

A Two Stepsize SQP Method for Nonlinear Equality Constrained Stochastic Optimization

Michael J. O'Neill

May 30, 2026

Abstract We develop a Sequential Quadratic Optimization (SQP) algorithm for minimizing a stochastic objective function subject to deterministic equality constraints. The method utilizes two different stepsizes, one which exclusively scales the component of the step corrupted by the variance of the stochastic gradient estimates and a second which scales the entire step. We prove that this stepsize splitting scheme has a worst-case complexity result which improves over the best known result for this class of problems. In terms of approximately satisfying the constraint violation, this complexity result matches that of deterministic SQP methods, up to constant factors, while matching the known optimal rate for stochastic SQP methods to approximately minimize the norm of the gradient of the Lagrangian. We also propose and analyze multiple variants of our algorithm. One of these variants is based upon popular adaptive gradient methods for unconstrained stochastic optimization while another incorporates a safeguarded line search along the constraint violation. Preliminary numerical experiments show competitive performance against state of the art algorithms. In addition, in these experiments, we observe an improved rate of convergence in terms of the constraint violation, as predicted by the theoretical results.

Keywords nonlinear optimization, stochastic optimization, sequential quadratic optimization, worst-case complexity, adaptive stepsizes

1 Introduction

We propose a new algorithm for solving equality constrained optimization problems in which the objective function is the expectation of a stochastic

Department of Statistics and Operations Research, University of North Carolina, Chapel Hill, NC, USA; Email: mikeoneill@unc.edu; Corresponding Author: Michael J. O'Neill

function. Formally, we consider the optimization problem

$$\min_{x \in \mathbb{R}^n} f(x) \quad \text{s.t.} \quad c(x) = 0 \quad \text{with} \quad f(x) = \mathbb{E}[F(x, \omega)], \quad (1)$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $c : \mathbb{R}^n \rightarrow \mathbb{R}^m$, ω is a random variable with associated probability space (Ω, \mathcal{F}, P) , $F : \mathbb{R}^n \times \Omega \rightarrow \mathbb{R}$, and $\mathbb{E}[\cdot]$ denotes the expectation taken with respect to P . Problems of this form arise in numerous applications, including optimal control [7], PDE-constrained optimization [21, 28], and resource allocation [6] as well as modern machine learning applications, such as physics informed neural networks [13, 27], constraining the output labels of deep neural networks [25] and neural network compression via constraints [12].

The method we design is based on Sequential Quadratic Optimization (SQP) methods, a popular class of algorithms that has seen significant interest in recent years for solving stochastic equality constrained optimization problems, beginning with the influential work of [4]. Numerous extensions of this work have been proposed, such as stochastic SQP methods for problems with rank-deficient Jacobians [3], algorithms for problems with nonlinear inequality constraints [15], worst-case complexity analysis for stochastic SQP methods [14], algorithms which incorporate variance reduction [5] or adaptive sampling [2], as well as stochastic SQP methods which utilize an exact augmented Lagrangian as a merit function [23, 24]. At each iteration, these algorithms generate a search direction by solving a quadratic optimization problem defined in terms of a stochastic gradient estimate subject to a linearization of the constraints and then produce a new iterate by moving along this search direction. For stochastic SQP methods, the chosen step length is generally scaled in such a way as to control the variance of the stochastic gradient estimates, in a manner similar to stepsizes for stochastic gradient methods in unconstrained optimization. Our algorithm takes a different approach and directly utilizes the orthogonal step decomposition of SQP methods¹. It is well known in the stochastic SQP literature that the normal component of the step decomposition is independent of the current stochastic gradient estimate. Therefore, it is unnecessary to rescale this component by the stepsize which controls the variance in the stochastic gradient estimates in order to ensure convergence. Using this observation, we propose a method which employs two different stepsizes: one which controls the variance of the stochastic gradient estimates and scales only the tangential component and a second stepsize which scales the entire search direction.

We demonstrate the effectiveness of this stepsize splitting approach by developing a worst-case complexity result for our proposed algorithm. We consider the worst-case complexity in terms of finding a point x which satisfies,

$$\mathbb{E}[\|\nabla f(x) + \nabla c(x)y\|] \leq \epsilon_\ell, \quad \mathbb{E}[\|c(x)\|_1] \leq \epsilon_c, \quad (2)$$

where $y \in \mathbb{R}^m$ is some Lagrange multiplier and ϵ_ℓ and ϵ_c are some small positive tolerances. Prior work has established a variety of complexity results for

¹ Computation of the orthogonal decomposition may be unnecessary in certain cases, see Remark 1 for details.

different SQP algorithms. In [8], convergence rates are established for a class of SQP methods for inequality constrained optimization under the assumption the Kurdyka-Lojasiewicz (KL) property, with the convergence rate dependent on the Lojasiewicz exponent. In another line of work, [18] devised an SQP method based on “ghost penalties” with diminishing stepsizes for inequality constrained optimization. The authors proved convergence results under a variety of scenarios, including when an initial feasible point is known, in which case a fixed stepsize can be employed to obtain a faster convergence rate. Other approaches with complexity results include [10], which adds a quadratic penalty term involving the constraints to a cubic subproblem at each iteration. When only a first-order method is employed and everywhere LICQ holds, the resulting complexity result is $\mathcal{O}(\epsilon_\ell^{-2})$ and $\mathcal{O}(\epsilon_c^{-2})$.

For first-order, equality deterministic constrained optimization, the best known result (without assuming the KL property), for an SQP method is given in [14], which proved a worst-case complexity result of $\mathcal{O}(\epsilon_\ell^{-2})$ and $\mathcal{O}(\epsilon_c^{-1})$ (this result holds deterministically, not just in expectation). This work also proved a result for the stochastic SQP method of [4], which was shown to have a worst-case complexity of $\mathcal{O}(\epsilon_\ell^{-4})$ and $\mathcal{O}(\epsilon_c^{-2})$ when a lower bound on the merit parameter is known a-priori and $\tilde{\mathcal{O}}(\epsilon_\ell^{-4})$ and $\tilde{\mathcal{O}}(\epsilon_c^{-2})$ otherwise, where $\tilde{\mathcal{O}}$ ignores logarithmic factors. In terms of ϵ_ℓ , this result is optimal, due to information theoretic lower bounds for stochastic gradient methods [1]. However, with respect to the constraint violation, it turns out that this result can be improved. We show that the worst-case complexity of the two stepsize stochastic SQP method proposed in this work has a worst-case complexity of $\mathcal{O}(\epsilon_\ell^{-4})$ and $\mathcal{O}(\epsilon_c^{-1})$. That is, in terms of convergence in the constraint violation, this result matches that of a **deterministic** SQP method, modulo the expectation and constant factors. Furthermore, we avoid unnecessary assumptions which were required to derive a complexity result in [14] by not estimating a merit parameter during the course of the algorithm. Previously this parameter was estimated using stochastic gradient information, which may be highly inaccurate on any given iteration and thus required additional assumptions in order to ensure convergence. In addition to these results, a number of other works have also proposed methods with known worst-case complexity results for solving (1), including augmented Lagrangian [20,29] and stochastic SQP methods [23]. A summary of these worst-case complexity results is given in Table 1.

Unfortunately, the complexity result we prove for our initial algorithm requires certain choices of the stepsizes based on the potential difficulty of estimating parameters of the problem (such as Lipschitz constants and a reasonable setting of the merit parameter). To remedy this, we propose a variant of our method which incorporates stepsizes inspired by adaptive gradient methods for unconstrained stochastic optimization [17,22,31]. Specifically, we build upon the methodology commonly known as Adagrad-Norm, which estimates a stepsize using the prior stochastic gradient estimates. We show that we can generate both of the stepsizes used by our algorithm under this framework and derive a worst-case complexity result for this variant of our method of the order $\tilde{\mathcal{O}}(\epsilon_\ell^{-4})$ and $\tilde{\mathcal{O}}(\epsilon_c^{-1})$, without requiring any knowledge of problem

Algorithm	Conditions	Stationarity	Feasibility
SPD [20]	N/A	$\mathcal{O}\left(\epsilon_\ell^{-6}\right)$	$\mathcal{O}\left(\epsilon_c^{-6}\right)$
SPD [20]	x_0 feasible	$\mathcal{O}\left(\epsilon_\ell^{-5}\right)$	$\mathcal{O}\left(\epsilon_c^{-5}\right)$
MLALM [29]	N/A	$\mathcal{O}\left(\epsilon_\ell^{-4}\right)$	$\mathcal{O}\left(\epsilon_c^{-4}\right)$
MLALM [29]	x_0 near feasible	$\mathcal{O}\left(\epsilon_\ell^{-3}\right)$	$\mathcal{O}\left(\epsilon_c^{-3}\right)$
SSQP-AL [23]	N/A	$\mathcal{O}\left(\epsilon_\ell^{-4}\right)$	$\mathcal{O}\left(\epsilon_c^{-4}\right)$
SSQP [14]	τ_{\min} known	$\mathcal{O}\left(\epsilon_\ell^{-4}\right)$	$\mathcal{O}\left(\epsilon_c^{-2}\right)$
SSQP [14]	τ_{\min} unknown	$\tilde{\mathcal{O}}\left(\epsilon_\ell^{-4}\right)$	$\tilde{\mathcal{O}}\left(\epsilon_c^{-2}\right)$
SSQP-AS [2]	N/A	$\mathcal{O}\left(\epsilon_\ell^{-4}\right)$	$\mathcal{O}\left(\epsilon_c^{-2}\right)$
Algorithm 1	non-adaptive	$\mathcal{O}\left(\epsilon_\ell^{-4}\right)$	$\mathcal{O}\left(\epsilon_c^{-1}\right)$
Algorithm 1	adaptive	$\tilde{\mathcal{O}}\left(\epsilon_\ell^{-4}\right)$	$\tilde{\mathcal{O}}\left(\epsilon_c^{-1}\right)$

Table 1: Sample complexity of algorithms for solving (1). Convergence of each algorithm is proven underneath slightly different conditions. All methods except MLALM assume that the Jacobian has full rank at each iteration, while MLALM assumes a certain constraint qualification as well as mean-squared smoothness of the stochastic gradients. SSQP and SSQP-AS also make additional assumptions on the behavior of the merit parameter.

specific constants. In addition, both versions of our algorithm guarantee convergence when the stepsizes to be relaxed to lie in a certain set, from which the actual stepsize can be chosen, as was originally proposed in [4]. In order to choose a stepsize from this set, we propose a safeguarded linesearch in terms of the constraint violation and show how this can be implemented when the safeguarding is done in terms of the adaptive stepsize rule based on Adagrad-Norm. Finally, we provide preliminary numerical experiments for our algorithm and show that it compares favorably with a state of the art methods. These numerical experiments also demonstrate faster convergence in constraint violation when compared with previously proposed stochastic SQP methods, providing confirmation of our theoretical results.

The rest of this work is organized as follows. In Section 2, we formally define and discuss our proposed algorithm and prove some basic properties. We provide a worst-case complexity analysis in Section 3 for two variants of our algorithm. A safeguarded linesearch procedure is developed in Section 4 and numerical experiments are presented in Section 5. We provide concluding remarks in Section 6.

1.1 Notation

We adopt the notation that $\|\cdot\|$ denotes the ℓ_2 -norm for vectors and the vector-induced ℓ_2 -norm for matrices. The set of nonnegative integers is denoted as $\mathbb{N} := \{0, 1, 2, \dots\}$ and we denote the positive real numbers by $\mathbb{R}_{>0}$.

Given $\phi : \mathbb{R} \rightarrow \mathbb{R}$ and $\varphi : \mathbb{R} \rightarrow [0, \infty)$, we write $\phi(\cdot) = \mathcal{O}(\varphi(\cdot))$ to indicate that $|\phi(\cdot)| \leq c\varphi(\cdot)$ for some $c \in (0, \infty)$. Similarly, we write $\phi(\cdot) = \tilde{\mathcal{O}}(\varphi(\cdot))$ to indicate that $|\phi(\cdot)| \leq c\varphi(\cdot)|\log^{\bar{c}}(\cdot)|$ for some $c \in (0, \infty)$ and $\bar{c} \in (0, \infty)$. In this manner, one finds that $\mathcal{O}(\varphi(\cdot)|\log^{\bar{c}}(\cdot)|) \equiv \tilde{\mathcal{O}}(\varphi(\cdot))$ for any $\bar{c} \in (0, \infty)$.

The algorithm that we analyze is iterative, generating in each realization a sequence $\{x_k\}$. We also append the iteration number to other quantities corresponding an iteration, e.g., $f_k := f(x_k)$ for all $k \in \mathbb{N}$.

1.2 Assumptions and Background

Throughout, we require the following assumptions on f and c :

Assumption 1 *The objective function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is continuously differentiable and bounded below by $f_{\text{low}} \in \mathbb{R}$ and the corresponding gradient function $\nabla f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is bounded and Lipschitz continuous with constant $L \in (0, \infty)$. The constraint function $c : \mathbb{R}^n \rightarrow \mathbb{R}^m$ (where $m \leq n$) and the corresponding Jacobian function $J := \nabla c^\top : \mathbb{R}^n \rightarrow \mathbb{R}^{m \times n}$ are bounded, each gradient function $\nabla c_i : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is Lipschitz continuous with constant γ_i for all $i \in \{1, \dots, m\}$, and the singular values of $J \equiv \nabla c^\top$ are bounded below and away from zero.*

Under this assumption both the gradient of f and the constraint violation are bounded in norm by constants. We denote these constants as $\|\nabla f(x)\| \leq \kappa_g$ and $\|c_k\|_1 \leq \kappa_c$. In addition, we denote the Lipschitz constant of the Jacobian as $\Gamma := \sum_{i=1}^m \gamma_i$, where γ_i is the Lipschitz constant of each ∇c_i^\top .

Defining the Lagrangian $\ell : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$ corresponding to (1) by $\ell(x, y) := f(x) + c(x)^\top y$, first-order primal-dual stationarity conditions for (1), which are necessary for optimality under Assumption 1, are given by

$$0 = \begin{bmatrix} \nabla_x \ell(x, y) \\ \nabla_y \ell(x, y) \end{bmatrix} = \begin{bmatrix} \nabla f(x) + \nabla c(x)y \\ c(x) \end{bmatrix}. \quad (3)$$

We note that the complexity measure (2) is simply an approximate version of these optimality conditions.

As stated above, our algorithm generates a search direction at iteration k by solving the following quadratic optimization problem:

$$\min_{p \in \mathbb{R}^n} f_k + g_k^\top p + \frac{1}{2} p^\top H_k p \quad \text{subject to} \quad c_k + J_k p = 0, \quad (4)$$

where g_k is the current stochastic gradient estimate. It is well known that this is equivalent to solving the ‘‘Newton SQP system’’:

$$\begin{bmatrix} H_k & J_k^\top \\ J_k & 0 \end{bmatrix} \begin{bmatrix} p_k \\ y_k \end{bmatrix} = - \begin{bmatrix} g_k \\ c_k \end{bmatrix}. \quad (5)$$

In order to ensure the solution of this sub-problem is unique, we require the following assumption on H_k .

Assumption 2 *The sequence $\{H_k\}$ 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 $k \in \mathbb{N}$, the matrix H_k has the property that $u^T H_k u \geq \zeta \|u\|^2$ for all $u \in \mathbb{R}^n$ such that $J_k u = 0$.*

In order to analyze our algorithm, we utilize the ℓ_1 merit function $\phi : \mathbb{R}^n \times \mathbb{R}_{>0} \rightarrow \mathbb{R}$:

$$\phi(x, \tau) = \tau f(x) + \|c(x)\|_1. \quad (6)$$

In the above equation, τ is the merit parameter which balances between the function value and constraint violation. For the analysis, we also use the following local model of the merit function $l : \mathbb{R}^n \times \mathbb{R}_{>0} \times \mathbb{R}^n \rightarrow \mathbb{R}$ defined as

$$l(x, \tau, d) = \tau(f(x) + \nabla f(x)^T d) + \|c(x) + \nabla c(x)^T d\|_1. \quad (7)$$

In addition, we consider the reduction in the model for a direction $d \in \mathbb{R}^n$ with $c(x) + \nabla c(x)^T d = 0$ which is $\Delta l : \mathbb{R}^n \times \mathbb{R}_{>0} \times \mathbb{R}^n \rightarrow \mathbb{R}$ defined as

$$\begin{aligned} \Delta l(x, \tau, d) &:= l(x, \tau, 0) - l(x, \tau, d) \\ &= -\tau \nabla f(x)^T d + \|c(x)\|_1. \end{aligned} \quad (8)$$

We wish to stress here that unlike previous works, we do not attempt to estimate a good value of τ . We choose to avoid this as previous work relied upon strong assumptions (such as uniformly bounded stochastic gradients [4], sub-Gaussian stochastic gradients [14], or direct assumptions on “good behavior” of the merit parameter [3, 15, 16], which can be implied by the prior assumptions) in order to prove the existence of a lower bound on a stochastically estimated merit parameter sequence. By choosing to relegate the merit function and parameter exclusively to the analysis, we are able to avoid overcomplicating the analysis and adding unnecessary assumptions. We note that other work on complexity of SQP methods has taken a similar approach, notably [18].

2 Algorithm and Basic Properties

Recall that at each iteration, a search direction p_k is computed as the solution of (5). We assume that this step computation is performed in such a way that the orthogonal decomposition

$$p_k = u_k + v_k \text{ where } u_k \in \text{Null}(J_k) \text{ and } v_k \in \text{Range}(J_k^T), \quad (9)$$

is known². One important consequence of this decomposition is that the normal component, v_k , does not depend on the current stochastic gradient estimate g_k . Unlike prior work, we do not directly use p_k as our search direction. Instead, we rescale the tangential component, u_k , in order to generate our search direction d_k as follows,

$$d_k = \beta_k u_k + v_k, \quad (10)$$

where $\beta_k \in \mathbb{R}_{>0}$. Then, we find the next iterate x_{k+1} by setting $x_{k+1} = x_k + \alpha_k d_k$ for some $\alpha_k \in \mathbb{R}_{>0}$. This procedure is formalized in Algorithm 1.

² This is not necessary in certain circumstances, see Remark 1 for details.

Algorithm 1 Generic Two Stepsize Stochastic SQP Algorithm

Require: $x_0 \in \mathbb{R}^n$;
1: **for** $k = 0, 1, \dots$ **do**
2: Compute stochastic gradient g_k .
3: Compute (p_k, y_k) as the solution of (5).
4: Choose $\beta_k \in \mathbb{R}_{>0}$.
5: Set $d_k \leftarrow v_k + \beta_k u_k$, where $v_k \in \text{Range}(J_k^T)$ and $u_k \in \text{Null}(J_k)$ are the orthogonal decomposition of p_k .
6: Choose $\alpha_k \in \mathbb{R}_{>0}$.
7: Set $x_{k+1} \leftarrow x_k + \alpha_k d_k$.
8: **end for**

The choice of β_k is crucial to ensure convergence of our algorithm, as it controls the variance of the stochastic gradient estimates and plays a similar role as the stepsize in stochastic gradient methods. As such, it is natural to consider β_k to be quite small. Indeed, to ensure our complexity result, we set $\beta_k = O(1/\sqrt{K})$, where K is the total number of iterations we intend to perform. On the other hand, α_k does **not** need to control the error in the stochastic gradients and thus may be set independent of K . Thus, v_k , which is the component of d_k that drives the algorithm towards constraint satisfaction, is only scaled by a stepsize which is independent of K . This is the key insight that leads to our improved complexity.

Algorithm 1 is written generically, without specifying how to choose the stepsizes α_k and β_k . We consider two variants for choosing these stepsizes in Section 3 and analyze their behavior. First, in Section 3.1, we consider the case where β_k is defined by a pre-specified sequence and

$$\alpha_k \in [\nu, \nu + \theta\beta_k], \quad (11)$$

where $\nu \in \mathbb{R}_{>0}$ and $\theta \in \mathbb{R}_{>0}$. This case is essentially equivalent to the standard stochastic gradient regime with a pre-specified stepsize sequence (modulo the relaxation of α_k into a range, which was originally suggested for stochastic SQP methods in [4]). For this method, we prove the complexity result foreshadowed in Section 1. However, this result only holds under certain conditions on ν which depend on the Lipschitz constants of the gradient of f and Jacobian of c as well as a good estimate of the merit parameter τ . Unfortunately, it may not be reasonable to estimate these parameters a-priori.

To remedy this, in Section 3.2 we analyze a version of Algorithm 1 which utilizes adaptive stepsizes based on Adagrad-Norm, which is a popular approach for choosing stepsizes in the stochastic gradient literature [17, 22, 31]. In this case, some additional logarithmic factors appear in the final complexity result, but this approach does not require any knowledge of the Lipschitz constants or the merit parameter.

In addition, in both of the cases we analyze in Section 3, α_k may be chosen from a specific range. In Section 4, we describe a safeguarded line search procedure which can be used to determine α_k . We take advantage of the assumption that the constraints are deterministic and design a line search which only relies on the constraint violation and does not include stochastic gradient

information when computing an α_k . In addition, we provide a fully specified algorithm in Algorithm 2 that combines this linesearch procedure with an adaptive lower bound based on the Adagrad-Norm stepsizes developed in Section 3.2.

2.1 Properties of Algorithm 1

First, we restate a basic result from [4].

Lemma 1 ([4, Lemma 2.9]) *There exists $\kappa_v \in \mathbb{R}_{>0}$ such that, for all $k \in \mathbb{N}$, the normal component v_k satisfies $\max\{\|v_k\|, \|v_k\|^2\} \leq \kappa_v \|c_k\|$.*

During the analysis of our algorithm, it is often useful to consider the “true” step computation that would occur if the step was computed using the true gradient, $\nabla f(x_k)$, in place of the stochastic gradient estimate, g_k . Specifically, let $(p_k^{\text{true}}, y_k^{\text{true}})$ be the solution of the linear system:

$$\begin{bmatrix} H_k & J_k^T \\ J_k & 0 \end{bmatrix} \begin{bmatrix} p_k^{\text{true}} \\ y_k^{\text{true}} \end{bmatrix} = - \begin{bmatrix} \nabla f(x_k) \\ c_k \end{bmatrix}. \quad (12)$$

In addition, we define

$$d_k^{\text{true}} = \beta_k u_k^{\text{true}} + v_k, \quad (13)$$

where $p_k^{\text{true}} = u_k^{\text{true}} + v_k$ with $u_k^{\text{true}} \in \text{Null}(J_k)$ (we recall here that v_k is independent of g_k and $\nabla f(x_k)$ and thus is the same v_k as in (9)).

Lemma 2 *Let Assumptions 1 and 2 hold. Then,*

$$\|u_k^{\text{true}}\| \leq \zeta^{-1} \|\nabla f(x_k)\| + \zeta^{-1} \kappa_H \kappa_v \|c_k\| \leq \zeta^{-1} \kappa_g + \zeta^{-1} \kappa_H \kappa_v \kappa_c =: \kappa_u.$$

Proof By the first equation of (12) and the definition of u_k^{true} , we have $(u_k^{\text{true}})^T H_k (u_k^{\text{true}} + v_k) = -\nabla f(x_k)^T u_k^{\text{true}}$. Then, by Assumption 2 and Lemma 1,

$$\begin{aligned} \zeta \|u_k^{\text{true}}\|^2 &\leq (u_k^{\text{true}})^T H_k u_k^{\text{true}} \\ &= -\nabla f(x_k)^T u_k^{\text{true}} - v_k^T H_k u_k^{\text{true}} \\ &\leq \|\nabla f(x_k)\| \|u_k^{\text{true}}\| + \kappa_H \kappa_v \|c_k\| \|u_k^{\text{true}}\|. \end{aligned}$$

Dividing this inequality through by $\|u_k^{\text{true}}\|$ proves the first result. The final result follows by Assumption 1 and Lemma 1. \square

Now we state an important property about the merit parameter τ .

Lemma 3 *Let Assumptions 1 and 2 hold and let $\sigma \in (0, 1)$. Let $\beta_k \leq \kappa_\beta$ hold for all k and define*

$$\tau_{\min} := \frac{1 - \sigma}{\kappa_v (\kappa_\beta \kappa_H \kappa_u + \kappa_g)}. \quad (14)$$

Then,

$$\tau_{\min} (\nabla f(x_k)^T d_k^{\text{true}} + \beta_k (u_k^{\text{true}})^T H_k u_k^{\text{true}}) \leq (1 - \sigma) \|c_k\|_1. \quad (15)$$

Proof By (12) and the definition of u_k^{true} ,

$$\begin{aligned}\nabla f(x_k)^T d_k^{\text{true}} &= \nabla f(x_k)^T (\beta_k u_k^{\text{true}} + v_k) \\ &= -\beta_k (u_k^{\text{true}})^T H_k u_k^{\text{true}} - \beta_k v_k^T H_k u_k^{\text{true}} + \nabla f(x_k)^T v_k.\end{aligned}$$

Thus, by Assumptions 1 and 2, Lemma 1 and Lemma 2,

$$\begin{aligned}\nabla f(x_k)^T d_k^{\text{true}} + \beta_k (u_k^{\text{true}})^T H_k u_k^{\text{true}} &= -\beta_k v_k^T H_k u_k^{\text{true}} + \nabla f(x_k)^T v_k \\ &\leq (\beta_k \kappa_H \|u_k^{\text{true}}\| + \|\nabla f(x_k)\|) \|v_k\| \\ &\leq \kappa_v (\kappa_\beta \kappa_H \kappa_u + \kappa_g) \|c_k\|_1.\end{aligned}$$

Combining this with (14), proves (15). \square

Remark 1 Under the condition that H_k preserves the null space of J_k (i.e. for any $u \in \text{Null}(J_k)$, $H_k u \in \text{Null}(J_k)$), we can sidestep the requirement to compute the orthogonal decomposition of p_k by simply rescaling the matrix H_k by β_k^{-1} and directly use the computed direction as d_k . We note that this additional requirement is necessary when using rescaling in order to prove a result similar to Lemma 3, as otherwise the crossing term $v_k^T H_k u_k^{\text{true}}$ picks up a factor of β_k^{-1} . This, in turn, means that it is not possible to provide a bound on τ_{\min} that is independent of a **lower** bound on β_k , thus causing serious issues in the final complexity result. For the sake of generality, we don't consider this rescaling approach, though when H_k preserves the nullspace of J_k our results still hold, albeit with potentially different constant factors.

In cases where this condition does not hold, one can compute the step decomposition in a variety of ways. For example, one could directly solve the quadratic subproblem using the nullspace method (see [26, Section 16.2]), which directly solves for u_k and v_k . Alternatively, a simple approach is to compute p_k by (5) and then compute the projection of p_k onto $\text{Null}(J_k)$ to find u_k (and subsequently v_k). We take this approach throughout our numerical experiments and observe it incurs only a minor additional cost, see Appendix C for details. In addition, if one allows for inexact computation of v_k (which is outside the scope of the current manuscript), a stepsize decomposition such as the one employed in [3] could be used, where v_k is computed as the Cauchy step of a simple trust region subproblem while u_k is computed by solving a linear system comparable to (5). In such a case, the cost to compute u_k and v_k is comparable to the cost to compute p_k by (5).

A direct consequence of the previous lemma is

$$\Delta l(x_k, \tau_{\min}, d_k^{\text{true}}) \geq \tau_{\min} \beta_k (u_k^{\text{true}})^T H_k u_k^{\text{true}} + \sigma \|c_k\|_1, \quad (16)$$

which will be used to prove the final convergence result. Given this inequality, it should be clear that with an upper bound on Δl , we would expect convergence in the constraint violation. To see the connection between the quantities in (16) and first order stationarity, we prove the following lemma, which shows that the quadratic term can be lower bounded in terms of the gradient of the Lagrangian at x_k for a specific Lagrange multiplier.

Lemma 4 *Let Assumptions 1 and 2 hold. Then,*

$$(u_k^{\text{true}})^T H_k u_k^{\text{true}} \geq \zeta \kappa_H^{-2} \|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|^2 - (1 + 2\kappa_u) \zeta \kappa_v \|c_k\|_1.$$

Proof By Assumption 2,

$$(u_k^{\text{true}})^T H_k u_k^{\text{true}} \geq \zeta \|u_k^{\text{true}}\|^2 \geq \zeta \kappa_H^{-2} \|H_k u_k^{\text{true}}\|^2.$$

Then, by (12) and Lemmas 1 and 2,

$$\begin{aligned} & \|H_k u_k^{\text{true}} + H_k v_k - H_k v_k\|^2 \\ &= \|H_k u_k^{\text{true}} + H_k v_k\|^2 - 2v_k^T H_k H_k (u_k^{\text{true}} + v_k) + \|H_k v_k\|^2 \\ &\geq \|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|^2 - 2v_k^T H_k H_k u_k^{\text{true}} - \|H_k v_k\|^2 \\ &\geq \|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|^2 - (1 + 2\kappa_u) \kappa_H^2 \kappa_v \|c_k\|_1, \end{aligned}$$

which proves the result. \square

Thus, given these results, we can see that the convergence rate in terms of the gradient of the Lagrangian should be directly related to the choice of β_k while convergence in the constraint violation will be largely independent of this stepsize (provided it is chosen to sufficiently control the noise in g_k). This is in contrast to the results in [14], where the norm of the constraint violation is multiplied by β_k and is the root cause of the improvement in the complexity result for the constraint violation that we prove in the sequel.

We finish this subsection with the following generic descent lemma.

Lemma 5 *Let Assumptions 1 and 2 hold. Then, with τ_{\min} defined as in (14),*

$$\begin{aligned} & \phi(x_k + \alpha_k d_k, \tau_{\min}) - \phi(x_k, \tau_{\min}) \\ & \leq -\alpha_k \Delta l(x_k, \tau_{\min}, d_k^{\text{true}}) + \frac{\alpha_k^2 \beta_k^2}{2} (\tau_{\min} L + \Gamma) \|u_k\|^2 \\ & \quad + \frac{\alpha_k^2}{2} (\kappa_v (\tau_{\min} L + \Gamma) + 4) \|c_k\|_1 + \alpha_k \tau_{\min} \nabla f(x_k)^T (d_k - d_k^{\text{true}}). \end{aligned} \quad (17)$$

Proof By L -Lipschitz continuity of $\nabla f(x)$ and Γ -Lipschitz continuity of J_k , we have

$$\begin{aligned} & \phi(x_k + \alpha_k d_k, \tau_{\min}) - \phi(x_k, \tau_{\min}) \\ & \leq \alpha_k \tau_{\min} \nabla f(x_k)^T d_k + \|c_k + \alpha_k J_k d_k\|_1 - \|c_k\|_1 + \frac{\alpha_k^2}{2} (\tau_{\min} L + \Gamma) \|d_k\|^2 \\ & = \alpha_k \tau_{\min} \nabla f(x_k)^T d_k^{\text{true}} + |1 - \alpha_k| \|c_k\|_1 - \|c_k\|_1 + \frac{\alpha_k^2}{2} (\tau_{\min} L + \Gamma) \|d_k\|^2 \\ & \quad + \alpha_k \tau_{\min} \nabla f(x_k)^T (d_k - d_k^{\text{true}}), \end{aligned}$$

where the equality follows from $J_k d_k = -c_k$.

Using the fact that $|1 - \alpha_k| \leq 1 - \alpha_k + 2\alpha_k^2$, we have

$$\phi(x_k + \alpha_k d_k, \tau_{\min}) - \phi(x_k, \tau_{\min})$$

$$\begin{aligned} &\leq \alpha_k \tau_{\min} \nabla f(x_k)^T d_k^{\text{true}} - \alpha_k \|c_k\|_1 + 2\alpha_k^2 \|c_k\|_1 + \frac{\alpha_k^2}{2} (\tau_{\min} L + \Gamma) \|d_k\|^2 \\ &\quad + \alpha_k \tau_{\min} \nabla f(x_k)^T (d_k - d_k^{\text{true}}). \end{aligned}$$

Then, using the orthogonal decomposition $d_k = \beta_k u_k + v_k$, Lemma 1, and (8), we have

$$\begin{aligned} &\phi(x_k + \alpha_k d_k, \tau_{\min}) - \phi(x_k, \tau_{\min}) \\ &\leq \alpha_k \tau_{\min} \nabla f(x_k)^T d_k^{\text{true}} - \alpha_k \|c_k\|_1 + 2\alpha_k^2 \|c_k\|_1 \\ &\quad + \frac{\alpha_k^2}{2} (\tau_{\min} L + \Gamma) (\beta_k^2 \|u_k\|^2 + \|v_k\|^2) + \alpha_k \tau_{\min} \nabla f(x_k)^T (d_k - d_k^{\text{true}}) \\ &\leq -\alpha_k \Delta l(x_k, \tau_{\min}, d_k^{\text{true}}) + \frac{\alpha_k^2 \beta_k^2}{2} (\tau_{\min} L + \Gamma) \|u_k\|^2 \\ &\quad + \frac{\alpha_k^2}{2} (\kappa_v (\tau_{\min} L + \Gamma) + 4) \|c_k\|_1 + \alpha_k \tau_{\min} \nabla f(x_k)^T (d_k - d_k^{\text{true}}), \end{aligned}$$

proving the result. \square

2.2 Stochastic Assumptions and Properties

In order to analyze the convergence of our algorithm, let \mathcal{F}_k denote the natural filtration adapted to Algorithm 1 and let $\mathbb{E}_k[\cdot] = \mathbb{E}[\cdot | \mathcal{F}_k]$. Under these definitions, we have the following assumption on our stochastic gradient estimates, g_k .

Assumption 3 *There exists $M \in \mathbb{R}_{>0}$ such that, for all k , one finds*

$$\mathbb{E}_k[g_k] = \nabla f(x_k) \quad \text{and} \quad \mathbb{E}_k[\|g_k - \nabla f(x_k)\|_2^2] \leq M. \quad (18)$$

This assumption is largely standard in the stochastic gradient literature. We note that relaxing the uniformly bounded variance assumption to an assumption which allows the variance to grow with the norm of the gradient of f (such as in [9]) is no more general under Assumption 1 since $\|\nabla f(x)\| \leq \kappa_g$.

Under Assumption 3, we have the following properties.

Lemma 6 *Let Assumptions 1, 2, and 3 hold. Then, $\mathbb{E}_k[u_k] = u_k^{\text{true}}$, $\mathbb{E}_k[y_k] = y_k^{\text{true}}$,*

$$\mathbb{E}_k[\|u_k - u_k^{\text{true}}\|^2] \leq \zeta^{-2} M,$$

and

$$\mathbb{E}_k[\|u_k\|^2] \leq \|u_k^{\text{true}}\|^2 + \zeta^{-2} M \leq \zeta^{-1} (u_k^{\text{true}})^T H_k u_k^{\text{true}} + \zeta^{-2} M.$$

Proof The first two claims follow directly by the statement of [4, Lemma 3.8]. To prove the third result, let Z_k be an orthogonal basis for the null space of J_k (which, by Assumption 1 is a matrix in $\mathbb{R}^{n \times (n-m)}$) and let $u_k = Z_k w_k$ and $u_k^{\text{true}} = Z_k w_k^{\text{true}}$. Then, by (5), it follows that

$$Z_k w_k = -Z_k (Z_k^T H_k Z_k)^{-1} Z_k^T (g_k + H_k v_k).$$

Similarly,

$$Z_k w_k^{\text{true}} = -Z_k (Z_k^T H_k Z_k)^{-1} Z_k^T (\nabla f(x_k) + H_k v_k),$$

so that

$$u_k - u_k^{\text{true}} = Z_k (Z_k^T H_k Z_k)^{-1} Z_k^T (\nabla f(x_k) - g_k)$$

and thus, by Assumptions 2 and 3,

$$\mathbb{E}_k[\|u_k - u_k^{\text{true}}\|^2] \leq \mathbb{E}_k[\|Z_k (Z_k^T H_k Z_k)^{-1} Z_k^T\|^2 \|\nabla f(x_k) - g_k\|^2] \leq \zeta^{-2} M.$$

The final set of inequalities follows directly by the first and third claim, with final inequality due to Assumption 2. \square

In addition, we have the following property under the stronger assumption that β_k is measurable to \mathcal{F}_k .

Lemma 7 *Let Assumptions 1, 2, and 3 hold and let β_k be measurable to \mathcal{F}_k . Then, $\mathbb{E}_k[d_k] = d_k^{\text{true}}$ and*

$$\mathbb{E}_k[\|d_k - d_k^{\text{true}}\|] \leq \beta_k \zeta^{-1} \sqrt{M}.$$

Proof Since β_k is measurable to \mathcal{F}_k , it follows that

$$\mathbb{E}_k[d_k] = \beta_k \mathbb{E}_k[u_k] + v_k = \beta_k u_k^{\text{true}} + v_k = d_k^{\text{true}}.$$

For the second result, we have that

$$d_k - d_k^{\text{true}} = \beta_k (u_k - u_k^{\text{true}}) = \beta_k Z_k (Z_k^T H_k Z_k)^{-1} Z_k^T (\nabla f(x_k) - g_k).$$

Thus, by Assumptions 2 and 3, as well as Jensen's inequality,

$$\mathbb{E}_k[\|d_k - d_k^{\text{true}}\|] \leq \beta_k \mathbb{E}_k[\|Z_k (Z_k^T H_k Z_k)^{-1} Z_k^T\| \|\nabla f(x_k) - g_k\|] \leq \beta_k \zeta^{-1} \sqrt{M}.$$

\square

3 Convergence Analysis

In this section, we derive our main convergence results for two variants of Algorithm 1, which differ on how α_k and β_k are chosen at each iteration.

3.1 Coverage with Pre-specified Stepsize Sequences

Throughout this subsection, we analyze Algorithm 1 when $\{\beta_k\}$ is a pre-specified sequence and α_k lies in a pre-specified range, i.e.,

$$\{\beta_k\} \subset \mathbb{R}_{>0}, \quad \alpha_k \in [\nu, \nu + \theta\beta_k], \quad \forall k, \quad (19)$$

for some $\nu \in \mathbb{R}_{>0}$ and $\theta \in \mathbb{R}_{>0}$.

Under this stepsize scheme, we prove a preliminary result about the final term that appears in Lemma 5.

Lemma 8 *Let Assumptions 1, 2, and 3 hold. Then,*

$$\mathbb{E}_k[\alpha_k \tau_{\min} \nabla f(x_k)^T (d_k - d_k^{\text{true}})] \leq \beta_k^2 \theta \tau_{\min} \kappa_g \zeta^{-1} \sqrt{M}.$$

Proof Let $\xi_k \in [0, 1]$ be the random variable such that $\alpha_k = \nu + \xi_k \theta \beta_k$. Then, by Lemma 7 and the fact that ν and β_k are measurable to \mathcal{F}_k ,

$$\begin{aligned} \mathbb{E}_k[\alpha_k \tau_{\min} \nabla f(x_k)^T (d_k - d_k^{\text{true}})] &= \mathbb{E}_k[(\nu + \xi_k \theta \beta_k) \tau_{\min} \nabla f(x_k)^T (d_k - d_k^{\text{true}})] \\ &= \mathbb{E}_k[\xi_k \theta \beta_k \tau_{\min} \nabla f(x_k)^T (d_k - d_k^{\text{true}})] \\ &\leq \mathbb{E}_k[\theta \beta_k \tau_{\min} \|\nabla f(x_k)\| \|d_k - d_k^{\text{true}}\|] \\ &\leq \beta_k^2 \theta \tau_{\min} \kappa_g \zeta^{-1} \sqrt{M}. \end{aligned}$$

□

Now, we are ready to derive our first main result.

Theorem 1 *Let Assumptions 1, 2, and 3 hold. Let $\sigma \in (0, 1)$, let $\{\beta_k\} \subset \mathbb{R}_{>0}$ be a pre-specified sequence such that $\beta_k \leq \kappa_\beta$ holds for all k , let $\alpha_k \in [\nu, \nu + \theta\beta_k]$, for some $\theta \in \mathbb{R}_{>0}$, $\nu \in (0, \sigma/(2\kappa_v(\tau_{\min}L + \Gamma) + 4))$, and let τ_{\min} be defined as in Lemma 3. Let*

$$\begin{aligned} \kappa_1 := & \frac{(\nu + \theta\kappa_\beta)^2 (\tau_{\min}L + \Gamma) (\kappa_u^2 + \zeta^{-2}M)}{2} \\ & + \theta(\theta\kappa_c(\kappa_v(\tau_{\min}L + \Gamma) + 4) + \tau_{\min}\kappa_g\zeta^{-1}\sqrt{M}). \end{aligned} \quad (20)$$

Then, for any $K \in \mathbb{N}$,

$$\begin{aligned} & \sum_{k=0}^{K-1} \mathbb{E}[\alpha_k \beta_k \tau_{\min} (u_k^{\text{true}})^T H_k u_k^{\text{true}} + \frac{\alpha_k \sigma}{2} \|c_k\|_1] \\ & \leq \tau_{\min}(f(x_0) - f_{\text{low}}) + \|c_0\|_1 + \kappa_1 \sum_{k=0}^{K-1} \beta_k^2. \end{aligned} \quad (21)$$

Proof Taking the conditional expectation on both sides of (17) and applying the results of Lemma 2, Lemma 7, and Lemma 8 (noting that β_k is measurable to \mathcal{F}_k),

$$\mathbb{E}_k[\phi(x_k + \alpha_k d_k, \tau_{\min})] - \phi(x_k, \tau_{\min})$$

$$\begin{aligned}
&\leq -\mathbb{E}_k[\alpha_k \Delta l(x_k, \tau_{\min}, \nabla f(x_k), d_k^{\text{true}})] + \mathbb{E}_k \left[\frac{\alpha_k^2 \beta_k^2}{2} (\tau_{\min} L + \Gamma) \|u_k\|^2 \right] \\
&\quad + \mathbb{E}_k \left[\frac{\alpha_k^2}{2} (\kappa_v (\tau_{\min} L + \Gamma) + 4) \|c_k\|_1 \right] + \mathbb{E}_k[\alpha_k \tau_{\min} \nabla f(x_k)^T (d_k - d_k^{\text{true}})] \\
&\leq -\alpha_k \Delta l(x_k, \tau_{\min}, \nabla f(x_k), d_k^{\text{true}}) + (\nu^2 + \theta^2 \beta_k^2) (\kappa_v (\tau_{\min} L + \Gamma) + 4) \|c_k\|_1 \\
&\quad + \frac{\alpha_k^2 \beta_k^2}{2} (\tau_{\min} L + \Gamma) (\|u_k^{\text{true}}\|^2 + \zeta^{-2} M) + \theta \beta_k^2 \tau_{\min} \kappa_g \zeta^{-1} \sqrt{M} \\
&\leq -\alpha_k \Delta l(x_k, \tau_{\min}, \nabla f(x_k), d_k^{\text{true}}) + \frac{\alpha_k^2 \beta_k^2}{2} (\tau_{\min} L + \Gamma) (\kappa_u^2 + \zeta^{-2} M) \\
&\quad + \frac{\alpha_k \sigma}{2} \|c_k\|_1 + \beta_k^2 \theta (\theta \kappa_c (\kappa_v (\tau_{\min} L + \Gamma) + 4) + \tau_{\min} \kappa_g \zeta^{-1} \sqrt{M}) \\
&= -\alpha_k \Delta l(x_k, \tau_{\min}, \nabla f(x_k), d_k^{\text{true}}) + \frac{\alpha_k \sigma}{2} \|c_k\|_1 + \beta_k^2 \kappa_1,
\end{aligned}$$

where the final inequality follows by $\nu \leq \alpha_k$.

Now, by (16),

$$\begin{aligned}
&\mathbb{E}_k[\phi(x_k + \alpha_k d_k, \tau_{\min})] - \phi(x_k, \tau_{\min}) \\
&\leq -\alpha_k \Delta l(x_k, \tau_{\min}, \nabla f(x_k), d_k^{\text{true}}) + \frac{\alpha_k \sigma}{2} \|c_k\|_1 + \beta_k^2 \kappa_1 \\
&\leq -\alpha_k (\beta_k \tau_{\min} (u_k^{\text{true}})^T H_k u_k^{\text{true}} + \sigma \|c_k\|_1) + \frac{\alpha_k \sigma}{2} \|c_k\|_1 + \beta_k^2 \kappa_1 \\
&= -\alpha_k \beta_k \tau_{\min} (u_k^{\text{true}})^T H_k u_k^{\text{true}} - \frac{\alpha_k \sigma}{2} \|c_k\|_1 + \beta_k^2 \kappa_1.
\end{aligned}$$

Taking the total expectation of this inequality, rearranging and summing from $k = 0, \dots, K - 1$,

$$\begin{aligned}
&\sum_{k=0}^{K-1} \mathbb{E}[\alpha_k \beta_k \tau_{\min} (u_k^{\text{true}})^T H_k u_k^{\text{true}} + \frac{\alpha_k \sigma}{2} \|c_k\|_1] \\
&\leq \phi(x_0, \tau_{\min}) - \mathbb{E}[\phi(x_K, \tau_{\min})] + \kappa_1 \sum_{k=0}^{K-1} \beta_k^2.
\end{aligned}$$

Due to Assumption 1, we have,

$$-\mathbb{E}[\phi(x_K, \tau_{\min})] = -\mathbb{E}[\tau_{\min} f(x_K) + \|c_K\|_1] \leq -\tau_{\min} f_{\text{low}},$$

so that

$$\begin{aligned}
&\sum_{k=0}^{K-1} \mathbb{E}[\alpha_k \beta_k \tau_{\min} (u_k^{\text{true}})^T H_k u_k^{\text{true}} + \frac{\alpha_k \sigma}{2} \|c_k\|_1] \\
&\leq \phi(x_0, \tau_{\min}) - \tau_{\min} f_{\text{low}} + \kappa_1 \sum_{k=0}^{K-1} \beta_k^2,
\end{aligned}$$

which proves the result. \square

Next, we establish some convergence results under different choices of β_k .

Corollary 1 *Let the assumptions of Theorem 1 hold. Then, if $\beta_k = \beta > 0$ for all $K \in \mathbb{N}$,*

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|c_k\|_1] \leq \frac{2(\tau_{\min}(f(x_0) - f_{\text{low}}) + \|c_0\|_1)}{\nu\sigma K} + \frac{\beta^2 \kappa_1}{\nu\sigma} \xrightarrow{K \rightarrow \infty} \frac{\beta^2 \kappa_1}{\nu\sigma}, \quad (22)$$

where κ_1 is defined in (20), and

$$\begin{aligned} \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|^2] &\leq \frac{\kappa_H^2(\tau_{\min}(f(x_0) - f_{\text{low}}) + \|c_0\|_1)}{\beta\nu\eta\tau_{\min}\zeta K} + \frac{\kappa_H^2 \kappa_1 \beta}{\nu\eta\tau_{\min}\zeta} \\ &+ \frac{2(1 + 2\kappa_u)\kappa_v \kappa_H^2(\tau_{\min}(f(x_0) - f_{\text{low}}) + \|c_0\|_1)}{\nu\sigma K} + \frac{2(1 + 2\kappa_u)\kappa_v \kappa_H^2 \beta^2 \kappa_1}{\nu\sigma} \\ &\xrightarrow{K \rightarrow \infty} \frac{\kappa_H^2 \kappa_1 \beta}{\nu\eta\tau_{\min}\zeta} + \frac{2(1 + 2\kappa_u)\kappa_v \kappa_H^2 \beta^2 \kappa_1}{\nu\sigma}. \end{aligned} \quad (23)$$

If $\sum_{k=0}^{\infty} \beta_k^2 < \infty$, then,

$$\lim_{k \rightarrow \infty} \mathbb{E}[\|c_k\|_1] = 0. \quad (24)$$

If, in addition, $\sum_{k=0}^{\infty} \beta_k = \infty$, then,

$$\liminf_{k \rightarrow \infty} \mathbb{E}[\|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|^2] = 0. \quad (25)$$

Proof By Theorem 1, the definition of β_k , and $\nu \leq \alpha_k$, it follows that

$$\sum_{k=0}^{K-1} \mathbb{E}[\|c_k\|_1] \leq \frac{2(\tau_{\min}(f(x_0) - f_{\text{low}}) + \|c_0\|_1)}{\nu\sigma} + \frac{2K\kappa_1\beta^2}{\nu\sigma} \quad (26)$$

Dividing both sides of this inequality by K yields the first result.

Now, by Theorem 1 and Lemma 4 as well as $\alpha_k \geq \nu$, we have

$$\begin{aligned} &\sum_{k=0}^{K-1} \mathbb{E}[\nu\beta_k\tau_{\min}\zeta\kappa_H^{-2}\|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|^2 + \frac{\nu\sigma}{2}\|c_k\|_1 \\ &\quad - \nu\tau_{\min}\beta_k(1 + 2\kappa_u)\zeta\kappa_v\|c_k\|_1] \\ &\leq \sum_{k=0}^{K-1} \mathbb{E}[\alpha_k\beta_k(u_k^{\text{true}})^T H_k(u_k^{\text{true}}) + \frac{\alpha_k\sigma}{2}\|c_k\|_1] \\ &\leq \tau_{\min}(f(x_0) - f_{\text{low}}) + \|c_0\|_1 + \kappa_1 \sum_{k=0}^{K-1} \beta_k^2. \end{aligned} \quad (27)$$

Rearranging this inequality and using $\beta_k = \beta$,

$$\begin{aligned} \sum_{k=0}^{K-1} \mathbb{E}[\|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|^2] &\leq \frac{\kappa_H^2(\tau_{\min}(f(x_0) - f_{\text{low}}) + \|c_0\|_1)}{\beta\nu\eta\tau_{\min}\zeta} + \frac{\kappa_H^2 \kappa_1 K \beta}{\nu\eta\tau_{\min}\zeta} \\ &\quad + (1 + 2\kappa_u)\kappa_v \kappa_H^2 \sum_{k=0}^{K-1} \mathbb{E}[\|c_k\|_1]. \end{aligned}$$

Dividing through by K and applying (22), it follows that

$$\begin{aligned} \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|^2] &\leq \frac{\kappa_H^2(\tau_{\min}(f(x_0) - f_{\text{low}}) + \|c_0\|_1)}{\beta\nu\eta\tau_{\min}\zeta K} \\ &\quad + \frac{2(1 + 2\kappa_u)\kappa_v\kappa_H^2(\tau_{\min}(f(x_0) - f_{\text{low}}) + \|c_0\|_1)}{\nu\sigma K} \\ &\quad + \frac{\kappa_H^2\kappa_1\beta}{\nu\eta\tau_{\min}\zeta} + \frac{2(1 + 2\kappa_u)\kappa_v\kappa_H^2\beta^2\kappa_1}{\nu\sigma}. \end{aligned}$$

which proves the second result.

By Theorem 1, and $\nu \leq \alpha_k$, it follows that

$$\sum_{k=0}^{\infty} \mathbb{E}[\|c_k\|_1] \leq \frac{2(\tau_{\min}(f(x_0) - f_{\text{low}}) + \|c_0\|_1)}{\nu\sigma} + \frac{2\kappa_1 \sum_{k=0}^{\infty} \beta_k^2}{\nu\sigma}. \quad (28)$$

By assumption, $\sum_{k=0}^{\infty} \beta_k^2 < \infty$, so the series converges, which directly implies the desired result.

For the final result, by (27) and $\beta_k \leq \kappa_\beta$ for all k , we have

$$\begin{aligned} &\sum_{k=0}^{\infty} \mathbb{E}[\nu\beta_k\tau_{\min}\zeta\kappa_H^{-2}\|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|^2] \\ &\leq \tau_{\min}(f(x_0) - f_{\text{low}}) + \|c_0\|_1 + \kappa_1 \sum_{k=0}^{\infty} \beta_k^2 + \nu\tau_{\min}(1 + 2\kappa_u)\zeta\kappa_v\kappa_\beta \sum_{k=0}^{\infty} \mathbb{E}[\|c_k\|_1]. \end{aligned}$$

By (28) and $\sum_{k=0}^{\infty} \beta_k^2 < \infty$, it follows that

$$\sum_{k=0}^{\infty} \mathbb{E}[\beta_k\|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|^2] = 0.$$

Recalling that $\sum_{k=0}^{\infty} \beta_k = \infty$, the desired result follows directly. \square

When comparing these results to those of [4, Corollary 3.14], we obtain the same radius of convergence for fixed β with respect to the gradient of the Lagrangian, which is proportional to β . In addition, we establish the same limit inferior result for this stationarity measure with a decaying β_k sequence. However, our results differ significantly with respect to the constraint violation. In particular, we establish a tighter radius of convergence, proportional to β^2 , for fixed β as well as convergence in the limit when β_k is a decaying sequence. These properties suggest superior convergence in the constraint violation, due to the two stepsize scheme we employ.

Now, we present our main complexity result of this section.

Corollary 2 *For any $K \in \mathbb{N}_{>0}$, let $\beta_k := \beta = \eta/\sqrt{K}$ for all $k \in [0, K - 1]$ where $\eta \in \mathbb{R}_{>0}$, let κ_1 be defined in (20) and let*

$$\kappa_2 := \tau_{\min}(f(x_0) - f_{\text{low}}) + \|c_0\|_1 + \eta^2\kappa_1. \quad (29)$$

Then, under the conditions of Theorem 1, we have

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|c_k\|_1] \leq \frac{2\kappa_2}{\nu\sigma K}, \quad (30)$$

and

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|^2] \leq \frac{\kappa_H^2 \kappa_2}{\tau_{\min} \zeta \nu \eta \sqrt{K}} + \frac{2\zeta(1+2\kappa_u)\kappa_v \kappa_H^2 \kappa_2}{\nu\sigma K}. \quad (31)$$

Finally, with probability at least $1 - \delta$,

$$\begin{aligned} \min_{k \in [0, K-1]} \tau_{\min} \zeta \kappa_H^{-2} \|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|^2 + \frac{\sigma \sqrt{K}}{2\eta} \|c_k\|_1 \\ \leq \frac{\kappa_2}{\nu \eta \delta \sqrt{K}} + \frac{2(1+2\kappa_u)\zeta \tau_{\min} \kappa_v \kappa_2}{\sigma \delta K}. \end{aligned} \quad (32)$$

Proof The first two results follow directly from Corollary 1 with this specific choice of β .

To prove the final result, by (27),

$$\begin{aligned} \sum_{k=0}^{K-1} \mathbb{E}[\nu \beta_k \tau_{\min} \zeta \kappa_H^{-2} \|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|^2 + \frac{\nu\sigma}{2} \|c_k\|_1] \\ \leq \tau_{\min}(f(x_0) - f_{\text{low}}) + \|c_0\|_1 + \kappa_1 \sum_{k=0}^{K-1} \beta_k^2 \\ + \sum_{k=0}^{K-1} \mathbb{E}[\nu \beta_k \tau_{\min}(1+2\kappa_u)\zeta \kappa_v \|c_k\|_1]. \end{aligned}$$

Applying the definition of β , multiplying through by $\frac{1}{\nu\beta K}$, and using (30),

$$\begin{aligned} \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\tau_{\min} \zeta \kappa_H^{-2} \|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|^2 + \frac{\sigma \sqrt{K}}{2\eta} \|c_k\|_1] \\ \leq \frac{\kappa_2}{\nu \eta \sqrt{K}} + \frac{2(1+2\kappa_u)\zeta \tau_{\min} \kappa_v \kappa_2}{\sigma K} \end{aligned}$$

and thus

$$\begin{aligned} \min_{k \in [0, K-1]} \mathbb{E}[\tau_{\min} \zeta \kappa_H^{-2} \|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|^2 + \frac{\sigma \sqrt{K}}{2\eta} \|c_k\|_1] \\ \leq \frac{\kappa_2}{\nu \sqrt{K}} + \frac{2(1+2\kappa_u)\zeta \tau_{\min} \kappa_v \kappa_2}{\nu\sigma K}. \end{aligned}$$

Applying Markov's inequality and Jensen's inequality, it follows that with probability at least $1 - \delta$ that

$$\begin{aligned} \min_{k \in [0, K-1]} \tau_{\min} \zeta \kappa_H^{-2} \|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|^2 + \frac{\sigma \sqrt{K}}{2\eta} \|c_k\|_1 \\ \leq \frac{\kappa_2}{\nu \delta \sqrt{K}} + \frac{2(1 + 2\kappa_u) \zeta \tau_{\min} \kappa_v \kappa_2}{\sigma \delta K}, \end{aligned}$$

which proves the final result. \square

From the result of Corollary 2, we can easily derive our worst-case complexity results, as promised in Section 1. It should be clear that in terms of the constraint violation, by (30), the maximum number of iterations until $\mathbb{E}[\|c_k\|_1]$ falls below ϵ_c is at most $\mathcal{O}(\epsilon_c^{-1})$. Similarly, by Jensen's inequality and (31), the maximum number of iterations until $\mathbb{E}[\|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|] \leq \epsilon_\ell$ is $\mathcal{O}(\epsilon_\ell^{-4})$. Finally, if one is interested in a combined result, we obtain the same $\mathcal{O}(K^{-1/2})$ convergence rate as [14], however, our convergence is in terms of a much stronger measure with respect to the constraint violation $\|c_k\|_1$, which is scaled by an additional factor of \sqrt{K} . Thus, we expect much faster convergence with respect to the constraint violation than the algorithm in [14] without harming the convergence rate in terms of the gradient of the Lagrangian.

3.2 Convergence with Adaptive Stepsizes

Now, we analyze the case where β_k and α_k are set adaptively, in a manner inspired by Adagrad-Norm [31]. Specifically, at each iteration k , let

$$b_k^2 = b_{k-1}^2 + \|u_k\|^2, \quad q_k^2 = q_{k-1}^2 + \|c_k\|_1, \quad (33)$$

and

$$\beta_k = \frac{\eta}{b_k}, \quad \alpha_k \in \left[\frac{\nu}{q_k}, \frac{\nu}{q_k} + \theta \min \left\{ \frac{1}{b_k}, \frac{1}{q_k} \right\} \right], \quad (34)$$

for some constants $\eta > 0$ and $\nu > 0$. We note here that the additional term at the upper end of the range for α_k is due to our adaptive setting of β_k using b_k , which is sufficient to control the stochasticity in g_k , but may be insufficient to control second order terms involving the constraint violation. We remedy this situation via the inclusion of the θ/q_k term. In addition, we remark that q_k can be set in many different ways, such as using $\|v_k\|^2$ or $\|c_k\|_2$ in place of $\|c_k\|_1$. In principle, one simply needs to choose quantities which are upper bounds on $\|v_k\|^2$ that are computable during the course of the algorithm, which is a relatively flexible condition given Lemma 1 and Assumption 1. These other strategies may lead to longer stepsizes, which could have important practical implications, however, we choose to use $\|c_k\|_1$ as it obtains the best constant factors in the convergence analysis among the relevant choices. Indeed, one can actually take the minimum over multiple quantities and the subsequent

analysis follows with only minor changes; in all of our numerical experiments, we use $q_k^2 = q_{k-1}^k + \min\{\|c_k\|_1, \|v_k\|, \|v_k\|^2\}$ which significantly improves the numerical results. We wish to highlight in the case where $c_k = 0$ for all k and $\theta = 0$, then α_k remains fixed at ν/q_{-1} for all k . In addition, if $H_k = I$, then $u_k = P_k g_k$, where P_k is the orthogonal projection matrix onto $\text{Null}(J_k)$, so that the step d_k is equivalent to the one generated by applying Adagrad-Norm [31] directly to the projected gradient mapping.

Throughout this section, since β_k is dependent on g_k , we redefine d_k^{true} as

$$d_k^{\text{true}} := v_k + \beta_{k-1} u_k^{\text{true}}, \quad (35)$$

so that it remains measurable to \mathcal{F}_k . We note that under this re-definition, the results of Lemmas 3 and 5 still hold.

Our subsequent analysis relies on the following lemma, which we give without proof as it is a well-known result in the adaptive gradient literature (see for example, [30, Lemma 10]).

Lemma 9 *Let $\{a_i\}_{i=0}^\infty$ be a series of non-negative real numbers with $a_0 \in \mathbb{R}_{>0}$. Then,*

$$\sum_{k=1}^T \frac{a_k}{\sum_{i=0}^k a_i} \leq \log \left(\sum_{k=0}^T a_k \right) - \log(a_0) \quad (36)$$

In order to prove convergence of our algorithm, the key issue posed by the adaptive stepsizes is the final term in (17), which requires a more detailed analysis than in Lemma 8 as β_k is no longer measurable to \mathcal{F}_k and d_k^{true} has been redefined in (35). We give a bound on this term in the following lemma.

Lemma 10 *Let Assumptions 1, 2, and 3 hold and let*

$$\kappa_3 := \kappa_u^2 + \zeta^{-2} M \quad (37)$$

and

$$\kappa_4 := \max \left\{ \zeta^{-1} \kappa_H^2, \beta_{-1} (1 + 2\kappa_u) \kappa_H^2 \kappa_v \tau_{\min} / \sigma \right\} \quad (38)$$

Then,

$$\begin{aligned} \mathbb{E}_k \left[\alpha_k \nabla f(x_k)^T (d_k - d_k^{\text{true}}) \right] &\leq \mathbb{E}_k \left[\frac{\alpha_k \beta_{k-1}}{2} (u_k^{\text{true}})^T H_k u_k^{\text{true}} + \frac{\alpha_k \sigma}{2\tau_{\min}} \|c_k\|_1 \right. \\ &\quad \left. + \left(\frac{3\eta^2 \kappa_3 \kappa_4 (\nu + \theta)^2}{2q_{-1} b_{-1}} + \frac{3\kappa_4 \theta^2 (\eta^2 + \beta_{-1}^2 \kappa_3)}{2\eta\nu} + \frac{3\zeta^{-2} M \kappa_4 \theta^2 \beta_{-1}}{2q_{-1}} \right) \frac{\|u_k\|^2}{b_k^2} \right]. \end{aligned} \quad (39)$$

Proof By the definition of d_k^{true} , we have

$$\mathbb{E}_k \left[\alpha_k \nabla f(x_k)^T (d_k - d_k^{\text{true}}) \right] = \mathbb{E}_k \left[\alpha_k \nabla f(x_k)^T (\beta_k u_k - \beta_{k-1} u_k^{\text{true}}) \right].$$

Let $\xi_k \in [0, 1]$ be the random variable such that $\alpha_k = \frac{\nu}{q_k} + \xi_k \min\{\frac{\theta}{b_k}, \frac{\theta}{q_k}\}$. Then, by Lemma 7 and the fact that β_{k-1} , ν , and q_k are measurable to \mathcal{F}_k ,

$$\mathbb{E}_k \left[\alpha_k \nabla f(x_k)^T (\beta_k u_k - \beta_{k-1} u_k^{\text{true}}) \right]$$

$$\begin{aligned}
&= \mathbb{E}_k \left[\left(\frac{\nu}{q_k} + \xi_k \min \left\{ \frac{\theta}{b_k}, \frac{\theta}{q_k} \right\} \right) \nabla f(x_k)^T (\beta_k u_k - \beta_{k-1} u_k^{\text{true}}) \right] \\
&= \mathbb{E}_k \left[\frac{\nu}{q_k} (\beta_k - \beta_{k-1}) \nabla f(x_k)^T u_k \right. \\
&\quad \left. + \xi_k \min \left\{ \frac{\theta}{b_k}, \frac{\theta}{q_k} \right\} \nabla f(x_k)^T (\beta_k u_k - \beta_{k-1} u_k^{\text{true}}) \right] \\
&= \mathbb{E}_k \left[\frac{\nu}{q_k} (\beta_k - \beta_{k-1}) (\nabla f(x_k) + J_k^T y_k^{\text{true}})^T u_k \right. \\
&\quad \left. + \xi_k \min \left\{ \frac{\theta}{b_k}, \frac{\theta}{q_k} \right\} (\nabla f(x_k) + J_k^T y_k^{\text{true}})^T (\beta_k u_k - \beta_{k-1} u_k^{\text{true}}) \right] \\
&= \mathbb{E}_k \left[\frac{\nu}{q_k} (\beta_k - \beta_{k-1}) (\nabla f(x_k) + J_k^T y_k^{\text{true}})^T u_k \right. \\
&\quad \left. + \xi_k \min \left\{ \frac{\theta}{b_k}, \frac{\theta}{q_k} \right\} \beta_k (\nabla f(x_k) + J_k^T y_k^{\text{true}})^T (u_k - u_k^{\text{true}}) \right. \\
&\quad \left. + \xi_k \min \left\{ \frac{\theta}{b_k}, \frac{\theta}{q_k} \right\} (\beta_k - \beta_{k-1}) (\nabla f(x_k) + J_k^T y_k^{\text{true}})^T u_k^{\text{true}} \right] \\
&= \mathbb{E}_k \left[\left(\frac{\nu}{q_k} + \xi_k \min \left\{ \frac{\theta}{b_k}, \frac{\theta}{q_k} \right\} \right) (\beta_k - \beta_{k-1}) (\nabla f(x_k) + J_k^T y_k^{\text{true}})^T u_k \right. \\
&\quad \left. + \xi_k \min \left\{ \frac{\theta}{b_k}, \frac{\theta}{q_k} \right\} \beta_k (\nabla f(x_k) + J_k^T y_k^{\text{true}})^T (u_k - u_k^{\text{true}}) \right. \\
&\quad \left. + \xi_k \min \left\{ \frac{\theta}{b_k}, \frac{\theta}{q_k} \right\} (\beta_k - \beta_{k-1}) (\nabla f(x_k) + J_k^T y_k^{\text{true}})^T (u_k^{\text{true}} - u_k) \right] \\
&\leq \mathbb{E}_k \left[\frac{\nu + \theta}{q_k} |\beta_k - \beta_{k-1}| \|\nabla f(x_k) + J_k^T y_k^{\text{true}}\| \|u_k\| \right. \\
&\quad \left. + \min \left\{ \frac{\theta}{b_k}, \frac{\theta}{q_k} \right\} \beta_k \|\nabla f(x_k) + J_k^T y_k^{\text{true}}\| \|u_k - u_k^{\text{true}}\| \right. \\
&\quad \left. + \frac{\theta}{q_k} |\beta_k - \beta_{k-1}| \|\nabla f(x_k) + J_k^T y_k^{\text{true}}\| \|u_k^{\text{true}} - u_k\| \right], \tag{40}
\end{aligned}$$

where the third equality follows by $u_k, u_k^{\text{true}} \in \text{Null}(J_k)$ and the inequality by the Cauchy-Schwarz inequality and $\xi_k \leq 1$.

Now, we focus on the first term in (40),

$$|\beta_k - \beta_{k-1}| = \frac{\eta}{b_{k-1}} - \frac{\eta}{b_k} = \frac{\eta \|u_k\|^2}{b_{k-1} b_k (b_k + b_{k-1})} \leq \frac{\eta \|u_k\|}{b_{k-1} b_k}, \tag{41}$$

where the inequality follows by $\|u_k\| \leq b_k$. Therefore, applying Young's inequality, we have, for any $\lambda_1 > 0$ measurable to \mathcal{F}_k ,

$$\begin{aligned}
&\mathbb{E}_k \left[\frac{\nu + \theta}{q_k} |\beta_{k-1} - \beta_k| \|\nabla f(x_k) + J_k^T y_k^{\text{true}}\| \|u_k\| \right] \\
&\leq \eta \mathbb{E}_k \left[\frac{(\nu + \theta) \|\nabla f(x_k) + J_k^T y_k^{\text{true}}\| \|u_k\|^2}{q_k b_{k-1} b_k} \right]
\end{aligned}$$

$$\begin{aligned}
&\leq \frac{\eta \|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|^2}{2b_{k-1}q_k\lambda_1} \mathbb{E}_k[\|u_k\|^2] + \mathbb{E}_k \left[\frac{\eta(\nu + \theta)^2 \lambda_1 \|u_k\|^2}{2q_k b_{k-1} b_k^2} \right] \\
&\leq \frac{\eta \kappa_3 \beta_{k-1} \|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|^2}{2q_k \lambda_1} + \mathbb{E}_k \left[\frac{\eta(\nu + \theta)^2 \lambda_1 \|u_k\|^2}{2q_k b_{k-1} b_k^2} \right] \quad (42)
\end{aligned}$$

where the final inequality follows by Assumption 2 as well as the results of Lemma 2 and Lemma 6.

Now, for the second term in (40), by Young's inequality, for any $\lambda_2 > 0$ measurable to \mathcal{F}_k ,

$$\begin{aligned}
&\mathbb{E}_k \left[\beta_k \min \left\{ \frac{\theta}{b_k}, \frac{\theta}{q_k} \right\} \|\nabla f(x_k) + J_k^T y_k^{\text{true}}\| \|u_k - u_k^{\text{true}}\| \right] \\
&\leq \frac{1}{2q_k b_k \lambda_2} \|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|^2 + \mathbb{E}_k \left[\frac{\lambda_2 \theta^2 \beta_k^2}{2} \|u_k - u_k^{\text{true}}\|^2 \right]. \quad (43)
\end{aligned}$$

Working with the last term in this inequality, since q_k and β_{k-1} are measurable to \mathcal{F}_k , by Lemma 6,

$$\begin{aligned}
&\mathbb{E}_k \left[\frac{\lambda_2 \theta^2 \beta_k^2}{2} \|u_k - u_k^{\text{true}}\|^2 \right] \\
&= \frac{\lambda_2 \theta^2}{2} \mathbb{E}_k [\beta_k^2 (\|u_k\|^2 + \|u_k^{\text{true}}\|^2 - 2u_k^T u_k^{\text{true}})] \\
&= \frac{\lambda_2 \theta^2}{2} \mathbb{E}_k [\beta_k^2 (\|u_k\|^2 + \|u_k^{\text{true}}\|^2 - 2u_k^T u_k^{\text{true}}) + \beta_{k-1}^2 (2u_k^T u_k^{\text{true}} - 2\|u_k^{\text{true}}\|^2)] \\
&\leq \frac{\lambda_2 \theta^2}{2} \mathbb{E}_k [\beta_k^2 \|u_k\|^2 + 2|\beta_k^2 - \beta_{k-1}^2| \|u_k\| \|u_k^{\text{true}}\| - \beta_{k-1}^2 \|u_k^{\text{true}}\|^2] \\
&= \frac{\lambda_2 \theta^2}{2} \mathbb{E}_k \left[\beta_k^2 \|u_k\|^2 + 2\eta^2 \frac{\|u_k\|^2}{b_{k-1}^2 b_k^2} \|u_k\| \|u_k^{\text{true}}\| - \beta_{k-1}^2 \|u_k^{\text{true}}\|^2 \right] \\
&\leq \frac{\lambda_2 \theta^2}{2} \mathbb{E}_k \left[\beta_k^2 \|u_k\|^2 + 2\eta^2 \frac{\|u_k\|^2}{b_{k-1}^2 b_k} \|u_k^{\text{true}}\| - \beta_{k-1}^2 \|u_k^{\text{true}}\|^2 \right] \quad (44)
\end{aligned}$$

where the first inequality follows by $\beta_k \leq \beta_{k-1}$ and the second inequality follows by $\|u_k\| \leq b_k$. Dealing with the second term in this inequality, again applying Young's inequality and using $\beta_{k-1} = \eta/b_{k-1}$, by Lemmas 2 and 6 as well as Assumption 2,

$$\begin{aligned}
&\frac{\lambda_2 \theta^2}{2} \mathbb{E}_k \left[\frac{2\beta_{k-1}^2 \|u_k\|^2}{b_k} \|u_k^{\text{true}}\| \right] \\
&\leq \frac{\lambda_2 \theta^2 \eta^2}{2} \mathbb{E}_k \left[\frac{\beta_{k-1}^2 \|u_k\|^2 \|u_k^{\text{true}}\|^2}{\kappa_3} + \frac{\beta_{k-1}^2 \kappa_3 \|u_k\|^2}{b_k^2} \right] \\
&\leq \frac{\lambda_2 \theta^2}{2} \mathbb{E}_k \left[\frac{\beta_{k-1}^2 (\kappa_u^2 + \zeta^{-2} M) \|u_k^{\text{true}}\|^2}{\kappa_3} + \frac{\beta_{k-1}^2 \kappa_3 \|u_k\|^2}{b_k^2} \right] \\
&= \frac{\lambda_2 \theta^2}{2} \mathbb{E}_k \left[\beta_{k-1}^2 \|u_k^{\text{true}}\|^2 + \frac{\beta_{k-1}^2 \kappa_3 \|u_k\|^2}{b_k^2} \right],
\end{aligned}$$

so that the first term cancels with the last in (44).

Now, for the final term in (40), by (41), applying Young's inequality for some $\lambda_3 > 0$ that is measurable to \mathcal{F}_k , by Lemma 6,

$$\begin{aligned} & \mathbb{E}_k \left[\frac{\theta}{q_k} |\beta_k - \beta_{k-1}| \|\nabla f(x_k) + J_k^T y_k^{\text{true}}\| \|u_k^{\text{true}} - u_k\| \right] \\ & \leq \mathbb{E}_k \left[\frac{\theta}{q_k} \frac{\beta_{k-1} \|u_k\|}{b_k} \|\nabla f(x_k) + J_k^T y_k^{\text{true}}\| \|u_k^{\text{true}} - u_k\| \right] \\ & \leq \mathbb{E}_k \left[\frac{\beta_{k-1}}{2\lambda_3 q_k} \|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|^2 \|u_k^{\text{true}} - u_k\|^2 + \frac{\theta^2 \lambda_3 \beta_{k-1} \|u_k\|^2}{2q_k b_k^2} \right] \\ & \leq \frac{\zeta^{-2} M \beta_{k-1}}{2\lambda_3 q_k} \|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|^2 + \mathbb{E}_k \left[\frac{\theta^2 \lambda_3 \beta_{k-1} \|u_k\|^2}{2q_k b_k^2} \right] \end{aligned}$$

Therefore, combining (40), (42), (43), and (44) we have

$$\begin{aligned} & \mathbb{E}_k [\alpha_k \nabla f(x_k)^T (\beta_k u_k - \beta_{k-1} u_k^{\text{true}})] \\ & \leq \mathbb{E}_k \left[\left(\frac{\eta \kappa_3 \beta_{k-1}}{2q_k \lambda_1} + \frac{1}{2q_k b_k \lambda_2} + \frac{\zeta^{-2} M \beta_{k-1}}{2q_k \lambda_3} \right) \|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|^2 \right. \\ & \quad \left. + \left(\frac{\lambda_1 \eta (\nu + \theta)^2}{2q_k b_{k-1}} + \frac{\lambda_2 \theta^2 (\eta^2 + \beta_{k-1}^2 \kappa_3)}{2} + \frac{\lambda_3 \theta^2 \beta_{k-1}}{2q_k} \right) \frac{\|u_k^{\text{true}}\|^2}{b_k^2} \right] \end{aligned}$$

Applying Lemma 4,

$$\begin{aligned} & \mathbb{E}_k [\alpha_k \nabla f(x_k)^T (\beta_k u_k - \beta_{k-1} u_k^{\text{true}})] \\ & \leq \mathbb{E}_k \left[\left(\frac{\eta \kappa_3 \beta_{k-1}}{2q_k \lambda_1} + \frac{1}{2q_k b_k \lambda_2} + \frac{\zeta^{-2} M \beta_{k-1}}{2q_k \lambda_3} \right) (\zeta^{-1} \kappa_H^2 (u_k^{\text{true}})^T H_k u_k^{\text{true}} + (1 + 2\kappa_u) \kappa_H^2 \kappa_v \|c_k\|) \right. \\ & \quad \left. + \left(\frac{\lambda_1 \eta (\nu + \theta)^2}{2q_k b_{k-1}} + \frac{\lambda_2 \theta^2 (\eta^2 + \beta_{k-1}^2 \kappa_3)}{2} + \frac{\lambda_3 \theta^2 \beta_{k-1}}{2q_k} \right) \frac{\|u_k\|^2}{b_k^2} \right]. \end{aligned}$$

Choosing $\lambda_1 = \frac{3\eta \kappa_3 \kappa_4}{\nu}$, $\lambda_2 = \frac{3\kappa_4}{\nu \eta}$, and $\lambda_3 = \frac{3\zeta^{-2} M \kappa_4}{\nu}$ and using $\nu/q_k \leq \alpha_k$, $q_k \geq q_{-1}$, and $b_{k-1} \geq b_{-1}$ proves the result. \square

Now, we are prepared to present the first main result of this subsection.

Theorem 2 *Let Assumptions 1, 2, and 3 hold. Let*

$$\kappa_5 := \frac{(\nu + \theta)^2 (\kappa_v (\tau_{\min} L + \Gamma) + 4)}{2}. \quad (45)$$

$$\begin{aligned} \kappa_6 := & \frac{\eta^2 (\nu + \theta)^2 (\tau_{\min} L + \Gamma)}{2q_{-1}^2} + \frac{3\eta^2 \kappa_3 \kappa_4 (\nu + \theta)^2}{2q_{-1} b_{-1}} \\ & + \frac{3\kappa_4 \theta^2 (\eta^2 + \beta_{-1}^2 \kappa_3)}{2\eta \nu} + \frac{3\zeta^{-2} M \kappa_4 \theta^2 \beta_{-1}}{2q_{-1} b_{-1}^2}. \end{aligned} \quad (46)$$

Then,

$$\begin{aligned} & \mathbb{E} \left[\sum_{k=0}^{K-1} \frac{\alpha_k \tau_{\min} \beta_{k-1}}{2} (u_k^{\text{true}})^T H_k u_k^{\text{true}} + \frac{\alpha_k \sigma}{2} \|c_k\|_1 \right] \\ & \leq \tau_{\min}(f_{-1} - f_{\min}) + \|c_{-1}\|_1 + \kappa_5 \log(1 + \kappa_c K/q_{-1}^2) \\ & \quad + \kappa_6 \log(1 + (\kappa_u^2 + \zeta^{-2} M)K/b_{-1}^2). \end{aligned} \quad (47)$$

In addition,

$$\begin{aligned} \mathbb{E}[q_{K-1}] & \leq q_{-1} + \frac{2}{\nu\sigma} (\tau_{\min}(f_{-1} - f_{\min}) + \|c_{-1}\|_1 + \kappa_5 \log(1 + \kappa_c K/q_{-1}^2)) \\ & \quad + \kappa_6 \log(1 + (\kappa_u^2 + \zeta^{-2} M)K/b_{-1}^2). \end{aligned}$$

Proof By Lemma 5, we have

$$\begin{aligned} & \mathbb{E}_k[\phi(x_k + \alpha_k d_k, \tau_{\min})] - \phi(x_k, \tau_{\min}) \\ & \leq -\mathbb{E}_k[\alpha_k \Delta l(x_k, \tau_{\min}, d_k^{\text{true}})] + \mathbb{E}_k \left[\frac{\alpha_k^2 \beta_k^2}{2} (\tau_{\min} L + \Gamma) \|u_k\|^2 \right] \\ & \quad + \mathbb{E}_k \left[\frac{\alpha_k^2}{2} (\kappa_v (\tau_{\min} L + \Gamma) + 4) \|c_k\|_1 \right] + \mathbb{E}_k[\alpha_k \tau_{\min} \nabla f(x_k)^T (d_k - d_k^{\text{true}})]. \end{aligned}$$

To prove the result, we need to bound the final three terms. Starting with the first of these, we have that

$$\mathbb{E}_k \left[\frac{\alpha_k^2 \beta_k^2}{2} (\tau_{\min} L + \Gamma) \|u_k\|^2 \right] \leq \frac{\eta^2 (\nu + \theta)^2 (\tau_{\min} L + \Gamma)}{2q_{-1}^2} \mathbb{E}_k \left[\frac{\|u_k\|^2}{b_k^2} \right],$$

where the inequality follows due to the definition of α_k and $q_k \geq q_{-1}$. For the next term,

$$\frac{\alpha_k^2 (\kappa_v (\tau_{\min} L + \Gamma) + 4)}{2} \|c_k\|_1 \leq \frac{(\nu + \theta)^2 (\kappa_v (\tau_{\min} L + \Gamma) + 4)}{2q_k^2} \|c_k\|_1 = \frac{\kappa_5 \|c_k\|_1}{q_k^2}.$$

Now, applying the result of Lemma 10, we have

$$\begin{aligned} & \mathbb{E}_k[\phi(x_k + \alpha_k d_k, \tau_{\min})] - \phi(x_k, \tau_{\min}) \\ & \leq -\mathbb{E}_k[\alpha_k \Delta l(x_k, \tau_{\min}, d_k^{\text{true}})] + \mathbb{E}_k \left[\frac{\alpha_k \tau_{\min} \beta_{k-1}}{2} (u_k^{\text{true}})^T H_k u_k^{\text{true}} \right] \\ & \quad + \mathbb{E}_k \left[\frac{\alpha_k \sigma}{2} \|c_k\|_1 \right] + \frac{\kappa_5 \|c_k\|_1}{q_k^2} + \kappa_6 \mathbb{E}_k \left[\frac{\|u_k\|^2}{b_k^2} \right]. \end{aligned}$$

Then, applying (16) (where we note that under the re-definition of d_k^{true} in (35), β_k is replaced by β_{k-1}), it follows that

$$\begin{aligned} & \mathbb{E}_k[\phi(x_k + \alpha_k d_k, \tau_{\min})] - \phi(x_k, \tau_{\min}) \\ & \leq -\mathbb{E}_k \left[\frac{\alpha_k \tau_{\min} \beta_{k-1}}{2} (u_k^{\text{true}})^T H_k u_k^{\text{true}} \right] - \mathbb{E}_k \left[\frac{\alpha_k \sigma}{2} \|c_k\|_1 \right] \end{aligned}$$

$$+ \frac{\kappa_5 \|c_k\|_1}{q_k^2} + \kappa_6 \mathbb{E}_k \left[\frac{\|u_k\|^2}{b_k^2} \right].$$

Next, taking the total expectation of this inequality and summing for all $k = 0, \dots, K-1$,

$$\begin{aligned} & \mathbb{E} \left[\sum_{k=0}^{K-1} \frac{\alpha_k \tau_{\min} \beta_{k-1}}{2} (u_k^{\text{true}})^T H_k u_k^{\text{true}} + \frac{\alpha_k \sigma}{2} \|c_k\|_1 \right] \\ & \leq \phi(x_{-1}, \tau_{\min}) - \mathbb{E}[\phi(x_K, \tau_{\min})] + \mathbb{E} \left[\kappa_5 \sum_{k=0}^{K-1} \frac{\|c_k\|_1}{q_k^2} \right] + \mathbb{E} \left[\kappa_6 \sum_{k=0}^{K-1} \frac{\|u_k\|^2}{b_k^2} \right]. \end{aligned}$$

By the definition of ϕ and Assumption 1, it follows that

$$\begin{aligned} \phi(x_{-1}, \tau_{\min}) - \mathbb{E}[\phi(x_K, \tau_{\min})] &= \tau_{\min} f_{-1} + \|c_{-1}\|_1 - \mathbb{E}[\tau_{\min} f_K - \|c_K\|_1] \\ &\leq \tau_{\min}(f_{-1} - f_{\min}) + \|c_{-1}\|_1. \end{aligned}$$

Now, applying Lemma 9 twice, by Assumption 1, it follows that

$$\begin{aligned} & \mathbb{E} \left[\sum_{k=0}^{K-1} \frac{\alpha_k \tau_{\min} \beta_{k-1}}{2} (u_k^{\text{true}})^T H_k u_k^{\text{true}} + \frac{\alpha_k \sigma}{2} \|c_k\|_1 \right] \\ & \leq \tau_{\min}(f_{-1} - f_{\min}) + \|c_{-1}\|_1 + \kappa_5 \log(1 + \kappa_c K / q_{-1}^2) \\ & \quad + \kappa_6 \mathbb{E} \left[\log \left(\frac{b_{-1}^2 + \sum_{k=0}^{K-1} \|u_k\|^2}{b_{-1}^2} \right) \right]. \end{aligned}$$

Using Jensen's inequality, the tower rule, and the results of Lemma 2 and Lemma 6,

$$\mathbb{E} \left[\log \left(\frac{b_{-1}^2 + \sum_{k=0}^{K-1} \|u_k\|^2}{b_{-1}^2} \right) \right] \leq \log(1 + (\kappa_u^2 + \zeta^{-2} M) K / b_{-1}^2) \quad (48)$$

and thus,

$$\begin{aligned} & \mathbb{E} \left[\sum_{k=0}^{K-1} \frac{\alpha_k \tau_{\min} \beta_{k-1}}{2} (u_k^{\text{true}})^T H_k u_k^{\text{true}} + \frac{\alpha_k \sigma}{2} \|c_k\|_1 \right] \\ & \leq \tau_{\min}(f_{-1} - f_{\min}) + \|c_{-1}\|_1 + \kappa_5 \log(1 + \kappa_c K / q_{-1}^2) \\ & \quad + \kappa_6 \log(1 + (\kappa_u^2 + \zeta^{-2} M) K / b_{-1}^2), \end{aligned}$$

proving the first result.

To prove the second result, note that

$$q_{K-1} = \frac{q_{-1}^2 + \sum_{k=0}^{K-1} \|c_k\|_1}{q_{K-1}} \leq q_{-1} + \sum_{k=0}^{K-1} \frac{\|c_k\|_1}{q_k} \leq q_{-1} + \frac{1}{\nu} \sum_{k=0}^{K-1} \alpha_k \|c_k\|_1,$$

and therefore, by (47),

$$\begin{aligned} \mathbb{E}[q_{K-1}] &\leq q_{-1} + \frac{2}{\nu\sigma}(\tau_{\min}(f_{-1} - f_{\min}) + \|c_{-1}\|_1 + \kappa_5 \log(1 + \kappa_c K/q_{-1}^2)) \\ &\quad + \kappa_6 \log(1 + (\kappa_u^2 + \zeta^{-2}M)K/b_{-1}^2). \end{aligned}$$

□

Next, we derive the following corollary, from which our complexity results for this subsection will follow directly.

Corollary 3 *Let the assumptions of Theorem 2 hold. Let*

$$\begin{aligned} \kappa_7(K) &:= \tau_{\min}(f_{-1} - f_{\min}) + \|c_{-1}\|_1 + \kappa_5 \log(1 + \kappa_c K/q_{-1}) \\ &\quad + \kappa_6 \log(1 + (\kappa_u^2 + \zeta^{-2}M)K/b_{-1}) \end{aligned} \quad (49)$$

and

$$\kappa_8(K) := \sqrt{b_{-1}^2 + (\kappa_u^2 + \zeta^{-2}M)K}. \quad (50)$$

Then, with probability at least $1 - \delta_1$,

$$\mathbb{E} \left[\frac{1}{K} \sum_{k=0}^{K-1} \|c_k\|_1 \right] \leq \frac{2(\nu\sigma q_{-1} + 2\kappa_7(K))\kappa_7(K)}{\nu^2\sigma^2\delta_1 K}, \quad (51)$$

with probability at least $1 - \delta_2$,

$$\begin{aligned} &\mathbb{E} \left[\frac{1}{K} \sum_{k=0}^{K-1} \|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|^2 \right] \\ &\leq \frac{8\kappa_H^2 \kappa_8(K)(\nu\sigma q_{-1} + 2\kappa_7(K))\kappa_7(K)}{\tau_{\min}\nu^2\eta\zeta\sigma\delta_2^2 K} \\ &\quad + \frac{4(\nu\sigma q_{-1} + 2\kappa_7(K))\kappa_H^2(1 + 2\kappa_u)\kappa_v\kappa_7(K)}{\nu^2\sigma^2\delta_2 K}, \end{aligned} \quad (52)$$

and with probability at least $1 - \delta_3$,

$$\begin{aligned} &\min_{k \in [0, K-1]} \tau_{\min}\zeta\kappa_H^{-2} \|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|^2 + \frac{\sigma\kappa_8(K)}{\eta} \|c_k\|_1 \\ &\leq \frac{54\kappa_8(K)(\nu\sigma q_{-1} + 2\kappa_7(K))\kappa_7(K)}{\nu^2\eta\sigma\delta_3^3 K} \\ &\quad + \frac{18(\nu\sigma q_{-1} + 2\kappa_7(K))\kappa_7(K)\zeta\tau_{\min}(1 + 2\kappa_u)\kappa_v}{\nu^2\sigma^2\delta_3^2 K}. \end{aligned} \quad (53)$$

Proof Note that for all $k \in [0, K-1]$, we have that

$$\alpha_k \geq \frac{\nu}{q_k} \geq \frac{\nu}{q_{K-1}}.$$

Additionally, by Theorem 2 and Markov's inequality, it follows that with probability at least $1 - \delta_1$,

$$\frac{\nu}{q_{K-1}} \geq \frac{\nu^2 \sigma \delta_1}{\nu \sigma q_{-1} + 2\kappa_7(K)}. \quad (54)$$

Therefore, by (47) and Assumption 2, it follows that with probability at least $1 - \delta_1$,

$$\mathbb{E} \left[\sum_{k=0}^{K-1} \frac{\nu^2 \sigma^2 \delta_1}{2(\nu \sigma q_{-1} + 2\kappa_7(K))} \|c_k\|_1 \right] \leq \mathbb{E} \left[\sum_{k=0}^{K-1} \frac{\alpha_k \sigma}{2} \|c_k\|_1 \right] \leq \kappa_7(K),$$

and thus

$$\frac{1}{K} \mathbb{E} \left[\sum_{k=0}^{K-1} \|c_k\|_1 \right] \leq \frac{2(\nu \sigma q_{-1} + 2\kappa_7(K)) \kappa_7(K)}{\nu^2 \sigma^2 \delta_1 K},$$

which proves the first result.

Next, by the law of iterated expectation, Jensen's inequality, and the results of Lemmas 2 and 6,

$$\mathbb{E}[b_{K-1}] = \mathbb{E} \left[\sqrt{b_{-1}^2 + \sum_{k=0}^{K-1} \|u_k\|^2} \right] \leq \sqrt{b_{-1}^2 + (\kappa_u^2 + \zeta^{-2} M) K} = \kappa_8(K). \quad (55)$$

Therefore, using (54) with $\delta_1 = \delta_2/2$ and Markov's inequality with (55), with probability at least $1 - \delta_2$, by Assumption 2, (47), and the union bound,

$$\mathbb{E} \left[\sum_{k=0}^{K-1} (u_k^{\text{true}})^T H_k u_k^{\text{true}} \right] \leq \frac{8\kappa_8(K)(\nu \sigma q_{-1} + 2\kappa_7(K)) \kappa_7(K)}{\tau_{\min} \nu^2 \eta \sigma \delta_2^2}.$$

Next, applying Lemma 4,

$$\begin{aligned} \mathbb{E} \left[\sum_{k=0}^{K-1} \|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|^2 \right] &\leq \frac{8\kappa_H^2 \kappa_8(K)(\nu \sigma q_{-1} + 2\kappa_7(K)) \kappa_7(K)}{\tau_{\min} \nu^2 \eta \zeta \sigma \delta_2^2} \\ &\quad + \kappa_H^2 (1 + 2\kappa_u) \kappa_v \mathbb{E} \left[\sum_{k=0}^{K-1} \|c_k\|_1 \right]. \end{aligned}$$

Noting that this result holds under the same event as in (51) (with $\delta_1 = \delta_2/2$), it follows that with probability at least $1 - \delta_2$,

$$\begin{aligned} \mathbb{E} \left[\frac{1}{K} \sum_{k=0}^{K-1} \|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|^2 \right] &\leq \frac{8\kappa_H^2 \kappa_8(K)(\nu \sigma q_{-1} + 2\kappa_7(K)) \kappa_7(K)}{\tau_{\min} \nu^2 \eta \zeta \sigma \delta_2^2 K} \\ &\quad + \frac{4(\nu \sigma q_{-1} + 2\kappa_7(K)) \kappa_H^2 (1 + 2\kappa_u) \kappa_v \kappa_7(K)}{\nu^2 \sigma^2 \delta_2 K}. \end{aligned}$$

Finally, using (47), (54), (55), Markov's inequality and the union bound, with probability at least $1 - \frac{2}{3}\delta_3$,

$$\begin{aligned} & \mathbb{E} \left[\sum_{k=0}^{K-1} \tau_{\min}(u_k^{\text{true}})^T H_k u_k^{\text{true}} + \frac{\sigma \kappa_8(K)}{\eta} \|c_k\|_1 \right] \\ & \leq \frac{18\kappa_8(K)(\nu\sigma q_{-1} + 2\kappa_7(K))\kappa_7(K)}{\nu^2\eta\sigma\delta_3^2}. \end{aligned}$$

Thus, by Lemma 4

$$\begin{aligned} & \mathbb{E} \left[\frac{1}{K} \sum_{k=0}^{K-1} \tau_{\min} \zeta \kappa_H^{-2} \|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|^2 + \frac{\sigma \kappa_8(K)}{\eta} \|c_k\|_1 \right] \\ & \leq \frac{18\kappa_8(K)(\nu\sigma q_{-1} + 2\kappa_7(K))\kappa_7(K)}{\nu^2\eta\sigma\delta_3^2 K} + \zeta \tau_{\min}(1 + 2\kappa_u)\kappa_v \mathbb{E} \left[\frac{1}{K} \sum_{k=0}^{K-1} \|c_k\|_1 \right] \\ & \leq \frac{18\kappa_8(K)(\nu\sigma q_{-1} + 2\kappa_7(K))\kappa_7(K)}{\nu^2\eta\sigma\delta_3^2 K} \\ & \quad + \frac{6(\nu\sigma q_{-1} + 2\kappa_7(K))\kappa_7(K)\zeta\tau_{\min}(1 + 2\kappa_u)\kappa_v}{\nu^2\sigma^2\delta_3 K}. \end{aligned}$$

Therefore, applying Markov's inequality, Jensen's inequality, and the union bound, it follows that with probability at least $1 - \delta_3$,

$$\begin{aligned} & \min_{k \in [0, K-1]} \tau_{\min} \zeta \kappa_H^{-2} \|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|^2 + \frac{\sigma \kappa_8(K)}{\eta} \|c_k\|_1 \\ & \leq \frac{54\kappa_8(K)(\nu\sigma q_{-1} + 2\kappa_7(K))\kappa_7(K)}{\nu^2\eta\sigma\delta_3^3 K} \\ & \quad + \frac{18(\nu\sigma q_{-1} + 2\kappa_7(K))\kappa_7(K)\zeta\tau_{\min}(1 + 2\kappa_u)\kappa_v}{\nu^2\sigma^2\delta_3^2 K}. \end{aligned}$$

□

By the definitions of $\kappa_7(K) = \mathcal{O}(\log(K))$ and $\kappa_8(K) = \mathcal{O}(\sqrt{K})$, it follows that the results of Corollary 3 match, up to log factors, those we derived in Section 3.1 for the pre-specified stepsize setting. Thus, in terms of the complexity measures (2), this variant of Algorithm 1 has a worst-case complexity of $\tilde{\mathcal{O}}(\epsilon_\ell^{-4})$ and $\tilde{\mathcal{O}}(\epsilon_c^{-1})$.

4 Safeguarded Line Search

The convergence analysis in Section 3 specifies proper ranges for α_k in Algorithm 1 in order to ensure convergence, but does not provide any recommendations on how to choose α_k in this range. Commonly, in other stochastic SQP methods, the procedure used to set α_k incorporates the merit parameter τ_k , which is adaptively estimated at each iteration. However, the estimation of τ_k

may be highly inaccurate and noisy due to only having stochastic access to the gradient of f . For this reason, we do not attempt to rely on the stochastic gradient information in order to choose α_k and instead solely utilize the constraints.

Consider first the case where α_k satisfies $\alpha_k \in [\nu, \nu + \theta\beta_k]$ as it does in the analysis in Section 3.1. Then, we can find an α_k in this range through a safeguarded backtracking procedure. Starting from $\hat{\alpha}_k = \nu + \theta\beta_k$, we backtrack until

$$\|c(x_k + \hat{\alpha}_k d_k)\|_1 \leq (1 - \xi\hat{\alpha}_k)\|c_k\|_1, \quad (56)$$

holds for some $\xi \in (0, 1)$ where, when (56) fails to hold for $\hat{\alpha}_k$, we set $\hat{\alpha}_k = \rho\hat{\alpha}_k$ for some $\rho \in (0, 1)$. However, as we cannot guarantee termination, we safeguard this linesearch by ceasing the search procedure if $\hat{\alpha}_k$ ever falls below ν . When (56) holds for some $\hat{\alpha}_k \geq \nu$, we set $\alpha_k = \hat{\alpha}_k$. On the other hand, if (56) fails to hold prior to $\hat{\alpha}_k < \nu$, we instead set $\alpha_k = \nu$. Thus, this procedure is guaranteed to output an α_k in the specified range and therefore the convergence results of Section 3.1 hold. In addition, on any step where (56) is satisfied for $\hat{\alpha}_k \geq \nu$, we have confirmation of sufficient decrease in the constraint violation. Finally, we note that the number of backtracking steps at any iteration k is at most $\log(\nu/(\nu + \theta\beta_k))/\log(\rho)$ due to terminating the backtracking as soon as $\hat{\alpha}_k < \nu$.

Unfortunately, the convergence theory only holds for the previous procedure under certain conditions on ν . To relax these conditions, we once again turn to the one of the adaptive stepsize rules of Section 3.2. In particular, we consider the case where β_k is chosen as a pre-specified sequence and the lower bound for α_k is chosen in a manner similar to that of (34). As we use a slight modification of this stepsize, we give the full procedure (which is the algorithm used in the computational results of Section 5) in Algorithm 2.

The backtracking procedure in Algorithm 2 is very similar to the one described above, with a few minor differences. In particular, we set the lower bound adaptively, using the stepsize rule in Section 3.2. In addition, unlike in Section 3.2, we only update the lower bound when the backtracking procedure fails to satisfy the sufficient decrease condition prior to reducing $\hat{\alpha}_k$ below the lower bound. The logic for this is simple; if the lower bound was reached, then it is probably too large and should be reduced. On the other hand, when the sufficient decrease condition is satisfied at iteration k , we keep the lower bound as it was at the start of this iteration, since it is already sufficiently small to find a good steplength in terms of reducing the constraint violation. We numerically explore the impact of only updating the lower bound when the line search fails to satisfy the sufficient decrease condition in Appendix B.

While this is a relatively simple variant of Algorithm 1, the analysis in Section 3 does not directly translate. We provide the following lemma which provides a starting point for the analysis that can then easily be combined with the techniques in Section 3 to obtain a worst-case complexity result.

Lemma 11 *Let Assumptions 1, 2, and 3 hold and let x_k be generated by Algorithm 2. Let $\beta_k = \eta/\sqrt{K}$ hold for all k . Let*

$$\kappa_9 := \xi^{-1}(\|c_0\|_1 - (2 + \Gamma\kappa_v/2)\nu^2 \log(q_{-1}^2) + \nu^2\eta^2\Gamma(\kappa_u^2 + \zeta^{-2}M)/(2q_{-1}^2)) \quad (57)$$

Algorithm 2 Two Stepsize Stochastic SQP with Adaptive Backtracking

Require: $x_0 \in \mathbb{R}^n$, $\{\beta_k\} \subset \mathbb{R}_{>0}$, $\nu \in \mathbb{R}_{>0}$, $q_{-1} \in \mathbb{R}_{>0}$, $\theta \in \mathbb{R}_{>0}$, $\xi \in (0, 1)$, $\rho \in (0, 1)$;

- 1: **for** $k = 0, 1, \dots$ **do**
- 2: Compute stochastic gradient g_k .
- 3: Compute (p_k, y_k) as the solution of (5).
- 4: Set $d_k \leftarrow v_k + \beta_k u_k$, where $v_k \in \text{Range}(J_k^T)$ and $u_k \in \text{Null}(J_k)$ are the orthogonal decomposition of p_k .
- 5: Set $\hat{q}_k^2 \leftarrow q_{k-1}^2 + \|c_k\|_1$ and $\hat{\alpha}_k \leftarrow \frac{\nu}{\hat{q}_k} + \theta\beta_k$.
- 6: **while** $\|c(x_k + \hat{\alpha}_k d_k)\|_1 > (1 - \xi\hat{\alpha}_k)\|c_k\|_1$ and $\hat{\alpha}_k \geq \frac{\nu}{\hat{q}_k}$ **do**
- 7: Set $\hat{\alpha}_k \leftarrow \rho\hat{\alpha}_k$.
- 8: **end while**
- 9: **if** $\hat{\alpha}_k > \frac{\nu}{\hat{q}_k}$ **then**
- 10: Set $\alpha_k \leftarrow \hat{\alpha}_k$ and $q_k = q_{k-1}$.
- 11: **else**
- 12: Set $\alpha_k \leftarrow \frac{\nu}{\hat{q}_k}$ and $q_k = \hat{q}_k$.
- 13: **end if**
- 14: Set $x_{k+1} \leftarrow x_k + \alpha_k d_k$.
- 15: **end for**

and

$$\kappa_{10} := 2(q_{-1} + \kappa_9/\nu) + 8(4 + \Gamma\kappa_\nu)\xi^{-1}\nu \log(e + (4 + \Gamma\kappa_\nu)\xi^{-1}\nu). \quad (58)$$

Then, $\mathbb{E}[q_{K-1}] \leq \kappa_{10}$ and

$$\sum_{k=0}^{K-1} \mathbb{E}[\alpha_k \|c_k\|_1] \leq \kappa_9 + (4 + \Gamma\kappa_\nu)\xi^{-1}\nu \log(\kappa_{10}) \quad (59)$$

Proof Let \mathcal{K}_α denote the index set of iterations k such that $\alpha_k = \frac{\nu}{\hat{q}_k}$. Then, for any $k \in \mathcal{K}_\alpha$, by Γ -Lipschitz continuity of the Jacobian of c ,

$$\begin{aligned} \|c(x_k + \alpha_k d_k)\|_1 - \|c_k\|_1 &\leq \|c_k + \alpha_k J_k d_k\|_1 - \|c_k\|_1 + \frac{\alpha_k^2 \Gamma}{2} \|d_k\|^2 \\ &= |1 - \alpha_k| \|c_k\|_1 - \|c_k\|_1 + \frac{\alpha_k^2 \Gamma}{2} \|d_k\|^2, \end{aligned}$$

where the equality follows from $J_k d_k = -c_k$. Using the fact that $|1 - \alpha_k| \leq 1 - \alpha_k + 2\alpha_k^2$, whenever $k \in \mathcal{K}_\alpha$, we have

$$\|c(x_k + \alpha_k d_k)\|_1 - \|c_k\|_1 \leq -\xi\alpha_k \|c_k\|_1 + 2\alpha_k^2 \|c_k\|_1 + \frac{\alpha_k^2 \Gamma}{2} \|d_k\|^2,$$

where we used $\xi \in (0, 1)$.

Next, for any iteration where $k \in \mathcal{K}_\alpha^c$, it follows that

$$\|c(x_k + \alpha_k d_k)\|_1 - \|c_k\|_1 \leq (1 - \xi\alpha_k) \|c_k\|_1 - \|c_k\|_1 = -\xi\alpha_k \|c_k\|_1.$$

Combining these cases and summing this inequality for $k = 0, \dots, K-1$, it follows that

$$\|c(x_K)\|_1 - \|c_0\|_1 \leq -\sum_{k=0}^{K-1} \xi\alpha_k \|c_k\|_1 + \sum_{j \in \mathcal{K}_\alpha} 2\alpha_j^2 \|c_j\|_1 + \frac{\alpha_j^2 \Gamma}{2} \|d_j\|^2.$$

Next, by the orthogonal decomposition $d_k = v_k + \beta_k u_k$ and Lemma 1, we have

$$\begin{aligned}
& \|c(x_K)\|_1 - \|c_0\|_1 \\
& \leq - \sum_{k=0}^{K-1} \xi \alpha_k \|c_k\|_1 + \sum_{j \in \mathcal{K}_\alpha} 2\alpha_j^2 \|c_j\|_1 + \frac{\alpha_j^2 \Gamma}{2} (\|v_j\|^2 + \beta_j^2 \|u_j\|^2) \\
& \leq - \sum_{k=0}^{K-1} \xi \alpha_k \|c_k\|_1 + \sum_{j \in \mathcal{K}_\alpha} (2 + \Gamma \kappa_v / 2) \alpha_j^2 \|c_j\|_1 + \frac{\alpha_j^2 \beta_j^2 \Gamma}{2} \|u_j\|^2 \\
& \leq - \sum_{k=0}^{K-1} \xi \alpha_k \|c_k\|_1 + \sum_{j \in \mathcal{K}_\alpha} \frac{(2 + \Gamma \kappa_v / 2) \nu^2}{q_{-1}^2 + \sum_{\ell \in \mathcal{K}_\alpha, \ell \leq j} \|c_\ell\|_1} \|c_j\|_1 + \frac{\nu^2 \beta_j^2 \Gamma}{2q_{-1}^2} \|u_j\|^2, \quad (60)
\end{aligned}$$

where the final inequality follows by the definition of α_k for any $k \in \mathcal{K}_\alpha$. Next, taking the expectation of both sides of this inequality, using the definition of β , rearranging, and using the law of iterated expectation with the result of Lemma 6, we have

$$\begin{aligned}
& \sum_{k=0}^{K-1} \mathbb{E} [\alpha_k \|c_k\|_1] \\
& \leq \xi^{-1} \|c_0\|_1 + \mathbb{E} \left[\sum_{j \in \mathcal{K}_\alpha} \frac{\xi^{-1} (2 + \Gamma \kappa_v / 2) \nu^2}{q_{-1}^2 + \sum_{\ell \in \mathcal{K}_\alpha, \ell \leq j} \|c_\ell\|_1} \|c_j\|_1 + \frac{\nu^2 \eta^2 \Gamma}{2K q_{-1}^2} \|u_j\|^2 \right] \\
& \leq \xi^{-1} \|c_0\|_1 + \xi^{-1} (2 + \Gamma \kappa_v / 2) \nu^2 \mathbb{E} \left[\log(q_{-1}^2 + \sum_{j \in \mathcal{K}_\alpha} \|c_j\|_1) - \log(q_{-1}^2) \right] \\
& \quad + \frac{\nu^2 \eta^2 \xi^{-1} \Gamma (\kappa_u^2 + \zeta^{-2} M)}{2q_{-1}^2} \\
& \leq \xi^{-1} \|c_0\|_1 + \xi^{-1} (2 + \Gamma \kappa_v / 2) \nu^2 (2 \log(\mathbb{E}[q_{K-1}]) - \log(q_{-1}^2)) \quad (61) \\
& \quad + \frac{\nu^2 \eta^2 \xi^{-1} \Gamma (\kappa_u^2 + \zeta^{-2} M)}{2q_{-1}^2},
\end{aligned}$$

where the second inequality follows by Lemma 9 and the final inequality follows by Jensen's inequality and the concavity of $\log(x)$.

Therefore, since $\alpha_k \geq \nu / q_{K-1}$ for all $k \leq K-1$, it follows that

$$\mathbb{E} \left[\sum_{k \in \mathcal{K}_\alpha} \frac{\|c_k\|_1}{q_{K-1}} \right] \leq \kappa_9 / \nu + (4 + \Gamma \kappa_v) \xi^{-1} \nu \log(\mathbb{E}[q_{K-1}]).$$

Now, by the definition of q_{K-1} and the fact that $x \leq a + b \log(x)$ implies $x \leq 2a + 8b \log(e + b)$ for any $a > 0$ and $b > 0$ [30],

$$\mathbb{E}[q_{K-1}] = \mathbb{E} \left[\frac{q_{-1}^2 + \sum_{k \in \mathcal{K}_\alpha} \|c_k\|_1}{q_{K-1}} \right]$$

$$\begin{aligned} &\leq q_{-1} + \kappa_9/\nu + (4 + \Gamma\kappa_v)\xi^{-1}\nu \log(\mathbb{E}[q_{K-1}]) \\ &\leq 2(q_{-1} + \kappa_9/\nu) + 8(4 + \Gamma\kappa_v)\xi^{-1}\nu \log(e + (4 + \Gamma\kappa_v)\xi^{-1}\nu), \end{aligned}$$

proving the first result. Thus, by (65), it follows that

$$\sum_{k=0}^{K-1} \mathbb{E}[\alpha_k \|c_k\|_1] \leq \kappa_9 + (4 + \Gamma\kappa_v)\xi^{-1}\nu \log(\mathbb{E}[q_{K-1}]) \leq \kappa_9 + (4 + \Gamma\kappa_v)\xi^{-1}\nu \log(\kappa_{10}).$$

□

From this proof, we can see that we still obtain a convergence rate of $\mathcal{O}(1/K)$ in terms of the average constraint violation. In addition, as in Section 3.1, any second order terms in the convergence analysis can be split into either terms involving β_k^2 or $\alpha_k^2 \|c_k\|_1$ terms. Since α_k is bounded from above by a constant, it should be clear by the prior lemma that the sum of the $\alpha_k^2 \|c_k\|_1$ terms are bounded, in expectation, by a constant factor. In addition, given the bound on $\mathbb{E}[q_{K-1}]$, we can combine the analysis of Sections 3.1 and 3.2 to derive a convergence result with a worst-case complexity of $\mathcal{O}(\epsilon_\ell^{-4})$ and $\mathcal{O}(\epsilon_c^{-1})$, matching the results of Section 3.1. We leave the full complexity analysis as an exercise to the reader.

Finally, we repeat a similar lemma in the case where β_k is chosen adaptively by (33) and (34).

Lemma 12 *Let Assumptions 1, 2, and 3 hold and let x_k be generated by Algorithm 2. Let β_k be defined as in (33) and (34) for all k . Let*

$$\begin{aligned} \kappa_{11}(K) &:= \xi^{-1}(\|c_0\|_1 - (2 + \Gamma\kappa_v/2)\nu^2 \log(q_{-1}^2)) \\ &\quad + \nu^2 \eta^2 \Gamma \log(1 + (\kappa_u^2 + \zeta^{-2}M)K/b_{-1}^2) / (2q_{-1}^2) \end{aligned} \quad (62)$$

and

$$\kappa_{12}(K) := 2(q_{-1} + \kappa_{11}(K))/\nu + 8(4 + \Gamma\kappa_v)\xi^{-1}\nu \log(e + (4 + \Gamma\kappa_v)\xi^{-1}\nu). \quad (63)$$

Then, $\mathbb{E}[q_{K-1}] \leq \kappa_{11}$ and

$$\sum_{k=0}^{K-1} \mathbb{E}[\alpha_k \|c_k\|_1] \leq \kappa_9 + (4 + \Gamma\kappa_v)\xi^{-1}\nu \log(\kappa_{10}) \quad (64)$$

Proof By the proof of Lemma 11, (60) holds in this case as well. Thus, taking the expectation of this inequality, rearranging, and applying Lemma 9 twice, we have,

$$\begin{aligned} &\sum_{k=0}^{K-1} \mathbb{E}[\alpha_k \|c_k\|_1] \\ &\leq \xi^{-1}\|c_0\|_1 + \mathbb{E} \left[\sum_{j \in \mathcal{K}_\alpha} \frac{\xi^{-1}(2 + \Gamma\kappa_v/2)\nu^2}{q_{-1}^2 + \sum_{\ell \in \mathcal{K}_\alpha, \ell \leq j} \|c_\ell\|_1} \|c_j\|_1 + \frac{\xi^{-1}\nu^2 \beta_j^2 \Gamma}{2q_{-1}^2} \|u_j\|^2 \right] \end{aligned}$$

$$\begin{aligned}
&\leq \xi^{-1} \|c_0\|_1 + \xi^{-1} (2 + \Gamma \kappa_v / 2) \nu^2 \mathbb{E} \left[\log(q_{-1}^2 + \sum_{j \in \mathcal{K}_\alpha} \|c_j\|_1) - \log(q_{-1}^2) \right] \\
&\quad + \frac{\xi^{-1} \nu^2 \eta^2 \Gamma}{2q_{-1}^2} \sum_{j=0}^{K-1} \mathbb{E} \left[\log \left(\frac{b_{-1}^2 + \sum_{k=0}^{K-1} \|u_k\|^2}{b_{-1}^2} \right) \right] \\
&\leq \xi^{-1} \|c_0\|_1 + \xi^{-1} (2 + \Gamma \kappa_v / 2) \nu^2 \eta^2 (2 \log(\mathbb{E}[q_{K-1}]) - \log(q_{-1}^2)) \quad (65) \\
&\quad + \frac{\xi^{-1} \nu^2 \Gamma \log(1 + (\kappa_u^2 + \zeta^{-2} M) K / b_{-1}^2)}{2q_{-1}^2},
\end{aligned}$$

where the final inequality follows by (48). From here, the proof follows by an identical argument to that of Lemma 11 with replacement of κ_9 by $\kappa_{11}(K)$ and κ_{10} by $\kappa_{12}(K)$, respectively. \square

Clearly, this proof shows that we still obtain an $\tilde{\mathcal{O}}(1/K)$ in terms of the average constraint violation, as was the case in Section 3.2. Similarly to the previous case, one can directly use this lemma with the proofs in Section 3.2 to prove the analagous $\tilde{\mathcal{O}}(\epsilon_\ell^{-4})$ and $\tilde{\mathcal{O}}(\epsilon_c^{-1})$ complexity result.

5 Numerical Experiments

In this section, we numerically validate the performance of our proposed algorithm. We focus our attention on Algorithm 2, as it is a fully specified version of the generic Algorithm 1. We consider two settings, equality constrained problems from the CUTEst collection [19] with simulated noise and constrained logistic regression problems using datasets from the LIBSVM collection [11].

In both cases, we compare Algorithm 2 with the Github implementation of Algorithm 3 (which we refer to as “SSQP” throughout this section) in [4] and use the parameter settings provided in [4]³, with the exception of θ , which we set as $\theta = 10^4$ for fair comparison. We use a variant of this algorithm which finds α_k using the procedure described in Section 5.1 of [3]. In addition, for the constrained logistic regression problems, we implemented the Momentum-based Linearized Augmented Lagrangian Method (MLALM) of [29]. This algorithm was not implemented for the CUTEst problems as it requires access to specific mini-batches, which is not available given the structure of the noise. We describe the parameter settings for this algorithm in more detail in Section 5.2.

We consider three variants of Algorithm 1, Algorithm 2 with a fixed stepsize β , which refer to as “TSSQP” throughout this section, Algorithm 2 with β_k chosen adaptively by (33) and (34) (“TSSQPU”), and Algorithm 1 with α_k and β_k chosen by (33) and (34), where we always select the stepsize $\alpha_k = \frac{\nu}{q_k}$ (“TSSQPUV”). As mentioned in Section 3.2, we use $q_k^2 = q_{k-1}^k +$

³ <https://github.com/frankecurtis/StochasticSQP>

$\min\{\|c_k\|_1, \|v_k\|, \|v_k\|^2\}$ when updating the lower bound for α_k as it significantly improves the numerical performance of the algorithms, without harming the complexity results, beyond constant factors. In order to not overwhelm the main body of the paper with numerical results, the findings for the fully adaptive methods (TSSQPU and TSSQPUV) are presented in Appendix A.

For all of the TSSQP algorithms, we use the following default parameter settings $\nu = 1$, $q_{-1} = 10^{-9}$, $\theta = 10^4$, $\xi = 10^{-3}$ and $\rho = 1/2$, while β_k is specified below based on the specific problem and algorithm variant in use. In all of our experiments, we choose H_k to be the identity matrix and compute the decomposition by first computing p_k directly by (5) and then find u_k by projecting p_k on $\text{Null}(J_k)$ and v_k by setting $v_k = p_k - u_k$.

Additional experiments involving variants of Algorithm 1 and 2 applied to these problems can be found in Appendix A and Appendix B. Information regarding the wallclock time required to solve these problems is given in C.

5.1 CUTEst Experiments

We consider the performance of Algorithm 2 on a subset of the equality constrained problems from the CUTEst collection [19]. We follow the experimental setup of [4] and select equality constrained optimization problems for which (i) f is not a constant function, (ii) $n+m \leq 1000$ and (iii) the Jacobian of c was non-singular at every iteration performed in our experiments. This selection resulted in a total of 60 problems, each of which has specified initial point, which we used in our experiments. We consider these problems at six different noise levels of $\epsilon_N \in \{10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1\}$. At iteration k , a stochastic gradient is generated such that $g_k \sim \mathcal{N}(\nabla f(x_k), \epsilon_N I)$. For each problem and noise level, we ran a total of 20 instances for each algorithm. For each instance, all algorithms were given a total budget of 1000 iterations or constraint evaluations, in order to control for the potentially additional effort imposed by the backtracking routine in Algorithm 2.

For every trial performed, we computed a resulting feasibility and optimality error. If a trial produced a sufficiently feasible iterate in the sense that $\|c_k\|_\infty \leq 10^{-6}$ for some k , then, we report the feasibility error as $\|c_k\|_\infty$ and the optimality error was reported as $\|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|_\infty$, where y_k^{true} was computed as a least-squares multiplier using the true gradient $\nabla f(x_k)$ and J_k . (This ensures that the reported optimality error is not based on a stochastic gradient and is instead an accurate measure of optimality corresponding to the iterate x_k .) On the other hand, if no sufficiently feasible iterate was produced on a given run, then the feasibility error and optimality error were computed using the same measures at the least infeasible iterate computed. In addition to terminating when the maximum iteration limit is reached, the algorithms were terminated if they ever computed a point x_k which was both sufficiently feasible and the stationarity error was smaller than 10^{-4} . Algorithms SSQP and TSSQP were tuned over 5 stepsizes, $\beta \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1\}$. For each problem and noise level, the best performing β was found by choosing the

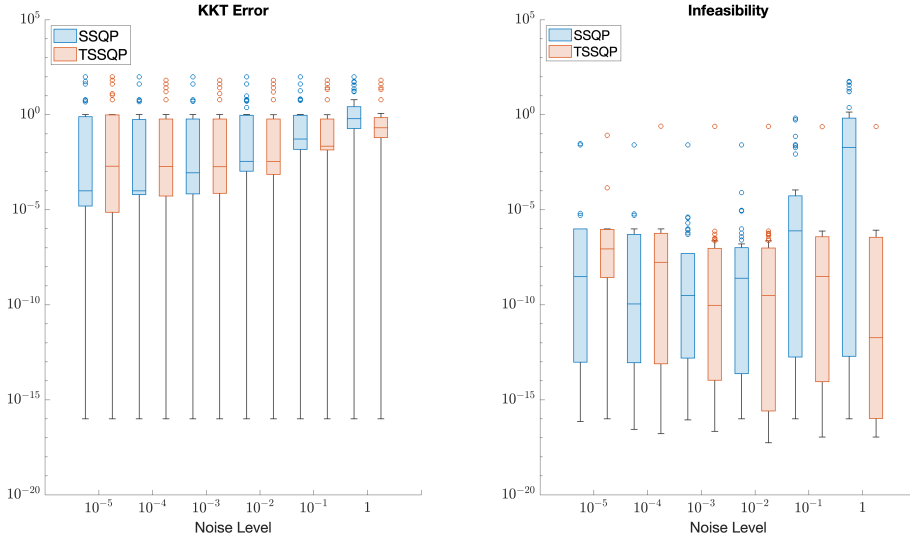


Fig. 1: Box plots of optimality error (left) and feasibility error (right) across various noise levels on CUTEst problems. SSQP is Algorithm 3 of [4] and TSSQP is Algorithm 2.

parameter setting which yielded the lowest average infeasibility error, when no parameter achieved sufficient feasibility on average, and the best average optimality error among sufficiently feasible solutions, when sufficient feasibility was achieved. The results of this experiment are presented in Figure 1.

As we can see from this plot, the computed stationarity error are relatively similar between these algorithms across all noise levels, with SSQP slightly outperforming TSSQP in stationarity error when the noise level is lower. This may be attributed to SSQP using an estimate of the merit parameter τ , which is more likely to be well-behaved in a low noise setting. However, once the noise level increases to $\epsilon_N = 10^{-2}$, the gap between these algorithms vanishes for the stationarity error and is eventually in favor of TSSQP. On the other hand, when the noise level is low, these algorithms perform similarly in terms of the infeasibility error, with SSQP slightly outperforming TSSQP. However, as the noise level increases, the performance of SSQP degrades significantly with respect to infeasibility, while the performance of TSSQP is largely unchanged. We view this as confirmation of our theoretical results as it demonstrates the superior ability of TSSQP to converge with respect to constraint violation while having minimal to no impact on its ability converge with respect to the KKT error.

5.2 Constrained Logistic Regression

In this subsection, we consider equality constrained logistic regression problems of the form

$$\min_{x \in \mathbb{R}^n} f(x) = \frac{1}{N} \sum_{i=1}^N \log \left(1 + e^{-y_i (X_i^T x)} \right) \quad \text{s.t. } Ax = b, \|x\|_2^2 = 1, \quad (66)$$

where $X \in \mathbb{R}^{n \times N}$ contains feature data of N data points (with X_i representing the i -th column of X), $y_i \in \{-1, 1\}^N$ contains label data, and $A \in \mathbb{R}^{m \times n}$ and $b \in \mathbb{R}^m$. For the datasets (X, y) , we consider all binary classification datasets from the LIBSVM collection [11] for which $22 \leq n \leq 1000$ and $202 \leq N \leq 100000$, which resulted in 9 datasets. For datasets with multiple versions, e.g., the {a1a, ..., a9a} datasets, we consider only the largest version.) The names and sizes of the datasets are given in Table 2. For the linear constraints, we chose $m = 10$ across all problems and generated A and b randomly with entries drawn from the standard normal distribution. The initial vector was chosen to be randomly distributed on the ball of norm 10^{-4} , a small, random initialization.

Table 2: Names and sizes of datasets. (Source: [11].)

Dataset	Dimension (n)	Datapoints (N)
a9a	123	32,561
ijcnn1	22	49,990
ionosphere	34	351
madelon	500	2,000
mushrooms	112	8,124
phising	68	11,055
sonar	60	208
splice	60	1,000
w8a	300	49,749

For each dataset, we considered two noise levels, where the level is dictated by the mini-batch size of each stochastic gradient estimate. For all problems, we used mini-batch sizes of 16 and 128. For each dataset and mini-batch size, we ran 20 instances with different random seeds. A budget of 10 epochs (i.e., number of effective passes over the dataset) was used for all SQP methods. For fair comparison, MLALM was given a budget of 300 epochs for each problem and batchsize, as it does not require linear system solves to generate search directions. This additional budget made MLALM the most expensive method in terms of time, see Appendix C for more details.

Similarly to the CUTEst problems, we considered 5 stepsizes for SSQP and TSSQP $\{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1\}$ while the adaptive methods were run with $\eta = 1$. For MLALM, following the prescriptions in [29], we set $\beta = \rho = K^{1/4}$, where K is total number of iterations to be performed. We consider 6 possible values for $\alpha = \{0, 0.2, 0.4, 0.6, 0.8, 1\}$ and four possible stepsize choices

Table 3: Average feasibility error, along with 95% confidence intervals, of MLALM, SSQP, and TSSQP. The results for the best-performing algorithm are shown in bold.

dataset	batch	MLALM	SSQP	TSSQP
		Feasibility	Feasibility	Feasibility
a9a	16	$1.46e-07 \pm 7.19e-08$	$1.32e-05 \pm 2.99e-12$	$3.95e-08 \pm 2.77e-16$
a9a	128	$1.51e-07 \pm 2.81e-08$	$1.48e-05 \pm 1.12e-10$	$5.25e-07 \pm 3.07e-15$
ijccn1	16	$5.02e-07 \pm 1.10e-07$	$1.20e-06 \pm 1.17e-13$	$5.17e-09 \pm 1.22e-16$
ijccn1	128	$8.82e-07 \pm 5.55e-08$	$1.20e-06 \pm 8.90e-14$	$1.03e-07 \pm 1.33e-16$
ionosphere	16	$3.67e-07 \pm 1.49e-07$	$1.20e-03 \pm 5.75e-07$	$6.90e-08 \pm 9.06e-16$
ionosphere	128	$2.39e-03 \pm 2.39e-04$	$1.79e-03 \pm 3.04e-09$	$4.49e-08 \pm 9.46e-17$
madelon	16	$3.73e-07 \pm 1.56e-07$	$9.76e-01 \pm 3.02e-03$	$2.79e-04 \pm 1.44e-04$
madelon	128	$7.36e-02 \pm 5.84e-02$	$9.91e-01 \pm 1.03e-03$	$1.31e-04 \pm 7.12e-08$
mushrooms	16	$8.32e-07 \pm 5.35e-08$	$1.00e-06 \pm 1.89e-12$	$7.99e-08 \pm 3.06e-15$
mushrooms	128	$8.59e-07 \pm 5.47e-08$	$1.52e-06 \pm 3.62e-12$	$8.67e-07 \pm 4.00e-14$
phishing	16	$1.01e-07 \pm 1.31e-07$	$1.16e-06 \pm 1.14e-12$	$6.01e-08 \pm 2.66e-16$
phishing	128	$5.10e-07 \pm 1.50e-07$	$4.56e-06 \pm 3.67e-12$	$9.37e-07 \pm 2.33e-14$
sonar	16	$3.84e-03 \pm 4.68e-04$	$3.35e-03 \pm 4.86e-09$	$8.59e-10 \pm 1.51e-16$
sonar	128	$3.41e-01 \pm 1.22e-04$	$2.35e-03 \pm 4.54e-09$	$2.60e-06 \pm 1.04e-16$
splice	16	$6.78e-07 \pm 1.01e-07$	$3.24e-03 \pm 2.06e-07$	$6.82e-07 \pm 2.01e-14$
splice	128	$1.61e-02 \pm 1.95e-04$	$3.91e-03 \pm 3.24e-09$	$5.46e-07 \pm 1.85e-14$
w8a	16	$9.07e-07 \pm 4.67e-08$	$4.87e-06 \pm 4.79e-13$	$2.43e-08 \pm 1.94e-16$
w8a	128	$6.65e-07 \pm 9.64e-08$	$5.31e-06 \pm 1.06e-12$	$3.04e-07 \pm 3.92e-16$

Table 4: Average stationarity error, along with 95% confidence intervals, of MLALM, SSQP, and TSSQP. The results for the best-performing algorithm are shown in bold.

dataset	batch	MLALM	SSQP	TSSQP
		Stationarity	Stationarity	Stationarity
a9a	16	$7.98e-02 \pm 3.70e-02$	$4.30e-02 \pm 3.61e-08$	$2.92e-02 \pm 2.25e-10$
a9a	128	$6.19e-02 \pm 1.13e-03$	$1.07e-02 \pm 5.98e-07$	$1.16e-02 \pm 3.80e-10$
ijccn1	16	$1.94e-01 \pm 2.20e-03$	$4.82e-02 \pm 2.68e-08$	$4.88e-02 \pm 9.18e-15$
ijccn1	128	$2.03e-01 \pm 7.47e-04$	$1.59e-02 \pm 6.46e-11$	$9.10e-03 \pm 3.62e-13$
ionosphere	16	$1.77e-01 \pm 4.16e-02$	$7.22e-02 \pm 5.00e-05$	$1.03e-01 \pm 7.93e-10$
ionosphere	128	$1.44e-01 \pm 2.00e-02$	$5.20e-02 \pm 5.88e-08$	$6.92e-02 \pm 4.21e-10$
madelon	16	$5.38e+03 \pm 8.23e+00$	$1.88e+02 \pm 3.69e+01$	$1.57e+02 \pm 4.59e+01$
madelon	128	$4.76e+03 \pm 6.53e+02$	$2.19e+02 \pm 5.13e+01$	$8.79e+01 \pm 3.02e-03$
mushrooms	16	$1.27e-01 \pm 1.47e-04$	$3.83e-03 \pm 1.45e-08$	$3.81e-03 \pm 1.56e-10$
mushrooms	128	$1.32e-01 \pm 4.39e-04$	$2.84e-03 \pm 9.11e-09$	$2.48e-03 \pm 5.33e-10$
phishing	16	$1.76e-01 \pm 9.98e-05$	$1.10e-02 \pm 1.07e-08$	$7.30e-03 \pm 7.03e-11$
phishing	128	$1.76e-01 \pm 4.04e-05$	$8.20e-03 \pm 4.23e-09$	$4.04e-03 \pm 8.70e-11$
sonar	16	$3.16e-01 \pm 6.76e-02$	$1.03e-01 \pm 1.25e-07$	$1.17e-01 \pm 1.37e-10$
sonar	128	$8.52e-01 \pm 1.11e-02$	$2.80e-02 \pm 3.29e-08$	$1.68e-01 \pm 1.38e-09$
splice	16	$4.28e+00 \pm 5.94e-01$	$1.06e+00 \pm 3.09e-06$	$4.15e-01 \pm 1.11e-08$
splice	128	$4.08e-01 \pm 1.03e-02$	$1.92e-01 \pm 7.93e-09$	$3.97e-01 \pm 5.03e-09$
w8a	16	$2.07e-01 \pm 1.49e-03$	$8.90e-03 \pm 2.49e-08$	$5.46e-03 \pm 1.05e-11$
w8a	128	$2.02e-01 \pm 1.33e-03$	$3.07e-03 \pm 3.28e-09$	$3.09e-03 \pm 7.04e-11$

$\eta = \{10^{-5}K^{-1/4}, 10^{-4}K^{-1/4}, 10^{-3}K^{-1/4}, 10^{-2}K^{-1/4}\}$, given a total of 24 different parameter settings. We follow a similar procedure as in the CUTEst experiments for reporting infeasibility and optimality error, with the caveat that we only record points at the end of each epoch. As before, if a sufficiently feasible iterate is found ($\|c_k\|_\infty \leq 10^{-6}$), then we report the sufficiently feasible point with the lowest optimality error, $\|\nabla f(x_k) + J_k^T y_k^{\text{true}}\|_\infty$. Otherwise, we report the least infeasible point. We chose the best parameter setting for each algorithm and batch size via the same procedure as in Section 5.1.

The results of these experiments can be seen in Tables 3 and 4. With respect to convergence in feasibility, TSSQP and MLALM perform the best, with TSSQP being the top performing algorithm more frequently than MLALM (13 times vs 4 times, with 1 tie). While not far behind in most instance, SSQP is never the top performing method. On the other hand, with respect to stationarity, SSQP and TSSQP are the top two performing methods, with TSSQP being the best algorithm 10 times versus 7 times for SSQP (with 1

tie). The MLALM algorithm is only the top performing algorithm with respect to feasibility on one instance. From these results, we conclude that the TSSQP method is the most effective overall as it performs well with respect to both measures the most frequently.

6 Conclusion

In this paper, we propose and analyze a new SQP method for equality constrained optimization with a stochastic objective function. The algorithm uses a stepsize splitting scheme in order to improve upon the worst-case complexity of recently proposed stochastic SQP methods. We show that the proposed method matches the rate of convergence of a deterministic SQP method in terms of constraint violation and obtains the optimal rate for a stochastic method in terms of the gradient of the Lagrangian.

There are number of possible directions of future research. Fundamentally, this stepsize splitting scheme can be incorporated into any of the previously proposed stochastic SQP methods in the literature, including those for rank deficient Jacobians, inequality constraints, and inexact subproblem solutions. Extending the algorithm to incorporate inexact subproblem solutions could likely be done by enforcing conditions similar to those derived in [16] and handling the additional error terms. Such an extension is outside the scope of the current manuscript.

References

1. Y. ARJEVANI, Y. CARMON, J. C. DUCHI, D. J. FOSTER, N. SREBRO, AND V. WOODWORTH, *Lower bounds for non-convex stochastic optimization*, *Mathematical Programming*, 199 (2023), pp. 165–214.
2. A. S. BERAHAS, R. BOLLAPRAGADA, AND B. ZHOU, *An adaptive sampling sequential quadratic programming method for equality constrained stochastic optimization*, arXiv preprint arXiv:2206.00712, (2022).
3. A. S. BERAHAS, F. E. CURTIS, M. J. O’NEILL, AND D. P. ROBINSON, *A stochastic sequential quadratic optimization algorithm for nonlinear-equality-constrained optimization with rank-deficient jacobians*, *Mathematics of Operations Research*, (2023).
4. A. S. BERAHAS, F. E. CURTIS, D. P. ROBINSON, AND B. ZHOU, *Sequential quadratic optimization for nonlinear equality constrained stochastic optimization*, *SIAM Journal on Optimization*, 31 (2021), pp. 1352–1379.
5. A. S. BERAHAS, J. SHI, Z. YI, AND B. ZHOU, *Accelerating stochastic sequential quadratic programming for equality constrained optimization using predictive variance reduction*, *Computational Optimization and Applications*, 86 (2023), pp. 79–116.
6. D. P. BERTSEKAS, *Network optimization: continuous and discrete models*, Athena Scientific Belmont, 1998.
7. J. T. BETTS, *Practical methods for optimal control and estimation using nonlinear programming*, SIAM, 2010.
8. J. BOLTE AND E. PAUWELS, *Majorization-minimization procedures and convergence of sqp methods for semi-algebraic and tame programs*, *Mathematics of Operations Research*, 41 (2016), pp. 442–465.
9. L. BOTTOU, F. E. CURTIS, AND J. NOCEDAL, *Optimization methods for large-scale machine learning*, *SIAM review*, 60 (2018), pp. 223–311.

10. C. CARTIS, N. I. GOULD, AND P. L. TOINT, *On the evaluation complexity of cubic regularization methods for potentially rank-deficient nonlinear least-squares problems and its relevance to constrained nonlinear optimization*, SIAM Journal on Optimization, 23 (2013), pp. 1553–1574.
11. C.-C. CHANG AND C.-J. LIN, *Libsvm: A library for support vector machines*, ACM transactions on intelligent systems and technology (TIST), 2 (2011), pp. 1–27.
12. C. CHEN, F. TUNG, N. VEDULA, AND G. MORI, *Constraint-aware deep neural network compression*, in Proceedings of the European Conference on Computer Vision (ECCV), 2018, pp. 400–415.
13. S. CUOMO, V. DI COLA, F. GIAMPAOLO, G. ROZZA, M. RAISSI, AND F. PICCIALI, *Scientific machine learning through physics-informed neural networks: Where we are and what’s next*, Journal of Scientific Computing, 92 (2022), p. 88.
14. F. E. CURTIS, M. J. O’NEILL, AND D. P. ROBINSON, *Worst-case complexity of an sqp method for nonlinear equality constrained stochastic optimization*, Mathematical Programming, 205 (2024), pp. 431–483.
15. F. E. CURTIS, D. P. ROBINSON, AND B. ZHOU, *Sequential quadratic optimization for stochastic optimization with deterministic nonlinear inequality and equality constraints*, SIAM Journal on Optimization, 34 (2024), pp. 3592–3622.
16. F. E. CURTIS, D. P. ROBINSON, AND B. ZHOU, *A stochastic inexact sequential quadratic optimization algorithm for nonlinear equality-constrained optimization*, INFORMS Journal on Optimization, 6 (2024), pp. 173–195.
17. J. DUCHI, E. HAZAN, AND Y. SINGER, *Adaptive subgradient methods for online learning and stochastic optimization.*, Journal of machine learning research, 12 (2011).
18. F. FACCHINEL, V. KUNGURTSEV, L. LAMPARIELLO, AND G. SCUTARI, *Ghost penalties in nonconvex constrained optimization: Diminishing stepsizes and iteration complexity*, Mathematics of Operations Research, 46 (2021), pp. 595–627.
19. N. I. GOULD, D. ORBAN, AND P. L. TOINT, *Cutest: a constrained and unconstrained testing environment with safe threads for mathematical optimization*, Computational optimization and applications, 60 (2015), pp. 545–557.
20. L. JIN AND X. WANG, *A stochastic primal-dual method for a class of nonconvex constrained optimization*, Computational Optimization and Applications, 83 (2022), pp. 143–180.
21. F. KUPFER AND E. W. SACHS, *Numerical solution of a nonlinear parabolic control problem by a reduced sqp method*, Computational Optimization and Applications, 1 (1992), pp. 113–135.
22. B. H. MCMAHAN AND M. STREETER, *Adaptive bound optimization for online convex optimization*, arXiv preprint arXiv:1002.4908, (2010).
23. S. NA, M. ANITESCU, AND M. KOLAR, *An adaptive stochastic sequential quadratic programming with differentiable exact augmented lagrangians*, Mathematical Programming, 199 (2023), pp. 721–791.
24. S. NA, M. ANITESCU, AND M. KOLAR, *Inequality constrained stochastic nonlinear optimization via active-set sequential quadratic programming*, Mathematical Programming, 202 (2023), pp. 279–353.
25. Y. NANDWANI, A. PATHAK, AND P. SINGLA, *A primal dual formulation for deep learning with constraints*, Advances in neural information processing systems, 32 (2019).
26. J. NOCEDAL AND S. J. WRIGHT, *Numerical optimization*, Springer Science & Business Media, 2006.
27. M. RAISSI, P. PERDIKARIS, AND G. E. KARNIADAKIS, *Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations*, Journal of Computational physics, 378 (2019), pp. 686–707.
28. T. REES, H. S. DOLLAR, AND A. J. WATHEN, *Optimal solvers for pde-constrained optimization*, SIAM Journal on Scientific Computing, 32 (2010), pp. 271–298.
29. Q. SHI, X. WANG, AND H. WANG, *A momentum-based linearized augmented lagrangian method for nonconvex constrained stochastic optimization*, Mathematics of Operations Research, 51 (2026), pp. 92–133.
30. B. WANG, H. ZHANG, Z. MA, AND W. CHEN, *Convergence of adagrad for non-convex objectives: Simple proofs and relaxed assumptions*, in The Thirty Sixth Annual Conference on Learning Theory, PMLR, 2023, pp. 161–190.

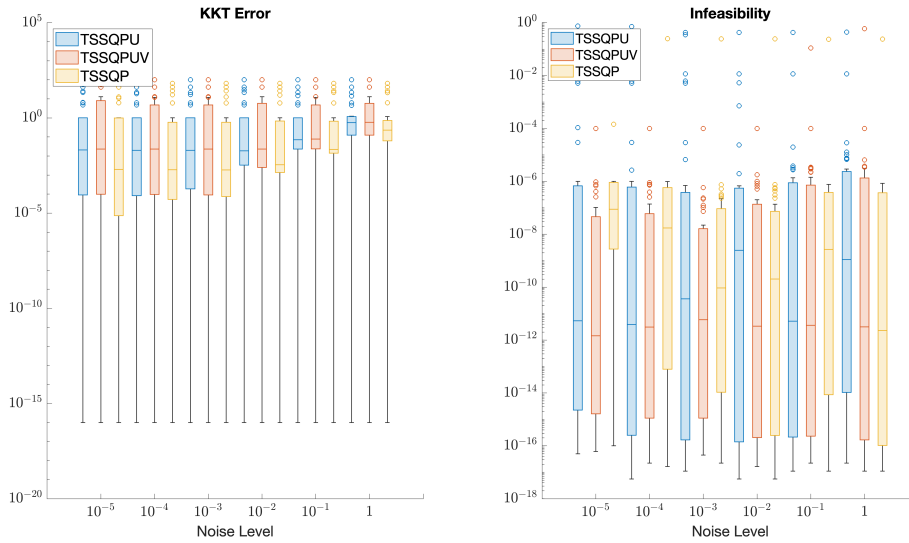


Fig. 2: Box plots of optimality error (left) and feasibility error (right) across various noise levels on CUTEst problems for TSSQPU and TSSQPUV. TSSQP is included for reference.

31. R. WARD, X. WU, AND L. BOTTOU, *Adagrad stepsizes: Sharp convergence over non-convex landscapes*, Journal of Machine Learning Research, 21 (2020), pp. 1–30.

A Numerical Performance of Adaptive Methods

In this appendix, we present our numerical results for the adaptive methods TSSQPU and TSSQPUV. For reference, we also include TSSQP in the results. We note that TSSQPU and TSSQPUV are fully adaptive in the sense that only default parameter settings were used ($\eta = 1$, $\nu = 1$), while the presented results for TSSQP were tuned for the best fixed stepsize β , using the same procedure as in Section 5.

A.1 Adaptive Methods on CUTEst Problems

In this subsection, we present the performance of TSSQPU and TSSQPUV on the CUTEst test set. For reference, we plot TSSQP along with TSSQPU and TSSQPUV in Figure 2.

We observe that, with respect to the infeasibility measure, the three methods perform largely the same, with the fully adaptive methods performing slightly better than TSSQP. On the other hand, both adaptive methods struggle with respect to the stationarity error when compared with TSSQP. Considering that both adaptive methods do not use any tuning, while TSSQP does, we find the results promising for the adaptive methods, which may be useful in scenarios where parameter tuning is difficult or expensive. In addition, TSSQPUV may be useful in contexts where computation of $c(x)$ is expensive, as it does not rely on the backtracking procedure in Algorithm 2 but still manages to converge to approximate feasibility consistently.

Table 5: Average feasibility error, along with 95% confidence intervals, of TSSQPU and TSSQPUV. For reference, we also include TSSQP. The results for the best-performing algorithm are shown in bold.

dataset	batch	TSSQPU	TSSQPUV	TSSQP
		Feasibility	Feasibility	Feasibility
a9a	16	$5.96e-06 \pm 1.14e-09$	$5.68e-02 \pm 5.49e-12$	$3.95e-08 \pm 2.77e-16$
a9a	128	$7.65e-05 \pm 1.21e-06$	$1.00e+00 \pm 2.13e-16$	$5.25e-07 \pm 3.07e-15$
ijccn1	16	$8.29e-07 \pm 2.57e-10$	$1.06e-06 \pm 3.92e-13$	$5.17e-09 \pm 1.22e-16$
ijccn1	128	$2.03e-05 \pm 4.30e-09$	$1.00e+00 \pm 2.13e-16$	$1.03e-07 \pm 1.33e-16$
ionosphere	16	$2.84e-04 \pm 4.45e-06$	$1.00e+00 \pm 2.13e-16$	$6.90e-08 \pm 9.06e-16$
ionosphere	128	$1.55e-03 \pm 1.39e-05$	$1.00e+00 \pm 2.13e-16$	$4.49e-08 \pm 9.46e-17$
madelon	16	$2.67e-06 \pm 2.50e-06$	$1.00e+00 \pm 2.13e-16$	$2.79e-04 \pm 1.44e-04$
madelon	128	$1.45e-04 \pm 8.61e-05$	$1.00e+00 \pm 2.13e-16$	$1.31e-04 \pm 7.12e-08$
mushrooms	16	$5.57e-06 \pm 1.37e-08$	$1.00e+00 \pm 2.13e-16$	$7.99e-08 \pm 3.06e-15$
mushrooms	128	$2.07e-06 \pm 2.88e-07$	$1.00e+00 \pm 2.13e-16$	$8.67e-07 \pm 4.00e-14$
phishing	16	$2.21e-05 \pm 1.63e-08$	$1.00e+00 \pm 2.13e-16$	$6.01e-08 \pm 2.66e-16$
phishing	128	$1.71e-04 \pm 4.64e-07$	$1.00e+00 \pm 2.13e-16$	$9.37e-07 \pm 2.33e-14$
sonar	16	$4.76e-04 \pm 7.92e-07$	$1.00e+00 \pm 2.13e-16$	$8.59e-10 \pm 1.51e-16$
sonar	128	$5.94e-02 \pm 4.50e-04$	$1.00e+00 \pm 2.13e-16$	$2.60e-06 \pm 1.04e-16$
splice	16	$4.70e-04 \pm 7.15e-07$	$1.00e+00 \pm 2.13e-16$	$6.82e-07 \pm 2.01e-14$
splice	128	$1.36e-03 \pm 2.16e-06$	$1.00e+00 \pm 2.13e-16$	$5.46e-07 \pm 1.85e-14$
w8a	16	$9.97e-06 \pm 8.54e-10$	$1.04e-04 \pm 3.26e-13$	$2.43e-08 \pm 1.94e-16$
w8a	128	$9.83e-05 \pm 1.67e-07$	$1.00e+00 \pm 2.13e-16$	$3.04e-07 \pm 3.92e-16$

Table 6: Average stationarity error, along with 95% confidence intervals, of TSSQPU and TSSQPUV. For reference, we also include TSSQP. The results for the best-performing algorithm are shown in bold.

dataset	batch	TSSQPU	TSSQPUV	TSSQP
		Stationarity	Stationarity	Stationarity
a9a	16	$3.75e-02 \pm 1.74e-07$	$1.59e-01 \pm 1.05e-05$	$2.92e-02 \pm 2.25e-10$
a9a	128	$1.39e-02 \pm 1.15e-06$	$1.47e+03 \pm 1.81e-03$	$1.16e-02 \pm 3.80e-10$
ijccn1	16	$4.06e-02 \pm 3.82e-09$	$1.24e-01 \pm 6.21e-07$	$4.88e-02 \pm 9.18e-15$
ijccn1	128	$8.69e-03 \pm 2.67e-08$	$2.00e+03 \pm 2.75e-03$	$9.10e-03 \pm 3.62e-13$
ionosphere	16	$1.65e-01 \pm 6.96e-05$	$2.70e+03 \pm 2.07e-02$	$1.03e-01 \pm 7.93e-10$
ionosphere	128	$1.88e-02 \pm 2.76e-05$	$2.70e+03 \pm 4.95e-03$	$6.92e-02 \pm 4.21e-10$
madelon	16	$1.03e+02 \pm 2.91e+01$	$7.76e+02 \pm 1.71e+01$	$1.57e+02 \pm 4.59e+01$
madelon	128	$8.83e+01 \pm 1.55e+01$	$7.62e+02 \pm 1.08e+01$	$8.79e+01 \pm 3.02e-03$
mushrooms	16	$1.25e-02 \pm 2.55e-07$	$1.37e+03 \pm 1.03e-02$	$3.81e-03 \pm 1.56e-10$
mushrooms	128	$1.95e-03 \pm 1.53e-04$	$1.37e+03 \pm 2.59e-03$	$2.48e-03 \pm 5.33e-10$
phishing	16	$5.13e-03 \pm 9.04e-08$	$1.56e+03 \pm 2.13e-03$	$7.30e-03 \pm 7.03e-11$
phishing	128	$9.24e-03 \pm 1.20e-07$	$1.56e+03 \pm 7.63e-04$	$4.04e-03 \pm 8.70e-11$
sonar	16	$8.57e-02 \pm 3.50e-06$	$1.59e+03 \pm 1.63e-02$	$1.17e-01 \pm 1.37e-10$
sonar	128	$5.94e-02 \pm 4.50e-04$	$1.59e+03 \pm 7.26e-03$	$1.68e-01 \pm 1.38e-09$
splice	16	$2.88e-01 \pm 2.03e-05$	$1.59e+03 \pm 1.53e-01$	$4.15e-01 \pm 1.11e-08$
splice	128	$4.27e-02 \pm 1.07e-05$	$1.59e+03 \pm 6.15e-02$	$3.97e-01 \pm 5.03e-09$
w8a	16	$2.81e-02 \pm 5.57e-08$	$1.48e-01 \pm 7.84e-07$	$5.46e-03 \pm 1.05e-11$
w8a	128	$9.25e-03 \pm 2.08e-06$	$9.52e+02 \pm 1.65e-03$	$3.09e-03 \pm 7.04e-11$

A.2 Adaptive Methods on Logistic Regression Problems

The results for the adaptive methods on the logistic regression problems from Subsection 5.2 can be found in Tables 5 and 6, where TSSQP is included for reference.

From these tables, we observe that TSSQPU performs reasonably well, especially with respect to stationarity error, where it is the best algorithm more often than TSSQP. In terms of feasibility error, TSSQP is more often the superior algorithm, but TSSQPU still performs reasonably well, even without tuning to the problem instance. We view this as confirmation of the algorithm’s ability to adapt to the problem structure effectively. On the other hand, TSSQPUV completely fails to converge on the majority of problems, both in terms of stationarity and feasibility (the initial infeasibility is $\|c(x_0)\|_\infty \approx 1$, so the instances where TSSQPUV records 1 for feasibility suggests that it never finds a more feasible point than the initial one). These results suggest that the linesearch procedure in Algorithm 2 may be very useful for certain problem instances, such as the logistic regression problems explored here, as it is the only difference between TSSQPU and TSSQPUV, while the gap between their results is significant.

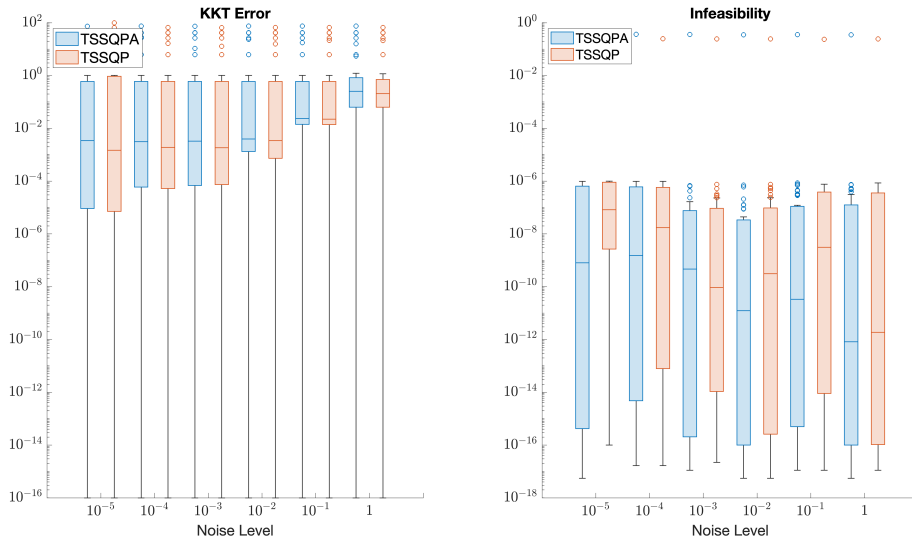


Fig. 3: Box plots of optimality error (left) and feasibility error (right) across various noise levels on CUTEst problems of TSSQPA vs TSSQP.

B Numerical Performance of Linesearch Variants and Stepsize Plots

In this subsection, we report results for a variant of Algorithm 2 in which the term q_k always accumulates (i.e., in line 10 of Algorithm 2, we set $q_k = \hat{q}_k$). The experiments were carried out in the same manner as those in Section 5 and we refer to this variant as “TSSQPA”, for always accumulate, in this section. To see the effect of this modification, we only compare directly with TSSQP. In addition, at the end of this appendix, we include plots of the behavior of the adaptive stepsizes computed by TSSQP, TSSQPA and TSSQPU on a logistic regression problem.

B.1 CUTEst Results

First, we report our results on the CUTEst problem set, which can be seen in Figure 3. We see relatively similar performance between the methods across all noise levels, with a slight advantage in KKT error for TSSQP and a slightly larger advantage in infeasibility error for TSSQPA. This aligns with the idea that TSSQPA always accumulates the $\|c_k\|_1$ in q_k , thus leading to overall shorter steps, which may hinder progress in achieving optimality.

B.2 Logistic Regression Results

Next, we report the results for the logistic regression problems, which can be seen below in Table 7 and 8. We see that TSSQPA and TSSQP perform very similarly in terms of both feasibility and stationarity, with TSSQPA having a slight edge in feasibility and TSSQP having a slight edge in stationarity. On a whole, these results suggest that both methods are effective and are reasonable choices for these types of problems.

Table 7: Average feasibility error, along with 95% confidence intervals, of TSSQP and TSSQPA. The results for the best-performing algorithm are shown in bold.

dataset	batch	TSSQPA		TSSQP	
		Feasibility		Feasibility	
a9a	16	2.15e-07	± 2.93e-15	3.95e-08	± 2.77e-16
a9a	128	1.96e-07	± 1.70e-14	5.25e-07	± 3.07e-15
ijccn1	16	2.55e-08	± 1.43e-16	5.17e-09	± 1.22e-16
ijccn1	128	5.12e-07	± 1.63e-16	1.03e-07	± 1.33e-16
ionosphere	16	3.74e-07	± 5.02e-15	6.90e-08	± 9.06e-16
ionosphere	128	3.90e-08	± 2.28e-16	4.49e-08	± 9.46e-17
madelon	16	1.91e-05	± 1.18e-06	2.79e-04	± 1.44e-04
madelon	128	3.58e-06	± 1.95e-12	1.31e-04	± 7.12e-08
mushrooms	16	4.38e-07	± 2.35e-13	7.99e-08	± 3.06e-15
mushrooms	128	2.40e-07	± 3.17e-13	8.67e-07	± 4.00e-14
phishing	16	4.71e-07	± 2.39e-15	6.01e-08	± 2.66e-16
phishing	128	2.72e-07	± 3.23e-14	9.37e-07	± 2.33e-14
sonar	16	1.03e-10	± 9.84e-14	8.59e-10	± 1.51e-16
sonar	128	7.76e-10	± 1.63e-16	2.60e-06	± 1.04e-16
splice	16	7.74e-08	± 3.29e-15	6.82e-07	± 2.01e-14
splice	128	1.58e-07	± 5.32e-15	5.46e-07	± 1.85e-14
w8a	16	1.04e-07	± 2.55e-15	2.43e-08	± 1.94e-16
w8a	128	1.61e-07	± 1.31e-14	3.04e-07	± 3.92e-16

Table 8: Average stationarity error, along with 95% confidence intervals, of TSSQP and TSSQPA. The results for the best-performing algorithm are shown in bold.

dataset	batch	TSSQPA		TSSQP	
		Stationarity		Stationarity	
a9a	16	3.76e-02	± 4.05e-09	2.92e-02	± 2.25e-10
a9a	128	1.13e-01	± 6.33e-09	1.16e-02	± 3.80e-10
ijccn1	16	4.89e-02	± 1.01e-11	4.88e-02	± 9.18e-15
ijccn1	128	8.91e-03	± 1.45e-11	9.10e-03	± 3.62e-13
ionosphere	16	1.11e-01	± 9.27e-10	1.03e-01	± 7.93e-10
ionosphere	128	1.35e-01	± 3.34e-10	6.92e-02	± 4.21e-10
madelon	16	1.20e+02	± 1.01e+01	1.57e+02	± 4.59e+01
madelon	128	1.33e+01	± 2.32e-05	8.79e+01	± 3.02e-03
mushrooms	16	3.79e-03	± 5.55e-09	3.81e-03	± 1.56e-10
mushrooms	128	6.26e-03	± 1.18e-08	2.48e-03	± 5.33e-10
phishing	16	1.12e-02	± 2.24e-10	7.30e-03	± 7.03e-11
phishing	128	8.18e-03	± 7.24e-10	4.04e-03	± 8.70e-11
sonar	16	1.17e-01	± 1.66e-08	1.17e-01	± 1.37e-10
sonar	128	8.13e-02	± 2.05e-10	1.68e-01	± 1.38e-09
splice	16	4.53e-01	± 1.58e-08	4.15e-01	± 1.11e-08
splice	128	4.01e-01	± 5.01e-09	3.97e-01	± 5.03e-09
w8a	16	2.04e-02	± 3.46e-09	5.46e-03	± 1.05e-11
w8a	128	8.76e-02	± 5.26e-09	3.09e-03	± 7.04e-11

B.3 Stepsize Behavior of Different Adaptive Methods

In this subsection, we provide plots of the behavior of the adaptive stepsizes computed by three methods, TSSQP, TSSQPA, and TSSQPU on the logistic regression problem w8a with a batch size of 16, as it is one of the largest problems in the collection. For this problem, the best settings for β was 10^{-3} for TSSQP and 10^{-4} for TSSQPA. We used these settings when generating the following plots.

The plots for TSSQP and TSSQPA are given in Figures 4 and 5. We focus on the first 30 iterations of TSSQPA, as the stepsize drops rapidly in that interval then remains largely unchanged for the rest of the procedure. We see that in both plots, the stepsize lower bound ($\nu/\sqrt{q_k}$) quickly decreases before becoming essentially flat, due to the constraint violation becoming very small. However, in the case of TSSQP, the lower bound stops at a much higher value (≈ 0.52) while TSSQPA stops at a lower value (≈ 0.02). This smaller stepsize impacts the ability of TSSQPA to achieve as accurate of a stationarity measure, as seen in Table 8. Either way, we see that the adaptive stepsize strategy quickly falls to a level such that the algorithm can make significant progress towards feasibility (even if it does increase initially, such as in the case of TSSQPA).

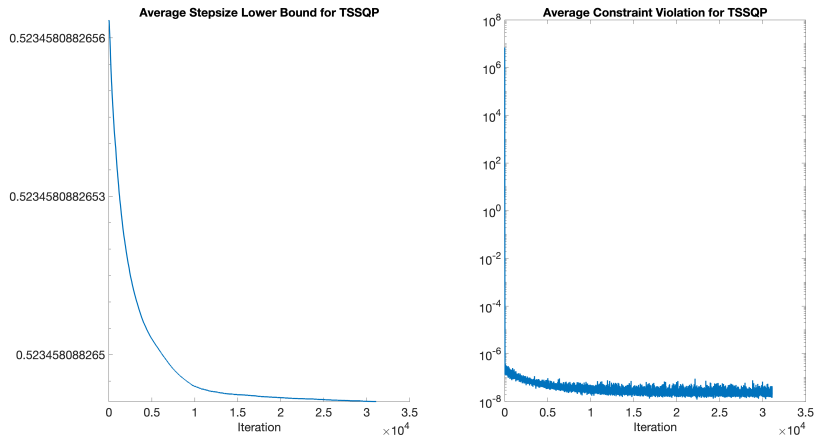


Fig. 4: Stepsize lower bound of TSSQP on w8a with a batchsize of 16 and the constraint violation.

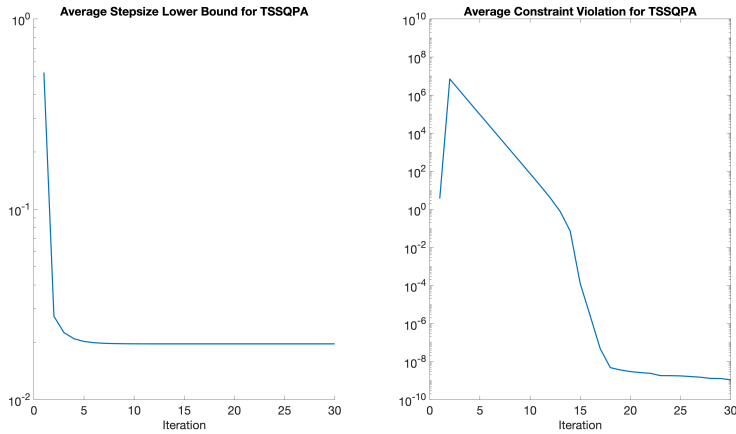


Fig. 5: Stepsize lower bound of TSSQPA on w8a with a batchsize of 16 and the constraint violation. We focus on the first 30 iterations, as the lower bound changes rapidly then remains largely flat afterwards.

Next, in Figure 6, we plot the behavior of β_k when it is chosen adaptively by (33) and (34). Due to the way β_k is constructed, the expected tail behavior is $\beta_k \sim 1/\sqrt{k}$. As such, we fit a curve of the form a/\sqrt{k} and plot this as well. We see that this fit is relatively accurate, suggesting that after an initial adaptive period where β falls rapidly to $\beta_k \approx 10^{-1}$, it settles into a decaying stepsize regime with a rate of decay proportional to $1/\sqrt{k}$. We note that this is essentially the behavior predicted by the quantity $\kappa_{\mathcal{G}}(K)$ in Corollary 3 (see (55)).

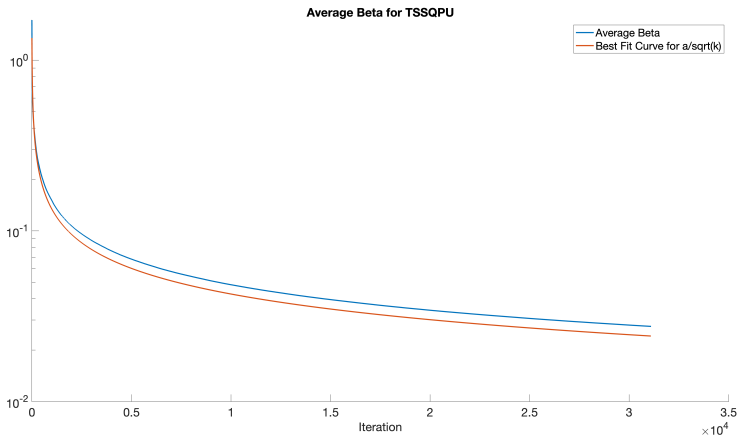


Fig. 6: Adaptive value of β_k for TSSQPU on w8a. We also include a curve which is fit to a/\sqrt{k} , with $a = 4.266$, to compare with the expected tail behavior.

C Wallclock Performance of Methods

In this appendix, we provide information on the wallclock performance of the numerical experiments. Caution should be exercised when considering the results of this appendix as the numerical experiments were run in a cluster environment, so there is no guarantee all experiments were run on identical hardware. We present two main pieces of information in this appendix: overall wallclock time of the methods and, for TSSQP, the relative wallclock time required to compute the decomposition u and v given a search direction d .

C.1 Time Comparison for CUTEst Problems

In this subsection, we compare the time for SSQP and TSSQP to solve the CUTEst problems from Subsection 5.1. To compare these quantities, for each problem, we compute the following measure:

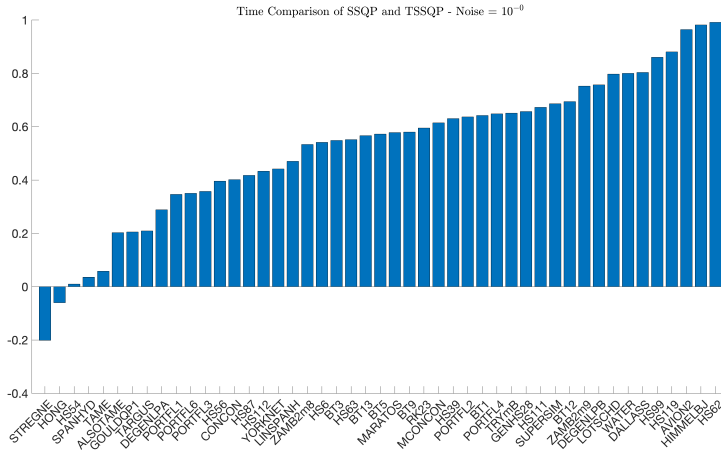
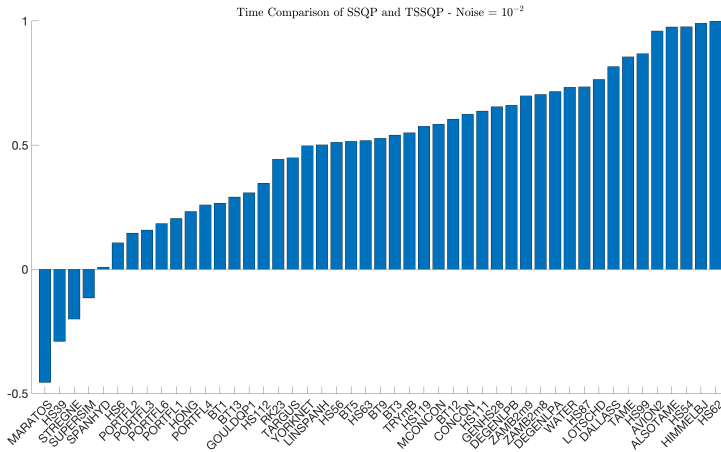
$$\frac{T_{SSQP}(i) - T_{TSSQP}(i)}{\max\{T_{SSQP}(i), T_{TSSQP}(i)\}}, \quad (67)$$

where $T_{SSQP}(i)$ and $T_{TSSQP}(i)$ are the average wallclock times of SSQP and TSSQP on problem i , respectively. This measure is always in $[-1, 1]$, with numbers closer to -1 representing SSQP taking less wallclock time than TSSQP and numbers closer to 1 representing TSSQP taking less time than SSQP. These results are plotted in Figures 7, 8, and 9 for noise levels of 1 , 10^{-2} , and 10^{-4} , respectively. We see that TSSQP consistently takes less time than SSQP, which may be a product of the early stopping done on these problems, as outlined in Section 5.1. However, this clearly shows that the additional cost of computing the decomposition and backtracking does not increase the computational costs so as to outweigh their benefits.

In addition to these plots, we also provide plots which calculate the relative time to compute the decomposition u_k and v_k , i.e.

$$T_{decomp}(i)/T_{TSSQP}(i), \quad (68)$$

where $T_{decomp}(i)$ is the time spent computing u_k and v_k after obtaining p_k on problem i . These results are given in Figures 10, 11, and 12. We see that for most problems, the relative

Fig. 7: Plots related to the measure (67) with $\epsilon_N = 1$.Fig. 8: Plots related to the measure (67) with $\epsilon_N = 10^{-2}$.

amount of time spent finding the decomposition u_k and v_k is no more that 15% of the total runtime. Thus, while computing this decomposition requires additional work, it does not overwhelm the cost of the full algorithm and appears to be beneficial enough to be worth the cost of this additional overhead. There are, however, a few instances, for which this cost is significantly higher, up to 60% of the time for one instance when the noise level was 10^{-4} . In such cases, other methods which do not rely on this decomposition may be preferred.

C.2 Logistic Regression Results

We conclude this appendix with wallclock times of the logistic regression experiments in Section 5.2. We include MLALM, SSQP, and TSSQP for comparison. Recall that MLALM

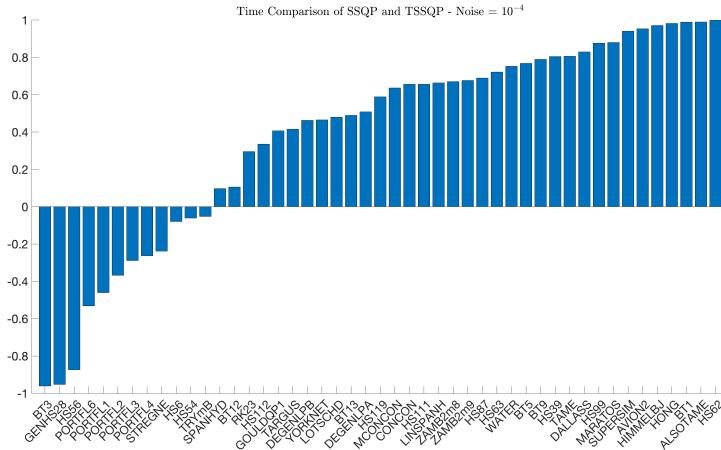


Fig. 9: Plots related to the measure (67) with $\epsilon_N = 10^{-4}$.

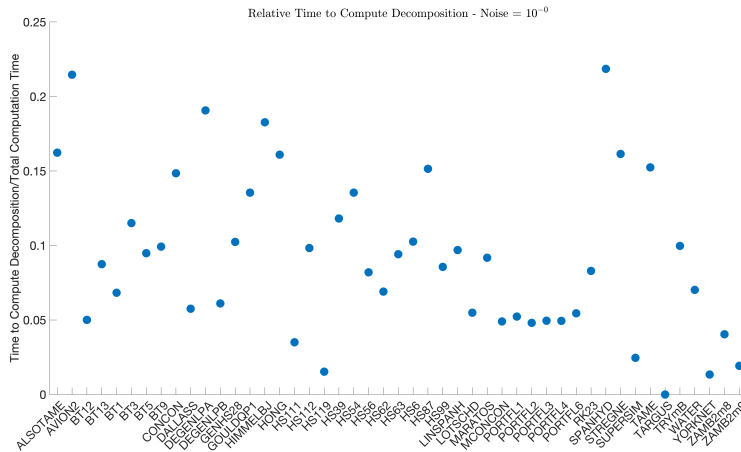


Fig. 10: Relative time to compute decomposition u and v in TSSQP, as measured by (68) with noise level $\epsilon_N = 1$.

was given a significantly larger budget (300 epochs) than the SQP methods (10 epochs), to account for the difference in step computation costs. In addition, for TSSQP we include a column, called relative decomposition time, which records (68) for the logistic regression experiments. The values can be found in Table 9. We see that MLALM almost always takes the most amount of time of the algorithms, validating that the number of additional epochs was sufficient for fair comparison. In addition, we see that TSSQP usually takes less time to complete than SSQP, validating that computing the decomposition and using the backtracking procedure in Algorithm 2 are only modest costs that do not detract from effectiveness the algorithm. Lastly, the final column reports the relative decomposition time, which shows that the computing the decomposition is less than 10% across all problem instances and under 5% for most.

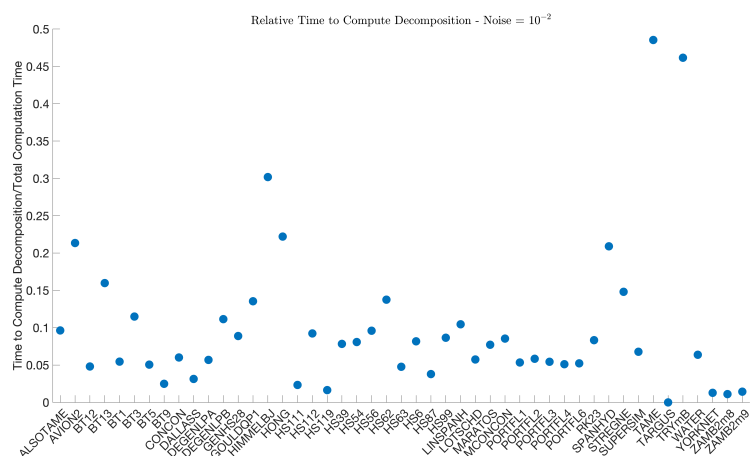


Fig. 11: Relative time to compute decomposition u and v in TSSQP, as measured by (68) with noise level $\epsilon_N = 10^{-2}$.

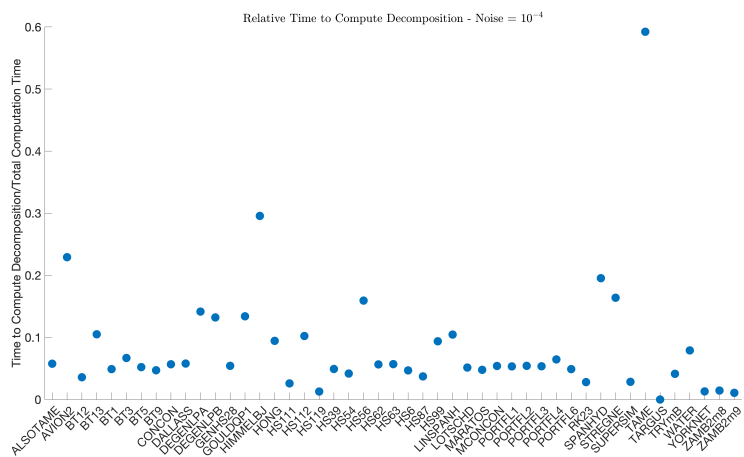


Fig. 12: Relative time to compute decomposition u and v in TSSQP, as measured by (68) with noise level $\epsilon_N = 10^{-4}$.

Table 9: Wallclock times for MLALM, SSQP, and TSSQP on logistic regression problems. The relative decomposition, computed by (68) is the final column for TSSQP.

dataset	batch	MLALM	SSQP	TSSQP	
		Total Time	Total Time	Total Time	Decomp Time
a9a	16	1.19e + 02	4.76e + 01	3.09e + 01	0.05
a9a	128	6.64e + 01	7.43e + 00	5.35e + 00	0.04
ijccn1	16	7.21e + 01	7.87e + 01	4.11e + 01	0.04
ijccn1	128	2.36e + 01	1.19e + 01	8.13e + 00	0.03
ionosphere	16	6.42e - 01	1.39e + 00	7.75e - 01	0.04
ionosphere	128	9.80e - 02	1.52e - 01	4.51e - 02	0.03
madelon	16	1.51e + 01	2.64e + 00	1.29e + 01	0.05
madelon	128	1.11e + 01	3.93e - 01	8.50e - 01	0.03
mushrooms	16	2.67e + 01	6.90e + 00	9.79e + 00	0.04
mushrooms	128	2.56e + 01	1.10e + 00	1.69e + 00	0.03
phishing	16	1.95e + 01	8.28e + 00	1.26e + 01	0.04
phishing	128	8.57e + 00	1.26e + 00	2.44e + 00	0.03
sonar	16	2.46e - 01	1.63e - 01	1.53e + 00	0.09
sonar	128	6.34e - 02	3.08e - 02	2.26e - 02	0.05
splice	16	1.68e + 00	9.33e - 01	1.19e + 00	0.05
splice	128	3.41e - 01	1.14e - 01	9.44e - 02	0.05
w8a	16	3.24e + 02	7.27e + 01	7.27e + 01	0.05
w8a	128	4.89e + 02	9.17e + 00	1.26e + 01	0.03