{k,0}\|_2^2.\end{aligned}$$

Taking total expectation conditioned on $E_{\tau_{\min}}$, for all $k \geq \bar{k}_\tau$ and $s \in [\bar{S}]$, combining

the results above and the definitions of λ_s, Λ_s , it follows that

$$\begin{aligned}
\mathbb{E}_{\tau_{\min}}[R_{k,s+1}] &= \mathbb{E}_{\tau_{\min}}[\phi(x_{k,s+1}, \bar{\tau}_{k,s}) + \lambda_{s+1}\|x_{k,s+1} - x_{k,0}\|^2] \\
&\leq \mathbb{E}_{\tau_{\min}}[\phi(x_{k,s}, \bar{\tau}_{k,s})] + \left[\lambda_{s+1}(1+z)\left(1 + \frac{a_u^2\beta^2L^2}{\kappa_l z b \zeta}\right) \right. \\
&\quad \left. + \frac{\alpha_u \bar{\tau}_{\min} L^2}{2b}(\kappa_H \kappa_d + \frac{\beta^2}{\zeta}) \right] \mathbb{E}_{\tau_{\min}}[\|x_{k,s} - x_{k,0}\|^2] \\
&\quad - \left(a_l \beta - \frac{1}{2} a_u \beta^2 - \frac{\alpha_u \kappa_H \kappa_d \beta^2}{2\kappa_l} - \lambda_{s+1}(1+z) \frac{a_u^2 \beta^2}{\bar{\tau}_{\min} \kappa_l z} \right) \\
&\quad \mathbb{E}_{\tau_{\min}}[\Delta l(x_{k,s}, \bar{\tau}_{\min}, g_{k,s}, d_{k,s})] \\
&\leq \mathbb{E}_{\tau_{\min}}[R_{k,s}] - \Lambda_s \mathbb{E}_{\tau_{\min}}[\Delta l(x_{k,s}, \bar{\tau}_{\min}, g_{k,s}, d_{k,s})] \\
&\leq \mathbb{E}_{\tau_{\min}}[R_{k,s}] - \Lambda_{\min} \mathbb{E}_{\tau_{\min}}[\Delta l(x_{k,s}, \bar{\tau}_{\min}, g_{k,s}, d_{k,s})].
\end{aligned}$$

Summing over all inner iterations ($s \in [\bar{S}]$), we have

$$\begin{aligned}
\sum_{s=0}^{S-1} \mathbb{E}_{\tau_{\min}}[\Delta l(x_{k,s}, \bar{\tau}_{\min}, g_{k,s}, d_{k,s})] &\leq \frac{\mathbb{E}_{\tau_{\min}}[R_{k,0} - R_{k,S}]}{\Lambda_{\min}} \\
&= \frac{\mathbb{E}_{\tau_{\min}}[\phi(x_{k,0}, \bar{\tau}_{\min}) - \phi(x_{k+1,0}, \bar{\tau}_{\min})]}{\Lambda_{\min}}.
\end{aligned}$$

The equality follows from the fact that $\lambda_S = 0$ and $x_{k,S} = x_{k+1,0}$. The desired conclusion (4.14) then follows by summing this inequality for $j \in \{\bar{k}_\tau, \bar{k}_\tau + 1, \dots, k\}$. \square

As a consequence of Lemma 4.18, in Theorem 4.19 we present the main convergence result of this subsection.

Theorem 4.19. *Suppose Assumption 4.8 holds. Let λ_s, Λ_s and Λ_{\min} be defined as in Lemma 4.18. Suppose $\beta = \frac{\mu_1 b}{(\bar{\tau}_{\min} L + \Gamma) N^\gamma} \in (0, 1]$ with $\mu_1 \in (0, 1]$, $z = \frac{\bar{\tau}_{\min} L + \Gamma}{N^{\gamma/2}}$, and $S \leq$*

$\left\lfloor \frac{N^{\gamma/2}}{(\bar{\tau}_{\min} L + \Gamma) + \frac{\mu_1^2 a_u^2 L^2 b}{(\bar{\tau}_{\min} L + \Gamma)^3 \kappa_l \zeta} + \frac{\mu_1^2 a_u^2 L^2 b}{(\bar{\tau}_{\min} L + \Gamma)^2 \kappa_l \zeta}} \right\rfloor$. Define the quantity $\Lambda_{\min} = \min_s \Lambda_s$. Then for $b < N^\gamma$, there exists universal constants μ_1, ν_1 such that: $\Lambda_{\min} \geq \frac{\nu_1 b}{(\bar{\tau}_{\min} L + \Gamma) N^\gamma}$ and

$$\begin{aligned}
\mathbb{E}_{\tau_{\min}} \left[\frac{1}{(k - \bar{k}_\tau + 1)S} \sum_{j=\bar{k}_\tau}^k \sum_{s=0}^{S-1} \Delta l(x_{j,s}, \bar{\tau}_{\min}, g_{j,s}, d_{j,s}) \right] \\
\leq \frac{(\bar{\tau}_{\min} L + \Gamma) N^\gamma (\mathbb{E}_{\tau_{\min}}[\phi(x_{\bar{k}_\tau, 0}, \bar{\tau}_{\min})] - \phi_{\inf})}{(k - \bar{k}_\tau + 1) S \nu_1 b}.
\end{aligned}$$

Proof. By the recursive definition of λ_s and the fact that $\lambda_S = 0$, we have that

$$\lambda_0 = (\kappa_H \kappa_d + \frac{\beta^2}{\zeta}) \frac{\alpha_u \bar{\tau}_{\min} L^2}{2b} \frac{(1+\rho)^S - 1}{\rho}, \tag{4.15}$$

with

$$\begin{aligned}
\rho &= z + (1+z) \frac{a_u^2 \beta^2 L^2}{\kappa_l z b \zeta} \\
&= (\bar{\tau}_{\min} L + \Gamma) N^{-\gamma/2} + \left(\frac{N^{\gamma/2}}{\bar{\tau}_{\min} L + \Gamma} + 1 \right) \frac{\mu_1^2 a_u^2 L^2 b}{(\bar{\tau}_{\min} L + \Gamma)^2 \kappa_l N^{2\gamma} \zeta} \\
&\leq \left((\bar{\tau}_{\min} L + \Gamma) + \frac{\mu_1^2 a_u^2 L^2 b}{(\bar{\tau}_{\min} L + \Gamma)^3 \kappa_l \zeta} + \frac{\mu_1^2 a_u^2 L^2 b}{(\bar{\tau}_{\min} L + \Gamma)^2 \kappa_l \zeta} \right) N^{-\gamma/2}
\end{aligned}$$

It follows that

$$\lambda_0 \leq \frac{\alpha_u \bar{\tau}_{\min} L^2}{2b} \left(\kappa_H \kappa_d + \frac{\mu_1^2 b^2}{(\bar{\tau}_{\min} L + \Gamma)^2 N^{2\gamma} \zeta} \right) \frac{e-1}{(\bar{\tau}_{\min} L + \Gamma) N^{-\gamma/2}}$$

where the inequality is obtained by noticing that for $l > 0$, $(1 + \frac{1}{l})^l$ is an increasing function and $(1 + \frac{1}{l})^l \rightarrow e$ as $l \rightarrow \infty$. Hence, $(1 + \rho)^S \leq e$ by the definition of S . Now, with the upper bound of λ_0 , the fact that λ_s is decreasing as s increases from 0 to S , and $\mu_0 \in (0, 1]$ and $N \geq 1$, we can lower bound Λ_{\min} as

$$\begin{aligned}
\Lambda_{\min} &= \min_{0 \leq s \leq S-1} \left\{ \alpha_l \beta - \frac{1}{2} \alpha_u \beta^2 - \frac{\alpha_u \kappa_H \kappa_d \beta^2}{2\kappa_l} - \lambda_{s+1} (1+z) \frac{\alpha_u^2 \beta^2}{\bar{\tau}_{\min} \kappa_l z} \right\} \\
&> \alpha_l \beta - \frac{1}{2} \alpha_u \beta^2 - \frac{\alpha_u \kappa_H \kappa_d \beta^2}{2\kappa_l} - \lambda_0 (1+z) \frac{\alpha_u^2 \beta^2}{\bar{\tau}_{\min} \kappa_l z} \\
&\geq \beta \left[\alpha_l - \frac{\alpha_u b \mu_1}{2(\bar{\tau}_{\min} L + \Gamma) N^\gamma} - \frac{\alpha_u \kappa_H \kappa_d \mu_1}{2\kappa_l (\bar{\tau}_{\min} L + \Gamma) N^\gamma} - \frac{\alpha_u^3 L^2 \kappa_H \kappa_d (e-1) \mu_1}{2\kappa_l (\bar{\tau}_{\min} L + \Gamma)^3} \right. \\
&\quad \left. - \frac{\alpha_u^3 L^2 \kappa_H \kappa_d (e-1) \mu_1}{2\kappa_l (\bar{\tau}_{\min} L + \Gamma)^2 N^{\gamma/2}} - \frac{\alpha_u^3 \bar{\tau}_{\min} L^2 b^2 (e-1) \mu_1}{2(\bar{\tau}_{\min} L + \Gamma)^5 N^{2\gamma} \zeta \kappa_l} - \frac{\alpha_u^3 \bar{\tau}_{\min} L^2 b^2 (e-1) \mu_1}{2(\bar{\tau}_{\min} L + \Gamma)^4 N^{5\gamma/2} \zeta \kappa_l} \right]
\end{aligned}$$

Let $\nu_1 = \alpha_l - \frac{\alpha_u b \mu_1}{2(\bar{\tau}_{\min} L + \Gamma) N^\gamma} - \frac{\alpha_u \kappa_H \kappa_d \mu_1}{2\kappa_l (\bar{\tau}_{\min} L + \Gamma) N^\gamma} - \frac{\alpha_u^3 L^2 \kappa_H \kappa_d (e-1) \mu_1}{2\kappa_l (\bar{\tau}_{\min} L + \Gamma)^3} - \frac{\alpha_u^3 L^2 \kappa_H \kappa_d (e-1) \mu_1}{2\kappa_l (\bar{\tau}_{\min} L + \Gamma)^2 N^{\gamma/2}} - \frac{\alpha_u^3 \bar{\tau}_{\min} L^2 b^2 (e-1) \mu_1}{2(\bar{\tau}_{\min} L + \Gamma)^5 N^{2\gamma} \zeta \kappa_l} - \frac{\alpha_u^3 \bar{\tau}_{\min} L^2 b^2 (e-1) \mu_1}{2(\bar{\tau}_{\min} L + \Gamma)^4 N^{5\gamma/2} \zeta \kappa_l}$. By choosing μ_1 (independent of N) such that $\nu_1 > 0$, it follows that $\Lambda_{\min} \geq \frac{b \nu_1}{(\bar{\tau}_{\min} L + \Gamma) N^\gamma}$. Combining this lower bound with Lemma 4.11 yields the desired result. \square

We conclude this section by noting that an analogue of Corollary 4.13 can be proven for the case in which adaptive step sizes are utilized. For brevity we omit this corollary since it is identical to Corollary 4.13 up to constants.

5 Numerical Results

In this section, we demonstrate the empirical performance of a Matlab implementation of Algorithm 1, with both **Options I** and **II**, for solving equality constrained binary classification machine learning problems. Specifically, we consider constrained logistic regression problems (datasets from the LIBSVM collection [8]) with linear equality constraints or an ℓ_2 norm squared constraint. All experiments were run in Matlab R2021b on macOS 12.2 with an Apple M1 Pro chip and 16GB memory.

In order to illustrate the merits of our proposed algorithm, we compared two variants of the SVR-SQP method (constant step sizes SVR-SQP-C and adaptive steps sizes SVR-SQP-A) with the stochastic SQP method from [4] (Sto-SQP) and a Stochastic Subgradient method that utilizes SVRG-type variance reduced gradient approximations (Sto-Subgrad-VR). The goals of this section can be summarized as follows: (1) illustrate the power and robustness of the adaptive step size variant of the SVR-SQP method; (2) show the advantages of utilizing variance reduced gradient approximations; (3) demonstrate the advantage of the SQP paradigm over a simple stochastic subgradient method; and, (4) show the robustness of the SVR-SQP method to user-defined parameters such as the inner iteration length (S) and the adaptive step size parameter (β).

5.1 Problem Specification, Implementation Details and Evaluation Metrics

Throughout this section we consider the following two constrained binary classification problems:

$$\min_{x \in \mathbb{R}^n} f(x) = \frac{1}{N} \sum_{i=1}^N \log \left(1 + e^{-y_i (X_i^T x)} \right) \quad \text{s.t.} \quad Ax = a_1 \quad (5.1)$$

$$\min_{x \in \mathbb{R}^n} f(x) = \frac{1}{N} \sum_{i=1}^N \log \left(1 + e^{-y_i (X_i^T x)} \right) \quad \text{s.t.} \quad \|x\|_2^2 = a_2 \quad (5.2)$$

where $X \in \mathbb{R}^{N \times n}$ is the data matrix (containing feature data for N data points; X_i representing the i th column of X) and $y \in \{-1, 1\}^N$ are the labels (for each data point), and $A \in \mathbb{R}^{m \times n}$, $a_1 \in \mathbb{R}^m$ and $a_2 = 1 \in \mathbb{R}_{>0}$ define the constraints. We consider 10 datasets, listed in Table 2, from the LIBSVM collection [8]. For the linear constraints in (5.1), we generated normal random A and a_1 for each problem with $m = 10$.

Table 2: Binary classification data set details. For more information see [8].

dataset	dimension (n)	datapoints (N)
a9a	123	32,561
australian	14	621
heart	13	270
ijcnn1	22	35,000
ionosphere	34	351
mushroom	112	5,500
phising	68	11,055
sonar	60	208
splice	60	3,175
w8a	300	49,749

A budget of 30 epochs (i.e., number of effective passes over the dataset; equivalent to the number of gradient evaluations of the objective function) was used for all methods. For all problems and algorithms, the initial primal iterate (x_0) was set to a normal random vector scaled to have norm 0.1, and the multipliers were initialized as $y_0 = \arg \min_{y \in \mathbb{R}^m} \|g_0 + J_0^T y\|^2$. For each method, we considered two batch sizes $b = 16$ (small batch) and $b = 128$ (large batch). For each problem, dataset, algorithm and batch size, we ran 10 instances with different random seeds. With regards to the constraint Lipschitz constant estimates, we used the true constants $\Gamma_{k,s} = 0$ (for (5.1)) and $\Gamma_{k,s} = 2$ (for (5.2)) for all $k \in \mathbb{N}$ and $s \in [\bar{S}]$ for all algorithms. We set $L_{k,s} = L$ for all $k \in \mathbb{N}$ and $s \in [\bar{S}]$ for all algorithms, where L was estimated by differences of gradients around the initial iterate. The details of all parameter settings are given below.

- **SVR-SQP-C** and **SVR-SQP-A**: $\sigma = 0.5$, $\theta = 10^4$, $\bar{\tau}_{-1,0} = 0.1$, and $\epsilon_\tau = 10^{-6}$
 - **SVR-SQP-C** [Algorithm 1; constant step sizes]:
 $\bar{\alpha}_{k,s} = \alpha \in \{10^{-3}, 10^{-2}, 10^{-1}, 10^0, 10^1\}$ for all $(k, s) \in \mathbb{N} \times [\bar{S}]$;
 - **SVR-SQP-A** [Algorithm 1; adaptive step sizes]: $\bar{\alpha}_{k,s}$ computed via (3.11)–(3.13) for all $(k, s) \in \mathbb{N} \times [\bar{S}]$, $\alpha_u = 10^6$, $\beta = 1$.
- **Sto-SQP** [4, Algorithm 3.1]: $\theta = 10^4$, $\bar{\tau}_{-1} = 0.1$, $\epsilon_\tau = 10^{-6}$, $\bar{\xi}_{-1} = 0.1$, $\epsilon_\xi = 10^{-2}$, $\sigma = 0.5$, and $\beta_k = \beta \in \{10^{-3}, 10^{-2}, 10^{-1}, 10^0, 10^1\}$ for all $k \in \mathbb{N}$.
- **Sto-Subgrad-VR**: $\bar{\alpha}_{k,s} = \frac{\alpha}{\tau L + \Gamma}$ with $\alpha \in \{10^{-3}, 10^{-2}, 10^{-1}, 10^0, 10^1\}$ and $\tau_{k,s} = \tau \in \{10^{-10}, 10^{-9}, \dots, 10^0\}$ for all $(k, s) \in \mathbb{N} \times [\bar{S}]$. We should note that **Sto-Subgrad-VR** applies the SVRG algorithm [16] to directly minimize the nonsmooth merit function (3.4); SVRG directly applied to the smooth part of merit function with the subgradient of the nonsmooth part added.

For all algorithms with inner outer iterations, the inner iteration length was set as $S = \lfloor \frac{N}{2b} \rfloor$, unless otherwise specified.

In all of our experiments, results are given in terms of feasibility and stationarity errors discussed below. We present the evolution of these measures with respect to epochs in Figures 1, 3, 4 and 5. Moreover, in Figure 2 and Tables 3 and 4, we report the error metrics at the best iterate found within the budget defined as follows. Given a fixed epoch budget, assume we have $k \in \{0, \dots, K\}$ for some $K \in \mathbb{N}$. If

$$\min\{\|c_k\|_\infty : k = 0, \dots, K\} > 10^{-6}, \quad \text{then } x_{\text{best}} \leftarrow \arg \min_{x_k \in \{x_0, \dots, x_K\}} \|c_k\|_\infty.$$

Otherwise, if $\min\{\|c_k\|_\infty : k = 0, \dots, K\} \leq 10^{-6}$, then we set

$$x_{\text{best}} \leftarrow \arg \min_{x_k \in \{x_0, \dots, x_K\}} \|\nabla f_k + J_k^T y_{k,1s}\|_\infty \quad \text{s.t. } \|c_k\|_\infty \leq 10^{-6},$$

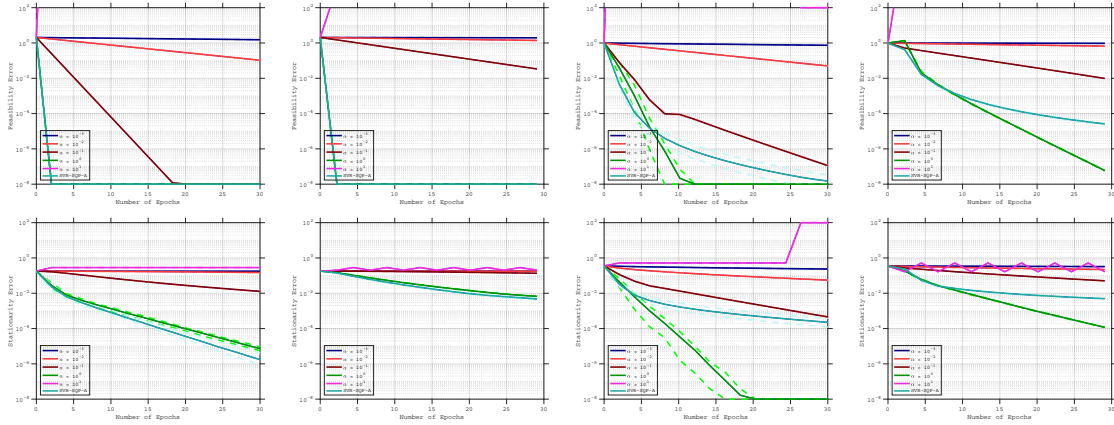
where $y_{k,1s}$ is the least-squares multiplier at x_k . Given x_{best} and the corresponding dual variables $y_{\text{best},1s}$, we report feasibility error ($\|c(x_{\text{best}})\|_\infty$) and stationarity error ($\|\nabla f(x_{\text{best}}) + \nabla c(x_{\text{best}})y_{\text{best},1s}\|_\infty$).

5.2 Comparison: SVR-SQP-C and SVR-SQP-A

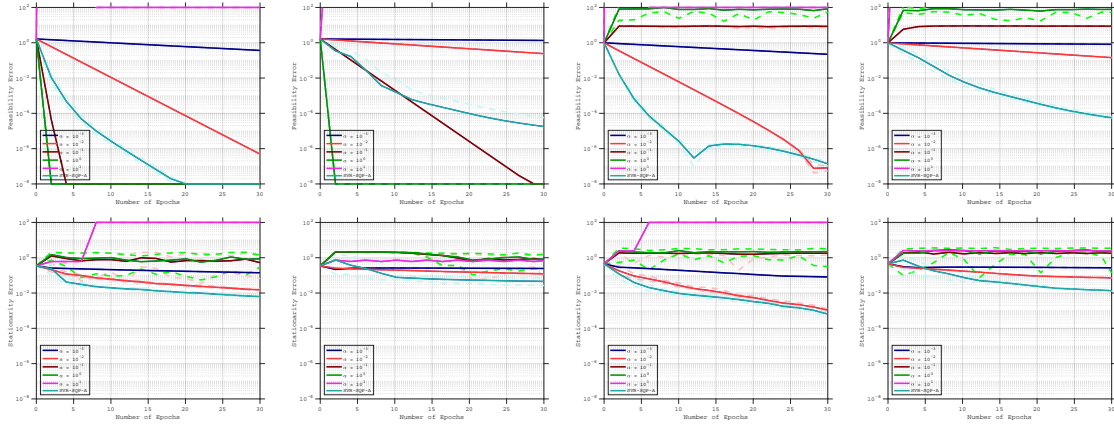
In this section, we compare the performance of SVR-SQP-C and SVR-SQP-A on (5.1) and (5.2). We ran all methods for 30 epochs with two different batch sizes. For SVR-SQP-C we tuned the step size $\bar{\alpha}_{k,s} = \alpha \in \{10^{-3}, 10^{-2}, 10^{-1}, 10^0, 10^1\}$ for all $(k, s) \in \mathbb{N} \times [\bar{S}]$. For the SVR-SQP-A we set $\beta = 1$. For both methods we used $S = \lfloor \frac{N}{2b} \rfloor$.

Figs. 1a and 1b show the stationarity and feasibility errors versus epochs for two datasets (`australian` and `splice`) for the SVR-SQP-C (different values of α) and SVR-SQP-A ($\beta = 1$) methods with different batch sizes. For each method, the figure shows the average trajectory (solid line) over the 10 random seeds of the measures with respect to epochs, and the 95% confidence interval (dashed lines). As is clear, SVR-SQP-A appears to be competitive with the best tuned version of the SVR-SQP-C method.

Similar behavior was observed on other datasets. Fig. 2 presents feasibility and stationarity errors for all datasets in Table 2 for the best iterates found by four variants of SVR-SQP-C and SVR-SQP-A. For each problem, we report as boxplots the feasibility and stationarity errors for the best iterate found by each method for two different batch sizes (4 box plots per problem per method). From Fig. 2, we observe that for both batch size options and for both constraints types, SVR-SQP-A performs as good as (if not better than) SVR-SQP-C with the best tuned step size in terms of stationarity and feasibility.



(a) **australian** dataset. Top row: feasibility error; Bottom row: stationarity error.



(b) **splice** dataset. Top row: feasibility error; Bottom row: stationarity error.

Figure 1: Performance of SVR-SQP-C with different step sizes $\alpha \in \{10^{-3}, 10^{-2}, 10^{-1}, 10^0, 10^1\}$ and SVR-SQP-C on logistic regression problems with linear (columns 1 and 2) and ℓ_2 norm (columns 3 and 4) constraints. First and third columns: batch size 16; Second and fourth columns: batch size 128.

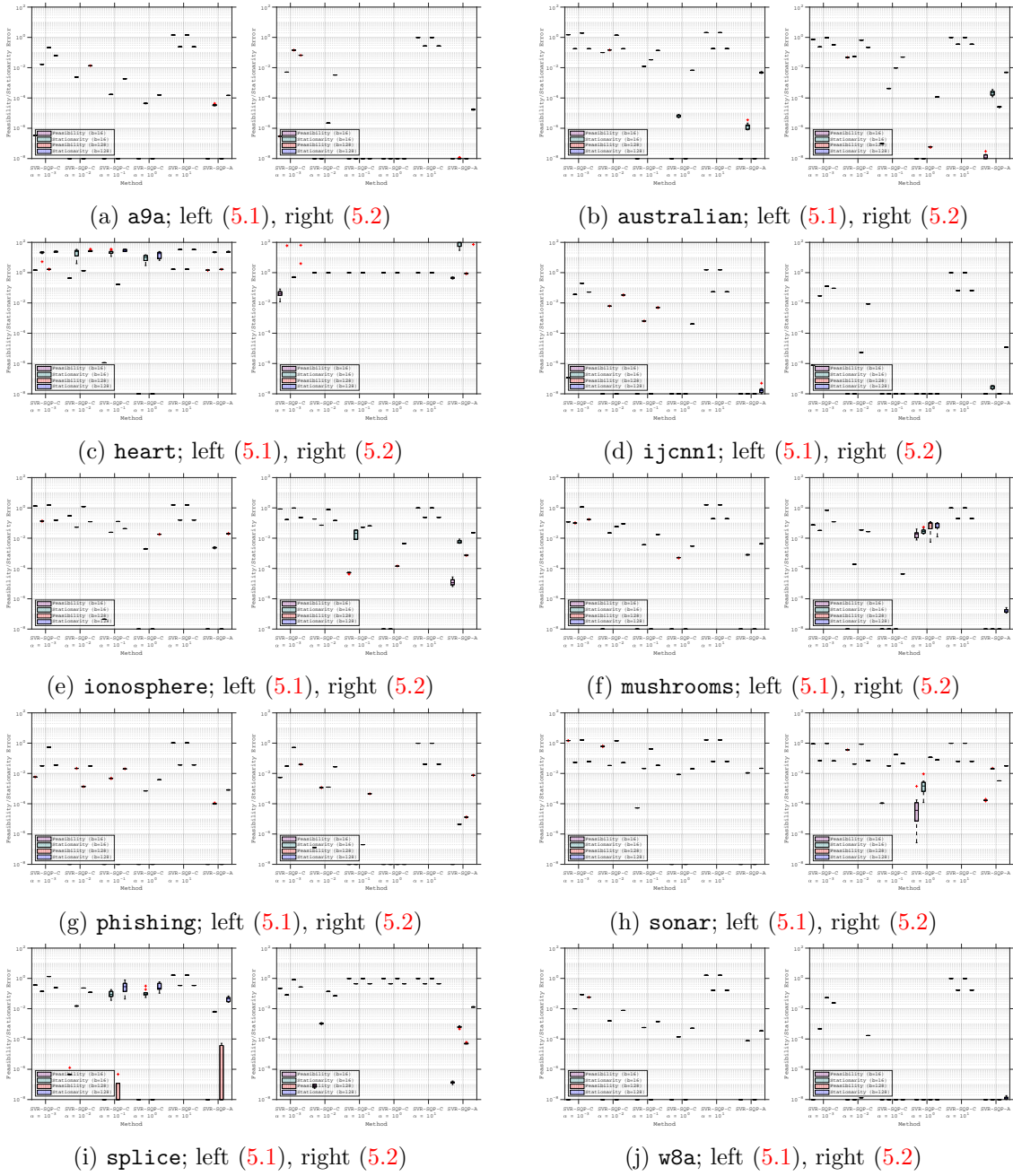
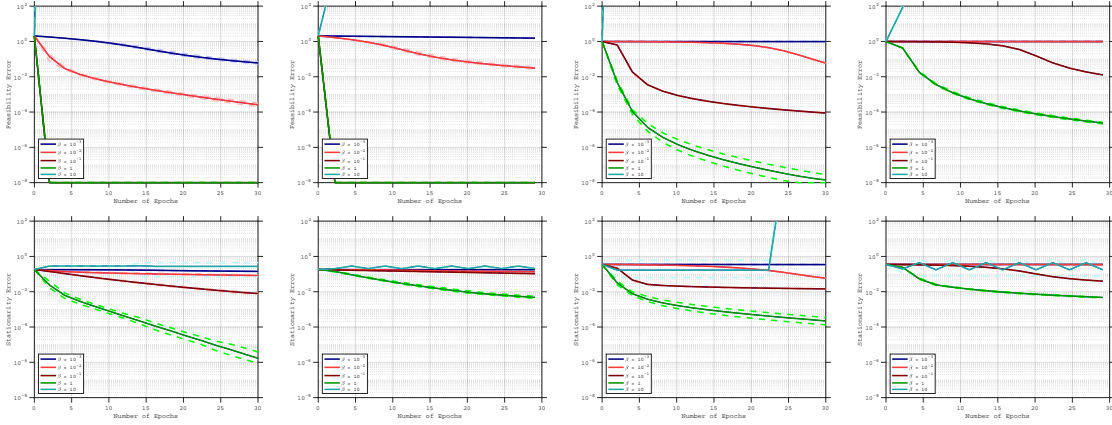


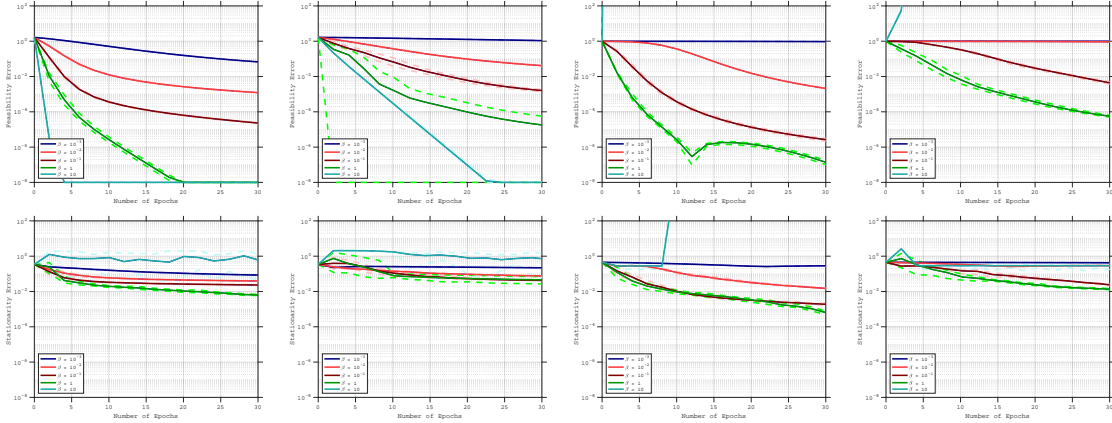
Figure 2: Best feasibility and stationarity errors, for SVR-SQP-C and SVR-SQP-A on (5.1) and (5.2).

5.3 Sensitivity to user-defined parameters

Given the encouraging numerical results for SVR-SQP-A (Section 5.2), in this subsection we investigate the robustness of SVR-SQP-A to two user-defined parameters: (1) the step size parameter $\beta \in \{10^{-3}, 10^{-2}, 10^{-1}, 10^0, 10^1\}$ (Fig. 3), and (2) the number of inner iterations $S \in \{\lfloor \frac{N}{b} \rfloor, \lfloor \frac{N}{2b} \rfloor, \lfloor \frac{N}{4b} \rfloor\}$ (Fig. 4) for two datasets (australian and splice). Overall, the results on these two datasets suggest that $\beta = 1$ is often the best choice. Moreover, our results in Fig. 4 illustrate the robustness of SVR-SQP-A to the choice of the number of inner iterations.



(a) `australian` dataset. Top row: feasibility error; Bottom row: stationarity error.



(b) `splice` dataset. Top row: feasibility error; Bottom row: stationarity error.

Figure 3: Performance of SVR-SQP-A with different step sizes parameter values $\beta \in \{10^{-3}, 10^{-2}, 10^{-1}, 10^0, 10^1\}$ on logistic regression problems with linear (columns 1 and 2) and ℓ_2 norm (columns 3 and 4) constraints. First and third columns: batch size 16; Second and fourth columns: batch size 128.

5.4 Comparison: SVR-SQP-A, Sto-SQP and Sto-Subgrad-VR

In this final subsection, we compare the performance of **SVR-SQP-A** to that of **Sto-SQP** [4, Algorithm 2] and **Sto-Subgrad-VR**. A budget of 30 epochs was used for all methods. For all methods, the inner iteration length was set to $S = \lfloor \frac{N}{2b} \rfloor$. For the **Sto-SQP** method the step size parameter was tuned $\beta \in \{10^{-3}, 10^{-2}, 10^{-1}, 10^0, 10^1\}$ for all $k \in \mathbb{N}$, and for the **Sto-Subgrad-VR** method the step size parameter and the merit parameter were tuned $\bar{\alpha}_{k,s} = \frac{\alpha}{\tau L + \Gamma}$ for all $(k, s) \in \mathbb{N} \times [\bar{S}]$ where $\alpha \in \{10^{-3}, 10^{-2}, 10^{-1}, 10^0, 10^1\}$, $\tau_{k,s} = \tau \in \{10^{-10}, 10^{-9}, \dots, 10^0\}$. For the **SVR-SQP-A**, we set $\beta = 1$. Overall, this meant that the **Sto-SQP** and **Sto-Subgrad-VR** methods were effectively run for 5 and 55 times the number of epochs, respectively, than were allowed for our method.

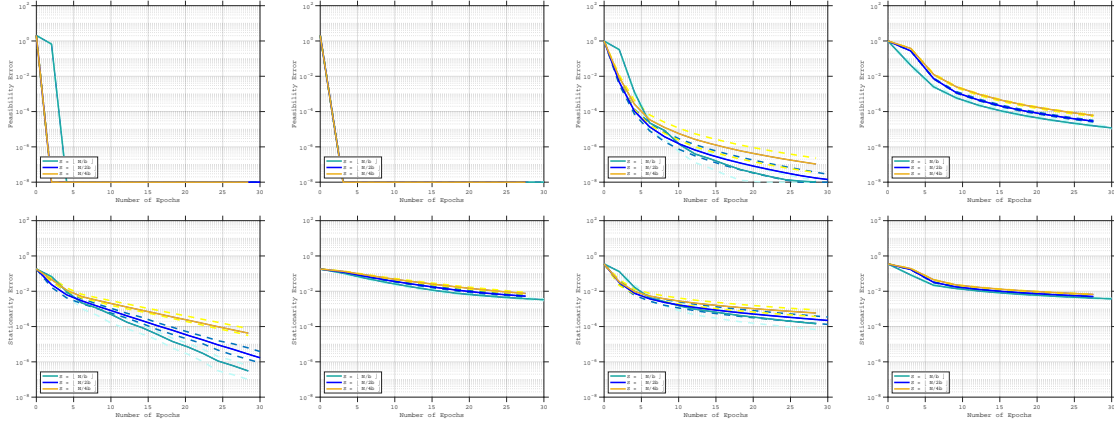
The results of these experiments are reported in Figs. 5a and 5b and in Tables 3 and 4. For each batch size and dataset, we report the average feasibility and stationarity errors for the best iterates generated (defined in Section 5.1) for the best hyper-parameter settings for each method in Tables 3 and 4. The results suggest that, when small batch sizes are employed (i.e., $b = 16$), **SVR-SQP-A** consistently outperforms the other methods for both sets of constraints. When large batch sizes are used (i.e., $b = 128$), **SVR-SQP-A** is competitive with **Sto-SQP**, even though the adaptive step size parameter β is well-tuned for **Sto-SQP** whereas for **SVR-SQP-A** we simply set $\beta = 1$. We should note again that 5 and 55 times the tuning effort was allocated to **Sto-SQP** and **Sto-Subgrad-VR**, respectively, as compared to **SVR-SQP-A**.

Table 3: Average feasibility and stationarity errors over 10 independent runs for each experiment, along with 95% confidence intervals represented by ‘ \pm ’, of **best tuned** variants of Sto-Subgrad-VR and Sto-SQP, and SVR-SQP-A with $\beta = 1$ and $S = \lfloor \frac{N}{2b} \rfloor$ on logistic regression problems with linear constraints (5.1). The results for the best-performing algorithm for each batch size are shown in **bold**. The symbol \star indicates that all the runs for a given method converged to $\min\{\|c_k\|_\infty : k = 0, \dots, K\} \leq 10^{-6}$.

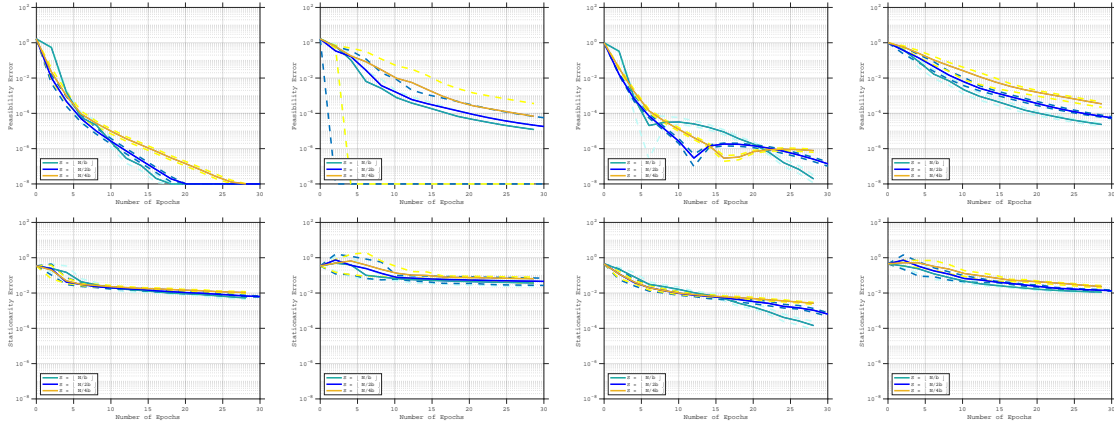
Dataset	Batch	Sto-Subgrad-VR		Sto-SQP		SVR-SQP-A	
		Feasibility	Stationarity	Feasibility	Stationarity	Feasibility	Stationarity
a9a	16	$1.2 \times 10^{-1} \pm 7.2 \times 10^{-3}$	$2.6 \times 10^{-2} \pm 1.1 \times 10^{-2}$	\star	$8.4 \times 10^{-4} \pm 2.9 \times 10^{-5}$	\star	$3.5 \times 10^{-5} \pm 4.5 \times 10^{-6}$
a9a	128	$1.1 \times 10^{-1} \pm 8.4 \times 10^{-3}$	$8.7 \times 10^{-2} \pm 1.6 \times 10^{-2}$	\star	$7.0 \times 10^{-4} \pm 2.8 \times 10^{-5}$	\star	$1.5 \times 10^{-4} \pm 1.1 \times 10^{-6}$
australian	16	$4.1 \times 10^{-2} \pm 1.1 \times 10^{-2}$	$2.0 \times 10^{-1} \pm 1.5 \times 10^{-3}$	\star	$1.2 \times 10^{-3} \pm 2.2 \times 10^{-4}$	\star	$1.5 \times 10^{-6} \pm 5.6 \times 10^{-7}$
australian	128	$3.1 \times 10^{-1} \pm 3.5 \times 10^{-2}$	$2.0 \times 10^{-1} \pm 4.7 \times 10^{-3}$	\star	$2.3 \times 10^{-3} \pm 7.3 \times 10^{-4}$	\star	$4.8 \times 10^{-3} \pm 2.7 \times 10^{-4}$
heart	16	$2.3 \times 10^{-1} \pm 4.8 \times 10^{-2}$	$1.2 \times 10^1 \pm 3.9 \times 10^0$	$7.9 \times 10^{-2} \pm 1.2 \times 10^{-2}$	$2.3 \times 10^1 \pm 5.1 \times 10^0$	$1.5 \times 10^0 \pm 1.6 \times 10^{-2}$	$2.2 \times 10^1 \pm 1.3 \times 10^0$
heart	128	$8.7 \times 10^{-1} \pm 1.6 \times 10^{-2}$	$2.3 \times 10^1 \pm 3.0 \times 10^0$	$1.2 \times 10^0 \pm 4.1 \times 10^{-2}$	$2.4 \times 10^1 \pm 2.0 \times 10^0$	$1.6 \times 10^0 \pm 2.3 \times 10^{-3}$	$2.6 \times 10^1 \pm 1.8 \times 10^0$
ijcnn1	16	$9.9 \times 10^{-1} \pm 1.1 \times 10^{-1}$	$6.6 \times 10^{-3} \pm 5.5 \times 10^{-3}$	\star	$1.8 \times 10^{-4} \pm 1.9 \times 10^{-5}$	\star	$1.3 \times 10^{-8} \pm 7.8 \times 10^{-9}$
ijcnn1	128	$9.7 \times 10^{-1} \pm 1.8 \times 10^{-1}$	$1.6 \times 10^{-2} \pm 6.5 \times 10^{-3}$	\star	$2.1 \times 10^{-4} \pm 2.2 \times 10^{-5}$	\star	$1.0 \times 10^{-8} \pm 1.8 \times 10^{-10}$
ionosphere	16	$3.2 \times 10^{-2} \pm 4.9 \times 10^{-3}$	$1.4 \times 10^{-1} \pm 5.3 \times 10^{-3}$	\star	$4.2 \times 10^{-3} \pm 4.3 \times 10^{-4}$	\star	$2.4 \times 10^{-3} \pm 1.8 \times 10^{-4}$
ionosphere	128	$4.4 \times 10^{-1} \pm 5.4 \times 10^{-2}$	$1.4 \times 10^{-1} \pm 1.1 \times 10^{-2}$	\star	$1.2 \times 10^{-2} \pm 9.8 \times 10^{-4}$	\star	$2.0 \times 10^{-2} \pm 4.0 \times 10^{-4}$
mushrooms	16	$4.8 \times 10^{-2} \pm 6.0 \times 10^{-3}$	$1.3 \times 10^{-1} \pm 2.0 \times 10^{-2}$	\star	$4.2 \times 10^{-4} \pm 9.0 \times 10^{-6}$	\star	$8.0 \times 10^{-4} \pm 2.4 \times 10^{-5}$
mushrooms	128	$5.7 \times 10^{-2} \pm 5.5 \times 10^{-3}$	$1.8 \times 10^{-1} \pm 4.5 \times 10^{-3}$	\star	$2.5 \times 10^{-3} \pm 9.2 \times 10^{-5}$	\star	$4.2 \times 10^{-3} \pm 7.7 \times 10^{-5}$
phising	16	$3.7 \times 10^{-1} \pm 5.4 \times 10^{-2}$	$2.6 \times 10^{-2} \pm 2.5 \times 10^{-4}$	\star	$4.6 \times 10^{-4} \pm 1.2 \times 10^{-5}$	\star	$9.9 \times 10^{-5} \pm 3.2 \times 10^{-6}$
phising	128	$6.0 \times 10^{-1} \pm 2.2 \times 10^{-2}$	$3.6 \times 10^{-2} \pm 1.4 \times 10^{-4}$	\star	$2.7 \times 10^{-3} \pm 3.6 \times 10^{-5}$	\star	$8.1 \times 10^{-4} \pm 1.3 \times 10^{-5}$
sonar	16	$4.1 \times 10^{-2} \pm 2.8 \times 10^{-3}$	$3.9 \times 10^{-2} \pm 1.2 \times 10^{-3}$	\star	$7.5 \times 10^{-3} \pm 4.4 \times 10^{-4}$	\star	$1.1 \times 10^{-2} \pm 3.1 \times 10^{-4}$
sonar	128	$4.1 \times 10^{-1} \pm 3.2 \times 10^{-2}$	$4.1 \times 10^{-2} \pm 3.0 \times 10^{-3}$	\star	$1.9 \times 10^{-2} \pm 2.9 \times 10^{-4}$	\star	$2.2 \times 10^{-2} \pm 7.8 \times 10^{-5}$
splice	16	$6.8 \times 10^{-4} \pm 6.2 \times 10^{-5}$	$2.6 \times 10^{-1} \pm 2.2 \times 10^{-4}$	\star	$4.1 \times 10^{-3} \pm 3.7 \times 10^{-4}$	\star	$7.6 \times 10^{-3} \pm 2.0 \times 10^{-4}$
splice	128	$7.0 \times 10^{-3} \pm 3.3 \times 10^{-4}$	$2.6 \times 10^{-1} \pm 7.6 \times 10^{-4}$	\star	$1.8 \times 10^{-2} \pm 8.0 \times 10^{-4}$	$1.9 \times 10^{-3} \pm 9.4 \times 10^{-5}$	$4.3 \times 10^{-2} \pm 1.2 \times 10^{-3}$
w8a	16	$5.1 \times 10^{-1} \pm 8.6 \times 10^{-3}$	$4.8 \times 10^{-3} \pm 1.1 \times 10^{-4}$	\star	$2.6 \times 10^{-4} \pm 1.9 \times 10^{-5}$	\star	$7.5 \times 10^{-5} \pm 7.9 \times 10^{-7}$
w8a	128	$6.9 \times 10^{-1} \pm 2.5 \times 10^{-3}$	$2.6 \times 10^{-2} \pm 1.8 \times 10^{-4}$	\star	$1.8 \times 10^{-4} \pm 5.0 \times 10^{-6}$	\star	$3.4 \times 10^{-4} \pm 2.1 \times 10^{-6}$

Table 4: Average feasibility and stationarity errors over 10 independent runs for each experiment, along with 95% confidence intervals represented by ‘ \pm ’, of **best tuned** variants of Sto-Subgrad-VR and Sto-SQP, and SVR-SQP-A with $\beta = 1$ and $S = \lfloor \frac{N}{2b} \rfloor$ on logistic regression problems with ℓ_2 constraint (5.2). The results for the best-performing algorithm for each batch size are shown in **bold**. The symbol \star indicates that all the runs for a given method converged to $\min\{\|c_k\|_\infty : k = 0, \dots, K\} \leq 10^{-6}$.

Dataset	Batch	Sto-Subgrad-VR		Sto-SQP		SVR-SQP-A	
		Feasibility	Stationarity	Feasibility	Stationarity	Feasibility	Stationarity
a9a	16	$3.0 \times 10^{-6} \pm 9.7 \times 10^{-7}$	$2.5 \times 10^{-1} \pm 7.4 \times 10^{-7}$	\star	$1.1 \times 10^{-2} \pm 7.8 \times 10^{-4}$	\star	$1.5 \times 10^{-8} \pm 1.9 \times 10^{-9}$
a9a	128	$2.9 \times 10^{-5} \pm 3.8 \times 10^{-6}$	$2.5 \times 10^{-1} \pm 3.6 \times 10^{-8}$	\star	$8.1 \times 10^{-3} \pm 2.2 \times 10^{-4}$	\star	$1.7 \times 10^{-5} \pm 1.2 \times 10^{-6}$
australian	16	$2.9 \times 10^{-4} \pm 5.8 \times 10^{-5}$	$3.1 \times 10^{-1} \pm 1.3 \times 10^{-6}$	$2.9 \times 10^{-4} \pm 1.4 \times 10^{-4}$	$2.3 \times 10^{-2} \pm 1.9 \times 10^{-2}$	\star	$2.1 \times 10^{-4} \pm 4.9 \times 10^{-5}$
australian	128	$2.6 \times 10^{-3} \pm 1.1 \times 10^{-3}$	$2.5 \times 10^{-1} \pm 3.2 \times 10^{-3}$	$1.3 \times 10^{-4} \pm 4.1 \times 10^{-6}$	$9.1 \times 10^{-3} \pm 1.1 \times 10^{-3}$	$2.4 \times 10^{-5} \pm 9.6 \times 10^{-7}$	$4.7 \times 10^{-3} \pm 9.2 \times 10^{-5}$
heart	16	$1.2 \times 10^{-3} \pm 5.2 \times 10^{-4}$	$2.5 \times 10^0 \pm 2.3 \times 10^{-1}$	$3.3 \times 10^{-1} \pm 6.1 \times 10^{-2}$	$9.1 \times 10^1 \pm 2.2 \times 10^1$	$4.7 \times 10^{-1} \pm 3.8 \times 10^{-2}$	$8.8 \times 10^1 \pm 7.6 \times 10^0$
heart	128	$2.9 \times 10^{-2} \pm 1.5 \times 10^{-2}$	$1.1 \times 10^2 \pm 2.5 \times 10^1$	$9.9 \times 10^{-1} \pm 5.7 \times 10^{-5}$	$5.9 \times 10^0 \pm 1.1 \times 10^0$	$8.8 \times 10^{-1} \pm 8.3 \times 10^{-3}$	$1.1 \times 10^2 \pm 3.4 \times 10^1$
ijcnn1	16	$1.2 \times 10^{-5} \pm 8.1 \times 10^{-6}$	$6.8 \times 10^{-2} \pm 8.3 \times 10^{-5}$	$1.0 \times 10^{-6} \pm 2.6 \times 10^{-9}$	$3.8 \times 10^{-2} \pm 7.9 \times 10^{-4}$	\star	$3.1 \times 10^{-8} \pm 3.8 \times 10^{-9}$
ijcnn1	128	$4.4 \times 10^{-6} \pm 1.8 \times 10^{-7}$	$6.8 \times 10^{-2} \pm 2.1 \times 10^{-10}$	$1.0 \times 10^{-6} \pm 3.2 \times 10^{-10}$	$3.1 \times 10^{-2} \pm 5.8 \times 10^{-4}$	\star	$1.2 \times 10^{-5} \pm 1.1 \times 10^{-7}$
ionosphere	16	$1.4 \times 10^{-3} \pm 5.4 \times 10^{-4}$	$1.1 \times 10^{-1} \pm 6.4 \times 10^{-4}$	$4.3 \times 10^{-4} \pm 4.0 \times 10^{-4}$	$5.2 \times 10^{-2} \pm 1.7 \times 10^{-2}$	$1.4 \times 10^{-5} \pm 3.2 \times 10^{-6}$	$6.1 \times 10^{-3} \pm 7.0 \times 10^{-4}$
ionosphere	128	$8.0 \times 10^{-3} \pm 7.1 \times 10^{-3}$	$1.1 \times 10^{-1} \pm 1.5 \times 10^{-3}$	$5.8 \times 10^{-4} \pm 1.9 \times 10^{-5}$	$2.0 \times 10^{-2} \pm 1.2 \times 10^{-3}$	$7.6 \times 10^{-4} \pm 1.4 \times 10^{-5}$	$2.3 \times 10^{-2} \pm 3.9 \times 10^{-4}$
mushrooms	16	$3.4 \times 10^{-5} \pm 9.1 \times 10^{-6}$	$1.8 \times 10^{-1} \pm 6.4 \times 10^{-8}$	\star	$7.1 \times 10^{-3} \pm 1.2 \times 10^{-4}$	\star	$1.0 \times 10^{-8} \pm 2.1 \times 10^{-11}$
mushrooms	128	$4.0 \times 10^{-4} \pm 2.4 \times 10^{-4}$	$2.9 \times 10^{-2} \pm 1.6 \times 10^{-3}$	\star	$2.2 \times 10^{-2} \pm 5.8 \times 10^{-4}$	\star	$1.6 \times 10^{-7} \pm 2.7 \times 10^{-8}$
phising	16	$3.9 \times 10^{-5} \pm 2.1 \times 10^{-5}$	$4.0 \times 10^{-2} \pm 4.1 \times 10^{-6}$	$8.4 \times 10^{-5} \pm 1.2 \times 10^{-6}$	$1.0 \times 10^{-3} \pm 4.0 \times 10^{-4}$	\star	$4.4 \times 10^{-6} \pm 5.2 \times 10^{-8}$
phising	128	$7.1 \times 10^{-4} \pm 3.5 \times 10^{-4}$	$2.9 \times 10^{-2} \pm 1.7 \times 10^{-3}$	$1.2 \times 10^{-5} \pm 3.4 \times 10^{-7}$	$3.0 \times 10^{-3} \pm 3.4 \times 10^{-4}$	$1.3 \times 10^{-5} \pm 6.2 \times 10^{-8}$	$7.5 \times 10^{-3} \pm 2.0 \times 10^{-5}$
sonar	16	$2.5 \times 10^{-3} \pm 6.5 \times 10^{-4}$	$1.4 \times 10^{-1} \pm 2.0 \times 10^{-5}$	$7.4 \times 10^{-4} \pm 1.9 \times 10^{-4}$	$2.3 \times 10^{-2} \pm 3.9 \times 10^{-3}$	$1.7 \times 10^{-4} \pm 1.1 \times 10^{-5}$	$2.0 \times 10^{-2} \pm 5.6 \times 10^{-4}$
sonar	128	$3.3 \times 10^{-2} \pm 2.1 \times 10^{-2}$	$2.9 \times 10^{-2} \pm 2.5 \times 10^{-3}$	$8.9 \times 10^{-4} \pm 6.0 \times 10^{-5}$	$2.7 \times 10^{-2} \pm 1.8 \times 10^{-3}$	$3.2 \times 10^{-3} \pm 3.9 \times 10^{-7}$	$3.2 \times 10^{-2} \pm 3.5 \times 10^{-6}$
splice	16	$3.0 \times 10^{-5} \pm 3.9 \times 10^{-4}$	$1.7 \times 10^0 \pm 1.6 \times 10^{-4}$	$1.3 \times 10^{-5} \pm 8.4 \times 10^{-6}$	$1.0 \times 10^{-1} \pm 3.8 \times 10^{-2}$	\star	$6.5 \times 10^{-4} \pm 9.0 \times 10^{-5}$
splice	128	$7.1 \times 10^{-4} \pm 2.8 \times 10^{-4}$	$1.6 \times 10^0 \pm 2.1 \times 10^{-3}$	$1.0 \times 10^{-5} \pm 5.4 \times 10^{-7}$	$4.5 \times 10^{-2} \pm 2.2 \times 10^{-2}$	$5.2 \times 10^{-5} \pm 4.1 \times 10^{-6}$	$1.4 \times 10^{-2} \pm 5.2 \times 10^{-4}$
w8a	16	$3.5 \times 10^{-5} \pm 2.2 \times 10^{-5}$	$1.6 \times 10^{-1} \pm 8.4 \times 10^{-5}$	$1.3 \times 10^{-6} \pm 5.4 \times 10^{-8}$	$3.6 \times 10^{-3} \pm 4.6 \times 10^{-5}$	\star	$1.4 \times 10^{-8} \pm 2.1 \times 10^{-9}$
w8a	128	$1.5 \times 10^{-5} \pm 1.0 \times 10^{-6}$	$1.6 \times 10^{-1} \pm 2.2 \times 10^{-9}$	$1.0 \times 10^{-6} \pm 3.6 \times 10^{-10}$	$3.4 \times 10^{-3} \pm 8.9 \times 10^{-5}$	\star	$2.4 \times 10^{-8} \pm 2.4 \times 10^{-9}$

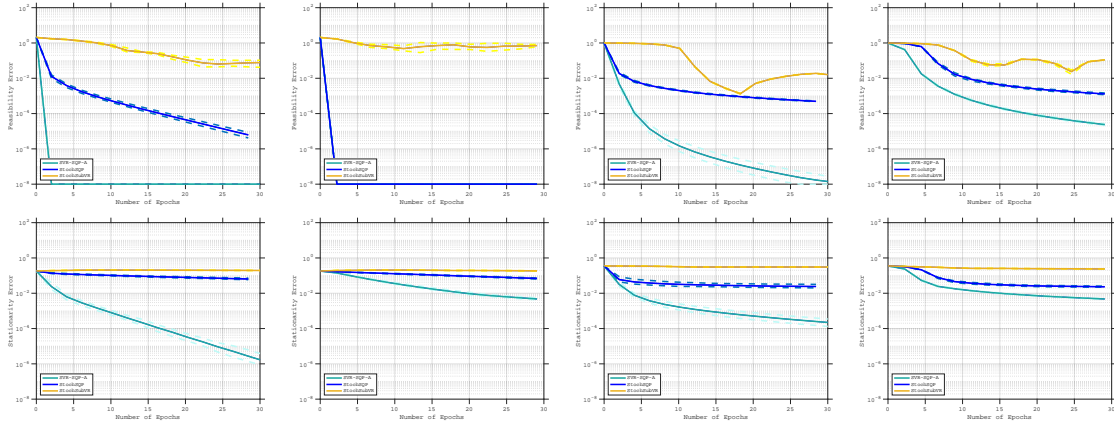


(a) **australian** dataset. Top row: feasibility error; Bottom row: stationarity error.

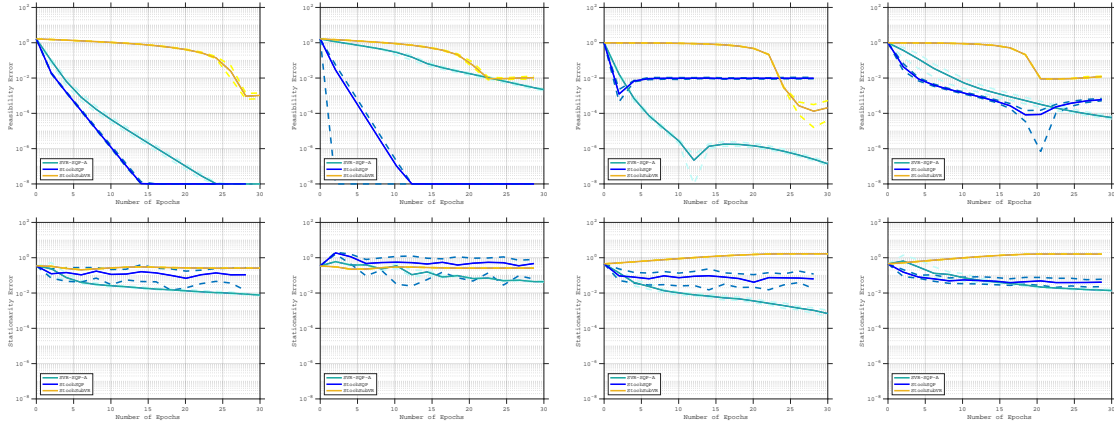


(b) **splice** dataset. Top row: feasibility error; Bottom row: stationarity error.

Figure 4: Performance of SVR-SQP-A with different step sizes parameter values $S \in \{\lfloor \frac{N}{6} \rfloor, \lfloor \frac{N}{12} \rfloor, \lfloor \frac{N}{24} \rfloor\}$ on logistic regression problems with linear (columns 1 and 2) and ℓ_2 norm (columns 3 and 4) constraints. First and third columns: batch size 16; Second and fourth columns: batch size 128.



(a) **australian** dataset. Top row: feasibility error; Bottom row: stationarity error.



(b) **splice** dataset. Top row: feasibility error; Bottom row: stationarity error.

Figure 5: Performance of **best tuned** variants of **Sto-Subgrad-VR** and **Sto-SQP**, and **SVR-SQP-A** with $\beta = 1$ and $S = \lfloor \frac{N}{2b} \rfloor$ on logistic regression problems with linear (columns 1 and 2) and ℓ_2 norm (columns 3 and 4) constraints. First and third columns: batch size 16; Second and fourth columns: batch size 128.

6 Final Remarks

We have designed and analyzed an adaptive variance reduced SQP method for minimizing general smooth finite-sum optimization problems with deterministic nonlinear equality constraints. Under common assumptions, with constant or adaptive (non-diminishing) step sizes, we presented comprehensive convergence guarantees for our proposed method. Specifically, we proved that the **SVR-SQP** method generates a sequence of iterates whose first-order stationarity measure converges to zero in expectation. Our theoretical results can be viewed as analogues of those of the SVRG method on general unconstrained non-convex finite-sum optimization problems [31]. The numerical experiments presented on classification problems from the LIBSVM collection [8] demonstrated the efficiency, efficacy and robustness of the proposed method.

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] Jianchao Bai, William W Hager, and Hongchao Zhang. An inexact accelerated stochastic admm for separable convex optimization. *Computational Optimization and Applications*, 81(2):479–518, 2022.
- [3] 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-deficient Jacobians. *arXiv preprint arXiv:2106.13015*, 2021.
- [4] 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.
- [5] Fengmiao Bian, Jingwei Liang, and Xiaoqun Zhang. A stochastic alternating direction method of multipliers for non-smooth and non-convex optimization. *Inverse Problems*, 37(7):075009, 2021.
- [6] Léon Bottou, Frank E Curtis, and Jorge Nocedal. Optimization methods for large-scale machine learning. *Siam Review*, 60(2):223–311, 2018.
- [7] 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.

- [8] 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.
- [9] Nilanjan Chatterjee, Yi-Hau Chen, Paige Maas, and Raymond J Carroll. Constrained maximum likelihood estimation for model calibration using summary-level information from external big data sources. *Journal of the American Statistical Association*, 111(513):107–117, 2016.
- [10] Changan Chen, Frederick Tung, Naveen Vedula, and Greg Mori. Constraint-aware deep neural network compression. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 400–415, 2018.
- [11] 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.
- [12] 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.
- [13] Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives. In *Advances in neural information processing systems*, pages 1646–1654, 2014.
- [14] Charles J Geyer. Constrained maximum likelihood exemplified by isotonic convex logistic regression. *Journal of the American Statistical Association*, 86(415):717–724, 1991.
- [15] Saeed Ghadimi, Guanghui Lan, and Hongchao Zhang. Mini-batch stochastic approximation methods for nonconvex stochastic composite optimization. *Mathematical Programming*, 155(1-2):267–305, 2016.
- [16] Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance reduction. *Advances in neural information processing systems*, 26:315–323, 2013.
- [17] Guanghui Lan. An optimal method for stochastic composite optimization. *Mathematical Programming*, 133(1):365–397, 2012.
- [18] Guanghui Lan. *First-order and Stochastic Optimization Methods for Machine Learning*. Springer Nature, 2020.
- [19] Rudolf Lioutikov, Alexandros Paraschos, Jan Peters, and Gerhard Neumann. Sample-based information-theoretic stochastic optimal control. In *2014 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3896–3902. IEEE, 2014.

- [20] Andreas A Malikopoulos. Stochastic optimal control for series hybrid electric vehicles. In *2013 American Control Conference*, pages 1189–1194. IEEE, 2013.
- [21] Pablo Márquez-Neila, Mathieu Salzmann, and Pascal Fua. Imposing hard constraints on deep networks: Promises and limitations. *arXiv: 1706.02025*, 2017.
- [22] Sen Na, Mihai Anitescu, and Mladen Kolar. An Adaptive Stochastic Sequential Quadratic Programming with Differentiable Exact Augmented Lagrangians. *arXiv preprint arXiv:2102.05320*, 2021.
- [23] Sen Na, Mihai Anitescu, and Mladen Kolar. Inequality Constrained Stochastic Non-linear Optimization via Active-Set Sequential Quadratic Programming. *arXiv preprint arXiv:2109.11502*, 2021.
- [24] Yatin Nandwani, Abhishek Pathak, and Parag Singla. A primal-dual formulation for deep learning with constraints. In *Proceedings of Neural Information Processing Systems (NeurIPS)*, pages 12157–12168, 2019.
- [25] Geoffrey Négiar, Gideon Dresdner, Alicia Tsai, Laurent El Ghaoui, Francesco Locatello, Robert Freund, and Fabian Pedregosa. Stochastic frank-wolfe for constrained finite-sum minimization. In *International Conference on Machine Learning*, pages 7253–7262. PMLR, 2020.
- [26] Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on optimization*, 19(4):1574–1609, 2009.
- [27] Lam M Nguyen, Jie Liu, Katya Scheinberg, and Martin Takáč. SARAH: A novel method for machine learning problems using stochastic recursive gradient. In *International Conference on Machine Learning*, pages 2613–2621, 2017.
- [28] Jorge Nocedal and Stephen Wright. *Numerical optimization*. Springer Series in Operations Research and Financial Engineering. Springer-Verlag New York, New York, 2006.
- [29] Hua Ouyang, Niao He, Long Tran, and Alexander Gray. Stochastic alternating direction method of multipliers. In *International conference on machine learning*, pages 80–88. PMLR, 2013.
- [30] 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.

- [31] Sashank J Reddi, Ahmed Hefny, Suvrit Sra, Barnabas Poczos, and Alex Smola. Stochastic variance reduction for nonconvex optimization. In *International conference on machine learning*, pages 314–323. PMLR, 2016.
- [32] Sashank J Reddi, Suvrit Sra, Barnabás Póczos, and Alex Smola. Stochastic frank-wolfe methods for nonconvex optimization. In *2016 54th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, pages 1244–1251. IEEE, 2016.
- [33] Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathematical statistics*, pages 400–407, 1951.
- [34] Sheldon M Ross. *Simulation*. Academic Press, Amsterdam, 2013.
- [35] Soumava Kumar Roy, Zakaria Mhammedi, and Mehrtash Harandi. Geometry aware constrained optimization techniques for deep learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4460–4469, 2018.
- [36] Mark Schmidt, Nicolas Le Roux, and Francis Bach. Minimizing finite sums with the stochastic average gradient. *Mathematical Programming*, 162(1-2):83–112, 2017.
- [37] Shai Shalev-Shwartz and Tong Zhang. Stochastic dual coordinate ascent methods for regularized loss minimization. *Journal of Machine Learning Research*, 14(2), 2013.
- [38] Alexander Shapiro, Darinka Dentcheva, and Andrzej Ruszczyński. *Lectures on stochastic programming: modeling and theory*. SIAM, 2021.
- [39] Jiahao Shi and James C Spall. SQP-based Projection SPSA Algorithm for Stochastic Optimization with Inequality Constraints. In *2021 American Control Conference (ACC)*, pages 1244–1249. IEEE, 2021.
- [40] Tyler Summers, Joseph Warrington, Manfred Morari, and John Lygeros. Stochastic optimal power flow based on conditional value at risk and distributional robustness. *International Journal of Electrical Power & Energy Systems*, 72:116–125, 2015.
- [41] Stanislav Uryasev and Panos M Pardalos. *Stochastic optimization: algorithms and applications*, volume 54. Springer Science & Business Media, 2013.
- [42] Maria Vrakopoulou, Johanna L Mathieu, and Göran Andersson. Stochastic optimal power flow with uncertain reserves from demand response. In *2014 47th Hawaii International Conference on System Sciences*, pages 2353–2362. IEEE, 2014.
- [43] Allen J Wood, Bruce F Wollenberg, and Gerald B Sheblé. *Power generation, operation, and control*. John Wiley & Sons, New Jersey, USA, 2013.

- [44] Wenliang Zhong and James Kwok. Fast stochastic alternating direction method of multipliers. In *International conference on machine learning*, pages 46–54. PMLR, 2014.
- [45] 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.
- [46] William T Ziemba and Raymond G Vickson. *Stochastic optimization models in finance*. Academic Press, 2014.