AN ASYMPTOTICALLY OPTIMAL METHOD FOR CONSTRAINED STOCHASTIC OPTIMIZATION*

BY SEN NA¹, YIHANG GAO², MICHAEL K. NG³, AND MICHAEL W. MAHONEY¹

¹ICSI and Department of Statistics, University of California, Berkeley

²Department of Mathematics, The University of Hong Kong

³Department of Mathematics, Hong Kong Baptist University

We perform statistical inference for the solution of stochastic optimization problems with equality and box inequality constraints. The considered problems are prevalent in statistics and machine learning, encompassing constrained *M*-estimation, PDE-constrained problems, physics-inspired networks, and algorithmic fairness. We introduce a stochastic sequential quadratic programming method (StoSQP) to solve these problems, where we determine the search direction by performing a quadratic approximation of the objective and a linear approximation of the constraints. Despite having access to unbiased estimates of population gradients, a key challenge in constrained problems lies in dealing with the bias in the search direction. To address this challenge, we introduce a novel gradient averaging technique to debias the direction step, leading to Debiased-StoSQP. Our method achieves global almost sure convergence and exhibits local asymptotic normality with an optimal limiting covariance matrix in Hájek and Le Cam's sense. Additionally, a plug-in estimator of the covariance matrix is provided for practical inference purposes. To our knowledge, Debiased-StoSQP is the first fully online method to achieve asymptotic minimax optimality without relying on projection operators to the constraint set, which are incomputable for nonlinear problems. Through extensive experiments on benchmark nonlinear problems in the CUTEst test set, as well as on constrained generalized linear models and portfolio allocation problems, with both synthetic and real data, we demonstrate the superior performance of the method.

1. Introduction. We consider stochastic optimization problems with equality and box inequality constraints, given by the form:

(1.1)
$$\min_{\boldsymbol{x} \in \mathbb{R}^d} \quad f(\boldsymbol{x}) = \mathbb{E}_{\zeta \sim \mathcal{P}} \left[F(\boldsymbol{x}; \zeta) \right],$$
s.t. $\boldsymbol{c}(\boldsymbol{x}) = \boldsymbol{0}, \quad \boldsymbol{\ell} < \boldsymbol{x} < \boldsymbol{u}.$

Here, the vectors ℓ and u denote the lower and upper bounds, respectively, with the symbol " \leq " representing element-wise comparison; and $\zeta \sim \mathcal{P}$ is a random variable. The function $F(\cdot;\zeta): \mathbb{R}^d \to \mathbb{R}$ denotes a realization of the stochastic objective f, and $c: \mathbb{R}^d \to \mathbb{R}^m$ encodes the deterministic equality constraints. Throughout this paper, we assume that f, c, and $F(\cdot;\zeta)$ for each realization ζ are twice continuously differentiable. We aim to develop a *practical*, fully online, and asymptotically optimal method to solve Problem (1.1).

Constraints are useful tools for integrating prior models information, ensuring models' identifiability, and reducing dimensionality. We will provide concrete motivating examples in Section 1.1. Given the ubiquity of Problem (1.1), it is of particular interest to estimate its (local) solution x^* with n samples. Arguably, the most primitive estimator is the classical

^{*}S. Na and Y. Gao contribute equally.

Keywords and phrases: Local asymptotic minimax optimality, Stochastic approximation, Sequential quadratic programming, Nonlinear nonconvex optimization.

M-estimator, where we generate samples $\zeta_1, \ldots, \zeta_n \stackrel{\text{iid}}{\sim} \mathcal{P}$ and solve the constrained problem by replacing the population loss f with the empirical loss \hat{f}_n :

$$\hat{\boldsymbol{x}}_n = rg\min_{\boldsymbol{x}\in\mathbb{R}^d} \hat{f}_n(\boldsymbol{x}) \coloneqq rac{1}{n} \sum_{i=1}^n F(\boldsymbol{x};\zeta_i),$$

s.t. $\boldsymbol{c}(\boldsymbol{x}) = \boldsymbol{0}, \quad \boldsymbol{\ell} \leq \boldsymbol{x} \leq \boldsymbol{u}.$

In fact, the above constrained *M*-estimator is optimal in Hájek and Le Cam's sense [41, 72]. That is, the asymptotic consistency and normality of the minimizer \hat{x}_n is given by

(1.2)
$$\sqrt{n} \left(\hat{\boldsymbol{x}}_n - \boldsymbol{x}^* \right) \xrightarrow{d} \mathcal{N} \left(\boldsymbol{0}, \boldsymbol{L}^{\dagger} \operatorname{Cov} \left(\nabla F(\boldsymbol{x}^*; \zeta) \right) \boldsymbol{L}^{\dagger} \right),$$

where $\boldsymbol{L} = \boldsymbol{P}_J \nabla^2 \mathcal{L}(\boldsymbol{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) \boldsymbol{P}_J$, $\boldsymbol{P}_J = \boldsymbol{I} - \boldsymbol{J}^\top (\boldsymbol{J} \boldsymbol{J}^\top)^\dagger \boldsymbol{J}$ is the projection matrix, \boldsymbol{J} is the Jacobian matrix of active constraints at \boldsymbol{x}^* , and $\mathcal{L}(\boldsymbol{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$ is the Lagrangian function at the optimal primal-dual points.

In this era of large-scale data, optimization problems such as Problem (1.1) have wideranging applications, including but not limited to signal processing, neural network learning, PDE-constrained optimization, and (often randomized) numerical linear algebra. By introducing auxiliary variables (also called slack variables), general equality- and inequalityconstrained problems can be transformed into the form of Problem (1.1). Given this equivalency, the focus of this paper is on developing stochastic optimization algorithms to solve Problem (1.1).

Stochastic optimization algorithms for optimizing an objective f(x) have a rich history and can be traced back (at least) to stochastic gradient descent (SGD), which solves Problem (1.1) in an unconstrained setting. While SGD is computationally and storage-efficient, subsequent research has developed and enhanced its global convergence and local asymptotic properties. For instance, Ruppert [62], Polyak and Juditsky[55] introduced the concept of Polyak-Ruppert averaging, achieving asymptotic normality for averaged iterates. Chen et al. [15] proposed the plug-in estimator and developed a more efficient batch-means estimator to approximate the covariance matrix and estimate the corresponding confidence intervals. Anastasiou et al. [2] developed non-asymptotic convergence rates for normal approximation of SGD with Polyak-Ruppert averaging. Leluc and Potier [42] extend the analysis to conditioned SGD, thereby encompassing a broader class of algorithms like Newton's methods and Quasi-Newton's methods.

Newton's methods are often favored over first-order methods like gradient descent, particularly for their faster convergence rates, which are made possible by incorporating (exact or approximate) Hessian information [38, 49, 79]. Beyond theoretical advantages, Newton's methods, in particular randomized versions [61, 75, 76, 78], have exceptional performance in practical applications. For example, Yao et al. [77] introduced AdaHessian, employng an adaptive Newton's methods to speed up deep neural network training; and Liu et al. [43] then used ideas very similar to AdaHessian to develop Sophia, which reduced the computational cost for training large language models. Although efforts have been made to enhance gradient descent-based algorithms by partially extracting Hessian information [1, 12, 13], the unique benefits of Newton's methods continue to make them a focal point of ongoing research.

Sequential Quadratic Programming (SQP) is recognized as a potent method for tackling constrained optimization problems, particularly when dealing with nonlinear constraints. As Nocedal and Wright emphasized in their seminal work [38], SQP stands as one of the most effective techniques for solving such problems in the deterministic setting. In contrast to deterministic SQP methods, which assume full access to the objective f(x) as described in [9, 38], our work considers a stochastic objective alongside deterministic constraints, as

formulated in Problem (1.1). This paradigm introduces challenges, as the exact values of the objective function, its gradients, and Hessian matrices are generally inaccessible. While recent research has extended SQP algorithms to stochastic settings [5, 18, 19, 20, 23, 27, 47, 48, 51], these works have focused predominantly on problems with only equality constraints. A more exhaustive literature review will be provided in Section 1.3.

Asymptotic analysis serves as a critical tool for a nuanced understanding of the local behavior of iterates in stochastic algorithms. In the context of constrained optimization, we define the primal-dual solution $(x^*, \lambda^*, \mu_1^*, \mu_2^*)$, especially the primal solution x^* , as the optimal solution of the Problem (1.1) with expected objective. The dual variable $(\lambda^*, \mu_1^*, \mu_2^*)$ corresponds to the equality constraints c(x) = 0, the lower-bound box constraints $\ell - x \leq 0$, and the upper-bound box constraints x - u < 0, respectively. While global convergence results offer a broad understanding of the algorithm's behavior, they often fall short in revealing detailed convergence characteristics, especially in the presence of noisy observations related to the objective f(x), gradients, and Hessians. Consider $\{(x_k, \lambda_k, \mu_{1,k}, \mu_{2,k})\}$ as the sequence of primal-dual iterates generated by an algorithm for solving Problem (1.1). The statistical inference drawn from $\{(\boldsymbol{x}_k - \boldsymbol{\lambda}^*, \boldsymbol{\lambda}_k - \boldsymbol{\lambda}^*, \boldsymbol{\mu}_{1,k} - \boldsymbol{\mu}_1^*, \boldsymbol{\mu}_{2,k} - \boldsymbol{\mu}_2^*)\}$ can provide more granular insights. Based on this, we can develop the local asymptotic distributions and estimate associated statistical properties, such as covariance matrices and confidence intervals. The stochastic nature of the iterates is reflected by the asymptotic distribution (especially the confidence interval), provided that the number of iterations is sufficiently large. Such statistical insights offer a quantified measure of confidence and a mechanism to manage uncertainty (and thus inference) in stochastic optimization.

Given these considerations, a natural question is:

• Can we develop an (asymptotically) optimal algorithm, in Hájek and Le Cam's sense [41, 72], and in a manner analogous to the classical M-estimator, for the constrained stochastic optimization Problem (1.1)?

In this paper, we answer this question in the affirmative. To accomplish this, we introduce a novel stochastic Sequential Quadratic Programming algorithm (Debiased-StoSQP and its refinement version Debiased-StoSQP-v2) to solve Problem (1.1), with global almost sure convergence guarantees. We also develop asymptotic normality results and practical estimators for covariance matrices of the generated iterates. The derived limiting covariance matrix matches that of the M-estimator, as in Equation (1.2), showing that our algorithm Debiased-StoSQP-v2 is asymptotically optimal.

Unlike previous analyses of SGD algorithms that focus on averaged iterates [15, 55], our statistical inference targets the last iterate, rendering our approach more aligned with practical applications. Importantly, the presence of inequality constraints in Problem (1.1) introduces a bias in the solutions of the quadratic subproblems for direction estimates, i.e., the obtained search direction is biased. This happens even when $\nabla f(x;\zeta)$ is an unbiased estimators of $\nabla f(x)$. This makes our problem formulation more challenging, compared to problems with only equality constraints [6, 51]. To mitigate the bias, we employ moving averaging techniques for gradient estimation. Our results on the asymptotic normality of iterates establish optimality in terms of the min-max lower bound on the covariance matrix in Hájek and Le Cam's sense [41, 72].

1.1. *Motivating examples.* We now present specific examples from machine learning and statistics that can be cast into the forms of Problem (1.1). We re-emphasize that the general constrained problem can be converted into the form of Problem (1.1) by introducing auxiliary variables, where both forms share the same KKT points (where the first-order optimality condition holds).

NA ET AL.

1.1.1. Constrained regression. In regression models, issues like multicollinearity can lead to unreliable inference results that conflict with both intuition and empirical evidence. One way to mitigate such issues is by incorporating prior information into the model via constraints on the model parameters. For instance, we observed such complexities while working with Poisson regression models for Chicago air pollution and death rate data; further details are discussed in Section 5.4. Constraints can also be an inherent part of the problem formulation itself. For example, in portfolio allocation problems, each entry of x denotes the weight assigned to an asset. Thus, it is common to constrain the estimation within the set $\{x \in \mathbb{R}^d : \mathbf{1}^\top x = 1, x \ge \mathbf{0}\}$. In certain contexts, alternative constraints are imposed for particular purposes, including box constraints $||x||_{\infty} \leq u$ and affine constraints Ax = b [25, 26] (e.g., a negative weight signifies shorting the asset). In semiparametric index models, we impose $\{x \in \mathbb{R}^d : \|x\|_2^2 = 1, x_1 > 0\}$ to make models identifiable [50, 52]. In factor analysis, constraints can prevent Heywood cases (i.e., a negative estimate for the variance) [66]. In algorithmic fairness, constraints can prevent classifiers from yielding disparate outcomes based on sensitive features like gender and ethnicity [80]. For a more comprehensive review of constrained regression models, including different types of constraints, we refer the reader to [22, 24, 26, 50, 52, 65, 67].

1.1.2. *Physics-informed machine learning.* In scientific machine learning, models must adhere to domain knowledge, often described by partial differential equations (PDEs) constraints [17, 33, 40, 53]. In specialized network architectures (like neural ODEs [14], physics-informed neural networks [39, 40, 58], and physics-informed DeepONets [73]), constraints derived from PDEs are applied to the network, i.e., the neural network is motivated from the following constrained optimization:

$$\min_{\boldsymbol{x}} \ \mathcal{L}_{data}(\boldsymbol{x}),$$
s.t. $\mathcal{C}_{PDE}(\boldsymbol{x}) = \mathbf{0},$

where x represents the neural network parameters; $\mathcal{L}_{data}(x)$ and $\mathcal{C}_{PDE}(x) = 0$ are the constructed data fitting loss to be minimized and the governing PDEs, respectively. (It is known that failure to enforce these constraints can lead to serious problems/instabilities in such models [33, 40, 53].)

1.1.3. Adversarial training. Constraints are also frequently employed in adversarial training scenarios. For example, in the training of Wasserstein generative adversarial networks, constraints on the norm of the network parameters are often imposed to ensure model effectiveness [3]. Various types of constraints have been found to improve the adversarial robustness of the model and reduce its sensitivity to perturbations in the input data [16]. In the context of adversarial attacks, constraints are formulated to ensure that the search space of adversarial examples remains close to the original data samples [30, 68, 81].

1.1.4. Constrained neural networks. With advances in computational power and storage, as well as the increasing complexity of problems to solve, the number of parameters in modern neural networks can range from millions to billions or more. This scale is often much larger than the size of available training samples. In such situations, constraints play a crucial role in mitigating overfitting by limiting the flexibility of the network space. One simple yet effective way to improve a neural network's generalization capability is to apply L_2 -norm (regularization) constraints on the network parameters. Well-known implementation tricks include batch normalization and layer normalization are practical tools to regularize neural network parameters to avoid gradient vanishing or explosion and accelerate convergence. Additionally, the importance of parameter constraints for neural networks is seen in generalization analysis [54]; and it is common to work with neural networks that are (implicitly) constrained to interpolate the data [35, 36]. 1.2. *Our contributions*. Our main contribution is to introduce a stochastic SQP algorithm with averaged gradients for solving the constrained optimization Problem (1.1). We summarize our primary contributions, described in more detail in Sections 2-4, here.

- (a) We revisit a standard SQP algorithm (namely, RelaxedSQP, Algorithm 1), which is applicable to deterministic objectives in the form of Problem (1.1), where a relaxation parameter is introduced for the feasibility of the quadratic subproblem. Significantly, we establish a connection between this relaxation scheme and constraint qualifications, providing a deeper understanding of constrained optimization problems.
- (b) Building on the RelaxedSQP, we introduce Debiased-StoSQP (Algorithm 2), a stochastic counterpart of RelaxedSQP. To address bias issues and challenges due to inequality constraints, we employ averaged gradients for debiasing. We introduce separate step sizes, denoted by α_k and β_k , for iterative updates and gradient averaging, respectively; and we prove that the KKT residual of the sequence of iterates $\{x_k\}$, along with least square estimates of dual variables, converges to zero almost surely.
- (c) We perform statistical inference on the iterates generated by Debiased-StoSQP-v2 algorithm (Algorithm 3), which is a refinement of Algorithm 2. To do so, an averaging scheme is introduced for the Hessian, along with an additional update for dual variables. Under mild conditions, we show almost sure convergence of the dual variables $\{(\lambda_k, \mu_{1,k}, \mu_{2,k})\}$ to their optimal values $\{(\lambda^*, \mu_1^*, \mu_2^*)\}$. We establish the asymptotic normality for the iterates,

$$1/\sqrt{\alpha_k^{\min}(\boldsymbol{x}_k-\boldsymbol{\lambda}^*,\boldsymbol{\lambda}_k-\boldsymbol{\lambda}^*,\boldsymbol{\mu}_{1,k}-\boldsymbol{\mu}_1^*,\boldsymbol{\mu}_{2,k}-\boldsymbol{\mu}_2^*)} \stackrel{d}{\longrightarrow} \mathcal{N}\left(\boldsymbol{0},\Theta\boldsymbol{\Omega}^*\right),$$

where $\Theta \Omega^*$ is the Fisher information matrix of the algorithm (more specifically, the covariance matrix revealing the uncertainty of iterates). We achieve asymptotically optimal normality in terms of the covariance matrix, according to the min-max lower bound by Duchi and Ruan [23], in Hájek and Le Cam's sense [41, 72]. We also provide a practical estimator for the unknown covariance matrix Ω^* and show that Ω_k is convergent to Ω^* almost surely, where Ω_k is the estimation of Ω^* (details can be found in Section 4). It is a surprising and novel result that the algorithm with averaged gradients can indeed achieve the asymptotic normality.

There are additional details upon which we would like to elaborate. The concept of using a relaxation technique in the SQP subproblem was initially proposed by Powell [56]. He provided an intuitive explanation that the constrained problem is difficult/challenging if the relaxation technique is invalid. In Section 2, we extend this intuition by providing a more rigorous treatment of the relaxation technique, examining it from the perspective of constraint qualification. Unlike the approach in [20], which relies on increasing sample sizes to reduce bias (a computationally expensive and often impractical process), we use moving averaging techniques. These techniques allow for fully stochastic and online settings, and they achieve convergence with just a single sample of $f(x; \zeta)$, for gradient and Hessian estimation.

A key novelty in our approach is that we introduce two different step sizes with different decay rates for updating the iterates x_k and the gradient, respectively. This can be regarded as a "competition" between the iterates and the gradients. Specifically, for the global almost sure convergence and local asymptotic normality to be achieved, it is essential that the gradients converge faster than the iterates. This ensures that the algorithm is driven by the most current and relevant gradients, contributing to its effective performance. Our work develops the almost sure "lim" results that the KKT residual of generated iterates converges to zero almost surely, with the help of the averaging gradient technique and the least squares estimation on dual variables. This provides a more complete analysis than existing work [20],

which only proved the almost sure "liminf" convergence. We also employ statistical inference techniques to gain insights into the locally asymptotic behavior of Debiased-StoSQP-v2 algorithm. Using recent advancements in martingale difference arrays, our asymptotic analysis is developed to deal with the algorithm with averaging gradients where the gradients are highly correlated. By setting $\alpha_k^{\min} = 1/(k+1)$, we can achieve the asymptotically optimal normality

$$\sqrt{k}(\boldsymbol{x}_k - \boldsymbol{\lambda}^*, \boldsymbol{\lambda}_k - \boldsymbol{\lambda}^*, \boldsymbol{\mu}_{1,k} - \boldsymbol{\mu}_1^*, \boldsymbol{\mu}_{2,k} - \boldsymbol{\mu}_2^*) \overset{d}{\longrightarrow} \mathcal{N}(\boldsymbol{0}, \boldsymbol{\Omega}^*),$$

in terms of the covariance matrix, as shown in [23]. Our estimator Ω_k for Ω^* is more like the plug-in estimator in [15] and we further show that Ω_k is convergent to Ω^* almost surely.

Note that, unlike SGD, SQP is a constrained Newton's method, where we calculate the noisy Hessian matrix $\nabla^2 f(x; \zeta)$ in each iteration. Therefore, the plug-in estimator for the covariance matrix does not significantly increase the computational complexity. Previous works studying the asymptotic behavior of algorithms typically rely on the independence of gradients, while we consider averaged gradients that are highly correlated. The technique we apply involves the use of two distinct step sizes with different rates of decay for iterate and gradient updates. This enables us to balance the convergence behavior of both the iterates and the gradients, thereby achieving asymptotic normality even under conditions of gradient averaging. This represents a significant difference from existing methods, and it brings new perspectives into the study of the asymptotic behavior of algorithms. More details can be found in Section 4.

1.3. Related works. Constrained optimization problems in stochastic settings have gained increasing attention in recent years. Berahas et al. [6] initiated the study of stochastic SQP algorithms, with a focus on equality-constrained problems. They incorporated an ℓ_1 -penalized merit function and adaptive selection mechanisms for both penalty parameters and step sizes to ensure the sufficient decrease of Newton step on the merit function, proving "liminf" convergence for the expectation of the KKT residual. Na and Mahoney [51] extended this line of work by developing an algorithm with inexact subproblem solutions, and they showed the almost sure convergence of the KKT residual based on the sufficient decrease of the exact augmented Lagrangian merit function. An alternative method is the stochastic line search SQP proposed by Na et al. [47], where they adaptively select batch size depending on the decrease of the exact augmented Lagrangian merit function. Their method is more adaptive and powerful than fully stochastic algorithms, due to the growing batch size. However, the stochastic line search method is usually more computationally expensive, and some safeguarded techniques are required in practice, as the batch size cannot grow arbitrarily. Curtis et al. [19] aimed to reduce computational overhead by allowing inexact solutions for the quadratic subproblems, subject to specific termination tests. This approach effectively reduces computational effort, especially in high-dimensional scenarios. Similarly, Na and Mahoney [51] considered the sketch-and-project method in stochastic SQP, a randomized iterative solver introduced in [32], to approximately solve the Newton system in each iteration and reduce the total computation. Berahas et al. [7] also explored variance reduction techniques in gradient approximations, adding robustness to the algorithms at the expense of requiring exact gradient estimations in the outer loops. However, their method still requires exact estimations of gradients, which may be intractable in some applications. Most of the aforementioned works focus on equality-constrained problems, leaving inequality constraints as an area open for further research.

Equality- and inequality-constrained problems pose more formidable challenges than their equality-constrained analogs, particularly due to complications like SQP subproblem infeasibility and solution bias. Recent advancements, such as the method of Na et al. [47], use exact

augmented Lagrangian merit functions, concentrating on identifying each iteration's active set of constraints. However, this approach may impose stringent requirements concerning the linear independence of Jacobians for the active constraints. Alternatively, Curtis et al. [20] introduce an innovative two-stage algorithm. The first stage is designed to improve feasibility by solving a box-constrained, strongly convex quadratic problem. The second stage then zeroes in on optimizing the objective function using quadratic expansion. A comparable two-stage algorithm has been introduced by Qiu and Kungurtsev [57]. The primary distinction between the two methods lies in their handling of stochasticity and step-size selection. Specifically, Curtis et al. mandate increasing the sample size for gradient estimations to ensure convergence and employ adaptive step sizes. In contrast, Qiu and Kungurtsev's approach [57] necessitates a lower bound for the batch size to control gradient uncertainty and uses stochastic line search techniques. Duchi and Ruan [23] have formulated a Riemannian stochastic gradient algorithm that employs dual averaging to address inequality-constrained problems. To guarantee the feasibility of the solution iterates, their method incorporates manifold projections, a technique that tends to be computationally demanding.

There exists an extensive body of work focusing on the statistical properties of SGD and its various adaptations [15, 55]. We begin by reviewing some seminal contributions to the area of SGD. For instance, Toulis and Airoldi [70] introduced the concept of implicit SGD, which achieves asymptotic normality accompanied by an optimal covariance matrix. Mou et al. [45] further contributed by investigating the asymptotic behavior of SGD when fixed step sizes and Polyak-Ruppert averaging are employed in solving linear systems. Duchi and Ruan [23] extended this line of research by developing projected Riemannian SGD and offering statistical inferences for inequality-constrained convex problems. However, the local statistical behavior of stochastic Newton's methods remains relatively unexplored. More recently, Boyer and Godichon-Baggioni [10] turned their focus to the asymptotic normality of an advanced stochastic Quasi-Newton method tailored for regression issues. Na and Mahoney [51] provided a particularly interesting insight by showing that the iterates generated by stochastic SQP in equality-constrained problems tend towards an asymptotic Gaussian distribution with a nearly optimal covariance matrix. The basis for this near-optimality is the min-max lower bounds on the covariance matrices, as proven by Duchi and Ruan [23]. Despite these strides, a significant gap persists in literature concerning local statistical analyses for stochastic algorithms applied to both equality and inequality-constrained problems. Our research bridges this critical lacuna.

1.4. Structure of the paper. Our paper is organized as follows. In Section 2, we revisit the concept of constraint relaxation and establish a link with constraint qualifications. Section 3 is devoted to the introduction of the proposed relaxed stochastic SQP method (Debiased-StoSQP, Algorithm 2), where we also derive its global almost-sure convergence properties. In Section 4, we delve into the algorithm's (Debiased-StoSQP-v2, Algorithm 3) asymptotic behavior, establishing both asymptotic normality and the convergence rate of the iterates generated by our method. We also introduce a practical estimator designed for statistical inference. Section 5 presents experimental results, focusing on applications to CUTEst benchmark problems and regression analyses. Throughout the paper, we provide sketches of proofs following the theorems to aid in comprehension, but we defer all detailed proofs to the Appendices.

1.5. Notations. Throughout the paper, we use $\|\cdot\|_2$ to denote the 2-norm (Euclidean norm) for vectors and the corresponding spectral norm for matrices. The ∞ -norm of a vector, representing the maximal absolute value among its elements, is symbolized as $\|\cdot\|_{\infty}$. We use boldface capital and lowercase letters (e.g., A and a) to denote matrices and vectors respectively. Given a positive integer m, the symbol [m] represents the set $\{1, 2, \dots, m\}$. If we

know without any ambiguity that $\mathcal{I} \subseteq [m]$, then we define $\mathcal{I}^- := [m] \setminus \mathcal{I}$, which implies that $\mathcal{I} \cup \mathcal{I}^- = [m]$ and $\mathcal{I} \cap \mathcal{I}^- = \emptyset$. Let $\mathcal{I} \subseteq [m]$ be the set of indices and A be an $m \times m$ matrix. Then the notation $A_{\mathcal{I}}$ indicates the submatrix composed by columns of A with corresponding columns indices in \mathcal{I} , i.e., $A_{\mathcal{I}} = [a_{i_1}, a_{i_2}, \cdots, a_{i_{|\mathcal{I}|}}]$, where $A = [a_1, a_2, \cdots, a_m]$, $|\mathcal{I}|$ denotes the number of elements in the set \mathcal{I} , and $\mathcal{I} = \{i_1, i_2, \cdots, i_{|\mathcal{I}|}\}$. For an m-dimensional vector a and a set of indices \mathcal{I} , we denote $[a]_{\mathcal{I}}$ as the subvector of the vector a with $[a]_{\mathcal{I}} = (a_{i_1}, a_{i_2}, \cdots, a_{i_{|\mathcal{I}|}})^{\top}$, where $a = (a_1, a_2, \cdots, a_m)^{\top}$ and $\mathcal{I} = \{i_1, i_2, \cdots, i_{|\mathcal{I}|}\}$. Individual elements of a vector a are expressed as either $(a)_i$ or a_i , depending on the context. Unimportant constants are subsumed within the big O notation, $\mathcal{O}(\cdot)$, implying that $f = \mathcal{O}(g)$ if $f \leq C \cdot g$ for some constant C > 0. We use \mathcal{F}_k to denote the σ -algebra defined by event $\{\zeta_i\}_{i=0}^k$. The conditional expectation $\mathbb{E}[\cdot|\mathcal{F}_{k-1}]$ on ζ_k is abbreviated as $\mathbb{E}_k[\cdot]$. We use \odot to denote the element-wise multiplication between two vectors.

2. Constraints Relaxation and Deterministic SQP Algorithm. In this section, we start by considering the deterministic constrained problem, defined as follows:

(2.1)
$$\begin{array}{c} \min_{\boldsymbol{x} \in \mathbb{R}^d} \quad f(\boldsymbol{x}), \\ \text{s.t.} \quad \boldsymbol{c}(\boldsymbol{x}) = \boldsymbol{0}, \qquad \boldsymbol{\ell} < \boldsymbol{x} < \boldsymbol{u}, \end{array}$$

where $f : \mathbb{R}^d \to \mathbb{R}$ is the objective whose derivatives and Hessian are fully accessible under the deterministic setting, and $c : \mathbb{R}^d \to \mathbb{R}^r$ is the equality constraint. Throughout the paper, we assume that the constraints c are second-order continuously differentiable. The vectors ℓ and u define the lower and upper bounds, respectively. Here, we require $-\infty < \ell < u < \infty$ and that the feasible region $\Omega := \{x : c(x) = 0, \ell \le x \le u\}$ is non-empty. At the current iterate x_k , the classical SQP algorithm obtains the search direction by solving the following subproblem:

(2.2)
$$\min_{\boldsymbol{p} \in \mathbb{R}^d} \quad \nabla f(\boldsymbol{x}_k)^\top \boldsymbol{p} + \frac{1}{2} \boldsymbol{p}^\top \boldsymbol{B}_k \boldsymbol{p},$$

s.t. $\boldsymbol{c}(\boldsymbol{x}_k) + \nabla \boldsymbol{c}(\boldsymbol{x}_k)^\top \boldsymbol{p} = \boldsymbol{0}, \quad \boldsymbol{\ell} \leq \boldsymbol{x}_k + \boldsymbol{p} \leq \boldsymbol{u}.$

This is a quadratic expansion of the objective with a linearization of the constraints.

The solution p_k of Problem (2.2) then serves as the direction for updating the variable x in Problem (2.1), i.e., $x_{k+1} = x_k + \alpha_k p_k$, where $\alpha_k > 0$ is the step size. Interestingly, even though the feasible region Ω of the original Problem (1.1) is non-empty, the Problem (2.2) may yield an infeasible region that

$$\Omega_k := \{ oldsymbol{p} : oldsymbol{c}(oldsymbol{x}_k) +
abla oldsymbol{c}(oldsymbol{x}_k)^ op oldsymbol{p} = oldsymbol{0} \} \cap \{oldsymbol{p} : oldsymbol{\ell} \leq oldsymbol{x}_k + oldsymbol{p} \leq oldsymbol{u} \} = \emptyset$$

We demonstrate this through the following example.

EXAMPLE 1. Consider the following constrained optimization problem

 $+3y^{2}$,

$$\min_{(x,y)\in\mathbb{R}^2} \quad x+2x^2$$

(2.3)

s.t.
$$c(\boldsymbol{x}) := x^2 + y^2 - 9 = 0, \quad \boldsymbol{\ell} := \begin{pmatrix} 0 \\ 0 \end{pmatrix} \leq \begin{pmatrix} x \\ y \end{pmatrix} \leq \begin{pmatrix} 3 \\ 2 \end{pmatrix} := \boldsymbol{u}.$$

Here, the feasible region $\Omega = \{(x,y) : x^2 + y^2 - 9 = 0, 0 \le x \le 3, 0 \le y \le 2\}$ is nonempty. If $(x_k, y_k) = (2, 1)$, then the region $\Omega_k = \{(\Delta x, \Delta y) : -4 + 4\Delta x + 2\Delta y = 0, 0 \le 2 + \Delta x \le 3, 0 \le 1 + \Delta y \le 2\}$ is non-empty; but the region Ω_k at $(x_k, y_k) = (1, 1)$ (i.e., $\Omega_k = \{(\Delta x, \Delta y) : -7 + 2\Delta x + 2\Delta y = 0, 0 \le 1 + \Delta x \le 3, 0 \le 1 + \Delta y \le 2\}$) is empty. Fortunately, constraint relaxation provides an effective approach to circumvent this issue of infeasibility in the SQP subproblem. Specifically, we can relax the constraints by introducing a factor $\theta_k \in (0, 1]$, resulting in a relaxed feasible region, defined as

$$\widetilde{\Omega}_k := \{ \boldsymbol{p} : \theta_k \boldsymbol{c}(\boldsymbol{x}_k) + \nabla \boldsymbol{c}(\boldsymbol{x}_k)^\top \boldsymbol{p} = \boldsymbol{0} \} \cap \{ \boldsymbol{p} : \boldsymbol{\ell} \leq \boldsymbol{x}_k + \boldsymbol{p} \leq \boldsymbol{u} \}.$$

For certain $\theta_k \in (0, 1]$, this relaxed feasible region $\widetilde{\Omega}_k$ can be non-empty, even when Ω_k is empty. To illustrate, recall Example 1, where the SQP subproblem yielded an empty region Ω_k at $(x_k, y_k) = (1, 1)$. Introducing the relaxation factor $\theta_k = \frac{1}{2}$ makes the relaxed feasible region $\widetilde{\Omega}_k$ non-empty. This constraint relaxation strategy was originally proposed by Powell [56], whose insight was that the absence of a suitable relaxation parameter signifies that the nonlinear constraints cannot be locally improved in a first-order sense (e.g., linearization).

Naturally, this leads us to investigate the conditions under which a relaxation parameter exists or fails to exist. We found that this is intimately tied to the extended generalized Mangasarian-Fromowitz constraint qualification (EGMFCQ, as defined in Definition 2.4 of [74]). Constraint qualifications serve as conditions that assess the compatibility between nonlinear constraints and their linear approximations. When these qualifications are not met, the linear approximations are inadequate to capture the local geometric properties of the nonlinear constraints. We demonstrate that if x_k moves away from points where EGMFCQ is violated, then $\tilde{\Omega}_k$ is feasible for some $\theta_k \in (0, 1]$. The relationship between constraint relaxation and constraint qualifications is elaborated further in Appendix A.

DEFINITION 1 (EGMFCQ, Definition 2.4 in [74]). The extended generalized Mangasarian-Fromowitz constraint qualification (EGMFCQ) is said to be satisfied at a point $\bar{x} \in \mathbb{R}^d$, with respect to the equality constraints c(x) = 0 and the box constraints $\ell \le x \le u$, if the following conditions are met:

• there is a vector $\boldsymbol{z} \in \mathbb{R}^d$ such that

(2.4)

$$c(\bar{x}) + \nabla c(\bar{x})^{\top} z = \mathbf{0},$$

$$(z)_i > 0, \text{ if } (\bar{x})_i = (\boldsymbol{\ell})_i,$$

$$(z)_i < 0, \text{ if } (\bar{x})_i = (\boldsymbol{u})_i;$$

• columns of $\nabla c(\bar{x})$ are linearly independent.

Remark. It is not difficult to verify that EGMFCQ is weaker than linear independence constraint qualification (LICQ, Definition 2 in Appendix A.1) [38, Definition 12.4]. Note that the point \bar{x} does not necessarily satisfy the equality constraints c(x) = 0, but it is required to lie within the box constraints $\ell \le x \le u$.

LEMMA 1. For Problem (1.1) and the current iterate x_k , if EGMFCQ is satisfied at x_k , then the relaxed feasible region $\widetilde{\Omega}_k$ is nonempty for some $\theta_k \in (0,1]$. Moreover, let θ_k be selected within the interval (0,1] such that $\widetilde{\Omega}_k$ is nonempty with this θ_k but becomes empty when the relaxation parameter θ_k is replaced by min $\{1.1\theta_k, 1\}$. This selection of θ_k can always be achieved. If $\liminf_{k\to\infty} \theta_k = 0$, then there exists an accumulation point x^* of the sequence $\{x_k\}$ where EGMFCQ fails to hold at x^* .

Remark. Here, θ_k is selected such that it approximates the maximal relaxation parameter to make the relaxed region $\tilde{\Omega}_k$ feasible. As indicated by Lemma 1, if the relaxation parameter is not uniformly lower-bounded, a subsequence of $\{x_k\}$ will converge to a non-EGMFCQ point. In light of this, we assume that the maximal relaxation parameter remains

NA ET AL.

lower-bounded throughout the iterative process, as in Assumption 1 (below). At the beginning of each iteration, we first examine the relaxation parameter to ensure that it exceeds a predefined threshold. Failing this, we deduce that the current point is approaching a non-EGMFCQ point, implying that the nonlinear constraints may not be effectively approximated by linearization.

ASSUMPTION 1. For the iterates $\{x_k\}$ generated by the algorithm, there exists $\hat{\theta} \in (0, 1]$ such that the relaxed feasible region $\widetilde{\Omega}_k$ with $\theta_k \leq \tilde{\theta}$ is always nonempty.

Remark. We would like to discuss more the search direction derived from Problem (2.2). It is possible to contemplate a more cost-effective algorithm, where the search direction is first obtained from the equality-constrained QP subproblem, followed by applying projections to the box-constrained regions. However, the following example illustrates that this approach can potentially lead to a stagnation at a specific point, despite being under MFCQ/LICQ conditions.

EXAMPLE 2. Consider the following QP problem

Input: $\ell \le x_0 \le u, \tau, \tilde{\tau} \in (0, 1), \sigma \in (0, 1), \rho_{-1} > 0, \epsilon > 0, \beta \in (0, 1).$

$$\min_{\boldsymbol{p} \in \mathbb{R}^3} \quad \boldsymbol{g}^\top \boldsymbol{p} + \frac{1}{2} \boldsymbol{p}^\top \boldsymbol{B} \boldsymbol{p},$$
s.t. $1 + \boldsymbol{J}^\top \boldsymbol{p} = 0, \quad \begin{pmatrix} 0 \\ 0 \end{pmatrix} \leq \boldsymbol{x} + \boldsymbol{p} \leq \begin{pmatrix} 2 \\ 2 \end{pmatrix},$

where $\boldsymbol{g} = (-1,0)^{\top}$, $\boldsymbol{B} = \boldsymbol{I}_2$, $\boldsymbol{J} = (1,1)^{\top}$ and $\boldsymbol{x} = (1,0)^{\top}$. If we first solve the equalityconstrained problem, the search direction is $\boldsymbol{p} = (0,-1)^{\top}$, then the projection $P(\boldsymbol{x} + \alpha \boldsymbol{p}) = \boldsymbol{x}$ gets stuck for any $\alpha > 0$.

Algorithm 1 RelaxedSQP

1: for $k = 0, 1, 2, \cdots$ do

 $\theta_k = 1;$ 2: 3: while $\hat{\Omega}_k$ with θ_k is empty do $\theta_k = \theta_k \cdot \tilde{\tau};$ 4: 5: end while Compute an positive definite Hessian matrix \boldsymbol{B}_k and the gradient $\nabla f(\boldsymbol{x}_k)$. 6: Solve the relaxed SQP Subproblem (2.7) with θ_k , where the solution is denoted as p_k ; 7: 8: Let $\rho_k^{\text{trial}} = \begin{cases} 0, & \text{if } -\nabla f(\boldsymbol{x}_k)^\top \boldsymbol{p}_k - \boldsymbol{p}_k^\top \boldsymbol{B}_k \boldsymbol{p}_k \ge 0, \\ \frac{\nabla f(\boldsymbol{x}_k)^\top \boldsymbol{p}_k + \boldsymbol{p}_k^\top \boldsymbol{B}_k \boldsymbol{p}_k}{(1-\sigma)\theta_k \|\boldsymbol{c}(\boldsymbol{x}_k)\|_2}, & \text{otherwise;} \end{cases}$ (2.5)and $\rho_{k} = \begin{cases} \rho_{k-1}, & \text{if } \rho_{k}^{\text{trial}} \leq \rho_{k-1}, \\ (1+\epsilon)\rho_{k}^{\text{trial}}, & \text{otherwise}; \end{cases}$ (2.6)9: $\alpha_k = 1;$ while $\phi(\boldsymbol{x}_k + \alpha_k \boldsymbol{p}_k, \rho_k) > \phi(\boldsymbol{x}_k, \rho_k) - \beta \alpha_k \Delta q(\boldsymbol{x}_k, \boldsymbol{p}_k, \nabla f(\boldsymbol{x}_k), \boldsymbol{B}_k; \rho_k)$ do 10: $\alpha_k = \alpha_k \cdot \tau;$ 11: 12: end while 13: $\boldsymbol{x}_{k+1} = \boldsymbol{x}_k + \alpha_k \boldsymbol{p}_k$ 14: end for

Now, we are ready to deal with the standard SQP Subproblem (2.2), which is probably infeasible as discussed, by including a relaxation parameter $\theta_k \in (0, 1]$ in each iteration. More specifically, throughout the paper, instead of solving Problem (2.2) for the search direction in each iteration, we alternatively focus on the following relaxed SQP subproblem:

(2.7)
$$\min_{\boldsymbol{p}\in\mathbb{R}^{n}} \quad \nabla f(\boldsymbol{x}_{k})^{\top}\boldsymbol{p} + \frac{1}{2}\boldsymbol{p}^{\top}\boldsymbol{B}_{k}\boldsymbol{p},$$
$$\boldsymbol{s.t.} \quad \theta_{k}\boldsymbol{c}(\boldsymbol{x}_{k}) + \nabla \boldsymbol{c}(\boldsymbol{x}_{k})^{\top}\boldsymbol{p} = \boldsymbol{0},$$
$$\boldsymbol{\ell} \leq \boldsymbol{x}_{k} + \boldsymbol{p} \leq \boldsymbol{u},$$

for some $\theta_k \in (0, 1]$. We assume that Assumption 1 consistently holds for iterates generated by Algorithms 1-3.

We next offer a concise overview of the line-search technique incorporated into SQP for constrained optimization. To do so, we adopt the ℓ_2 regularized merit function, defined as

(2.8)
$$\phi(\boldsymbol{x};\rho) := f(\boldsymbol{x}) + \rho \|\boldsymbol{c}(\boldsymbol{x})\|_2,$$

to perform the backtracking line search. We define the expanded merit function at x_k with step p_k as

$$q(\boldsymbol{x}_k, \boldsymbol{p}_k, \nabla f(\boldsymbol{x}_k), \boldsymbol{B}_k, \rho_k) := f(\boldsymbol{x}_k) + \nabla f(\boldsymbol{x}_k)^\top \boldsymbol{p}_k + \frac{1}{2} \boldsymbol{p}_k^\top \boldsymbol{B}_k \boldsymbol{p}_k + \rho_k \|\boldsymbol{c}(\boldsymbol{x}_k) + \nabla \boldsymbol{c}(\boldsymbol{x}_k)^\top \boldsymbol{p}_k\|_2,$$

which combines the second-order approximation of the objective with the first-order linearization of constraints, and the corresponding improvement

$$\Delta q(\boldsymbol{x}_{k}, \boldsymbol{p}_{k}, \nabla f(\boldsymbol{x}_{k}), \boldsymbol{B}_{k}, \rho_{k})$$

$$:=q(\boldsymbol{x}_{k}, \boldsymbol{0}, \nabla f(\boldsymbol{x}_{k}), \boldsymbol{B}_{k}, \rho_{k}) - q(\boldsymbol{x}_{k}, \boldsymbol{p}_{k}, \nabla f(\boldsymbol{x}_{k}), \boldsymbol{B}_{k}, \rho_{k})$$

$$(2.10) = -\nabla f(\boldsymbol{x}_{k})^{\top} \boldsymbol{p}_{k} - \frac{1}{2} \boldsymbol{p}_{k}^{\top} \boldsymbol{B}_{k} \boldsymbol{p}_{k} + \rho_{k} \left(\|\boldsymbol{c}(\boldsymbol{x}_{k})\|_{2} - \left\|\boldsymbol{c}(\boldsymbol{x}_{k}) + \nabla \boldsymbol{c}(\boldsymbol{x}_{k})^{\top} \boldsymbol{p}_{k} \right\|_{2} \right)$$

$$= -\nabla f(\boldsymbol{x}_{k})^{\top} \boldsymbol{p}_{k} - \frac{1}{2} \boldsymbol{p}_{k}^{\top} \boldsymbol{B}_{k} \boldsymbol{p}_{k} + \rho_{k} \theta_{k} \|\boldsymbol{c}(\boldsymbol{x}_{k})\|_{2},$$

where the last equality comes from the equality constraints of the relaxed SQP Subproblem, Problem (2.7). To make sufficient improvement, we let $\rho_k > 0$ to be large enough such that $\Delta q(\boldsymbol{x}_k, \boldsymbol{p}_k, \nabla f(\boldsymbol{x}_k), \boldsymbol{B}_k; \rho_k) \geq \frac{1}{2} \boldsymbol{p}_k^\top \boldsymbol{B}_k \boldsymbol{p}_k + \sigma \rho_k \theta_k \|\boldsymbol{c}(\boldsymbol{x}_k)\|_2$ for some $\sigma \in (0, 1)$, i.e., we introduce the following strategy

$$\rho_k^{\text{trial}} = \begin{cases} 0, & \text{if } -\nabla f(\boldsymbol{x}_k)^\top \boldsymbol{p}_k - \boldsymbol{p}_k^\top \boldsymbol{B}_k \boldsymbol{p}_k \ge 0, \\ \frac{\nabla f(\boldsymbol{x}_k)^\top \boldsymbol{p}_k + \boldsymbol{p}_k^\top \boldsymbol{B}_k \boldsymbol{p}_k}{(1-\sigma)\theta_k \| \boldsymbol{c}(\boldsymbol{x}_k) \|_2}, & \text{otherwise;} \end{cases}$$

and

$$\rho_k = \begin{cases} \rho_{k-1}, & \text{if } \rho_k^{\text{trial}} \le \rho_{k-1}, \\ (1+\epsilon)\rho_k^{\text{trial}}, & \text{otherwise;} \end{cases}$$

for some $\epsilon > 0$. Here, the strategy guarantees that the sequence $\{\rho_k\}$ is monotonically increasing and sufficient improvement is secured. We summarize this algorithm in Algorithm 1.

In addition to the feasibility assumption (Assumption 1) on the constraints, the smoothness and boundedness assumptions on the objective and constraints, as stated in Assumption 2, are standard for convergence analysis.

ASSUMPTION 2. The objective function f and the constraints c are second-order continuously differentiable. Then for all $\ell \leq x \leq u$, there exist $M_{\nabla f}, M_{\ell,u}, \kappa_{\nabla f}, \kappa_{\nabla c} > 0$ such that

$$\|\boldsymbol{u} - \boldsymbol{\ell}\|_2 = M_{\boldsymbol{\ell}, \boldsymbol{u}}, \|\nabla f(\boldsymbol{x})\|_2 \leq M_{\nabla f},$$

and for all $\ell \leq x, y \leq u$ one has

$$\left\|\nabla f(\boldsymbol{x}) - \nabla f(\boldsymbol{y})\right\|_{2} \leq \kappa_{\nabla f} \left\|\boldsymbol{x} - \boldsymbol{y}\right\|_{2}, \left\|\nabla \boldsymbol{c}(\boldsymbol{x}) - \nabla \boldsymbol{c}(\boldsymbol{y})\right\|_{2} \leq \kappa_{\nabla c} \left\|\boldsymbol{x} - \boldsymbol{y}\right\|_{2}.$$

The approximate Hessian matrix B_k is positive definite, i.e., $\kappa_1 \mathbf{I} \preceq B_k \preceq \kappa_2 \mathbf{I}$ for some $0 < \kappa_1 \leq \kappa_2$. The Lagrangian multipliers $\{(\boldsymbol{\lambda}_k^{sub}, \boldsymbol{\mu}_{1,k}^{sub}, \boldsymbol{\mu}_{2,k}^{sub})\}$ for the Problem (2.7) are bounded, i.e., there exist $M_{Lag} > 0$ such that

$$\max\{\|\boldsymbol{\lambda}_k^{sub}\|_2, \|\boldsymbol{\mu}_{1,k}^{sub}\|_2, \|\boldsymbol{\mu}_{2,k}^{sub}\|_2\} \le M_{Lag}$$

Here, the boundedness assumption for the Lagrangian multipliers guarantees that the penalty parameter ρ_k is upper bounded, as illustrated by existing literature [8, 11]. Potential concerns regarding the boundedness of Lagrange multipliers are addressed by invoking the EGMFCQ condition. This reveals that the Lagrange multipliers for the SQP subproblems are indeed bounded under EGMFCQ condition. A detailed exposition of this can be found in Appendix A.2, with references [11, 29].

Before delving into the properties of RelaxedSQP (Algorithm 1), it is imperative to articulate the Karush-Kuhn-Tucker (KKT) optimality conditions, as well as the associated KKT residual specific to Problem (1.1). The KKT condition and the corresponding KKT residual for Problem (1.1) at x is formalized as

(2.11)

$$\nabla f(\boldsymbol{x}) + \nabla c(\boldsymbol{x})\boldsymbol{\lambda} - \boldsymbol{\mu}_1 + \boldsymbol{\mu}_2 = \boldsymbol{0},$$

$$c(\boldsymbol{x}) = \boldsymbol{0}, \quad \boldsymbol{\ell} \leq \boldsymbol{x} \leq \boldsymbol{u},$$

$$\boldsymbol{\mu}_1^\top (\boldsymbol{x} - \boldsymbol{\ell}) = 0, \boldsymbol{\mu}_2^\top (\boldsymbol{x} - \boldsymbol{u}) = 0,$$

$$\boldsymbol{\mu}_1 \geq \boldsymbol{0}, \quad \boldsymbol{\mu}_2 \geq \boldsymbol{0},$$

$$(\boldsymbol{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2) = \begin{pmatrix} \nabla f(\boldsymbol{x}) + \nabla c(\boldsymbol{x})\boldsymbol{\lambda} - \boldsymbol{\mu}_1 + \boldsymbol{\mu}_2 \\ c(\boldsymbol{x}) \\ \boldsymbol{\mu}_1 \odot (\boldsymbol{x} - \boldsymbol{\ell}) \\ \boldsymbol{\mu}_2 \odot (\boldsymbol{x} - \boldsymbol{u}) \end{pmatrix}$$

for some dual variables $(\lambda, \mu_1, \mu_2) \in \mathbb{R}^r \times \mathbb{R}^n_+ \times \mathbb{R}^n_+$. Notably, we exclude the inequality constraints $\ell \leq x \leq u$ in the residual definition, as they are intrinsically satisfied by the sequences generated via the proposed RelaxedSQP algorithm. Consequently, if the sequence $\{x_k\}$ with the accompanying Lagrangian multipliers $\{(\lambda_k, \mu_{1,k}, \mu_{2,k})\}$ satisfy $R(x_k, \lambda_k, \mu_{1,k}, \mu_{2,k}) \rightarrow 0$, then any accumulation point $(x^*, \lambda^*, \mu_1^*, \mu_2^*)$ of $\{(x_k, \lambda_k, \mu_1, \mu_2)\}$ satisfies the KKT condition in Equation (2.11), rendering x^* as a KKT (first-order) optimal point.

THEOREM 1. Under Assumptions 1 and 2, there exist sufficiently large $\widetilde{K} \in \mathbb{Z}_+$ and $\tilde{\rho} > 0$, such that $\rho_k = \tilde{\rho}$ for all $k \geq \tilde{K}$ and

(2.12)
$$\lim_{k \to \infty} \|\boldsymbol{p}_k\|_2 = 0 \text{ and } \lim_{k \to \infty} \|\boldsymbol{c}(\boldsymbol{x}_k)\|_2 = 0.$$

Furthermore, if we let $(\lambda_k^{sub}, \mu_{1,k}^{sub}, \mu_{2,k}^{sub})$ be the Lagrangian multipliers of Problem (2.7) at x_k , then

(2.13)
$$\lim_{k\to\infty} \left\| \boldsymbol{R}(\boldsymbol{x}_k,\boldsymbol{\lambda}_k^{sub},\boldsymbol{\mu}_{1,k}^{sub},\boldsymbol{\mu}_{2,k}^{sub}) \right\|_2 = 0.$$

Sketch of the proof. We begin by establishing that the penalty parameter ρ_k stabilizes after a sufficient number of iterations. Specifically, we show that there exists a sufficiently

large integer $\tilde{K} \in \mathbb{Z}_+$ and a constant $\tilde{\rho} > 0$, such that $\rho_k = \tilde{\rho}$ for all $k \ge \tilde{K}$. Capitalizing on the construction of the penalty parameter, we achieve a sufficient improvement in the merit function, formalized as $\Delta q(\boldsymbol{x}_k, \boldsymbol{p}_k, \nabla f(\boldsymbol{x}_k), \boldsymbol{B}_k; \rho_k) \ge \frac{1}{2}\boldsymbol{p}_k^\top \boldsymbol{B}_k \boldsymbol{p}_k + \sigma \rho_k \theta_k \|\boldsymbol{c}(\boldsymbol{x}_k)\|_2$. The lower-boundedness of the merit function, in conjunction with this sufficient improvement, ensures global convergence. This rationale is analogous to the convergence analysis in gradient descent and Newton's methods in unconstrained problems. While unconstrained optimization techniques ensure convergence by achieving a sufficient reduction in the objective function, constrained optimization employs a merit function that amalgamates both the objective function and constraint violations to achieve a similar result. Complete details are provided in Appendix A.3.

The RelaxedSQP specializes to the conventional SQP if the unit relaxation parameter (i.e., $\theta_k = 1$) is accepted, under the condition that x_k is close enough to a feasible and EGMFCQ point. Remarkably, the RelaxedSQP algorithm achieves superlinear local convergence, akin to the classical SQP method, under mild conditions. We do not elaborate on the local superlinear convergence of classical SQP algorithms, which are well-studied and can be found in [9, 38, 44, 63, 71]. Our principal focus here is on the intriguing behavior associated with the acceptance of unit relaxation parameters, which is given in the following lemma. The detailed proof is included in Appendix A.4.

LEMMA 2. Suppose that Assumptions 1 and 2 hold, and $\{x_k\} \to x^*$, where $c(x^*) = 0$ and EGMFCQ condition holds at x^* . Then the unit relaxation parameter $\theta_k = 1$ will be accepted when k is sufficiently large.

3. The Stochastic SQP Algorithm. In this section, we developed a stochastic optimization algorithm with almost sure convergence guarantees (Debiased-StoSQP, Algorithm 2) for solving Problem (1.1). To accomplish this, we first modify the deterministic SQP algorithm (presented in Algorithm 1) into a fully stochastic algorithm (see Algorithm 2), where the averaging gradient is used to reduce the biasedness introduced by the uncertainty in the SQP subproblem. Initially, we will establish the global almost sure "lim inf" convergence for the iterates. Subsequently, we extend these results to achieve the almost sure "lim" convergence by incorporating the least squares estimates of dual variables. Before analyzing the convergence performances of Debiased-StoSQP (Algorithm 2), we describe the algorithm the following steps:

• Step 1: Selection of relaxation parameter. The relaxation parameter is initialized to be 1 for k-th iterate x_k, i.e., θ_k = 1. The feasibility of the region Ω̃_k with θ_k is then assessed. The relaxation parameter θ_k is adjusted iteratively by scaling it down by a factor τ̃ ∈ (0, 1), i.e., θ_k ← θ_k · τ̃, until Ω̃_k is confirmed to be feasible. Under Assumption 1, a suitable relaxation parameter θ_k can be found after at most ⌈log_{τ̃} θ̃⌉ steps. In practice, we include θ̃ ∈ (0, 1] as a tolerance in the algorithm for θ_k. There are various ways to verify the feasibility of Ω̃_k with θ_k. A direct and practical way is to solve the following convex quadratic problem:

$$egin{aligned} \min_{m{p}} & \left\| heta_k m{c}(m{x}_k) +
abla m{c}(m{x}_k)^\top m{p}
ight\|_2^2, \end{aligned}$$
s.t. $m{\ell} \leq m{x}_k + m{p} \leq m{u}, \end{aligned}$

and verify whether the minimal objective is zero. Projected gradient descent and projected Newton's methods are popular and efficient solvers for the box-constrained quadratic problem. If the algorithm terminates with $\theta_k < \tilde{\theta}$, where the small $\tilde{\theta} \in (0, 1]$ is the tolerance included in Algorithm 2, then the iterate x_k approaches an undesirable point, where EGM-FCQ does not hold.

Step 2: Derivative and Hessian estimation. Since the exact derivative and Hessian are inaccessible at x_k, we obtain an estimated derivative g_k := ∇f(x_k; ζ_k) and approximated Hessian B_k. Note that B_k here is an approximation to the Hessian of the Lagrangian L(x, λ, μ₁, μ₂) := f(x) + λ^Tc(x) + μ₁^T (ℓ - x) + μ₂^T (x - u) at the primal variable x_k and the estimated dual variable (λ_k, μ_{1,k}, μ_{2,k}). A practical and cheap way to calculate B_k follows

(3.1)
$$\boldsymbol{B}_{k} = \nabla^{2} f(\boldsymbol{x}_{k}; \zeta_{k}) + \sum_{j=1}^{r} (\boldsymbol{\lambda}_{k-1}^{\mathrm{sub}})_{j} \nabla^{2} c_{j}(\boldsymbol{x}_{k}) + \boldsymbol{\Delta}_{k},$$

where $\lambda_{k-1}^{\text{sub}}$ is the dual variable of the SQP subproblem at (k-1)-th iteration and Δ_k is a regularizer to make B_k positive definite. To achieve the first-order optimality convergence, B_k is not necessarily an accurate approximation to $\nabla^2 \mathcal{L}(x_k, \lambda_k, \mu_{1,k}, \mu_{2,k})$. In fact, B_k is required to be positive definite to achieve a sufficient decrease direction, i.e., $\kappa_1 \mathbf{I} \preceq \mathbf{B}_k \preceq \kappa_2 \mathbf{I}$ for some $0 < \kappa_1 \leq \kappa_2$. For example, the approximate Hessian \mathbf{B}_k is set to an identity matrix in [20]. To achieve the local "optimal" convergence, we expect B_k to be an accurate approximation to the Hessian of Lagrangian, and the averaging technique is employed to reduce the noise of the stochasticity. For more details, see Lemma 4, where we show the almost sure convergence of B_k to the exact Hessian of Lagrangian B^* (defined later). In this part for the first-order global convergence, only the positivedefiniteness of B_k is enforced. Different from [20], where the estimated derivative g_k is directly used in the SQP subproblem, we apply the averaging technique for reducing the noise, i.e., $\bar{g}_k = \bar{g}_{k-1} + \beta_k (g_k - \bar{g}_{k-1})$. The averaging of derivatives is essential to inequality-constrained problems, in both theoretical convergence and experimental performance. It is not difficult to verify that the averaged gradients converge to the exact gradient, at least in expectation, i.e., $\lim_{k\to\infty} \mathbb{E}\left[\|\bar{\boldsymbol{g}}_k - \nabla f(\boldsymbol{x}_k)\|_2^2\right] = 0$, under some mild conditions. However, the simple estimated gradient is not close to the exact gradient, i.e., $\mathbb{E}\left|\left\|\boldsymbol{g}_{k}-\nabla f(\boldsymbol{x}_{k})\right\|_{2}^{2}\right|=\mathcal{O}(1)$. Without the averaging of derivatives, [20] achieves global convergence by reducing the noise level manually, i.e., increasing the sample size during the iterations. However, our algorithms (Debiased-StoSQP and its variant Debiased-StoSQP-v2) are still fully stochastic, i.e., the derivative estimate is only required to have bounded variance, and the noise level is reduced by the imposed averaging.

• Step 3: Obtaining the direction from SQP subproblem. Equipped with the estimated derivative \bar{g}_k , the approximate Hessian B_k and the relaxation parameter θ_k , we acquire the search direction \bar{p}_k as a solution of the following SQP subproblem

(3.2)
$$\min_{\boldsymbol{p} \in \mathbb{R}^n} \quad \bar{\boldsymbol{g}}_k^\top \boldsymbol{p} + \frac{1}{2} \boldsymbol{p}^\top \boldsymbol{B}_k \boldsymbol{p}, \\ \text{s.t.} \quad \theta_k \boldsymbol{c}(\boldsymbol{x}_k) + \nabla \boldsymbol{c}(\boldsymbol{x}_k)^\top \boldsymbol{p} = \boldsymbol{0}, \quad \boldsymbol{\ell} \leq \boldsymbol{x}_k + \boldsymbol{p} \leq \boldsymbol{u}$$

Here, the direction \bar{p}_k is "descent" for the merit function, owing to the convergence of \bar{g}_k to $\nabla f(x_k)$ and the positive-definiteness of B_k .

• Step 4: Adaptive step size selection. We first require that the pre-defined step size {γ_k} decays (asymptotically) polynomially, i.e., γ_k = ι₀ (k + 1)^{-b₁} for some ι₀ > 0 and b₁ ∈ (0,1]. The strategy is similar to the adaptive strategy in the deterministic algorithm. We alternatively select ρ_k such that Δq(x_k, p

_k, g

_k, B_k, ρ_k) ≥ ½p^T_kB_kp_k + σρ_kθ_k ||c(x_k)||₂ for some σ ∈ (0,1). The adaptive parameter ξ_k ≤ ξ_k^{trial} := Δq(x_k, p

_k, B_k, ρ_k) measures the quality of the direction p

_k in reducing the merit function. If ξ_k is large, which implies that p

_k is probably a promising direction, then a more aggressive step size α_k ∝ ξ_kγ_k is preferred, and vise versa. We

select the step size $\alpha_k \in [\alpha_k^{\min}, \alpha_k^{\max}] := \left[\frac{\xi_k \gamma_k}{\kappa_{\nabla f} + \rho_k \kappa_{\nabla c}}, \frac{\xi_k \gamma_k}{\kappa_{\nabla f} + \rho_k \kappa_{\nabla c}} + \varrho \gamma_k^2\right]$, where $\kappa_{\nabla f}$ and $\kappa_{\nabla c}$ are Lipschitz constant for ∇f and ∇c , respectively. We may efficiently estimate the Lipschitz constants by finite differences, an idea quite similar to the Armijo condition.

Step 5: Updating the variable. The primal variable is updated as $x_{k+1} = x_k + \alpha_k \bar{p}_k$. Notably, the algorithm does not necessitate the explicit use of dual variables for updating the primal variable, thus omitting an update scheme for the dual variables in this section. However, the "optimal" local convergence requires an accurate approximation of B_k to the Hessian of the Lagrangian function. This, in turn, demands a satisfactory estimation of dual variables. An update scheme for these dual variables, which is crucial for examining the local convergence properties of the iterates, is provided in Equation (4.1) in Section 4.

Algorithm 2 Debiased-StoSQP

 $\textbf{Input: } \boldsymbol{\ell} \leq \boldsymbol{x}_0 \leq \boldsymbol{u}, \tau, \tilde{\tau} \in (0, 1), \, \sigma \in (0, 1), \, \rho_{-1} > 0, \, \epsilon_{\rho}, \epsilon_{\xi}, \beta \in (0, 1), \, \mu \in (0, 1), \, \varrho > 0, \, \{\beta_k\}_{k=0}^{\infty}, \, \{\gamma_k\}_{k=0}^{\infty}, \, \{\gamma_k\}_{k=0}^$

- 1: for $k = 0, 1, 2, \cdots$ do
- 2: (Step 1.) $\theta_k = 1;$
- while $\widetilde{\Omega}_k$ with θ_k is empty do 3:
- $\theta_k = \theta_k \cdot \tilde{\tau};$ 4:
- 5: end while
- 6: (Step 2.) Compute a positive definite approximate Hessian matrix B_k and the estimated gradient $g_k =$ $\nabla f(\boldsymbol{x}_k; \zeta_k);$
- 7: Let

$$\bar{\boldsymbol{g}}_k = \bar{\boldsymbol{g}}_{k-1} + \beta_k (\boldsymbol{g}_k - \bar{\boldsymbol{g}}_{k-1});$$

- 8: (Step 3.) Solve the relaxed SQP Subproblem (2.7) with θ_k , B_k and \bar{g}_k , where the solution is denoted as \bar{p}_k ;
- 9: (Step 4.) Let (3.3)

$$\rho_{k}^{\text{trial}} = \begin{cases} 0, & \text{if } -\bar{\boldsymbol{g}}_{k}^{\top} \bar{\boldsymbol{p}}_{k} - \bar{\boldsymbol{p}}_{k}^{\top} \boldsymbol{B}_{k} \bar{\boldsymbol{p}}_{k} \ge 0, \\ \frac{\bar{\boldsymbol{g}}_{k}^{\top} \bar{\boldsymbol{p}}_{k} + \bar{\boldsymbol{p}}_{k}^{\top} \boldsymbol{B}_{k} \bar{\boldsymbol{p}}_{k}}{(1 - \sigma) \theta_{k} \| \boldsymbol{c}(\boldsymbol{x}_{k}) \|_{2}}, & \text{otherwise;} \end{cases} \text{ and } \rho_{k} = \begin{cases} \rho_{k-1}, & \text{if } \rho_{k}^{\text{trial}} \le \rho_{k-1}, \\ (1 + \epsilon_{\rho}) \rho_{k}^{\text{trial}}, & \text{otherwise;} \end{cases}$$

10: Let

1

$$(3.4) \quad \xi_k^{\text{trial}} = \frac{\Delta q(\boldsymbol{x}_k, \bar{\boldsymbol{p}}_k, \bar{\boldsymbol{g}}_k, \boldsymbol{B}_k, \rho_k)}{\|\bar{\boldsymbol{p}}_k\|_2^2}, \text{ and } \quad \xi_k = \begin{cases} \xi_{k-1}, & \text{if } \xi_{k-1} \leq \xi_k^{\text{trial}} \\ \min\{(1 - \epsilon_{\xi})\xi_{k-1}, \xi_k^{\text{trial}}\}, & \text{otherwise}; \end{cases}$$

$$11: \quad \text{Select } \alpha_k \in \left[\alpha_k^{\min}, \alpha_k^{\max}\right] := \left[\frac{\xi_k \gamma_k}{\kappa_{\nabla f} + \rho_k \kappa_{\nabla c}}, \frac{\xi_k \gamma_k}{\kappa_{\nabla f} + \rho_k \kappa_{\nabla c}} + \varrho \gamma_k^2\right];$$

$$12: \quad \alpha_k = \min\{\alpha_k, 1/\theta_k\};$$

$$13: \quad (\text{Step 5.}) \ \boldsymbol{x}_{k+1} = \boldsymbol{x}_k + \alpha_k \bar{\boldsymbol{p}}_k.$$

$$14: \text{ end for}$$

We make the following two key assumptions regarding the gradient estimate: that it is unbiased; and that it has bounded variance. Similar to the deterministic algorithm, we assume that the penalty parameter becomes stable after a sufficient number of iterations, in line with existing literature [6, 57]. See Assumption 3 (below). In addition, in Assumption 4 (below), we impose an additional condition stipulating that this stabilized penalty parameter must be sufficiently large, when compared to the corresponding parameter in deterministic algorithms. A similar assumption also appears in [6] for the convergence of the fully stochastic algorithms.

ASSUMPTION 3. Suppose that $g_k := \nabla f(x_k; \zeta_k)$ is a unbiased estimate of $\nabla f(x_k)$, *i.e.*, $\mathbb{E}_k[\boldsymbol{g}_k] = \mathbb{E}_{\zeta}[\nabla f(\boldsymbol{x}_k; \zeta) | \mathcal{F}_{k-1}]$, and there exists a positive number $\sigma_g > 0$ such that $\mathbb{E}_k \|\boldsymbol{g}_k - \nabla f(\boldsymbol{x}_k)\|_2^2 \leq \sigma_g^2$. We further assume that the penalty parameter becomes stable after \bar{K} iterations, i.e., $\rho_k = \bar{\rho}, \forall k \geq \bar{K}$.

ASSUMPTION 4. Suppose that the stable penalty parameter $\bar{\rho}$ for the stochastic algorithm is sufficiently large such that $\Delta q(\boldsymbol{x}_k, \boldsymbol{p}_k, \nabla f(\boldsymbol{x}_k), \boldsymbol{B}_k; \bar{\rho}) \geq \frac{1}{2} \boldsymbol{p}_k^\top \boldsymbol{B}_k \boldsymbol{p}_k + \sigma \bar{\rho} \theta_k \|\boldsymbol{c}(\boldsymbol{x}_k)\|_2$ holds for some $\sigma \in (0, 1)$ and for all $k \geq \bar{K}$.

Under these two additional assumptions, we have the following theorem. This theorem is an intermediate result that the KKT residual of the iterates obtained by the proposed Debiased-StoSQP algorithm achieves the "liminf" convergence. The detailed proofs are available in Appendix B.

THEOREM 2. Under Assumptions 1, 2 and 3, if $\alpha_k^{min} = \iota_1(k+1)^{-b_1}$ and $\beta_k = \iota_2(k+1)^{-b_2}$ for some $\iota_1, \iota_2 > 0$ and some b_1, b_2 satisfying $b_1 \in (\frac{3}{4}, 1]$ and $b_2 \in (2-2b_1, 2b_1-1)$, then

$$\liminf_{k \to \infty} \Delta q(\boldsymbol{x}_k, \boldsymbol{p}_k, \nabla f(\boldsymbol{x}_k), \boldsymbol{B}_k; \bar{\rho}) = 0, \text{ almost surely.}$$

If we further assume that Assumption 4 holds, then

$$\liminf_{k \to \infty} \left[\|\boldsymbol{p}_k\|_2^2 + \|\boldsymbol{c}(\boldsymbol{x}_k)\|_2 \right] = 0, \text{ almost surely.}$$

Furthermore, let $(\lambda_k^{sub}, \mu_{1,k}^{sub}, \mu_{2,k}^{sub})$ be the Lagrangian multipliers of Problem (2.7) at x_k with full gradient $\nabla f(x_k)$, then

(3.5)
$$\liminf_{k \to \infty} \left\| \boldsymbol{R}(\boldsymbol{x}_k, \boldsymbol{\lambda}_k^{sub}, \boldsymbol{\mu}_{1,k}^{sub}, \boldsymbol{\mu}_{2,k}^{sub}) \right\|_2 = 0, \text{ almost surely.}$$

On top of this result, we also aim to enhance this "lim inf" convergence to "lim" convergence by employing the least squares estimates of dual variables, rather than Lagrangian multipliers obtained from subproblems. Given an iterate x_k , the Lagrangian multipliers can be determined from the following least-square optimization problem

$$\min_{\boldsymbol{\lambda}, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2} F(\boldsymbol{\lambda}, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2; \boldsymbol{x}) = \left\| \nabla f(\boldsymbol{x}) + \nabla c(\boldsymbol{x})^\top \boldsymbol{\lambda} - \boldsymbol{\mu}_1 + \boldsymbol{\mu}_2 \right\|_2^2 + \left\| \boldsymbol{\mu}_1 \odot (\boldsymbol{x} - \boldsymbol{\ell}) \right\|_2^2 + \left\| \boldsymbol{\mu}_2 \odot (\boldsymbol{x} - \boldsymbol{u}) \right\|_2^2,$$

s.t. $\boldsymbol{\mu}_1 \ge \mathbf{0}, \boldsymbol{\mu}_2 \ge \mathbf{0}.$

The estimated optimal Lagrangian multipliers $(\lambda_k^*, \mu_{1,k}^*, \mu_{2,k}^*)$ corresponding to x_k serve as one of the feasible solutions of Problem (3.6) evaluated at x_k . The detailed proofs can be found in Appendix B.

THEOREM 3. Under Assumptions 1, 2, 3 and 4, if $\alpha_k^{\min} = \iota_1(k+1)^{-b_1}$ and $\beta_k = \iota_2(k+1)^{-b_2}$ for some $\iota_1, \iota_2 > 0$ and some b_1, b_2 satisfying $b_1 \in (\frac{3}{4}, 1]$ and $b_2 \in (2-2b_1, 2b_1-1)$, then we have

(3.7)
$$\lim_{k \to \infty} \mathbf{R}(\mathbf{x}_k, \boldsymbol{\lambda}_k^*, \boldsymbol{\mu}_{1,k}^*, \boldsymbol{\mu}_{2,k}^*) = \mathbf{0}, \text{ almost surely}.$$

Practical step size selection. In line 10 of Algorithm 2, the step size α_k is chosen from the interval that $\alpha_k \in \left[\frac{\xi_k \gamma_k}{\kappa_{\nabla f} + \rho_k \kappa_{\nabla c}}, \frac{\xi_k \gamma_k}{\kappa_{\nabla f} + \rho_k \kappa_{\nabla c}} + \varrho \gamma_k^2\right]$. By the definition of the adaptivity parameter $\xi_k \leq \xi_k^{\text{trial}} = \frac{\Delta q(\boldsymbol{x}_k, \bar{\boldsymbol{p}}_k, \nabla \bar{\boldsymbol{f}}(\boldsymbol{x}_k), \boldsymbol{B}_k, \rho_k)}{\|\bar{\boldsymbol{p}}_k\|_2^2}$ and the step size $\alpha_k \geq \frac{\xi_k \gamma_k}{\kappa_{\nabla f} + \rho_k \kappa_{\nabla c}}$, there is a potential

risk of underestimation, i.e., $\xi_k \ll \frac{\Delta q(\boldsymbol{x}_k, \bar{\boldsymbol{p}}_k, \nabla \bar{f}(\boldsymbol{x}_k), \boldsymbol{B}_k, \rho_k)}{\|\bar{\boldsymbol{p}}_k\|_2^2}$. This occurs especially due to the non-increasing strategy outlined in Equation (3.4) for the construction of the sequence $\{\xi_k\}$. To address this, we introduce additional flexibility in the upper bound of the step size selection by incorporating $\rho \gamma_k^2$. Defining a more aggressive trial step size $\alpha_k^{\text{trial}} = \frac{\xi_k^{\text{trial}} \gamma_k}{\kappa_{\nabla f} + \rho_k \kappa_{\nabla c}}$, we project the trail step size into the predefined interval $\left[\frac{\xi_k \gamma_k}{\kappa_{\nabla f} + \rho_k \kappa_{\nabla c}}, \frac{\xi_k \gamma_k}{\kappa_{\nabla f} + \rho_k \kappa_{\nabla c}} + \rho \gamma_k^2\right]$, resulting in

(3.8)
$$\alpha_{k} = \begin{cases} \alpha_{k}^{\text{trial}}, & \text{if } \alpha_{k}^{\text{trial}} \leq \frac{\xi_{k}\gamma_{k}}{\kappa_{\nabla f} + \rho_{k}\kappa_{\nabla c}} + \varrho\gamma_{k}^{2}, \\ \frac{\xi_{k}\gamma_{k}}{\kappa_{\nabla f} + \rho_{k}\kappa_{\nabla c}} + \varrho\gamma_{k}^{2}, & \text{otherwise.} \end{cases}$$

4. Asymptotic Normality and Convergence Rate. In this section, we further refine Debiased-StoSQP (Algorithm 2) to Debiased-StoSQP-v2 (Algorithm 3) and show the asymptotic normality property, where Debiased-StoSQP-v2 (Algorithm 3) is asymptotically optimal in Hájek and Le Cam's sense, as shown in Equation (1.2). Compared with Debiased-StoSQP (Algorithm 2), in Debiased-StoSQP-v2 (Algorithm 3), we provide detailed update scheme for approximate Hessian matrix B_k (as shown in Step 2') and dual variables (as shown in Step 6). Note that Debiased-StoSQP-v2 is a special version of Debiased-StoSQP, therefore, they share all properties developed in Section 3.

Recall that in Algorithm 2, the approximate Hessian matrix can be any bounded and positive definite matrices. Here, to show the optimal local convergence of $\{x_k\}$ to a local minimizer x^* , the convergence of the approximate Hessian B_k to the exact Hessian matrix $\nabla^2 f(x^*) + \sum_{i=1}^r (\lambda^*)_i \nabla^2 c(x^*)$ is essential. Besides the update scheme for the primal variable x_k in Algorithm 2, we must include extra updates for dual variables which are only useful for the calculating the approximate Hessian matrix B_k . More specifically, for the current primal-dual variables $(x_k, \lambda_k, \mu_{1,k}, \mu_{2,k})$, the approximate Hessian matrix B_k is estimated by

$$\boldsymbol{B}_{k} = \frac{1}{k} \sum_{i=1}^{k} \left(\nabla^{2} f(\boldsymbol{x}_{i}; \zeta_{i}) + \sum_{j=1}^{r} (\boldsymbol{\lambda}_{i})_{j} \nabla^{2} c_{j}(\boldsymbol{x}_{i}) \right) + \boldsymbol{\Delta}_{k},$$

where Δ_k is a regularization matrix that guarantees the positive-definiteness of B_k . We also include the averaging for the approximate Hessian to reduce the stochasticity and achieve the almost sure convergence (see Lemma 4). We describe it as **Step 2'** in Algorithm 3, since it is a specific case of Step 2 in Algorithm 2. Let $(\bar{p}_k, \lambda_k^{\text{sub}}, \mu_{1,k}^{\text{sub}}, \mu_{2,k}^{\text{sub}})$ be the primal-dual solution of the SQP subproblem

$$\begin{split} \min_{\boldsymbol{p} \in \mathbb{R}^n} \quad \bar{\boldsymbol{g}}_k^\top \boldsymbol{p} + \frac{1}{2} \boldsymbol{p}^\top \boldsymbol{B}_k \boldsymbol{p}, \\ \text{s.t.} \quad \theta_k \boldsymbol{c}(\boldsymbol{x}_k) + \nabla \boldsymbol{c}(\boldsymbol{x}_k)^\top \boldsymbol{p} = \boldsymbol{0}, \quad \boldsymbol{\ell} \leq \boldsymbol{x}_k + \boldsymbol{p} \leq \boldsymbol{u}, \end{split}$$

where \bar{g}_k is the averaged gradient, as in Algorithm 2. Then the dual variables $\lambda_{k+1}, \mu_{1,k+1}, \mu_{2,k+1}$ are obtained by

(4.1)
$$\begin{pmatrix} \boldsymbol{\lambda}_{k+1} \\ \boldsymbol{\mu}_{1,k+1} \\ \boldsymbol{\mu}_{2,k+1} \end{pmatrix} = \begin{pmatrix} \boldsymbol{\lambda}_k \\ \boldsymbol{\mu}_{1,k} \\ \boldsymbol{\mu}_{2,k} \end{pmatrix} + \alpha_k \begin{pmatrix} \boldsymbol{\lambda}_k^{\text{sub}} - \boldsymbol{\lambda}_k \\ \boldsymbol{\mu}_{1,k}^{\text{sub}} - \boldsymbol{\mu}_{1,k} \\ \boldsymbol{\mu}_{2,k}^{\text{sub}} - \boldsymbol{\mu}_{2,k} \end{pmatrix},$$

as in Step 6 in Algorithm 3. Note that the regularization matrix Δ_k guarantees that the Algorithm 3 is a specific case of Algorithm 2. Consequently, the almost sure convergence results

established in Section 3 for Algorithm 2 are also applicable to Algorithm 3. The following assumption requires that the LICQ condition, strictly complementary slackness condition, and strongly convexity conditions, namely second-order sufficient conditions (SOSC), are satisfied at the local minimizer. This kind of local condition is commonly considered crucial for analyzing local convergence behaviors, both in unconstrained and constrained optimization problems [38].

Algorithm 3 Debiased-StoSQP-v2

Input: $\ell \le x_0 \le u, \tau, \tilde{\tau} \in (0, 1), \sigma \in (0, 1), \rho_{-1} > 0, \epsilon_{\rho}, \epsilon_{\xi}, \beta \in (0, 1), \mu \in (0, 1), \varrho > 0, \{\beta_k\}_{k=0}^{\infty}, \{\gamma_k\}_{k=0}^{\infty}, \{\gamma_k\}_{k=0}^{\infty$

- 1: for $k = 0, 1, 2, \cdots$ do
- 2: (Step 1.) $\theta_k = 1$;
- 3: while $\tilde{\Omega}_k$ with θ_k is empty do
- 4: $\theta_k = \theta_k \cdot \tilde{\tau};$
- 5: end while
- 6: (Step 2'). Compute the estimated gradient $\boldsymbol{g}_k = \nabla f(\boldsymbol{x}_k; \zeta_k)$ and let

$$\bar{\boldsymbol{g}}_k = \bar{\boldsymbol{g}}_{k-1} + \beta_k (\boldsymbol{g}_k - \bar{\boldsymbol{g}}_{k-1})$$

7: Compute the positive definite approximate Hessian matrix B_k , i.e.,

$$\boldsymbol{B}_{k} = \frac{1}{k} \sum_{i=1}^{k} \left(\nabla^{2} f(\boldsymbol{x}_{i}; \zeta_{i}) + \sum_{j=1}^{r} (\boldsymbol{\lambda}_{i})_{j} \nabla^{2} c_{j}(\boldsymbol{x}_{i}) \right) + \boldsymbol{\Delta}_{k},$$

where Δ_k is a regularization matrix that guarantees the positive-definiteness of B_k .

- 8: (Step 3.) Solve the relaxed SQP Subproblem (2.7) with θ_k , B_k and \bar{g}_k , where the primal-dual solution is denoted as $(\bar{p}_k, \lambda_k^{\text{sub}}, \mu_{1,k}^{\text{sub}}, \mu_{2,k}^{\text{sub}})$;
- 9: (Step 4.) Let

$$\rho_{k}^{\text{trial}} = \begin{cases} 0, & \text{if } -\bar{\boldsymbol{g}}_{k}^{\top} \bar{\boldsymbol{p}}_{k} - \bar{\boldsymbol{p}}_{k}^{\top} \boldsymbol{B}_{k} \bar{\boldsymbol{p}}_{k} \ge 0, \\ \frac{\bar{\boldsymbol{g}}_{k}^{\top} \bar{\boldsymbol{p}}_{k} + \bar{\boldsymbol{p}}_{k}^{\top} \boldsymbol{B}_{k} \bar{\boldsymbol{p}}_{k}}{(1 - \sigma) \theta_{k} \| \boldsymbol{c}(\boldsymbol{x}_{k}) \|_{2}}, & \text{otherwise;} \end{cases} \text{ and } \rho_{k} = \begin{cases} \rho_{k-1}, & \text{if } \rho_{k}^{\text{trial}} \le \rho_{k-1}, \\ (1 + \epsilon_{\rho}) \rho_{k}^{\text{trial}}, & \text{otherwise;} \end{cases}$$

10: Let

$$\boldsymbol{\xi}_{k}^{\text{trial}} = \frac{\Delta q(\boldsymbol{x}_{k}, \bar{\boldsymbol{p}}_{k}, \bar{\boldsymbol{g}}_{k}, \boldsymbol{B}_{k}, \rho_{k})}{\|\bar{\boldsymbol{p}}_{k}\|_{2}^{2}}, \text{ and } \boldsymbol{\xi}_{k} = \begin{cases} \boldsymbol{\xi}_{k-1}, & \text{if } \boldsymbol{\xi}_{k-1} \leq \boldsymbol{\xi}_{k}^{\text{trial}} \\ \min\{(1 - \boldsymbol{\epsilon}_{\xi})\boldsymbol{\xi}_{k-1}, \boldsymbol{\xi}_{k}^{\text{trial}}\}, & \text{otherwise;} \end{cases}$$

11: Select
$$\alpha_k \in \left[\alpha_k^{\min}, \alpha_k^{\max}\right] := \left[\frac{\xi_k \gamma_k}{\kappa_{\nabla f} + \rho_k \kappa_{\nabla c}}, \frac{\xi_k \gamma_k}{\kappa_{\nabla f} + \rho_k \kappa_{\nabla c}} + \varrho \gamma_k^2\right];$$

12: $\alpha_k = \min\{\alpha_k, 1/\theta_k\};$

13: **(Step 5.)**
$$x_{k+1} = x_k + \alpha_k \bar{p}_k$$
.

14: (Step 6.)
$$\begin{pmatrix} \boldsymbol{\lambda}_{k+1} \\ \boldsymbol{\mu}_{1,k+1} \\ \boldsymbol{\mu}_{2,k+1} \end{pmatrix} = \begin{pmatrix} \boldsymbol{\lambda}_k \\ \boldsymbol{\mu}_{1,k} \\ \boldsymbol{\mu}_{2,k} \end{pmatrix} + \alpha_k \begin{pmatrix} \Delta \boldsymbol{\lambda}_k \\ \Delta \boldsymbol{\mu}_{1,k} \\ \Delta \boldsymbol{\mu}_{2,k} \end{pmatrix}$$
, where $\begin{pmatrix} \Delta \boldsymbol{\lambda}_k \\ \Delta \boldsymbol{\mu}_{1,k} \\ \Delta \boldsymbol{\mu}_{2,k} \end{pmatrix} := \begin{pmatrix} \boldsymbol{\lambda}_k^{\text{sub}} - \boldsymbol{\lambda}_k \\ \boldsymbol{\mu}_{1,k}^{\text{sub}} - \boldsymbol{\mu}_{1,k} \\ \boldsymbol{\mu}_{2,k}^{\text{sub}} - \boldsymbol{\mu}_{2,k} \end{pmatrix}$.

15: end for

ASSUMPTION 5. We assume that the generated sequence $\{x_k\}$ is convergent almost surely to a strict local solution x^* , where (i) LICQ holds for active constraints at x^* ; (ii) strictly complementary slackness condition holds, i.e., $(\mu_1^*)_i > 0$ if $(x)_i = (\ell)_i$, and $(\mu_2^*)_i > 0$ if $(x)_i = (u)_i$; (iii) $\nabla^2 f(x^*) + \sum_{i=1}^r (\lambda^*)_i \nabla^2 c_i(x^*)$ is positive definite.

LEMMA 3. Under Assumptions 2 and 5, the followings hold:

1. $p_k \rightarrow 0$ almost surely, where p_k is the solution of the relaxed SQP subproblem at x_k with exact gradient $\nabla f(x_k)$ and the approximate Hessian matrix B_k ;

- 2. $\bar{g}_k \nabla f(\boldsymbol{x}_k) \rightarrow \mathbf{0}$, almost surely;
- 3. there exist sufficiently sufficiently large K^* , such that $\mathcal{I}(\boldsymbol{x}_k + \bar{\boldsymbol{p}}_k) = \mathcal{I}(\boldsymbol{x}_k + \boldsymbol{p}_k) = \mathcal{I}(\boldsymbol{x}^*)$ and $\mathcal{J}(\boldsymbol{x}_k + \bar{\boldsymbol{p}}_k) = \mathcal{J}(\boldsymbol{x}_k + \boldsymbol{p}_k) = \mathcal{J}(\boldsymbol{x}^*)$, for $k \ge K^*$;
- 4. $(\boldsymbol{x}_k, \boldsymbol{\lambda}_k, \boldsymbol{\mu}_{1,k}, \boldsymbol{\mu}_{2,k}) \rightarrow (\boldsymbol{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}_1^*, \boldsymbol{\mu}_2^*)$ almost surely.

We put the detailed proof in Appendix C.1. In Lemma 3, we show that: (1) the exact search direction p_k converges almost surely to zero, implying that the iterate x_k approximately satisfies KKT conditions; (2) the averaged gradients are arbitrarily close to the exact gradients after a sufficient number of iterations, due to the updating schemes for the averaged gradients and iterates; (3) the active and inactive sets of constraints can be correctly identified; and (4) the primal-dual iterates converge almost surely to the optimal solution. By the correct identification of active and inactive sets in Lemma 3, the KKT condition and the strong convexity of the SQP subproblem further imply that p_k and \bar{p}_k are the solution of the following equality-constrained problems, respectively:

$$\begin{aligned} \boldsymbol{p}_{k} &= \operatorname*{arg\,min}_{\boldsymbol{p} \in \mathbb{R}^{d}} \quad \nabla f(\boldsymbol{x}_{k})^{\top} \boldsymbol{p} + \frac{1}{2} \boldsymbol{p}^{\top} \boldsymbol{B}_{k} \boldsymbol{p}, \qquad \bar{\boldsymbol{p}}_{k} = \operatorname*{arg\,min}_{\boldsymbol{p} \in \mathbb{R}^{d}} \quad \bar{\boldsymbol{g}}_{k}^{\top} \boldsymbol{p} + \frac{1}{2} \boldsymbol{p}^{\top} \boldsymbol{B}_{k} \boldsymbol{p}, \\ \text{s.t.} \quad \boldsymbol{c}(\boldsymbol{x}_{k}) + \nabla \boldsymbol{c}(\boldsymbol{x}_{k})^{\top} \boldsymbol{p} = \boldsymbol{0}, \qquad \text{and} \quad \text{s.t.} \quad \boldsymbol{c}(\boldsymbol{x}_{k}) + \nabla \boldsymbol{c}(\boldsymbol{x}_{k})^{\top} \boldsymbol{p} = \boldsymbol{0}, \\ (\boldsymbol{x}_{k} + \boldsymbol{p})_{i} = (\boldsymbol{\ell})_{i}, \text{ for } i \in \mathcal{I}(\boldsymbol{x}^{*}), \qquad (\boldsymbol{x}_{k} + \boldsymbol{p})_{i} = (\boldsymbol{\ell})_{i}, \text{ for } i \in \mathcal{I}(\boldsymbol{x}^{*}), \\ (\boldsymbol{x}_{k} + \boldsymbol{p})_{i} = (\boldsymbol{u})_{i}, \text{ for } i \in \mathcal{J}(\boldsymbol{x}^{*}), \qquad (\boldsymbol{x}_{k} + \boldsymbol{p})_{i} = (\boldsymbol{u})_{i}, \text{ for } i \in \mathcal{J}(\boldsymbol{x}^{*}). \end{aligned}$$

Without the loss of generality, we assume that the relaxation parameter is unit according to Lemma 2. The LICQ condition at x^* also implies the LICQ at x_k when x_k is sufficiently close to x^* . Then the KKT system of Problem (4.2) shows that (4.3)

$$\begin{pmatrix} \bar{\boldsymbol{p}}_{k} \\ \Delta \boldsymbol{\lambda}_{k} \\ [\Delta \boldsymbol{\mu}_{1,k}]_{\mathcal{I}(\boldsymbol{x}^{*})} \\ [\Delta \boldsymbol{\mu}_{2,k}]_{\mathcal{J}(\boldsymbol{x}^{*})} \end{pmatrix} = \begin{pmatrix} \boldsymbol{B}_{k} \quad \nabla \boldsymbol{c}(\boldsymbol{x}_{k}) \ [-\boldsymbol{I}]_{\mathcal{I}(\boldsymbol{x}^{*})} \ [\boldsymbol{I}]_{\mathcal{J}(\boldsymbol{x}^{*})} \\ \nabla \boldsymbol{c}(\boldsymbol{x}_{k})^{\top} \quad \boldsymbol{0} \quad \boldsymbol{0} \quad \boldsymbol{0} \\ [-\boldsymbol{I}]_{\mathcal{I}(\boldsymbol{x}^{*})}^{\top} \quad \boldsymbol{0} \quad \boldsymbol{0} \quad \boldsymbol{0} \\ [\boldsymbol{I}]_{\mathcal{J}(\boldsymbol{x}^{*})}^{\top} \quad \boldsymbol{0} \quad \boldsymbol{0} \quad \boldsymbol{0} \quad \boldsymbol{0} \\ [\boldsymbol{I}]_{\mathcal{J}(\boldsymbol{x}^{*})}^{\top} \quad \boldsymbol{0} \quad \boldsymbol{0} \quad \boldsymbol{0} \quad \boldsymbol{0} \\ \boldsymbol{I}]_{\mathcal{J}(\boldsymbol{x}^{*})}^{\top} \quad \boldsymbol{0} \quad \boldsymbol{0} \quad \boldsymbol{0} \quad \boldsymbol{0} \quad \boldsymbol{0} \\ [\boldsymbol{I}]_{\mathcal{J}(\boldsymbol{x}^{*})}^{\top} \quad \boldsymbol{0} \quad \boldsymbol{0} \quad \boldsymbol{0} \quad \boldsymbol{0} \quad \boldsymbol{0} \\ \boldsymbol{I}]_{\mathcal{J}(\boldsymbol{x}^{*})}^{\top} \quad \boldsymbol{0} \quad \boldsymbol{0} \quad \boldsymbol{0} \quad \boldsymbol{0} \quad \boldsymbol{0} \quad \boldsymbol{0} \\ \boldsymbol{I}]_{\mathcal{J}(\boldsymbol{x}^{*})}^{\top} \quad \boldsymbol{0} \quad \boldsymbol{0}$$

and

$$\left[\Delta\boldsymbol{\mu}_{2,k}\right]_{\mathcal{J}^{-}(\boldsymbol{x}^{*})} = -\left[\boldsymbol{\mu}_{2,k}\right]_{\mathcal{J}^{-}(\boldsymbol{x}^{*})},$$

under the almost sure convergence of primal-dual iterates and conditions in Assumption 5. For the parameters and the step size, we only consider the case where the penalty parameter ρ_k and adaptivity parameter ξ_k become stable, and we let $\alpha_k^{\min} = \iota_1(k+1)^{-b_1}$ and $\alpha_k^{\max} = \iota_1(k+1)^{-b_1} + \iota_0(k+1)^{-2b_1}$, where $\alpha_k \in [\alpha_k^{\min}, \alpha_k^{\max}]$. We denote the Jacobian matrix of the (estimated) KKT system at \boldsymbol{x}_k and \boldsymbol{x}^* as

$$m{H}_k = egin{pmatrix} m{B}_k &
abla m{c}(m{x}_k) \left[-m{I}
ight]_{\mathcal{I}(m{x}^*)} \left[m{I}
ight]_{\mathcal{J}(m{x}^*)} \ m{
abla} & m{0} & m{0} \ \left[-m{I}
ight]_{\mathcal{I}(m{x}^*)} & m{0} & m{0} & m{0} \ \left[-m{I}
ight]_{\mathcal{I}(m{x}^*)} & m{0} & m{0} & m{0} \ \left[m{I}
ight]_{\mathcal{J}(m{x}^*)}^\top & m{0} & m{0} & m{0} \ m{0} \ m{0} \ m{0} \ m{J} \ m{J}_{\mathcal{I}(m{x}^*)} & m{0} & m{0} & m{0} \ m{$$

and

respectively. Let the core covariance (also the Fisher information matrix of the algorithm) at x^* be defined as

(4.4)

$$\boldsymbol{\Omega}^{*} = \boldsymbol{H}^{*-1} \begin{pmatrix} \mathbb{E} \left[\nabla f(\boldsymbol{x}^{*}; \zeta) \nabla f(\boldsymbol{x}^{*}; \zeta)^{\top} \right] - \nabla f(\boldsymbol{x}^{*}) \nabla f(\boldsymbol{x}^{*})^{\top} \mathbf{0} \mathbf{0} \mathbf{0} \\ \mathbf{0} & \mathbf{0} \mathbf{0} \mathbf{0} \\ \mathbf{0} & \mathbf{0} \mathbf{0} \mathbf{0} \\ \mathbf{0} & \mathbf{0} \mathbf{0} \mathbf{0} \end{pmatrix} \boldsymbol{H}^{*-1} := \boldsymbol{H}^{*-1} \boldsymbol{\Sigma} \boldsymbol{H}^{*-1}$$

LEMMA 4. Under Assumptions 2 and 5, we have $B_k \to B^*$ and $H_k \to H^*$ almost surely, where $B^* := \nabla^2 f(x^*) + \sum_{i=1}^r (\lambda^*)_i \nabla^2 c_i(x^*)$.

The proof for Lemma 4 can be found in Appendix C.2. According to the almost sure convergence in Assumption 5, and Lemmas 3 and 4, we deduce that there exists a sufficiently large integer K^* , such that the active set of box inequalities remains constant. Specifically, $\mathcal{I}(\boldsymbol{x}_k + \bar{\boldsymbol{p}}_k) = \mathcal{I}(\boldsymbol{x}_k + \boldsymbol{p}_k) = \mathcal{I}(\boldsymbol{x}^*)$ and $\mathcal{J}(\boldsymbol{x}_k + \bar{\boldsymbol{p}}_k) = \mathcal{J}(\boldsymbol{x}_k + \boldsymbol{p}_k) = \mathcal{J}(\boldsymbol{x}^*)$, for $k \ge K^*$. Therefore, we can equivalently consider the equality-constrained SQP subproblem as given by Problem (4.2) for $k \ge K^*$.

ASSUMPTION 6. Assume that the random gradient has finite third-moment, i.e., the conditioned expectation

$$\mathbb{E}\left[\left\|\boldsymbol{g}_{k}-\nabla f(\boldsymbol{x}_{k})\right\|_{2}^{3}|\mathcal{F}_{k-1}\right]=\mathbb{E}_{\zeta}\left[\left\|\nabla f(\boldsymbol{x}_{k};\zeta)-\nabla f(\boldsymbol{x}_{k})\right\|_{2}^{3}|\mathcal{F}_{k-1}\right]\leq M_{m}$$

for some $M_m > 0$ and all \boldsymbol{x}_k in the feasible region $\{\boldsymbol{x} : \boldsymbol{\ell} \leq \boldsymbol{x} \leq \boldsymbol{u}\}$. The Lipschitzness of stochastic gradient holds, i.e., $\mathbb{E}\left[\|\nabla f(\boldsymbol{x};\zeta) - \nabla f(\boldsymbol{y};\zeta)\|_2^2\right] \leq \kappa_{\nabla f}^2 \|\boldsymbol{x} - \boldsymbol{y}\|_2^2$ for some $\kappa_{\nabla f} > 0$.

THEOREM 4. Under Assumptions 2, 5 and 6, and suppose that $\iota_1 > b_2$ if $b_1 = 1$, and $b_2 > \frac{1}{2}b_1$, then

(4.5)
$$\frac{1}{\sqrt{\alpha_k^{min}}} \begin{pmatrix} \boldsymbol{x}_{k+1} - \boldsymbol{x}^* \\ \boldsymbol{\lambda}_{k+1} - \boldsymbol{\lambda}^* \\ [\boldsymbol{\mu}_{1,k+1} - \boldsymbol{\mu}_1^*]_{\mathcal{I}(\boldsymbol{x}^*)} \\ [\boldsymbol{\mu}_{2,k+1} - \boldsymbol{\mu}_2^*]_{\mathcal{J}(\boldsymbol{x}^*)} \end{pmatrix} \overset{d}{\longrightarrow} \mathcal{N}(\boldsymbol{0}, \Theta \boldsymbol{\Omega}^*),$$

and

(4.6)
$$\left\| \begin{pmatrix} [\boldsymbol{\mu}_{1,k+1} - \boldsymbol{\mu}_{1}^{*}]_{\mathcal{I}^{-}(\boldsymbol{x}^{*})} \\ [\boldsymbol{\mu}_{2,k+1} - \boldsymbol{\mu}_{2}^{*}]_{\mathcal{J}^{-}(\boldsymbol{x}^{*})} \end{pmatrix} \right\|_{2} = \begin{cases} o\left(\alpha_{k}^{min}\right), & \text{if } b_{1} < 1, \\ \mathcal{O}\left(\alpha_{k}^{min}\right), & \text{if } b_{1} = 1. \end{cases}$$

where

$$\Theta := \begin{cases} 1/2, & \text{if } b_1 < 1, \\ 1/\left(2 - \frac{1}{\iota_1}\right), & \text{if } b_1 = 1. \end{cases}$$

Sketch of proof: We start by decomposing the primal-dual variable

$$(x_{k+1} - x^*, \lambda_{k+1} - \lambda^*, [\mu_{1,k+1} - \mu_1^*]_{\mathcal{I}(x^*)}, [\mu_{2,k+1} - \mu_2^*]_{\mathcal{J}(x^*)})$$

into three terms $\mathcal{Q}_{1,k}$, $\mathcal{Q}_{2,k}$ and $\mathcal{Q}_{3,k}$. Using the central limit theorem for martingale difference array, we establish that $\frac{1}{\sqrt{\alpha_k^{\min}}}\mathcal{Q}_{1,k} \xrightarrow{d} \mathcal{N}(\mathbf{0},\Theta\mathbf{\Omega}^*)$. Under the given conditions, we

can further show the remaining two terms satisfying $\mathbb{E}[\mathcal{Q}_{2,k}] = o\left(\sqrt{\alpha_k^{\min}}\right)$ and $\mathbb{E}[\mathcal{Q}_{3,k}] = o\left(\sqrt{\alpha_k^{\min}}\right)$. Then, the result is obtained by Slutsky's theorem. Detailed proof can be found

 $o(\sqrt{\alpha_k^{\min}})$. Then, the result is obtained by Slutsky's theorem. Detailed proof can be found in Appendix C.3.

It is a surprising and novel result that the algorithm with averaged gradients can achieve asymptotic normality. Previous works [15, 42, 51, 70] studying the asymptotic normality of algorithms mostly rely on the independence of gradients. However, the averaged gradients in our algorithm are highly correlated. The key idea here is the introduction of two distinct step sizes, with different decay rates for iterates and gradient updates. This can be regarded as a "competition" between the iterates and the gradients. Specifically, for asymptotic normality to be achieved, it is essential that the gradients converge faster than the iterates. This ensures that the algorithm is driven by the most current and relevant gradients, contributing to its effective performance. To the best of our knowledge, it is the first work establishing the asymptotic normality for the algorithm with averaged gradients.

COROLLARY 1. Under Assumptions 2, 5 and 6, and let $\iota_1 = 1$, $b_1 = 1$, $\iota_2 > 0$ and $b_2 \in (\frac{1}{2}, 1)$, then

(4.7)
$$\sqrt{k} \begin{pmatrix} \boldsymbol{x}_{k+1} - \boldsymbol{x}^* \\ \boldsymbol{\lambda}_{k+1} - \boldsymbol{\lambda}^* \\ [\boldsymbol{\mu}_{1,k+1} - \boldsymbol{\mu}_1^*]_{\mathcal{I}(\boldsymbol{x}^*)} \\ [\boldsymbol{\mu}_{2,k+1} - \boldsymbol{\mu}_2^*]_{\mathcal{J}(\boldsymbol{x}^*)} \end{pmatrix} \stackrel{d}{\longrightarrow} \mathcal{N}(\boldsymbol{0}, \boldsymbol{\Omega}^*) \,.$$

The asymptotic normality for M-estimators in Equation (1.2) establishes a lower bound for stochastic algorithms in solving Problem (1.1), in Hájek and Le Cam's sense [41, 72]. Our result in Corollary 1 complements this by demonstrating that the optimal local convergence behavior of the primal variable is achieved, as characterized by the covariance matrix Ω^* . Here, considering only the primal variable x_k , our results in Equation (4.7) match Equation (1.2) by applying knowledge of block matrix inverse [69, Corollary 2.3].

A practical estimator of the covariance matrix. Here, we provide a practical plug-in estimator for the unknown covariance matrix Ω^* . We then show that the estimator is convergent to the exact one, and therefore, it can be adopted in analyzing the local behavior of algorithms and conducting statistical inference. Let

(4.8)
$$\boldsymbol{\Omega}_k = \boldsymbol{H}_k^{-1} \boldsymbol{\Sigma}_k \boldsymbol{H}_k^{-1},$$

where

$$\boldsymbol{\Sigma}_{k} = \left(\begin{array}{c} \frac{1}{k+1} \sum_{i=0}^{k} \boldsymbol{g}_{i} \boldsymbol{g}_{i}^{\top} - \left(\frac{1}{k+1} \sum_{i=0}^{k} \boldsymbol{g}_{i} \right) \left(\frac{1}{k+1} \sum_{i=0}^{k} \boldsymbol{g}_{i} \right)^{\top} \boldsymbol{0} \\ \boldsymbol{0} \end{array} \right).$$

Observe that Ω_k can be cheaply calculated in each iteration, without additional sampling/estimation of gradients, since it shares the gradient estimate g_i with averaged gradients \bar{g}_k . The estimator Ω_k for the covariance matrix Ω^* can be used as a surrogate to the exact matrix for analyzing the local behavior of algorithms and conducting statistical inference. The following theorem establishes the almost sure convergence of the practical plug-in estimator to the exact covariance matrix. The proof can be found in Appendix C.4.

THEOREM 5. Under Assumptions 2, 5 and 6, and suppose that $\iota_1 > b_2$ if $b_1 = 1$, and $b_2 > \frac{1}{2}b_1$, then $\Sigma_k \to \Sigma^*$ and $\Omega_k \to \Omega^*$, almost surely.

NA ET AL.

5. Experiments. In this section, we describe comprehensive experiments we have conducted to demonstrate the effectiveness of the Debiased-StoSQP-v2 (Algorithm 3). Specifically, we applied it to a variety of problems, including benchmark optimization problems from CUTEst library [28, 31] as well as constrained regression problems. For regression problems, we consider the linear, logistic, and Poisson models. Under the classical regression setup, we define $\zeta_k = (\zeta_{b_k}, \zeta_{a_k})$, where ζ_{b_k} represents the k-th response and ζ_{a_k} the corresponding observation (attributes). In linear regression, the response is generated according to:

$$\zeta_{b_k} = \zeta_{\boldsymbol{a}_k}^\top \boldsymbol{x}^* + \varepsilon_k,$$

where x^* is the true parameter and $\{\varepsilon_k\}_k$ are i.i.d. noise terms. For logistic regression models, we consider binary responses $\zeta_{b_k} \in \{-1, 1\}$ generated via:

$$\mathbb{P}\left(\zeta_{b_k}|\zeta_{\boldsymbol{a}_k}\right) = \frac{1}{1 + \exp\left(-\zeta_{b_k} \cdot \zeta_{\boldsymbol{a}_k}^\top \boldsymbol{x}^*\right)}$$

In the case of Poisson regression, the response follows a conditional Poisson distribution depending on the observation, i.e.,

$$\zeta_{b_k} \sim \operatorname{Pois}\left(\lambda(\zeta_{\boldsymbol{a}_k})\right), \text{ where } \log(\lambda(\zeta_{\boldsymbol{a}_k})) = \zeta_{\boldsymbol{a}_k}^\top \boldsymbol{x}^*.$$

For each of these models, we can define an objective function corresponding to the model parameter x:

linear models:
$$f(\boldsymbol{x};\zeta_k) = \frac{1}{2} \left(\zeta_{b_k} - \zeta_{\boldsymbol{a}_k}^{\top} \boldsymbol{x} \right),$$

logistic models: $f(\boldsymbol{x};\zeta_k) = \log \left(1 + \exp \left(-\zeta_{b_k} \cdot \zeta_{\boldsymbol{a}_k}^{\top} \boldsymbol{x} \right) \right),$
Poisson models: $f(\boldsymbol{x};\zeta_k) = \zeta_{b_k} \cdot \zeta_{\boldsymbol{a}_k}^{\top} \boldsymbol{x} - \exp \left(\zeta_{\boldsymbol{a}_k}^{\top} \boldsymbol{x} \right).$

It is straightforward to verify that x^* is the optimal solution of the stochastic objective $f(x) = \mathbb{E}[f(x; \zeta_k)]$. Constraints on the model parameters x may be incorporated based on prior knowledge or specific problem requirements. We also explore portfolio optimization problems featuring exponential and logarithmic utility functions as the objective. In terms of the hyperparameters, we fix them for all experiments. The step sizes are set as $\gamma_k = 1/(k+1)^{0.751}$ and $\beta_k = 1/(k+1)^{0.5}$, which satisfy conditions in Theorems 3 and 4. Quadratic subproblems are solved by the ProxQP solver [4].¹ Implementation details are available as our supplementary code, which can be accessed at https://github.com/yihang-gao/Debiased-StoSQP/tree/main/code.

5.1. *CUTEst benchmark problems*. The CUTEst library collects various types of constrained and unconstrained optimization problems for evaluating the performances of optimization algorithms. We select a subset of the constrained optimization problems (e.g., HS problems) from the library, and we artificially add noise to gradient and Hessian as follows:

• Gaussian noise: Let $g_k = \nabla f(\boldsymbol{x}_k, \zeta_k)$ be perturbed such that $g_k = \nabla f(\boldsymbol{x}_k, \zeta_k) \sim \mathcal{N}(\nabla f(\boldsymbol{x}_k), \epsilon(\boldsymbol{I} + \boldsymbol{e}\boldsymbol{e}^{\top}))$. Similarly, the Hessian is given by $\nabla^2 f(\boldsymbol{x}_k, \zeta_k) \sim \nabla^2 f(\boldsymbol{x}_k) + \boldsymbol{E}$, where $\boldsymbol{E}_{i,j} \sim \mathcal{N}(0, \epsilon)$. The noise level ϵ is chosen from $\{1, 10^{-1}, 10^{-2}, 10^{-4}\}$.

¹In our implementation, we import the 'qpsolvers' package from https://qpsolvers.github.io/qpsolvers/.

Student's t-distribution noise: Let g_k = ∇f(x_k, ζ_k) be perturbed by noise following a Student's t-distribution, i.e., g_k = ∇f(x_k, ζ) ~ ∇f(x_k) + s, where each entry s_i ~ t(m). Similarly, the Hessian is perturbed as ∇²f(x_k, ζ) ~ ∇²f(x_k) + E, where each element E_{i,j} ~ t(m). Here, t(m) denotes the t-distributional noise, m denotes the degrees of freedom, and m is selected from the set {3, 4, 5}.

We first conduct a comparison between the proposed Debiased-StoSQP-v2 (Algorithm 3), ActiveSet-SQP [48] and StoSQP [20], where we evaluate each method by the KKT residual in Equation (2.11) and feasibility error. For each algorithm and problem, we run 10^5 iterations. Our empirical results indicate that the Debiased-StoSQP-v2 algorithm consistently outperforms the StoSQP (without debiasing techniques), both of which adopt fully stochastic gradients and Hessian. This superior performance can be attributed to our algorithm's use of gradient averaging. After a sufficient number of iterations, the averaged gradient approaches the exact gradient, thereby approximating the behavior of deterministic algorithms.

In particular, as shown in Figure 1, we plot the difference between the averaged gradients and the exact gradients during iterations, and the results validate our expectations and intuitions. The averaged gradients (the solid lines) move closer to the exact gradients, compared with estimated gradients without averaging (the dashed lines). In contrast, StoSQP (without debiasing techniques) lacks this beneficial property, and it suffers from oscillations brought by the stochastic gradients. The ActiveSet-SQP, which uses a stochastic line search method, necessitates an increasing sample size, and it employs a safeguard technique to ensure the accuracy of the line search. Consequently, it requires a sufficiently large sample size to make the line search practically effective. In contrast, Debiased-StoSQP-v2 requires only a single sample to estimate both the gradient and the Hessian in each iteration. Therefore, it is unsurprising that ActiveSet-SQP performs better with higher noise levels. However, when the noise level is relatively low, Debiased-StoSQP-v2 can effectively mitigate the noise through averaging gradient, and it can achieve similar and even better performances than ActiveSet-SQP. The visualized results are shown in Figure 2.

We next test the local asymptotic normality behavior of the generated iterates. For each problem, we aim to estimate $\frac{1}{d} \mathbf{1}^{\top} \boldsymbol{x}^*$ and set the nominal coverage probability to 95%. Here, the confidence interval is constructed by

$$\left[\frac{1}{d}\mathbf{1}^{\top}\boldsymbol{x}_{k} - \frac{1.96}{d}\sqrt{\alpha_{k}^{\min}}\sqrt{\Theta\boldsymbol{e}_{[1:d]}^{\top}\boldsymbol{\Omega}_{t}\boldsymbol{e}_{[1:d]}}, \frac{1}{d}\mathbf{1}^{\top}\boldsymbol{x}_{k} + \frac{1.96}{d}\sqrt{\alpha_{k}^{\min}}\sqrt{\Theta\boldsymbol{e}_{[1:d]}^{\top}\boldsymbol{\Omega}_{t}\boldsymbol{e}_{[1:d]}}\right]$$

using the estimators and limiting normality results in Theorem 5. The performance of the method in terms of asymptotic normality is measured by the coverage rate (CovRate) of the confidence intervals and their average length (AvgLen) over 200 runs. The aggregated results are summarized in Table 1. We observe that the constructed (95%) confidence intervals by Debiased-StoSQP-v2 cover the true solution in probability closely aligned to 95%, thereby empirically validating our theoretical derivations on asymptotic normality. From the table, we note that the length of the confidence intervals tends to expand as the noise level increases, a behavior which is in line with our expectations, as the covariance matrix Ω^* is dependent on the Cov ($\nabla f(\boldsymbol{x}^*; \zeta)$).

5.2. Constrained regression problems. Here, we implement Debiased-StoSQP-v2 algorithm (Algorithm 3) on constrained regression problems, including both the linear and the logistic regression. The response ζ_{b_k} is generated based on observations $\zeta_{a_k} \sim \mathcal{N}(\mu_a, \Sigma_a)$, where the mean vector is set as $\mu_a = (1, \dots, 1, -1, \dots, -1)$. We explore three different choices of the covariance matrix, as in [15]: (i) Identity matrix, i.e., $\Sigma_1 = I_d$; (ii) Toeplitz matrix, i.e., $(\Sigma_a)_{i,j} = r^{|i-j|}$ for some r > 0; and (iii) Equicorrelation matrix, i.e., $(\Sigma_a)_{i,j} = r$ for all $i \neq j$ and $(\Sigma_a)_{i,i} = 1$, for some r > 0. The true parameter vector of



Fig 1: Difference between the averaged gradients and the exact gradients on HS32 and FCCU problems. Solid lines: trajectories of gradient difference between the averaged gradients and the exact gradients during iterations, i.e., $\|\bar{g}_k - \nabla f(x_k)\|_2$. Dashed lines: expected error without averaging, i.e., $\mathbb{E}_k [\|g_k - \nabla f(x_k)\|_2]$.



Fig 2: KKT residuals and feasibility errors of Debiased-StoSQP-v2, StoSQP, and ActiveSet-SQP on CUTEst problems.

both two regression models is configured as $\boldsymbol{x}^* = \left(\frac{3}{2d}, \cdots, \frac{3}{2d}, \frac{1}{2d}, \cdots, \frac{1}{2d}\right)^{\top}$. We consider the non-negativity constraints, denoted by $\Omega := \{\boldsymbol{x} : \mathbf{1}^\top \boldsymbol{x} = 1, \boldsymbol{x} \ge \mathbf{0}\}$. In the linear regression problem, the noise ε_k is sampled from $\varepsilon_k \sim \mathcal{N}(0, 1)$. We aim to estimate $\hat{\boldsymbol{e}}^\top \boldsymbol{x}^*$, where $\hat{\boldsymbol{e}} = (1, \cdots, 1, -1, \cdots, -1)^{\top}$, by constructing 95% confidence intervals.

We report the results in Tables 2 and 3, highlighting different settings for the Toeplitz matrix and Equicorrelation matrix with r = 0.4, 0.5, 0.6 and r = 0.1, 0.2, 0.3, respectively. In each experiment, we run 200 times with varying random seeds to calculate the coverage rate (CovRate) and the average length (AvgLen) of the confidence interval. Our results affirm that the constructed 95% confidence intervals closely achieve a 95% coverage rate, thus empirically validating our theoretical conclusions on asymptotic normality. Moreover, we also observe that the average length of confidence intervals are in order of 10^{-2} , matching the experimental results reported by Chen et al. [15] and Na et al. [51]. The low standard deviation of these intervals' length relative to their average length suggests robustness across different random seeds.

5.3. Portfolio optimization problems. Here, we investigate portfolio optimization problems using 30 portfolios selected from the Fama-French 100 Portfolios DataSet, subject to the well-known gross-exposure constraint [26]: $\Omega_c := \{ \boldsymbol{x} : \mathbf{1}^\top \boldsymbol{x} = 1, \|\boldsymbol{x}\|_1 \leq c \}$, where we set c = 3, and where \boldsymbol{x} denotes the weights for corresponding stocks. We consider four prevalent portfolio optimization models:

Problem Noise Level		Gaussian		Student t	
Noise Level	CovRate(%)	AvgLen	Freedom	CovRate(%)	AvgLen
1E+0	97.0	2.50E-2 (7.73E-4)	3	86.0	3.77E-2 (1.90E-3)
1E-1	97.5	7.59E-3 (7.03E-5)	4	93.0	3.06E-2 (1.28E-3)
1E-2	97.0	2.40E-3 (8.69E-6)	5	94.0	2.79E-2 (9.25E-4)
1E-4	97.5	2.40E-4 (5.95E-7)	9	97.0	2.45E-2 (7.10E-4)
1E+0	94.5	1.87E-3 (6.82E-6)	3	96.5	3.18E-3 (1.66E-4)
1E-1	94.5	5.92E-4 (1.59E-6)	4	95.0	2.59E-3 (1.72E-5)
1E-2	95.0	1.87E-4 (4.96E-7)	5	95.0	2.37E-3 (1.11E-5)
1E-4	94.5	1.87E-5 (4.97E-8)	9	94.5	2.08E-3 (8.36E-6)
1E+0	97.0	2.31E-1 (4.85E-2)	3	95.5	3.00E-1 (1.32E-1)
1E-1	98.0	5.09E-2 (2.33E-3)	4	94.5	2.08E-1 (6.09E-2)
1E-2	98.5	1.58E-2 (2.23E-4)	5	95.0	1.81E-1 (5.07E-2)
1E-4	95.5	1.58E-3 (4.56E-6)	9	94.5	1.48E-1 (3.52E-2)
1E+0	97.0	1.95E-3 (1.44E-5)	3	94.0	3.34E-3 (1.23E-4)
1E-1	96.5	6.17E-4 (1.93E-6)	4	96.0	2.74E-3 (6.79E-3)
1E-2	96.5	1.95E-4 (5.20E-7)	5	96.5	2.49E-3 (2.51E-5)
1E-4	98.5	1.95E-5 (5.08E-8)	9	95.0	2.19E-3 (2.12E-5)
1E+0	94.5	3.49E-2 (3.17E-3)	3	91.0	5.04E-2 (7.56E-3)
1E-1	97.0	1.13E-2 (4.77E-5)	4	94.0	4.21E-2 (3.51E-3)
1E-2	98.0	3.58E-3 (9.63E-6)	5	94.5	3.88E-2 (2.42E-3)
1E-4	98.0	3.59E-4 (9.22E-7)	9	95.0	3.43E-2 (2.10E-3)
	Noise Level 1E+0 1E-1 1E-2 1E-4 1E-1 1E-2 1E-4	Noise Level Gate (%) 1E+0 97.0 1E-1 97.5 1E-2 97.0 1E-2 97.0 1E-2 97.0 1E-4 97.5 1E-2 97.0 1E-4 94.5 1E-1 94.5 1E-2 95.0 1E-4 94.5 1E-2 98.0 1E-4 94.5 1E-4 94.5 1E-4 94.5 1E-4 95.5 1E-4 96.5 1E-4 98.5 1E-4 94.5 1E-4 94.5 1E-4 94.5 1E-1 97.0 1E-2 98.0 1E-4 98.0 1E-4 98.0	Noise Level GrowRate(%) AvgLen 1E+0 97.0 2.50E-2 (7.73E-4) 1E-1 97.5 7.59E-3 (7.03E-5) 1E-2 97.0 2.40E-3 (8.69E-6) 1E-4 97.5 2.40E-3 (8.69E-6) 1E-4 97.5 2.40E-4 (5.95E-7) 1E+0 94.5 1.87E-3 (6.82E-6) 1E+1 94.5 5.92E-4 (1.59E-6) 1E-2 95.0 1.87E-5 (4.97E-8) 1E-4 94.5 1.87E-5 (4.97E-8) 1E-4 94.5 1.87E-5 (4.97E-8) 1E+4 94.5 1.87E-5 (4.97E-8) 1E+4 94.5 1.87E-3 (4.45E-2) 1E+1 98.0 5.09E-2 (2.33E-3) 1E-2 98.5 1.58E-3 (4.56E-6) 1E+4 95.5 1.58E-3 (4.56E-6) 1E+4 95.5 1.95E-3 (1.44E-5) 1E+1 96.5 1.95E-4 (5.20E-7) 1E+2 96.5 1.95E-5 (5.08E-8) 1E+4 94.5 3.49E-2 (3.17E-3) 1E+4 94.5 3.49E-2 (3.17E-3	Noise LevelGaussian CovRate(%)AvgLenFreedom $1E+0$ 97.0 $2.50E-2$ (7.73E-4)3 $1E-1$ 97.5 $7.59E-3$ (7.03E-5)4 $1E-2$ 97.0 $2.40E-3$ (8.69E-6)5 $1E-4$ 97.5 $2.40E-4$ (5.95E-7)9 $1E+0$ 94.5 $1.87E-3$ (6.82E-6)3 $1E-1$ 94.5 $5.92E-4$ (1.59E-6)4 $1E-2$ 95.0 $1.87E-5$ (4.97E-8)9 $1E+4$ 94.5 $1.87E-5$ (4.97E-8)9 $1E+4$ 98.0 $5.09E-2$ (2.33E-3)4 $1E-2$ 98.5 $1.58E-2$ (2.23E-4)5 $1E-4$ 95.5 $1.58E-3$ (4.56E-6)9 $1E+4$ 95.5 $1.58E-3$ (4.56E-6)9 $1E-1$ 96.5 $6.17E-4$ (1.93E-6)4 $1E-2$ 96.5 $1.95E-4$ (5.20E-7)5 $1E-4$ 98.5 $1.95E-5$ (5.08E-8)9 $1E+0$ 94.5 $3.49E-2$ (3.17E-3)3 $1E-1$ 97.0 $1.13E-2$ (4.77E-5)4 $1E-2$ 98.0 $3.58E-3$ (9.63E-6)5 $1E-4$ 98.0 $3.59E-4$ (9.22E-7)9	Noise LevelGaussianFreedomSteedom $1E+0$ 97.0 $2.50E-2$ ($7.73E-4$) 3 86.0 $1E-1$ 97.5 $7.59E-3$ ($7.03E-5$) 4 93.0 $1E-2$ 97.0 $2.40E-3$ ($8.69E-6$) 5 94.0 $1E-4$ 97.5 $2.40E-4$ ($5.95E-7$) 9 97.0 $1E+4$ 94.5 $1.87E-3$ ($6.82E-6$) 3 96.5 $1E-1$ 94.5 $5.92E-4$ ($1.59E-6$) 4 95.0 $1E-2$ 95.0 $1.87E-5$ ($4.97E-8$) 9 94.5 $1E-4$ 94.5 $1.87E-5$ ($4.97E-8$) 9 94.5 $1E-4$ 94.5 $1.87E-5$ ($4.97E-8$) 9 94.5 $1E-4$ 95.0 $1.87E-5$ ($4.97E-8$) 9 94.5 $1E-4$ 95.5 $1.58E-2$ ($2.33E-3$) 4 94.5 $1E-4$ 98.0 $5.09E-2$ ($2.33E-3$) 4 94.5 $1E-4$ 95.5 $1.58E-3$ ($4.56E-6$) 9 94.5 $1E-4$ 95.5 $1.58E-3$ ($4.56E-6$) 9 94.5 $1E-4$ 95.5 $1.95E-3$ ($1.44E-5$) 3 94.0 $1E-1$ 96.5 $6.17E-4$ ($1.93E-6$) 4 96.0 $1E-2$ 96.5 $1.95E-5$ ($5.08E-8$) 9 95.0 $1E+4$ 98.5 $1.95E-5$ ($5.08E-8$) 9 95.0 $1E+4$ 98.0 $3.58E-3$ ($9.63E-6$) 5 94.5 $1E-4$ 98.0 $3.58E-3$ ($9.63E-6$) 5 94.5 $1E-4$ 98.0 $3.58E-3$ ($9.63E-6$) 5

 TABLE 1

 The coverage rate (CovRate) and length of confidence intervals (AvgLen) for some CUTEst (constrained) problems. The standard deviation of the interval length is also reported.

TABLE 2

The coverage rate (CovRate) and length of confidence intervals (AvgLen) for constrained linear regression problems. The standard deviation of the interval length is also reported.

Cov Matrix	Dimension	CovRate(%)	AvgLen	Dimension	CovRate(%)	AvgLen
T.J	5	93.5	3.73E-2 (1.74E-4)	20	92.5	4.00E-2 (1.33E-4)
Identity	10	96.5	3.91E-2 (1.47E-4)	30	92.5	4.03E-2 (1.53E-4)
Tooplitz (0.4)	5	94.0	3.71E-2 (1.68E-4)	20	96.0	3.93E-2 (1.38E-4)
10epntz (0.4)	10	94.5	3.82E-2 (1.62E-4)	30	93.0	3.98E-2 (1.52E-4)
Tooplitz (0.5)	5	94.0	3.74E-2 (1.67E-4)	20	96.0	3.91E-2 (1.38E-4)
Toephitz (0.5)	10	95.5	3.82E-2 (1.60E-4)	30	93.0	3.95E-2 (1.61E-4)
Toeplitz (0.6)	5	94.5	3.78E-2 (1.70E-4)	20	96.5	3.90E-2 (1.36E-4)
	10	94.5	3.83E-2 (1.68E-4)	30	93.5	3.94E-2 (1.60E-4)
EquiCorr (0.1)	5	93.5	3.76E-2 (1.58E-4)	20	94.0	4.01E-2 (1.35E-4)
	10	93.0	3.92E-2 (1.40E-4)	30	92.5	4.05E-2 (1.56E-4)
FauiCorr (0.2)	5	92.5	3.79E-2 (1.59E-4)	20	93.5	4.02E-2 (1.26E-4)
Equicol1 (0.2)	10	95.0	3.94E-2 (1.50E-4)	30	96.0	4.05E-2 (1.44E-4)
EquiCorr (0.3)	5	92.5	3.83E-2 (1.65E-4)	20	93.0	4.03E-2 (1.31E-4)
	10	95.0	3.96E-2 (1.46E-4)	30	93.5	4.05E-2 (1.49E-4)

• Global minimum variance (GMV)

$$\min_{\boldsymbol{x}\in\Omega_c}\boldsymbol{x}^{\top}\boldsymbol{\Sigma}\boldsymbol{x},$$

where Σ is the covariance matrix of target stocks.

• Mean-variance (MV)

$$\min_{\boldsymbol{x}\in\Omega_c} - \boldsymbol{x}^\top \boldsymbol{\mu} + \boldsymbol{x}^\top \boldsymbol{\Sigma} \boldsymbol{x},$$

where μ and Σ are the mean and the covariance matrix of target stocks.

• Exponential utility (EXP)

$$\min_{\boldsymbol{x}\in\Omega_{c}}\mathbb{E}\left[\exp\left(-\eta\left(\boldsymbol{x}^{\top}\zeta_{\boldsymbol{a}}\right)\right)\right],$$

TABLE	3
IADLL	~

The coverage rate (CovRate) and length of confidence intervals (AvgLen) for constrained logistic regression problems. The standard deviation of the interval length is also reported.

Cov Matrix	Dimension	CovRate(%)	AvgLen	Dimension	CovRate(%)	AvgLen
T1	5	96.5	4.46E-2 (7.97E-5)	20	94.5	5.87E-2 (7.13E-5)
Identity	10	94.5	5.87E-2 (7.13E-5)	30	93.0	7.34E-2 (7.90E-5)
Tooplitz (0.4)	5	94.5	4.46E-2 (9.06E-5)	20	92.5	6.86E-2 (1.01E-4)
10epntz (0.4)	10	95.5	5.83E-2 (8.59E-5)	30	93.5	7.30E-2 (1.13E-4)
Toeplitz (0.5)	5	95.0	4.46E-2 (8.91E-5)	20	94.0	6.84E-2 (1.08E-4)
	10	94.5	5.83E-2 (8.77E-5)	30	93.0	7.28E-2 (1.24E-4)
	5	94.5	4.47E-2 (9.63E-5)	20	92.5	6.82E-2 (1.19E-4)
Toephitz (0.0)	10	94.0	5.83E-2 (8.77E-5)	30	94.5	7.26E-2 (1.32E-4)
EquiCorr (0.1)	5	95.0	4.47E-2 (9.22E-5)	20	93.0	6.69E-2 (9.40E-5)
	10	94.0	5.89E-2 (7.81E-5)	30	93.5	7.40E-2 (9.27E-5)
EquiCorr (0.2)	5	96.0	4.47E-2 (8.86E-4)	20	95.0	7.00E-2 (1.05E-4)
	10	95.0	5.92E-2 (7.32E-5)	30	92.5	7.46E-2 (1.02E-4)
EquiCorr (0.3)	5	95.0	4.48E-2 (8.59E-5)	20	93.5	7.05E-2 (1.09E-4)
	10	96.0	5.95E-2 (7.94E-4)	30	94.5	7.52E-2 (1.09E-4)

where ζ_a is the observed price changes and $\eta > 0$ is a scaling parameter set to be $\eta = 0.1$. • Logarithmic utility (LOG)

$$\min_{\boldsymbol{x}\in\Omega_c} -\mathbb{E}\left[\log\left(\boldsymbol{x}^{\top}\zeta_{\boldsymbol{a}}+\eta\right)\right]$$

where ζ_a is the observed price changes and $\eta > 0$ serves as the regularization parameter to ensure the feasibility of the logarithm, where we set $\eta = 15$.

Model	Return (%)	Max Drawdown	Sharp Ratio	Sortino Ratio
EW	15.10	0.22	0.73	1.15
GMV (ours)	34.94	0.27	2.81	4.28
GMV (det)	33.43	0.27	2.71	4.14
MV (ours)	42.21	0.28	3.36	5.09
MV (det)	40.31	0.28	3.29	5.02
EXP (ours)	52.50	0.32	2.60	3.98
EXP (det)	51.85	0.31	2.55	3.86
LOG (ours)	54.86	0.33	2.45	3.59
LOG (det)	55.08	0.32	2.46	3.57

TABLE 4Fama-French 100 Portfolios DataSet, 2021-2023

The exact mean, covariance, and expectations are inaccessible in practice. Instead, we estimate the stochastic gradient and Hessian of the expected objective using available observations, and we apply Debiased-StoSQP-v2 to solve the problems. For the portfolio strategy x, we use historical data from the past year as training samples. We assess the performance of our portfolio strategies using four key metrics, calculated over the data from years 2021-2023: the accumulative return, maximum drawdown, sharp ratio, and Sortino ratio. The accumulative return captures the overall gain or loss of the portfolio strategy. The other three are related to the risk of the strategy: the maximum drawdown measures the maximum observed loss from a peak to a trough; the sharp ratio compares the portfolio's return to its risk, taking into account the standard deviation of the portfolio returns; and the Sortino ratio is a variation of the sharp ratio, considering the standard deviation of negative portfolio returns. The results are summarized in Table 4.

Interestingly, we observe that the model of logarithmic utility achieves the best accumulative return, consistent with the results reported by [22]. In terms of risk control, however, the mean-variance model is more favorable. We also perform a comparison between Debiased-StoSQP-v2 and the deterministic approach denoted as "det". We found that the performance metrics across the two methods are quite similar, suggesting that Debiased-StoSQP-v2 approaches deterministic algorithms because of the use of the averaging gradient and Hessian.

In Figure 3, we visualize the weights of two stocks as an example, evaluated by the exponential and the logarithmic utility models. The blue line traces the trajectory of the weight corresponding to the stock over time. This is accompanied by a blue band, which represents the estimated standard deviation of the weight, as evaluated by the developed asymptotic normality. The yellow line is the accumulative return of the stock. We observe a significant correlation between the weight adjustments and the stock's return trajectory. Notably, abrupt changes in the stock's return are promptly followed by widened blue bands, indicating a surge in the estimated variance of the weight. This behavior matches well with intuition and underscores the hypothesis that the variance of the weight may serve as an indicator of the stock's inherent risk.



Fig 3: Weights and returns. Blue lines: the predicted weight for a specific stock in different months. Blue bands: the estimated standard deviation of the weight, evaluated by the derived asymptotic normality results. Yellow lines: the accumulative return of the stock. The predicted stock weights are highly correlated with the stock's price. Abrupt stock price (return) changes shown by the yellow lines are followed by widened blue bands.

5.4. Poisson regression: Chicago air pollution and death rate data. Here, we study the relationship between different attributes related to air pollution (e.g., PM10, PM25, O3, SO2) and time, with the death rate, by using Poisson regression. Let $\zeta_a \in \mathbb{R}^d$ represent the vector of air pollution and time attributes, and let $\zeta_b \in \mathbb{N}$ denote the death rate. We model the conditional distribution of death ζ_b given ζ_a as a Poisson distribution: $\zeta_b | \zeta_a \sim \text{Pois}(\lambda(\zeta_a))$, where

 $\log \lambda(\zeta_a) = \zeta_a^\top x^*$ and x^* is the true, but unknown, parameter vector for the Poisson linear model. The unconstrained Poisson regression problem is formulated as

(5.1)
$$\min_{\boldsymbol{x}} \quad \mathbb{E}\left[f(\boldsymbol{x};\zeta)\right] := \mathbb{E}_{(\zeta_b,\zeta_a)}\left[\zeta_b \cdot \zeta_a^\top \boldsymbol{x} - \exp\left(\zeta_a^\top \boldsymbol{x}\right)\right].$$

However, based on prior knowledge that air pollution attributes are likely to contribute to an increase in the death rate, we impose non-negativity constraints on the corresponding weights \boldsymbol{x} . The constrained Poisson regression model is

(5.2)
$$\min_{\boldsymbol{x}} \quad \mathbb{E}\left[f(\boldsymbol{x};\zeta)\right] := \mathbb{E}_{(\zeta_b,\zeta_a)}\left[\zeta_b \cdot \zeta_a^\top \boldsymbol{x} - \exp\left(\zeta_a^\top \boldsymbol{x}\right)\right],$$

s.t. $\boldsymbol{x}_{\mathcal{B}} \ge \boldsymbol{0},$

where \mathcal{B} is the set of indices of weights corresponding to pollution attributes.

Summary of Poisson regression (Model 1) on Chicago air pollution and death rate data						
Variables	Model	Coefficient (10^{-2})	95 % CI (10 ⁻²)	p-Value		
DM10	Model 1	0.42	[-0.57, 1.42]	0.396		
1 11110	Ours	0.13	[-0.56,0.79]	0.371		
DM25	Model 1	0.72	[-0.08, 1.52]	0.103		
F 1 V125	Ours	0.65	[0.02, 1.28]	0.023		
03	Model 1	-2.97	[-3.70, -2.24]	0.000		
03	Ours	0.00	active			
SO2	Model 1	1.38	[0.58, 2.20]	0.001		
	Ours	2.08	[1.43, 2.73]	0.000		
Time	Model 1	0.95	[0.17, 1.74]	0.008		
	Ours	1.13	[0.64, 1.63]	0.000		
Intoro	Model 1	4.6968	[4.690, 4.704]	0.000		
merc	Ours	4.6974	[4.692, 4.703]	0.000		

TABLE 5

TABLE 6

Summary of Poisson regression (Model 2) on Chicago air pollution and death rate data

Variables	Model	Coefficient (10^{-2})	95 % CI (10 ⁻²)	p-Value
PM10	Model 2	-0.86	[-1.82, 0.10]	0.062
	Ours	0.11	[-0.52, 0.74]	0.362
PM25	Model 2	1.37	[0.51, 2.23]	0.001
	Ours	0.65	[0.01, 1.28]	0.022
SO2	Model 2	2.06	[1.28, 2.84]	0.000
	Ours	2.08	[1.42, 2.73]	0.000
Time	Model 2	1.21	[0.53, 1.89]	0.001
	Ours	1.13	[0.64, 1.63]	0.000
Interc	Model 1	4.6972	[4.690, 4.704]	0.000
	Ours	4.6973	[4.692, 4.703]	0.000

We first consider the unconstrained Poisson regression model in Equation (5.1), including all five attributes, denoted as Model 1. The estimated model coefficients and their confidence turn to the constrained Poisson model in Equation (5.2), solved by Debiased-StoSQP-v2 (Algorithm 3). We also estimate confidence intervals and p-values using the derived asymptotic normality. We list the results in Table 5. Remarkably, under the non-negativity constraints, the weight of O3 is recognized to be active with the constraints, i.e., it is equal to zero. The estimated model coefficients from the constrained Poisson model are more consistent with our prior beliefs.

Next, we consider a reduced model that excludes the O3 attribute, denoted as Model 2. We find that in this model, the weight for PM10 becomes negative, once again contradicting our prior beliefs. Similarly, employing the constrained Poisson model and solving it via Debiased-StoSQP-v2 (Algorithm 3), the weight for PM10 becomes positive, although the significance level revealed by the p-value is not particularly strong. The results are reported in Table 6. These findings highlight the importance of incorporating domain-specific constraints in statistical models. They also emphasize the effectiveness of our approach in addressing such constrained optimization problems.

6. Conclusion. In this work, we proposed Debiased-StoSQP and its refinement Debiased-StoSQP-v2, fully stochastic Newton's methods for solving constrained optimization problems. We include the averaging technique for both the gradient and Hessian, reducing the impact of stochastic noise and improving the algorithm's performance, compared to existing fully stochastic algorithms. We then established the almost sure global convergence in terms of the first-order optimality (KKT) conditions. Furthermore, under certain mild conditions, we developed the asymptotic normality for Debiased-StoSQP-v2 (Algorithm 3). This is a particularly surprising and novel result since the gradients in Debiased-StoSQP-v2 are highly correlated, in contrast with previous work that primarily relies on the independence of gradients. We also provided a practical plug-in estimator for the covariance matrix. With our results, we are capable of applying Debiased-StoSQP-v2 to perform online inference for constrained optimization problems, as we have demonstrated with our empirical results.

While our algorithm Debiased-StoSQP-v2 has demonstrated promising results, there is still potential for further investigation and improvement. Specifically, the current implementation and analysis require the exact solution of quadratic subproblems, which could be computationally expensive. A possible extension of this work would be to explore the use of inexact solutions for the quadratic subproblems, possibly implemented within the Rand-BLAS/RandLAPACK library [46]. A recent work by Na and Mahoney [51] employed sketching techniques to inexactly solve linear systems for equality-constrained subproblems. The asymptotic normality behavior still holds for the StoSQP algorithm with sketching. It remains an open question whether the global almost sure convergence and the local asymptotic normality properties of Debiased-StoSQP-v2 are preserved when fast and inexact solvers are adopted. Investigating these areas would help in developing a more efficient and versatile algorithm, with broader applicability in constrained optimization scenarios.

Acknowledgements. MWM would like to acknowledge the IARPA, NSF, and ONR for providing partial support of this work.

REFERENCES

[1] ALLEN-ZHU, Z. (2018). How to make the gradients small stochastically: Even faster convex and nonconvex SGD. *Advances in Neural Information Processing Systems* **31**.

- [2] ANASTASIOU, A., BALASUBRAMANIAN, K. and ERDOGDU, M. A. (2019). Normal approximation for stochastic gradient descent via non-asymptotic rates of martingale CLT. In *Conference on Learning Theory* 115–137. PMLR.
- [3] ARJOVSKY, M., CHINTALA, S. and BOTTOU, L. (2017). Wasserstein Generative Adversarial Networks. In Proceedings of the 34th International Conference on Machine Learning. Proceedings of Machine Learning Research 70 214–223. PMLR.
- [4] BAMBADE, A., EL-KAZDADI, S., TAYLOR, A. and CARPENTIER, J. (2022). Prox-QP: Yet another quadratic programming solver for robotics and beyond. In *RSS 2022-Robotics: Science and Systems*.
- [5] BERAHAS, A. S., CURTIS, F. E., O'NEILL, M. J. and ROBINSON, D. P. (2021). A stochastic sequential quadratic optimization algorithm for nonlinear equality constrained optimization with rank-deficient jacobians. arXiv preprint arXiv:2106.13015. https://doi.org/arXiv:2106.13015
- [6] BERAHAS, A. S., CURTIS, F. E., ROBINSON, D. and ZHOU, B. (2021). Sequential Quadratic Optimization for Nonlinear Equality Constrained Stochastic Optimization. *SIAM Journal on Optimization* 31 1352– 1379. https://doi.org/10.1137/20m1354556
- [7] BERAHAS, A. S., SHI, J., YI, Z. and ZHOU, B. (2023). Accelerating stochastic sequential quadratic programming for equality constrained optimization using predictive variance reduction. *Computational Optimization and Applications* 86 79–116. https://doi.org/10.1007/s10589-023-00483-2
- [8] BERTSEKAS, D. P. (1997). Nonlinear Programming. Journal of the Operational Research Society 48 334– 334. https://doi.org/10.1057/palgrave.jors.2600425
- BOGGS, P. T. and TOLLE, J. W. (1995). Sequential Quadratic Programming. Acta Numerica 4 1–51. https: //doi.org/10.1017/s0962492900002518
- [10] BOYER, C. and GODICHON-BAGGIONI, A. (2022). On the asymptotic rate of convergence of Stochastic Newton algorithms and their Weighted Averaged versions. *Computational Optimization and Applications* 84 921–972. https://doi.org/10.1007/s10589-022-00442-3
- BURKE, J. V. and HAN, S.-P. (1989). A robust sequential quadratic programming method. *Mathematical Programming* 43 277–303. https://doi.org/10.1007/bf01582294
- [12] CARMON, Y., DUCHI, J. C., HINDER, O. and SIDFORD, A. (2017). "Convex Until Proven Guilty": Dimension-Free Acceleration of Gradient Descent on Non-Convex Functions. In *International Conference on Machine Learning* 654–663. PMLR.
- [13] CARMON, Y., DUCHI, J. C., HINDER, O. and SIDFORD, A. (2018). Accelerated Methods for NonConvex Optimization. SIAM Journal on Optimization 28 1751–1772. https://doi.org/10.1137/17m1114296
- [14] CHEN, R. T., RUBANOVA, Y., BETTENCOURT, J. and DUVENAUD, D. K. (2018). Neural ordinary differential equations. Advances in Neural Information Processing Systems 31.
- [15] CHEN, X., LEE, J. D., TONG, X. T. and ZHANG, Y. (2020). Statistical inference for model parameters in stochastic gradient descent. *The Annals of Statistics* 48 251–273. https://doi.org/10.1214/18-aos1801
- [16] CISSE, M., BOJANOWSKI, P., GRAVE, E., DAUPHIN, Y. and USUNIER, N. (2017). Parseval networks: Improving robustness to adversarial examples. In *Proceedings of the 34th International Conference on Machine Learning* 854–863. PMLR.
- [17] CUOMO, S., COLA, V. S. D., GIAMPAOLO, F., ROZZA, G., RAISSI, M. and PICCIALLI, F. (2022). Scientific Machine Learning Through Physics–Informed Neural Networks: Where we are and What's Next. *Journal of Scientific Computing* 92 88. https://doi.org/10.1007/s10915-022-01939-z
- [18] CURTIS, F. E., O'NEILL, M. J. and ROBINSON, D. P. (2023). Worst-case complexity of an SQP method for nonlinear equality constrained stochastic optimization. *Mathematical Programming*. https://doi. org/10.1007/s10107-023-01981-1
- [19] CURTIS, F. E., ROBINSON, D. P. and ZHOU, B. (2021). Inexact sequential quadratic optimization for minimizing a stochastic objective function subject to deterministic nonlinear equality constraints. arXiv preprint arXiv:2107.03512. https://doi.org/arXiv:2107.03512
- [20] CURTIS, F. E., ROBINSON, D. P. and ZHOU, B. (2023). Sequential Quadratic Optimization for Stochastic Optimization with Deterministic Nonlinear Inequality and Equality Constraints. arXiv preprint arXiv:2302.14790. https://doi.org/arXiv:2302.14790
- [21] DANIEL, J. W. (1973). Stability of the solution of definite quadratic programs. *Mathematical Programming* 5 41–53. https://doi.org/10.1007/bf01580110
- [22] DU, J.-H., GUO, Y. and WANG, X. (2022). High-Dimensional Portfolio Selection with Cardinality Constraints. *Journal of the American Statistical Association* **118** 779–791. https://doi.org/10.1080/ 01621459.2022.2133718
- [23] DUCHI, J. C. and RUAN, F. (2021). Asymptotic optimality in stochastic optimization. *The Annals of Statistics* 49. https://doi.org/10.1214/19-aos1831
- [24] DUPACOVA, J. and WETS, R. (1988). Asymptotic Behavior of Statistical Estimators and of Optimal Solutions of Stochastic Optimization Problems. *The Annals of Statistics* 16. https://doi.org/10.1214/aos/ 1176351052

- [25] FAN, J. (2007). Variable screening in high-dimensional feature space. In Proceedings of the 4th International Congress of Chinese Mathematicians 2 735–747. Citeseer.
- [26] FAN, J., ZHANG, J. and YU, K. (2012). Vast Portfolio Selection With Gross-Exposure Constraints. Journal of the American Statistical Association 107 592–606. https://doi.org/10.1080/01621459.2012.682825
- [27] FANG, Y., NA, S., MAHONEY, M. W. and KOLAR, M. (2022). Fully stochastic trust-region sequential quadratic programming for equality-constrained optimization problems. arXiv preprint arXiv:2211.15943. https://doi.org/arXiv:2211.15943
- [28] FOWKES, J., ROBERTS, L. and BÜRMEN, Á. (2022). PyCUTEst: an open source Python package of optimization test problems. *Journal of Open Source Software* 7 4377. https://doi.org/10.21105/joss.04377
- [29] GAUVIN, J. (1977). A necessary and sufficient regularity condition to have bounded multipliers in nonconvex programming. *Mathematical Programming* 12 136–138. https://doi.org/10.1007/bf01593777
- [30] GOODFELLOW, I., SHLENS, J. and SZEGEDY, C. (2015). Explaining and Harnessing Adversarial Examples. In International Conference on Learning Representations.
- [31] GOULD, N. I. M., ORBAN, D. and TOINT, P. L. (2014). CUTEst: a Constrained and Unconstrained Testing Environment with safe threads for mathematical optimization. *Computational Optimization and Applications* 60 545–557. https://doi.org/10.1007/s10589-014-9687-3
- [32] GOWER, R. M. and RICHTÁRIK, P. (2015). Randomized Iterative Methods for Linear Systems. SIAM Journal on Matrix Analysis and Applications 36 1660–1690. https://doi.org/10.1137/15m1025487
- [33] HANSEN, D., MADDIX, D. C., ALIZADEH, S., GUPTA, G. and MAHONEY, M. W. (2023). Learning Physical Models that Can Respect Conservation Laws. In *Proceedings of International Conference on Machine Learning* 202 12469–12510. PMLR.
- [34] HAO, S. and LIU, Q. (2014). Convergence rates in the law of large numbers for arrays of martingale differences. *Journal of Mathematical Analysis and Applications* 417 733–773. https://doi.org/10.1016/j. jmaa.2014.03.049
- [35] HODGKINSON, L., VAN DER HEIDE, C., ROOSTA, F. and MAHONEY, M. W. (2022). Monotonicity and Double Descent in Uncertainty Estimation with Gaussian Processes Technical Report No. Preprint: arXiv:2210.07612.
- [36] HODGKINSON, L., VAN DER HEIDE, C., SALOMONE, R., ROOSTA, F. and MAHONEY, M. W. (2023). The Interpolating Information Criterion for Overparameterized Models. arXiv preprint arXiv:2307.07785.
- [37] HOFFMAN, A. J. (2003). On Approximate Solutions of Systems of Linear Inequalities. In Selected Papers of Alan J Hoffman 174–176. World Scientific. https://doi.org/10.1142/9789812796936_0018
- [38] JORGE, N. and STEPHEN, J. W. (2006). Numerical optimization. Spinger. https://doi.org/10.1007/ 0-387-22742-3_18
- [39] KARNIADAKIS, G. E., KEVREKIDIS, I. G., LU, L., PERDIKARIS, P., WANG, S. and YANG, L. (2021). Physics-informed machine learning. *Nature Reviews Physics* 3 422–440. https://doi.org/10.1038/ s42254-021-00314-5
- [40] KRISHNAPRIYAN, A., GHOLAMI, A., ZHE, S., KIRBY, R. and MAHONEY, M. W. (2021). Characterizing possible failure modes in physics-informed neural networks. *Advances in Neural Information Processing Systems* 34 26548–26560.
- [41] LE CAM, L. M. and YANG, G. L. (2000). Asymptotics in statistics: some basic concepts. Springer Science & Business Media. https://doi.org/10.1007/978-1-4684-0377-0
- [42] LELUC, R. and PORTIER, F. (2020). Asymptotic Analysis of Conditioned Stochastic Gradient Descent. arXiv preprint arXiv:2006.02745. https://doi.org/arXiv:2006.02745
- [43] LIU, H., LI, Z., HALL, D., LIANG, P. and MA, T. (2023). Sophia: A Scalable Stochastic Second-order Optimizer for Language Model Pre-training. arXiv preprint arXiv:2305.14342. https://doi.org/arXiv: 2305.14342
- [44] LIU, X. and YUAN, Y. (2011). A Sequential Quadratic Programming Method Without A Penalty Function or a Filter for Nonlinear Equality Constrained Optimization. SIAM Journal on Optimization 21 545– 571. https://doi.org/10.1137/080739884
- [45] MOU, W., LI, C. J., WAINWRIGHT, M. J., BARTLETT, P. L. and JORDAN, M. I. (2020). On linear stochastic approximation: Fine-grained Polyak-Ruppert and non-asymptotic concentration. In *Conference on Learning Theory* 2947–2997. PMLR.
- [46] MURRAY, R., DEMMEL, J., MAHONEY, M. W., ERICHSON, N. B., MELNICHENKO, M., MALIK, O. A., GRIGORI, L., LUSZCZEK, P., DEREZIŃSKI, M., LOPES, M. E., LIANG, T., LUO, H. and DON-GARRA, J. (2023). Randomized Numerical Linear Algebra – A Perspective on the Field with an Eye to Software Technical Report No. Preprint: arXiv:2302.11474v2.
- [47] NA, S., ANITESCU, M. and KOLAR, M. (2022). An adaptive stochastic sequential quadratic programming with differentiable exact augmented lagrangians. *Mathematical Programming* 199 721–791. https: //doi.org/10.1007/s10107-022-01846-z

- [48] NA, S., ANITESCU, M. and KOLAR, M. (2023). Inequality constrained stochastic nonlinear optimization via active-set sequential quadratic programming. *Mathematical Programming*. https://doi.org/10. 1007/s10107-023-01935-7
- [49] NA, S., DEREZIŃSKI, M. and MAHONEY, M. W. (2022). Hessian averaging in stochastic Newton methods achieves superlinear convergence. *Mathematical Programming* 201 473–520. https://doi.org/10.1007/ s10107-022-01913-5
- [50] NA, S. and KOLAR, M. (2021). High-dimensional index volatility models via Stein's identity. *Bernoulli* 27. https://doi.org/10.3150/20-bej1238
- [51] NA, S. and MAHONEY, M. W. (2022). Asymptotic convergence rate and statistical inference for stochastic sequential quadratic programming. arXiv preprint arXiv:2205.13687. https://doi.org/arXiv:2205. 13687
- [52] NA, S., YANG, Z., WANG, Z. and KOLAR, M. (2019). High-dimensional Varying Index Coefficient Models via Stein's Identity. *Journal of Machine Learning Research* 20 152–1.
- [53] NÉGIAR, G., MAHONEY, M. W. and KRISHNAPRIYAN, A. (2023). Learning differentiable solvers for systems with hard constraints. In *The Eleventh International Conference on Learning Representations*.
- [54] NEYSHABUR, B., BHOJANAPALLI, S., MCALLESTER, D. and SREBRO, N. (2017). Exploring generalization in deep learning. Advances in Neural Information Processing Systems 30.
- [55] POLYAK, B. T. and JUDITSKY, A. B. (1992). Acceleration of Stochastic Approximation by Averaging. SIAM Journal on Control and Optimization 30 838–855. https://doi.org/10.1137/0330046
- [56] POWELL, M. J. D. (1978). A fast algorithm for nonlinearly constrained optimization calculations. In Lecture Notes in Mathematics 144–157. Springer Berlin Heidelberg. https://doi.org/10.1007/bfb0067703
- [57] QIU, S. and KUNGURTSEV, V. (2023). A sequential quadratic programming method for optimization with stochastic objective functions, deterministic inequality constraints and robust subproblems. arXiv preprint arXiv:2302.07947. https://doi.org/arXiv:2302.07947
- [58] RAISSI, M., PERDIKARIS, P. and KARNIADAKIS, G. E. (2019). Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational Physics* 378 686–707. https://doi.org/10.1016/j.jcp.2018.10.045
- [59] ROBBINS, H. and SIEGMUND, D. (1971). A convergence theorem for non negative almost supermartingales and some applications. In *Optimizing Methods in Statistics* 233–257. Elsevier. https://doi.org/10.1016/ b978-0-12-604550-5.50015-8
- [60] ROBINSON, S. M. (1976). Stability Theory for Systems of Inequalities, Part II: Differentiable Nonlinear Systems. SIAM Journal on Numerical Analysis 13 497–513. https://doi.org/10.1137/0713043
- [61] ROOSTA-KHORASANI, F. and MAHONEY, M. W. (2019). Sub-Sampled Newton Methods. *Mathematical Programming* 174 293–326. https://doi.org/10.1007/s10107-018-1346-5
- [62] RUPPERT, D. (1988). Efficient estimations from a slowly convergent Robbins-Monro process Technical Report, Cornell University Operations Research and Industrial Engineering.
- [63] SCHITTKOWSKI, K. and YUAN, Y.-X. (2011). Sequential Quadratic Programming Methods. https://doi. org/10.1002/9780470400531.eorms0984
- [64] SEABOLD, S. and PERKTOLD, J. (2010). Statsmodels: Econometric and statistical modeling with python. In Proceedings of the 9th Python in Science Conference 57 10–25080. Austin, TX.
- [65] SEN, P. K. (1979). Asymptotic Properties of Maximum Likelihood Estimators Based on Conditional Specification. *The Annals of Statistics* 7. https://doi.org/10.1214/aos/1176344785
- [66] SHAPIRO, A. (1985). Asymptotic distribution of test statistics in the analysis of moment structures under inequality constraints. *Biometrika* 72 133–144. https://doi.org/10.1093/biomet/72.1.133
- [67] SHAPIRO, A. (2000). On the asymptotics of constrained local \$M\$-estimators. The Annals of Statistics 28. https://doi.org/10.1214/aos/1015952006
- [68] SU, J., VARGAS, D. V. and SAKURAI, K. (2019). One Pixel Attack for Fooling Deep Neural Networks. *IEEE Transactions on Evolutionary Computation* 23 828–841. https://doi.org/10.1109/tevc.2019. 2890858
- [69] TIAN, Y. and TAKANE, Y. (2009). The inverse of any two-by-two nonsingular partitioned matrix and three matrix inverse completion problems. *Computers & Mathematics with Applications* 57 1294–1304. https://doi.org/10.1016/j.camwa.2009.01.025
- [70] TOULIS, P. and AIROLDI, E. M. (2017). Asymptotic and finite-sample properties of estimators based on stochastic gradients. *The Annals of Statistics* 45. https://doi.org/10.1214/16-aos1506
- [71] ULBRICH, S. (2003). On the superlinear local convergence of a filter-SQP method. *Mathematical Programming* 100. https://doi.org/10.1007/s10107-003-0491-6
- [72] VAN DER VAART, A. W. (2000). Asymptotic statistics 3. Cambridge university press. https://doi.org/10. 1017/CBO9780511802256

- [73] WANG, S., WANG, H. and PERDIKARIS, P. (2021). Learning the solution operator of parametric partial differential equations with physics-informed DeepONets. *Science Advances* 7. https://doi.org/10.1126/ sciadv.abi8605
- [74] XU, M., YE, J. J. and ZHANG, L. (2015). Smoothing SQP Methods for Solving Degenerate Nonsmooth Constrained Optimization Problems with Applications to Bilevel Programs. SIAM Journal on Optimization 25 1388–1410. https://doi.org/10.1137/140971580
- [75] XU, P., ROOSTA, F. and MAHONEY, M. W. (2020). Newton-type methods for non-convex optimization under inexact Hessian information. *Mathematical Programming* 184 35–70. https://doi.org/10.1007/ s10107-019-01405-z
- [76] XU, P., ROOSTA, F. and MAHONEY, M. W. (2020). Second-order optimization for non-convex machine learning: An empirical study. In *Proceedings of the 2020 SIAM International Conference on Data Mining* 199–207. SIAM. https://doi.org/10.1137/1.9781611976236.23
- [77] YAO, Z., GHOLAMI, A., SHEN, S., MUSTAFA, M., KEUTZER, K. and MAHONEY, M. (2021). ADA-HESSIAN: An Adaptive Second Order Optimizer for Machine Learning. *Proceedings of the AAAI Conference on Artificial Intelligence* 35 10665–10673. https://doi.org/10.1609/aaai.v35i12.17275
- [78] YAO, Z., XU, P., ROOSTA, F. and MAHONEY, M. W. (2021). Inexact nonconvex newton-type methods. INFORMS Journal on Optimization 3 154–182. https://doi.org/10.1287/ijoo.2019.0043
- [79] YUE, M.-C., ZHOU, Z. and SO, A. M.-C. (2019). On the Quadratic Convergence of the Cubic Regularization Method under a Local Error Bound Condition. SIAM Journal on Optimization 29 904–932. https://doi.org/10.1137/18m1167498
- [80] ZAFAR, M. B., VALERA, I., GOMEZ-RODRIGUEZ, M. and GUMMADI, K. P. (2019). Fairness constraints: A flexible approach for fair classification. *The Journal of Machine Learning Research* 20 2737–2778.
- [81] ZHU, L., FENG, K., PU, Z. and MA, W. (2023). Adversarial Diffusion Attacks on Graph-based Traffic Prediction Models. *IEEE Internet of Things Journal* 1–1. https://doi.org/10.1109/jiot.2023.3290401

NA ET AL.

APPENDIX A: CONSTRAINTS RELAXATION AND DETERMINISTIC ALGORITHM

A.1. Proof for Proposition 1. We first prove the first part of the proposition that the relaxation parameter is non-zero if EGMFCQ holds at the iterate. Define $\mathcal{I}_k := \mathcal{I}(\boldsymbol{x}_k) := \{i \in [d] : (\boldsymbol{x}_k)_i = (\boldsymbol{\ell})_i\}$ and $\mathcal{J}_k := \mathcal{J}(\boldsymbol{x}_k) := \{i \in [d] : (\boldsymbol{x}_k)_i = (\boldsymbol{u})_i\}$, then $\mathcal{I}_k \cap \mathcal{J}_k = \emptyset$ and we denote $\mathcal{A}_k := \mathcal{A}(\boldsymbol{x}_k) := \mathcal{I}(\boldsymbol{x}_k) \cup \mathcal{J}(\boldsymbol{x}_k)$ as the active set of \boldsymbol{x}_k . Suppose that \boldsymbol{z}_k is the vector satisfying Equation (2.4) at \boldsymbol{x}_k and we simply let

$$\varepsilon := \min\left\{\left|(\boldsymbol{x}_k - \boldsymbol{\ell})_i\right|, \left|(\boldsymbol{u} - \boldsymbol{x}_k)_i\right|, \left|(\boldsymbol{u} - \boldsymbol{\ell})_j\right| : i \in \mathcal{A}_k^-, j \in \mathcal{A}_k\right\} > 0,$$

and

$$ar{oldsymbol{z}}_k = rac{arepsilon}{\|oldsymbol{z}_k\|_2} oldsymbol{z}_k.$$

Then, it is not difficult to verify that

$$\boldsymbol{\ell} \leq \boldsymbol{x}_k + \bar{\boldsymbol{z}}_k \leq \boldsymbol{u},$$

and

(A.1)
$$\frac{\varepsilon}{\|\boldsymbol{z}_k\|_2} \boldsymbol{c}(\boldsymbol{x}_k) + \nabla \boldsymbol{c}(\boldsymbol{x}_k) \bar{\boldsymbol{z}}_k = \boldsymbol{0},$$

which imply that $\bar{z}_k \in \widetilde{\Omega}_k$ with $\theta_k = \frac{\varepsilon}{\|z_k\|_2}$. The following lemma shows that if $\widetilde{\Omega}_k$ with $\bar{\theta} \in (0,1]$ is feasible and $0 < \bar{\theta} \le \bar{\theta}$, then $\widetilde{\Omega}_k$ with $\bar{\theta}$ is also feasible. It further indicates that Assumption 1 on the lower-boundedness of the relaxation parameter makes sense.

LEMMA 5. If $\{\boldsymbol{p}: \bar{\boldsymbol{\theta}}\boldsymbol{c}_k + \boldsymbol{J}_k^\top \boldsymbol{p} = \boldsymbol{0}\} \cap \{\boldsymbol{p}: \boldsymbol{\ell} \leq \boldsymbol{x}_k + \boldsymbol{p} \leq \boldsymbol{u}\} \neq \emptyset$ and $0 < \bar{\bar{\boldsymbol{\theta}}} \leq \bar{\boldsymbol{\theta}}$, then $\{\boldsymbol{p}: \bar{\bar{\boldsymbol{\theta}}}\boldsymbol{c}_k + \boldsymbol{J}_k^\top \boldsymbol{p} = \boldsymbol{0}\} \cap \{\boldsymbol{p}: \boldsymbol{\ell} \leq \boldsymbol{x}_k + \boldsymbol{p} \leq \boldsymbol{u}\} \neq \emptyset$. Therefore, Assumption 1 makes sense.

PROOF. Suppose that $\bar{p} \in \{p : \bar{\theta}c_k + J_k^\top p = 0\} \cap \{p : \ell \le x_k + p \le u\}$, then $\bar{\theta}c_k + J_k^\top \bar{p} = 0$ and thus, $\bar{\theta}c_k + J_k^\top \left(\bar{\bar{\theta}}/\bar{\theta} \cdot \bar{p}\right) = 0$. Let $\bar{\bar{p}} = \bar{\bar{\theta}}/\bar{\theta} \cdot \bar{p}$, then $\bar{\bar{\theta}}c_k + J_k^\top \bar{\bar{p}} = 0$ and $\ell \le x_k + \bar{\bar{p}} \le u$, which complete the proof.

LEMMA 6 (Theorem 3 in [60]). If \bar{x} satisfies EGMFCQ, then there exists a neighborhood $\mathcal{B}(\bar{x};\bar{r}) := \{x: ||x - \bar{x}||_2 \le \bar{r}\}$ with some sufficiently small radius $\bar{r} > 0$, such that all points in the neighborhood satisfy EGMFCQ.

LEMMA 7. Suppose that EGMFCQ holds at \bar{x} , then EGMFCQ also holds at \bar{x}_k when \bar{x}_k is sufficiently close to \bar{x} , for any sequence $\bar{x}_k \to \bar{x}$, by Lemma 6. Let \bar{z} be the vectors satisfying Condition (2.4) at \bar{x} . Then we can always find a sequence of vectors $\{\bar{z}_k\}$ with \bar{z}_k satisfying Condition (2.4) at \bar{x}_k such that $\|\bar{z}_k - \bar{z}\|_2 \to 0$ as $\|\bar{x}_k - \bar{x}\|_2 \to 0$.

PROOF. Since the vector \bar{z} satisfies Condition (2.4) at \bar{x} , i.e., $c(\bar{x}) + \nabla c(\bar{x})^{\top} \bar{z} = 0$, by the smoothness of c(x) and the linear independence of columns of $\nabla c(\bar{x})$, we can find \bar{z}_k such that $c(\bar{x}_k) + \nabla c(\bar{x}_k)^{\top} \bar{z}_k = 0$ and $\|\bar{z}_k - \bar{z}\|_2 \to 0$ as $\|\bar{x}_k - \bar{x}\|_2 \to 0$. Let $\varepsilon := \min\{|(\bar{z})_i| : (\bar{x})_i = (\ell)_i \text{ or } (\bar{x})_i = (u)_i\}$. Due to the fact that $\mathcal{A}(\bar{x}_k) \subseteq \mathcal{A}(\bar{x})$, we have $(\bar{z}_k)_i > 0$, if $(\bar{x}_k)_i = (\ell)_i$ and $(\bar{z}_k)_i < 0$, if $(\bar{x}_k)_i = (u)_i$, when $\|\bar{z}_k - \bar{z}\|_2 \leq \varepsilon$.

LEMMA 8. Let θ_k be selected in (0,1] such that the relaxed feasible region Ω_k is nonempty with θ_k but is empty with $\min\{1.1\theta_k, 1\}$, and we can always achieve it based on Lemma 5. If $\liminf_{k\to\infty} \theta_k = 0$, then there exists an accumulation point x^* of $\{x_k\}$ where EGMFCQ does not hold at x^* . PROOF. Without the loss of generality, we assume that $\lim_{k\to\infty} \theta_k = 0$ and $\lim_{k\to\infty} x_k = x^*$. Let $l_k := \inf\{\|z_k\|_2 : z_k \text{ satisfies Condition (2.4) at } x_k\}$. The construction of $\theta_k = \frac{\varepsilon}{\|z_k\|_2}$ in Equation (A.1) shows that $\limsup_{k\to\infty} l_k = \infty$. Suppose that EGMFCQ holds at x^* and let $l^* = \|z^*\|_2 < \infty$ for some z^* satisfying Condition (2.4) at x^* . It is a contradiction to Lemma 7 as $\infty = \limsup_{k\to\infty} l_k \leq l^* < \infty$. Therefore, EGMFCQ does not hold at x^* . \Box

EGMFCQ and its multiple variants are common in constrained optimization algorithms, i.e., [11, 74]. According to the above proposition, EGMFCQ makes the relaxed SQP subproblem feasible. Instead of assuming the EGMFCQ at all iterates $\{x_k\}$, which is difficult to verify in real applications, a weaker and more explicit assumption (Assumption 1) is made. Proposition 5 shows the reasonability of Assumption 1. To verify Assumption 1, as shown in the deterministic SQP (RelaxedSQP, Algorithm 1), we first validate and adopt a feasible $\tilde{\Omega}_k$ with proper relaxation parameters θ_k . If $\tilde{\Omega}_k$ is not feasible for small θ_k below the predefined tolerance, then x_k is close to a point where EGMFCQ does not hold, by Lemma 1. For completeness, we put the definition of LICQ here.

DEFINITION 2 (Linear independence constraint qualification (LICQ)). The linear independence constraint qualification (LICQ) is satisfied at a point \tilde{x} , if columns of $[\nabla c(\tilde{x}), I_{\mathcal{A}(\tilde{x})}]$ are linearly independent, where $\mathcal{A}(\tilde{x}) := \{i : (\tilde{x})_i = (\ell)_i \text{ or } (\tilde{x})_i = (u)_i\}$ is the active set of inequality constraints at \tilde{x} .

A.2. EGMFCQ and Boundedness of Lagrangian Multipliers. The following Lemma 9 shows that if the sequence $\{x_k\}$ generated by the algorithm is convergent to a feasible point x^* satisfying EGMFCQ (Definition 1), then the corresponding Lagrangian multipliers of the SQP subproblem are bounded.

LEMMA 9. If EGMFCQ is satisfied at \bar{x} which is feasible for both the equality and inequality constraints (i.e., $c(\bar{x}) = 0$ and $\ell \leq \bar{x} \leq u$), then there exists a neighborhood $\mathcal{B}(\bar{x};r_0) := \{x : ||x - \bar{x}||_2 \leq r_0\}$ with some $r_0 > 0$, such that the Lagrangian multipliers of the SQP subproblems are bounded for all points in $\mathcal{N}(\bar{x};r_0)$, under Assumptions 1 and 2.

PROOF. We prove it by contradiction. Suppose that there exist sequences $\{(\bar{x}_k, \bar{B}_k, \bar{\lambda}_k^{\text{sub}}, \bar{\mu}_{1,k}^{\text{sub}}, \bar{\mu}_{2,k}^{\text{sub}})\}$ with Assumptions 1 and 2, such that $\bar{x}_k \to \bar{x}$, $\left\|(\bar{\lambda}_k^{\text{sub}}, \bar{\mu}_{1,k}^{\text{sub}}, \bar{\mu}_{2,k}^{\text{sub}})\right\|_2 \to \infty$ and $\kappa_1 \mathbf{I} \preceq \bar{B}_k \preceq \kappa_2 \mathbf{I}$, where \bar{p}_k and $(\bar{\lambda}_k^{\text{sub}}, \bar{\mu}_{1,k}^{\text{sub}}, \bar{\mu}_{2,k}^{\text{sub}})$ are the solution and the Lagrangian multipliers of the SQP subproblem at \bar{x}_k with corresponding relaxing parameters $\bar{\theta}_k$ satisfying

(A.2)

$$\nabla f(\bar{\boldsymbol{x}}_{k}) + \boldsymbol{B}_{k}\bar{\boldsymbol{p}}_{k} + \nabla \boldsymbol{c}(\bar{\boldsymbol{x}}_{k})\boldsymbol{\lambda}_{k}^{\text{sub}} - \bar{\boldsymbol{\mu}}_{1,k}^{\text{sub}} + \bar{\boldsymbol{\mu}}_{2,k}^{\text{sub}} = \boldsymbol{0},$$

$$\bar{\boldsymbol{\theta}}_{k}\boldsymbol{c}(\bar{\boldsymbol{x}}_{k}) + \nabla \boldsymbol{c}(\bar{\boldsymbol{x}}_{k})^{\top}\bar{\boldsymbol{p}}_{k} = \boldsymbol{0}, \quad \boldsymbol{\ell} \leq \bar{\boldsymbol{x}}_{k} + \bar{\boldsymbol{p}}_{k} \leq \boldsymbol{u},$$

$$\bar{\boldsymbol{\mu}}_{1,k}^{\text{sub}\top}(\bar{\boldsymbol{x}}_{k} + \bar{\boldsymbol{p}}_{k} - \boldsymbol{\ell}) = 0,$$

$$\bar{\boldsymbol{\mu}}_{2,k}^{\text{sub}\top}(\bar{\boldsymbol{x}}_{k} + \bar{\boldsymbol{p}}_{k} - \boldsymbol{u}) = 0,$$

$$\bar{\boldsymbol{\mu}}_{1,k}^{\text{sub}} \geq \boldsymbol{0}, \quad \bar{\boldsymbol{\mu}}_{2,k}^{\text{sub}} \geq \boldsymbol{0}.$$

Note that the sequence $\{(\bar{\lambda}_{k}^{\text{sub}}, \bar{\mu}_{1,k}^{\text{sub}}, \bar{\mu}_{2,k}^{\text{sub}}) / \|(\bar{\lambda}_{k}^{\text{sub}}, \bar{\mu}_{1,k}^{\text{sub}}, \bar{\mu}_{2,k}^{\text{sub}})^{\top}\|_{2}\}$ is bounded. Without the loss of generality, we assume that $(\bar{\lambda}_{k}^{\text{sub}}, \bar{\mu}_{1,k}^{\text{sub}}, \bar{\mu}_{2,k}^{\text{sub}}) / \|(\bar{\lambda}_{k}^{\text{sub}}, \bar{\mu}_{1,k}^{\text{sub}}, \bar{\mu}_{2,k}^{\text{sub}})^{\top}\|_{2} \rightarrow (\bar{\lambda}, \bar{\mu}_{1}, \bar{\mu}_{2}),$ $\bar{p}_{k} \rightarrow \bar{p}$ and $\bar{\theta}_{k} = \bar{\theta}$ (due to line 4 in Algorithm 1). Then dividing both two sides of the first equality in Equation (A.2) by $\|(\bar{\lambda}_{k}, \bar{\mu}_{1,k}, \bar{\mu}_{2,k})^{\top}\|_{2}$ and taking the limit of $k \rightarrow \infty$, we have (A.3) $\nabla c(\bar{x})\bar{\lambda} - \bar{\mu}_{1} + \bar{\mu}_{2} = \mathbf{0}.$ NA ET AL.

Moreover, the second equality in Equation (A.2) implies that

$$ar{ heta}ar{m{\lambda}}^{ op}m{c}(ar{m{x}}) = -ar{m{\lambda}}^{ op}
ablam{c}(ar{m{x}})^{ op}m{m{p}}$$

The third and the fourth equality in Equation (A.2) further shows that

$$\bar{\boldsymbol{\mu}}_1^{\top}(\bar{\boldsymbol{x}}-\boldsymbol{\ell}) = -\bar{\boldsymbol{\mu}}_1^{\top}\bar{\boldsymbol{p}} \text{ and } \bar{\boldsymbol{\mu}}_2^{\top}(\bar{\boldsymbol{x}}-\boldsymbol{u}) = -\bar{\boldsymbol{\mu}}_2^{\top}\bar{\boldsymbol{p}}.$$

Combing with the above four equalities, we have

(A.4)
$$\bar{\theta}\bar{\boldsymbol{\lambda}}^{\top}\boldsymbol{c}(\bar{\boldsymbol{x}}) + \bar{\boldsymbol{\mu}}_{1}^{\top}(\boldsymbol{\ell}-\bar{\boldsymbol{x}}) + \bar{\boldsymbol{\mu}}_{2}^{\top}(\bar{\boldsymbol{x}}-\boldsymbol{u}) = 0.$$

Note that $\bar{\mu}_1 \geq \mathbf{0}$ and $\bar{\mu}_2 \geq \mathbf{0}$, we can deduce from Equation (A.4) and $c(\bar{x}) = \mathbf{0}$ that $(\bar{\mu}_1) > 0$ only if $(\bar{x})_i = (\ell)_i$ and $(\bar{\mu}_2) > 0$ only if $(\bar{x})_i = (u)_i$. The EGMFCQ condition at \bar{x} (Definition 1) implies that there exists $p \in \mathbb{R}^d$ such that $c(\bar{x}) + \nabla c(\bar{x})^\top p = \mathbf{0}$, $(p)_i > 0$ if $(\bar{x})_i = (\ell)_i$, and $(p)_i < 0$ if $(\bar{x})_i = (u)_i$. Then $-p^\top \bar{\mu}_1 + p^\top \bar{\mu}_2 < 0$ if \bar{x} is on the boundary of the box constraints. Multiplying both two sides of Equation (A.3) by $-\bar{\theta}p$, we have $0 = -\bar{\theta}p^\top (\nabla c(\bar{x})\bar{\lambda} - \bar{\mu}_1 + \bar{\mu}_2) = \bar{\theta}c(\bar{x})^\top \bar{\lambda} + \bar{\theta}p^\top \bar{\mu}_1 - \bar{\theta}p^\top \bar{\mu}_2$. It is a contradiction to Equation (A.4). On the other hand, if \bar{x} is in the interior of the box constraints, $\bar{\mu}_1 = \bar{\mu}_2 = \mathbf{0}$. Together with Equation (A.3), the linear independence of the columns of $\nabla c(\bar{x})$ shows $\bar{\lambda} = \mathbf{0}$, which is a contradiction to the fact that $\|(\bar{\lambda}, \bar{\mu}_1, \bar{\mu}_2)^\top\|_2 = 1$.

COROLLARY 2. If all accumulation points of the sequence $\{x_k\}$ are feasible and satisfy EGMFCQ, then the Lagrangian multipliers of the corresponding SQP subproblems are bounded.

PROOF. We first show that the Lagrangian multipliers of the corresponding SQP subproblems are bounded at all accumulation points of $\{x_k\}$, denoted as \mathcal{X} . Note that the set \mathcal{X} is closed, any accumulation point of \mathcal{X} is also an accumulation point of $\{x_k\}$.

Secondly, by Lemma 9, for a sufficiently large number $M_{\text{Lag}} > 0$ and any point $x_i^* \in \mathcal{X}$, there exists $r_i > 0$ such that the Lagrangian multipliers of the corresponding SQP subproblems are bounded at x for any $x \in \bigcup_{i=1}^{\infty} \mathcal{B}(x_i^*; r_i)$. There must be a finite subset of $\{x_k\}$, that is outside $\bigcup_{i=1}^{\infty} \mathcal{B}(x_i^*; r_i)$ (otherwise, we can still find an accumulation point). We complete the proof.

A.3. Proof for Theorem 1. The proof directly comes from the following lemmas. The first lemma here shows that the directional derivative of the merit function is controlled by the improvement $\Delta q(\boldsymbol{x}, \boldsymbol{p}, \nabla f(\boldsymbol{x}), \boldsymbol{B}, \rho)$.

LEMMA 10. Under Assumption 2, given $(\boldsymbol{x}, \rho, \theta, \boldsymbol{B}, \boldsymbol{p}) \in \mathbb{R}^n \times \mathbb{R}_{>0} \times (0, 1] \times \mathbb{S}^n_+ \times \mathbb{R}^n$ with $\theta \boldsymbol{c}(\boldsymbol{x}) + \nabla \boldsymbol{c}(\boldsymbol{x})^\top \boldsymbol{p} = \boldsymbol{0}$, then the directional derivative of $\phi(\boldsymbol{x}, \rho)$ along \boldsymbol{p} satisfies

(A.5)

$$\phi'(\boldsymbol{x}, \rho; \boldsymbol{p}) = \nabla f(\boldsymbol{x})^{\top} \boldsymbol{p} - \rho \theta \|\boldsymbol{c}(\boldsymbol{x})\|_{2}$$

$$\leq \nabla f(\boldsymbol{x})^{\top} \boldsymbol{p} + \frac{1}{2} \boldsymbol{p}^{\top} \boldsymbol{B} \boldsymbol{p} - \rho \theta \|\boldsymbol{c}(\boldsymbol{x})\|_{2}$$

$$= -\Delta q(\boldsymbol{x}, \boldsymbol{p}, \nabla f(\boldsymbol{x}), \boldsymbol{B}, \rho).$$

PROOF. We prove it by the definition of the directional derivative. Suppose that $\|\nabla^2 f(x)\|_2 \le M$ for some M > 0. First,

$$(A.6) \qquad \begin{aligned} \phi(\boldsymbol{x} + \alpha \boldsymbol{p}, \rho) - \phi(\boldsymbol{x}, \rho) \\ = f(\boldsymbol{x} + \alpha \boldsymbol{p}) + \rho \|\boldsymbol{c}(\boldsymbol{x} + \alpha \boldsymbol{p})\|_{2} - f(\boldsymbol{x}) - \rho \|\boldsymbol{c}(\boldsymbol{x})\|_{2} \\ \leq \alpha \nabla f(\boldsymbol{x})^{\top} \boldsymbol{p} + \frac{\kappa_{\nabla f}}{2} \alpha^{2} \|\boldsymbol{p}\|_{2}^{2} + \rho \|\boldsymbol{c}(\boldsymbol{x} + \alpha \boldsymbol{p})\|_{2} - \rho \|\boldsymbol{c}(\boldsymbol{x})\|_{2} \\ = \alpha \nabla f(\boldsymbol{x})^{\top} \boldsymbol{p} + \frac{\kappa_{\nabla f}}{2} \alpha^{2} \|\boldsymbol{p}\|_{2}^{2} + \rho (|1 - \alpha \theta| - 1) \|\boldsymbol{c}(\boldsymbol{x})\|_{2} + \frac{\kappa_{\nabla c}}{2} \alpha^{2} \|\boldsymbol{p}\|_{2}^{2} \\ = \alpha \left(\nabla f(\boldsymbol{x})^{\top} \boldsymbol{p} - \rho \theta \|\boldsymbol{c}(\boldsymbol{x})\|_{2} \right) + \frac{\kappa_{\nabla f} + \kappa_{\nabla c}}{2} \alpha^{2} \|\boldsymbol{p}\|_{2}^{2}. \end{aligned}$$

On the other side, similarly, we have $\phi(\boldsymbol{x} + \alpha \boldsymbol{p}, \rho) - \phi(\boldsymbol{x}, \rho) \ge \alpha \left(\nabla f(\boldsymbol{x})^{\top} \boldsymbol{p} - \rho \theta \| \boldsymbol{c}(\boldsymbol{x}) \|_2 \right) - \frac{\kappa_{\nabla f} + \kappa_{\nabla c}}{2} \alpha^2 \| \boldsymbol{p} \|_2^2$. Taking limits for $\alpha \to 0$ and the definition, we have $\phi'(\boldsymbol{x}, \rho; \boldsymbol{p}) = \nabla f(\boldsymbol{x})^{\top} \boldsymbol{p} - \rho \theta \| \boldsymbol{c}(\boldsymbol{x}) \|_2 \le -\Delta q(\boldsymbol{x}, \boldsymbol{p}, \nabla f(\boldsymbol{x}), \boldsymbol{B}, \rho)$.

We incorporate a backtracking line search in the algorithm while [6] adopted Lipschitz constant estimation for step size selection. We prove that under mild smoothness conditions, the line search condition will be met after a finite number of search steps. Specifically, the backtracking search loop is guaranteed to terminate within a bounded number of iterations.

LEMMA 11. The strategies in Equations (2.5) and (2.6) for ρ_k guarantee that $\Delta q(\boldsymbol{x}_k, \boldsymbol{p}_k, \boldsymbol{B}_k; \rho_k) \geq \frac{1}{2} \boldsymbol{p}_k^\top \boldsymbol{B}_k \boldsymbol{p}_k + \sigma \rho_k \theta_k \|\boldsymbol{c}(\boldsymbol{x}_k)\|_2$ for some $\sigma \in (0, 1)$. Therefore, combining it with Lemma 10, we have that the backtracking line search condition $\phi(\boldsymbol{x}_k + \alpha_k \boldsymbol{p}_k, \rho_k) \leq \phi(\boldsymbol{x}_k, \rho_k) - \beta \alpha_k \Delta q(\boldsymbol{x}_k, \boldsymbol{p}_k, \nabla f(\boldsymbol{x}_k), \boldsymbol{B}_k; \rho_k)$ always holds for $\alpha_k \leq \frac{(1-\beta)\kappa_1}{\kappa_{\nabla f} + \kappa_{\nabla c}}$.

PROOF. Equation (A.6) in Lemma 10 shows that $\phi(\boldsymbol{x}_k + \alpha_k \boldsymbol{p}_k, \rho_k) - \phi(\boldsymbol{x}_k, \rho_k) \leq \alpha_k \left(\nabla f(\boldsymbol{x}_k)^\top \boldsymbol{p}_k - \rho_k \theta_k \|\boldsymbol{c}(\boldsymbol{x}_k)\|_2\right) + \frac{\kappa_{\nabla f} + \kappa_{\nabla c}}{2} \alpha_k^2 \|\boldsymbol{p}_k\|_2^2$. Here, we let α_k to be small enough such that

$$\begin{aligned} &\alpha_k \left(\nabla f(\boldsymbol{x}_k)^\top \boldsymbol{p}_k - \rho_k \theta_k \| \boldsymbol{c}(\boldsymbol{x}_k) \|_2 \right) + \frac{\kappa_{\nabla f} + \kappa_{\nabla c}}{2} \alpha_k^2 \| \boldsymbol{p}_k \|_2^2 \\ &\leq -\alpha_k \Delta q(\boldsymbol{x}_k, \boldsymbol{p}_k, \nabla f(\boldsymbol{x}_k), \boldsymbol{B}_k; \rho_k) + \frac{\kappa_{\nabla f} + \kappa_{\nabla c}}{2} \alpha_k^2 \| \boldsymbol{p}_k \|_2^2 \\ &\leq -\beta \alpha_k \Delta q(\boldsymbol{x}_k, \boldsymbol{p}_k, \nabla f(\boldsymbol{x}_k), \boldsymbol{B}_k; \rho_k), \end{aligned}$$

i.e., $\frac{\kappa_{\nabla f} + \kappa_{\nabla c}}{2} \alpha_k \|\boldsymbol{p}_k\|_2^2 \leq (1 - \beta) \Delta q(\boldsymbol{x}_k, \boldsymbol{p}_k, \nabla f(\boldsymbol{x}_k), \boldsymbol{B}_k; \rho_k). \text{ Here, we let } \frac{\kappa_{\nabla f} + \kappa_{\nabla c}}{2} \alpha_k \|\boldsymbol{p}_k\|_2^2 \leq \frac{(1 - \beta)\kappa_1}{2} \|\boldsymbol{p}_k\|_2^2 \leq \frac{1 - \beta}{2} \boldsymbol{p}_k^\top \boldsymbol{B}_k \boldsymbol{p}_k, \text{ i.e., } \alpha_k \leq \frac{(1 - \beta)\kappa_1}{\kappa_{\nabla f} + \kappa_{\nabla c}}. \text{ In conclusion, the backtracking line search condition holds when } \alpha_k \leq \frac{(1 - \beta)\kappa_1}{\kappa_{\nabla f} + \kappa_{\nabla c}}.$

The next lemma demonstrates that if the Lagrange multipliers are bounded, then the penalty parameter will stabilize. This result is crucial for the global convergence of the algorithm, as convergence is only assured subsequent to the penalty parameter's stabilization. Specifically, once the penalty parameter stabilizes, the merit function's convergence naturally leads to the convergence of the iterates.

LEMMA 12. Under Assumption 1, $\theta_k \geq \tilde{\tau} \tilde{\theta}$ holds for all $k = 0, 1, \dots$. If we further assume that Assumption 2 holds, then the sequence $\{\rho_k\}$ is monotonically increasing and there exists a large enough $\tilde{K} \in \mathbb{Z}$, such that $\rho_k = \tilde{\rho} > 0$ for all $k \geq \tilde{K}$, where $\tilde{\rho} \leq \frac{(1+\epsilon)M_{Lag}}{(1-\sigma)\tilde{\tau}\tilde{\theta}}$.

PROOF. Under Assumption 1, it is obvious that $\theta_k \geq \tilde{\tau}\tilde{\theta}$ holds in our algorithm, for all $k = 0, 1, \cdots$. If there does not exist $\tilde{\rho} > 0$ and $\tilde{K} \in \mathbb{Z}$ such that $\rho_k = \tilde{\rho} > 0$ for $k \geq \tilde{K}$, according to Equation (2.6), then there is an infinite sequence $\{k_j\} \subseteq \mathbb{Z}_+$ where $\rho_{k_j}^{\text{trial}} > \rho_{k_j-1}$ and $\rho_{k_j} = (1+\epsilon)\rho_{k_j}^{\text{trial}}$. It further implies that $-\nabla f(\boldsymbol{x}_k)^\top \boldsymbol{p}_k - \boldsymbol{p}_k^\top \boldsymbol{B}_k \boldsymbol{p}_k < 0$ and $\rho_{k_j}^{\text{trial}} = \frac{\nabla f(\boldsymbol{x}_k)^\top \boldsymbol{p}_k + \boldsymbol{p}_k^\top \boldsymbol{B}_k \boldsymbol{p}_k}{(1-\sigma)\theta_k \|\boldsymbol{c}(\boldsymbol{x}_k)\|_2}$, by Equation (2.5). The KKT conditions for the relaxed SQP Subproblem (2.7) show that there exist some $(\boldsymbol{\lambda}_k^{\text{sub}}, \boldsymbol{\mu}_{1,k}^{\text{sub}}, \boldsymbol{\mu}_{2,k}^{\text{sub}})$ satisfying

(A.7)

$$\nabla f(\boldsymbol{x}_{k}) + \boldsymbol{B}_{k}\boldsymbol{p}_{k} + \nabla \boldsymbol{c}(\boldsymbol{x}_{k})\boldsymbol{\lambda}_{k}^{\text{sub}} - \boldsymbol{\mu}_{1,k}^{\text{sub}} + \boldsymbol{\mu}_{2,k}^{\text{sub}} = \boldsymbol{0},$$

$$\theta_{k}\boldsymbol{c}(\boldsymbol{x}_{k}) + \nabla \boldsymbol{c}(\boldsymbol{x}_{k})^{\top}\boldsymbol{p}_{k} = \boldsymbol{0},$$

$$\boldsymbol{\ell} \leq \boldsymbol{x}_{k} + \boldsymbol{p}_{k} \leq \boldsymbol{u},$$

$$\boldsymbol{\mu}_{1,k}^{\text{sub}\top}(\boldsymbol{x}_{k} + \boldsymbol{p}_{k} - \boldsymbol{\ell}) = 0,$$

$$\boldsymbol{\mu}_{2,k}^{\text{sub}\top}(\boldsymbol{x}_{k} + \boldsymbol{p}_{k} - \boldsymbol{u}) = 0,$$

$$\boldsymbol{\mu}_{1,k}^{\text{sub}} \geq \boldsymbol{0} \text{ and } \boldsymbol{\mu}_{2,k}^{\text{sub}} \geq \boldsymbol{0}.$$

Multiplying both two sides of the first equality by p_k , we have

$$\begin{split} \nabla f(\boldsymbol{x}_k)^\top \boldsymbol{p}_k + \boldsymbol{p}_k^\top \boldsymbol{B}_k \boldsymbol{p}_k &= -\boldsymbol{p}_k^\top \nabla \boldsymbol{c}(\boldsymbol{x}_k) \boldsymbol{\lambda}_k^{\text{sub}} + \boldsymbol{p}_k^\top \boldsymbol{\mu}_{1,k}^{\text{sub}} - \boldsymbol{p}_k^\top \boldsymbol{\mu}_{2,k}^{\text{sub}} \\ &= \theta_k \boldsymbol{\lambda}_k^{\text{sub}\top} \boldsymbol{c}(\boldsymbol{x}_k) - \boldsymbol{\mu}_{1,k}^{\text{sub}\top} (\boldsymbol{x}_k - \boldsymbol{\ell}) + \boldsymbol{\mu}_{2,k}^{\text{sub}\top} (\boldsymbol{x}_k - \boldsymbol{u}) \\ &\leq \theta_k \boldsymbol{\lambda}_k^{\text{sub}\top} \boldsymbol{c}(\boldsymbol{x}_k) \leq \|\boldsymbol{\lambda}_k^{\text{sub}}\|_2 \|\boldsymbol{c}(\boldsymbol{x}_k)\|_2 \leq M_{\text{Lag}} \|\boldsymbol{c}(\boldsymbol{x}_k)\|_2 \end{split}$$

where the first inequality comes from $\mu_{1,k}^{ ext{sub}} \geq 0$, $\mu_{2,k}^{ ext{sub}} \geq 0$, and $\ell \leq x_k \leq u$. Then,

(A.8)
$$\rho_{k_j-1} < \rho_{k_j}^{\text{trial}} = \frac{\nabla f(\boldsymbol{x}_{k_j})^\top \boldsymbol{p}_{k_j} + \boldsymbol{p}_{k_j}^\top \boldsymbol{B}_{k_j} \boldsymbol{p}_{k_j}}{(1-\sigma)\theta_{k_j} \|\boldsymbol{c}(\boldsymbol{x}_{k_j})\|_2} \le \frac{M_{\text{Lag}}}{(1-\sigma)\theta_{k_j}} \le \frac{M_{\text{Lag}}}{(1-\sigma)\tilde{\tau}\tilde{\theta}}.$$

However, $\rho_{k_j} = (1 + \epsilon)\rho_{k_j}^{\text{trial}} > (1 + \epsilon)\rho_{k_j-1}$ implies that $\rho_{k_j-1} \to \infty$ as $k_j \to \infty$. It is a contradiction. Therefore, there exist $\tilde{\rho} > 0$ and a large enough $\tilde{K} \in \mathbb{Z}$, such that $\rho_k = \tilde{\rho} > 0$ for all $k \ge \tilde{K}$. Here, we can also conclude from Equation (A.8) that $\tilde{\rho} \le \frac{(1+\epsilon)M_{\text{Lag}}}{(1-\sigma)\tilde{\tau}\tilde{\theta}}$.

PROPOSITION 1. If we suppose that all accumulation points of the generated sequence $\{x_k\}$ satisfies EGMFCQ, then $\lim_k \rho_k < \infty$.

PROOF. Suppose that $\lim_{k\to\infty} \rho_k = \infty$, then we can find a subsequence $\{k_j\} \subseteq \mathbb{Z}_+$ such that $\rho_{k_j} > \rho_{k_j-1}$ and $\rho_k = \rho_{k-1}$ for $k \notin \{k_j\}$. By the fact that

$$\rho_{k_j}^{\text{trial}} = \frac{\nabla f(\boldsymbol{x}_{k_j})^\top \boldsymbol{p}_{k_j} + \boldsymbol{p}_{k_j}^\top \boldsymbol{B}_{k_j} \boldsymbol{p}_{k_j}}{(1 - \sigma)\theta_{k_j} \|\boldsymbol{c}(\boldsymbol{x}_{k_j})\|_2} \le \frac{M_{\nabla f} M_{\boldsymbol{\ell}, \boldsymbol{u}} + \kappa_2 M_{\boldsymbol{\ell}, \boldsymbol{u}}^2}{(1 - \sigma)\tilde{\tau}\tilde{\theta} \|\boldsymbol{c}(\boldsymbol{x}_{k_j})\|_2},$$

we have $\lim_{j\to\infty} \|\boldsymbol{c}(\boldsymbol{x}_{k_j})\|_2 = 0$. By Lemmas 9 and 12, it is a contradiction.

Proposition 1 shows the boundedness of the penalty parameters from the constraint qualification perspective. More specifically, if EGMFCQ holds for all accumulation points of the sequence $\{x_k\}$, then the penalty parameter is guaranteed to be bounded. Given that we update the penalty parameter by multiplying it by a factor greater than one, it follows that the penalty parameter will eventually stabilize.

LEMMA 13. Under Assumptions 1 and 2, there exist sufficiently large $\widetilde{K} \in \mathbb{Z}_+$ and $\tilde{\rho} > 0$, such that $\rho_k = \tilde{\rho}$ for all $k \geq \tilde{K}$ and

(A.9)
$$\phi(\boldsymbol{x}_{k},\tilde{\rho}) - \phi(\boldsymbol{x}_{k+1},\tilde{\rho}) \geq \frac{\beta(1-\beta)\tau\kappa_{1}\tilde{\rho}\tilde{\tau}\theta\sigma}{\kappa_{\nabla f} + \kappa_{\nabla c}} \|\boldsymbol{c}(\boldsymbol{x}_{k})\|_{2} + \frac{\beta(1-\beta)\kappa_{1}^{2}}{2(\kappa_{\nabla f} + \kappa_{\nabla c})} \|\boldsymbol{p}_{k}\|_{2}^{2}$$

PROOF. By Lemma 12, the penalty parameter ρ_k becomes stable when $k \ge \tilde{K}$ for some sufficiently large $\tilde{K} \in \mathbb{Z}_+$, i.e., $\rho_k = \tilde{\rho}$ for $k \ge \tilde{K}$. Next, we only consider the iterates when ρ_k becomes stable. The backtracking line search guarantees that

$$\phi(\boldsymbol{x}_{k},\tilde{\rho}) - \phi(\boldsymbol{x}_{k+1},\tilde{\rho}) \geq \beta \alpha_{k} \Delta q(\boldsymbol{x}_{k},\boldsymbol{p}_{k},\nabla f(\boldsymbol{x}_{k}),\boldsymbol{B}_{k},\tilde{\rho}) \geq \beta \alpha_{k} \cdot \left(\frac{1}{2}\boldsymbol{p}_{k}^{\top}\boldsymbol{B}_{k}\boldsymbol{p}_{k} + \sigma \tilde{\rho} \theta_{k} \|\boldsymbol{c}(\boldsymbol{x}_{k})\|_{2}\right)$$

By Lemma 11, we have $\alpha_k \geq \frac{\tau(1-\beta)\kappa_1}{\kappa_{\nabla f}+\kappa_{\nabla c}}$ by the backtracking line search. Furthermore, by the positive-definiteness of B_k (i.e., $B_k \succeq \kappa_1 \mathbf{I}$) and the lower-boundedness of θ_k (i.e., $\theta_k \geq \tilde{\tau}\tilde{\theta}$), together with the stabilization of ρ_k (i.e., $\rho_k = \tilde{\rho}$) and the lower-boundedness of α_k (i.e., $\alpha_k \geq \frac{\tau(1-\beta)\kappa_1}{\kappa_{\nabla f}+\kappa_{\nabla c}}$), we complete the proof for Equation (A.9).

Proof for Theorem 1. It is a direct result of Lemma 13. Here, we only consider the case where the penalty parameter ρ_k becomes stable. By the boundedness of the feasible region (i.e., $\ell \le x \le u$) and smoothness of the objective and the constraints, we have that $\phi(x, \tilde{\rho})$ is (lower and upper) bounded. Then Equation (A.9) implies that

$$\frac{\beta(1-\beta)\tau\kappa_1\tilde{\rho}\tilde{\tau}\tilde{\theta}\sigma}{\kappa_{\nabla f}+\kappa_{\nabla c}}\sum_{k=\tilde{K}}^{\infty}\|\boldsymbol{c}(\boldsymbol{x}_k)\|_2+\frac{\beta(1-\beta)\kappa_1^2}{2(\kappa_{\nabla f}+\kappa_{\nabla c})}\sum_{k=\tilde{K}}^{\infty}\|\boldsymbol{p}_k\|_2^2<\infty,$$

which completes the proof for Equation (2.12). Conditions in Equation (A.7) show that $\|\nabla f(\boldsymbol{x}_k) + \nabla c(\boldsymbol{x}_k) \boldsymbol{\lambda}_k^{\text{sub}} - \boldsymbol{\mu}_{1,k}^{\text{sub}} + \boldsymbol{\mu}_{2,k}^{\text{sub}}\|_2 = \|\boldsymbol{B}_k \boldsymbol{p}_k\|_2 \le \kappa_2 \|\boldsymbol{p}_k\|_2, \|\boldsymbol{\mu}_{1,k}^{\text{sub}} \odot (\boldsymbol{x} - \boldsymbol{\ell})\|_2 \le \mu_{1,k}^{\text{sub}\top} (\boldsymbol{x} - \boldsymbol{\ell}) \le M_{\text{Lag}} \|\boldsymbol{p}_k\|_2$ and $\|\boldsymbol{\mu}_{2,k}^{\text{sub}} \odot (\boldsymbol{x} - \boldsymbol{u})\|_2 \le \mu_{1,k}^{\text{sub}\top} (\boldsymbol{u} - \boldsymbol{x}) \le M_{\text{Lag}} \|\boldsymbol{p}_k\|_2$, then Equation (2.13) is straightforward.

A.4. Proof for Lemma 2. Denote $\mathcal{A}^* = \mathcal{A}(\boldsymbol{x}^*) := \{i : (\boldsymbol{x}^*)_i = (\boldsymbol{\ell})_i \text{ or } (\boldsymbol{x}^*)_i = (\boldsymbol{\ell})_i\}$ the active set of inequality constraints at \boldsymbol{x}^* . Denote $\varepsilon = \min\{(\boldsymbol{x}^* - \boldsymbol{\ell})_i, (\boldsymbol{u} - \boldsymbol{x}^*)_i, i \in \mathcal{A}^{*-}\} > 0$. First, let \boldsymbol{x}_k be sufficiently close to \boldsymbol{x}^* such that $\|\boldsymbol{x}_k - \boldsymbol{x}^*\|_{\infty} \leq \frac{1}{4}\varepsilon$, then $\min\{(\boldsymbol{x}_k - \boldsymbol{\ell})_i, (\boldsymbol{u} - \boldsymbol{x}_k)_i, i \in \mathcal{A}^{*-}\} \geq \frac{3}{4}\varepsilon$. Since EGMFCQ holds at \boldsymbol{x}^* , there exists a vector $\boldsymbol{z}^* \in \mathbb{R}^d$ satisfying Condition (2.4) at \boldsymbol{x}^* . The fact that $\boldsymbol{c}(\boldsymbol{x}^*) = \boldsymbol{0}$ further implies that we can scale the vector \boldsymbol{z}^* by some constants such that

$$c(x^*) + \nabla c(x^*)^\top z^* = 0,$$

(z*)_i > 0, if (x*)_i = (ℓ)_i,
(z*)_i < 0, if (x*)_i = (u)_i,
 $||z^*||_{\infty} \le \varepsilon/2.$

By the smoothness of $c(\boldsymbol{x})$ and the linear independence of columns of $\nabla c(\boldsymbol{x}^*)$, we can find \boldsymbol{z}_k such that $c(\boldsymbol{x}_k) + \nabla c(\boldsymbol{x}_k)^\top \boldsymbol{z}_k = \boldsymbol{0}$ and $\|\boldsymbol{z}_k - \boldsymbol{z}^*\|_{\infty} \leq \frac{1}{2} \min\{|(\boldsymbol{z}^*)_i| : (\boldsymbol{x}^*)_i =$ $(\boldsymbol{\ell})_i$ or $(\boldsymbol{x}^*)_i = (\boldsymbol{u})_i\} \leq \frac{1}{4}\varepsilon$ as $\|\boldsymbol{x}_k - \boldsymbol{x}^*\|_2 \to 0$. Then $\|\boldsymbol{z}_k\|_{\infty} \leq \|\boldsymbol{z}_k - \boldsymbol{z}^*\|_{\infty} + \|\boldsymbol{z}^*\|_{\infty} \leq \frac{3}{4}\varepsilon;$ $(\boldsymbol{z}_k)_i > 0$, if $(\boldsymbol{x}^*)_i = (\boldsymbol{\ell})_i$ and $(\boldsymbol{z}_k)_i < 0$, if $(\boldsymbol{x}^*)_i = (\boldsymbol{u})_i$. Together with the fact that $\min\{(\boldsymbol{x}_k - \boldsymbol{\ell})_i, (\boldsymbol{u} - \boldsymbol{x}_k)_i, i \in \mathcal{A}^{*-}\} \geq \frac{3}{4}\varepsilon$, we show $\boldsymbol{\ell} \leq \boldsymbol{x}_k + \boldsymbol{z}_k \leq \boldsymbol{u}$. Therefore, $\theta_k = 1$ is always accepted if \boldsymbol{x}_k is sufficiently close to \boldsymbol{x}^* .

APPENDIX B: PROOF FOR THEOREM 2 AND 3

B.1. Some Technical Lemmas for Theorem 2. We first show that the adaptivity parameter will stabilize after sufficient iterations.

LEMMA 14. Under Assumption 3, there exists a constant $\bar{\xi} > 0$ such that $\xi_k = \bar{\xi}$ for all sufficiently large k.

PROOF. Observe that the sequence $\{\xi_k\}$ is monotonically decreasing and $\xi_k < \xi_{k-1}$ holds if and only if $\xi_k^{\text{trial}} < \xi_{k-1}$ and $\xi_k \le (1 - \epsilon_{\xi})\xi_{k-1}$. Suppose $\lim_{k\to\infty} \xi_k = 0$, then it follows that $\liminf_{k\to\infty} \xi_k^{\text{trial}} = 0$. However, the selection of ρ_k guarantees that $\Delta q(\boldsymbol{x}_k, \bar{\boldsymbol{p}}_k, \bar{\boldsymbol{g}}_k, \boldsymbol{B}_k, \rho_k) \ge \frac{1}{2}\bar{\boldsymbol{p}}_k^\top \boldsymbol{B}_k \bar{\boldsymbol{p}}_k \ge \frac{\kappa_1}{2} \|\bar{\boldsymbol{p}}_k\|_2^2$, implying that $\xi_k^{\text{trial}} \ge \frac{\kappa_1}{2}$. It is a contradiction. Therefore, we conclude that $\lim_{k\to\infty} \xi_k > 0$.

The following lemma is essential in our subsequent analysis and is extended from Lemma A.3 in [51]. The results investigate the competition and reveal the asymptotic behavior between the two sequences $\{\alpha_k\}$ and $\{\beta_k\}$. Importantly, we observe that when $\{\alpha_k\}$ decays faster than $\{\beta_k\}$, the asymptotic behavior of terms described in the lemma is dominated by the sequence $\{\alpha_k\}$, resulting in the asymptotic normality of the generated iterates with averaged gradient as studied in Section 4.

LEMMA 15 (Lemma A.3 in [51]). For two sequences $\{\alpha_k\}$ and $\{\beta_k\}$ satisfying $\alpha_k = \iota_1(k+1)^{-b_1}$ and $\beta_k = \iota_2(k+1)^{-b_2}$ with $\iota_1, \iota_2 > 0$ and $b_1, b_2 > 0$, the followings hold

1. Define $\chi = 0$ if $0 < b_2 < 1$ and $\chi = -\frac{b_1}{t_2}$ if $b_2 = 1$, then

$$\lim_{k \to \infty} \frac{1}{\alpha_k} \sum_{i=0}^k \prod_{j=i+1}^k \prod_{t=1}^l (1 - a_t \beta_j) \beta_i \alpha_i = \frac{1}{\sum_{t=1}^l a_t + \chi},$$

where we require that $\sum_{t=1}^{l} a_t + \chi > 0$. Moreover,

$$\lim_{k \to \infty} \left\{ \frac{1}{\alpha_k} \sum_{i=0}^k \prod_{j=i+1}^k \prod_{t=1}^l (1 - a_t \beta_j) \,\beta_i \alpha_i e_i + b \prod_{j=0}^k \prod_{t=1}^l (1 - a_t \beta_j) \right\} = 0,$$

for any $b \in \mathbb{R}$ and $e_i \to 0$.

2. If $0 < b_2 < b_1 \le 1$, then

$$\lim_{k \to \infty} \frac{1}{\alpha_k} \sum_{i=0}^k \prod_{j=i+1}^k (1 - \alpha_j) (1 - \beta_j) \alpha_i \beta_i = 1.$$

LEMMA 16. For two given sequence $\{\alpha_k\}$ and $\{\beta_k\}$ satisfying $\lim_{k\to\infty} \alpha_k = 0$, $\lim_{k\to\infty} \beta_k = 0$, and $\lim_{k\to\infty} \alpha_k/\beta_k = 0$, then

(B.1)
$$\lim_{k \to \infty} \mathbb{E}\left[\left\| \bar{\boldsymbol{g}}_k - \nabla f(\boldsymbol{x}_k) \right\|_2^2 \right] = 0.$$

Therefore, there exists a number $M_{\sigma} > 0$ such that

$$\mathbb{E}\left[\|\bar{\boldsymbol{g}}_k - \nabla f(\boldsymbol{x}_k)\|_2^2\right] \le M_{\sigma}^2,$$

under Assumption 3.

PROOF. By the update scheme of \bar{g}_k , we have

$$\begin{split} \bar{g}_{k} - \nabla f(\boldsymbol{x}_{k}) &= \beta_{k} \left(\boldsymbol{g}_{k} - \nabla f(\boldsymbol{x}_{k}) \right) + (1 - \beta_{k}) \left(\bar{g}_{k-1} - \nabla f(\boldsymbol{x}_{k-1}) \right) \\ &+ (1 - \beta_{k}) \left(\nabla f(\boldsymbol{x}_{k-1}) - \nabla f(\boldsymbol{x}_{k}) \right) \\ &= \beta_{k} \left(\boldsymbol{g}_{k} - \nabla f(\boldsymbol{x}_{k}) \right) + (1 - \beta_{k}) \left\{ \beta_{k-1} \left(\boldsymbol{g}_{k-1} - \nabla f(\boldsymbol{x}_{k-1}) \right) \right\} + (1 - \beta_{k-1}) \left(\bar{g}_{k-2} - \nabla f(\boldsymbol{x}_{k-2}) \right) \\ &+ (1 - \beta_{k-1}) \left(\nabla f(\boldsymbol{x}_{k-2}) - \nabla f(\boldsymbol{x}_{k-1}) \right) \right\} + (1 - \beta_{k}) \left(\nabla f(\boldsymbol{x}_{k-1}) - \nabla f(\boldsymbol{x}_{k}) \right) \\ &= \cdots \\ &= \sum_{i=0}^{k} \left(\prod_{j=i+1}^{k} (1 - \beta_{j}) \right) \beta_{i} \left(\boldsymbol{g}_{i} - \nabla f(\boldsymbol{x}_{i}) \right) \\ &+ \sum_{i=1}^{k} \left(\prod_{j=i}^{k} (1 - \beta_{j}) \right) \left(\nabla f(\boldsymbol{x}_{i-1}) - \nabla f(\boldsymbol{x}_{i}) \right) \\ &:= \mathcal{W}_{1} + \mathcal{W}_{2}. \end{split}$$

Here, both W_1 and W_2 are random variables. By Lemma 15 and the fact that

$$\|\mathcal{W}_2\|_2 \leq \sum_{i=1}^k \left(\prod_{j=i}^k (1-\beta_j)\right) \alpha_{i-1} M_{\ell,\boldsymbol{u}},$$

we have $\mathcal{W}_2 \to 0$ as $k \to \infty$ since $\lim_{k\to\infty} \alpha_{i-1}/\beta_i = 0$. It follows that $\lim_{k\to\infty} \mathbb{E}\left[\bar{g}_k - \nabla f(\boldsymbol{x}_k)\right] = \lim_{k\to\infty} \mathbb{E}\left[\mathcal{W}_1\right] = 0$, since

$$\mathbb{E}\left[\|\mathcal{W}_1\|_2^2\right]$$

= $\sum_{i=0}^k \left(\prod_{j=i+1}^k (1-\beta_j)\right)^2 \beta_i^2 \mathbb{E}\left[\|\boldsymbol{g}_i - \nabla f(\boldsymbol{x}_i)\|_2^2\right]$
 $\leq \sigma_g^2 \sum_{i=0}^k \left(\prod_{j=i+1}^k (1-\beta_j)\right)^2 \beta_i^2 \to 0 \text{ as } k \to \infty,$

where the last convergence result comes from Lemma 15. Therefore, $\lim_{k\to\infty} \mathbb{E}\left[\|\bar{g}_k - \nabla f(x_k)\|_2^2\right] \le 2\lim_{k\to\infty} \mathbb{E}\left[\|\mathcal{W}_1\|_2^2\right] + 2\lim_{k\to\infty} \|\mathcal{W}_2\|_2^2 = 0$, which completes the first part of the proof. The second result is straightforward since a convergent sequence must be bounded.

The above lemma establishes the convergence of the averaged gradient to the exact gradient in expectation, by utilizing the asymptotic behavior of two sequences in Lemma 15. The lemma not only assures us of the asymptotic validity of using \bar{g}_k as a surrogate for $\nabla f(x_k)$, but also offers a bound for their difference, lending confidence in the effectiveness of the algorithm. Following this, the next lemma studies the perturbation robustness property of the quadratic SQP subproblems and implies that the solutions are Lipschitz continuous with respect to the gradients. Consequently, the fact that the averaged gradient is asymptotically convergent to the exact gradient implies that the Debiased-StoSQP is arbitrarily close to the deterministic algorithm after sufficiently many iterations. This constitutes one of the most significant advantages of Debiased-StoSQP, employing averaged gradients, over other fully stochastic algorithms [6, 20].

NA ET AL.

 $\|-\boldsymbol{p}_{k}\|_{2} \leq \kappa_{1}^{-1} \|\bar{\boldsymbol{g}}_{k} - \nabla f(\boldsymbol{x}_{k})\|_{2},$

LEMMA 17. Suppose Assumptions 1, 2 and 3 hold, then

$$(B.2) \|\bar{\boldsymbol{p}}_k$$

and

$$\mathbb{E}_{k} \left\| \bar{\boldsymbol{p}}_{k} - \boldsymbol{p}_{k} \right\|_{2} \leq \kappa_{1}^{-1} \mathbb{E}_{k} \left\| \bar{\boldsymbol{g}}_{k} - \nabla f(\boldsymbol{x}_{k}) \right\|_{2}.$$

PROOF. The relaxed SQP subproblem at $m{x}_k$ with the averaged gradient $ar{m{g}}_k$ can be written as

$$\min_{\boldsymbol{p}\in\widetilde{\Omega}_{k}}\frac{1}{2}\left\|\boldsymbol{p}+\boldsymbol{B}_{k}^{-1}\bar{\boldsymbol{g}}_{k}\right\|_{\boldsymbol{B}_{k}}^{2},$$

which is a convex-constrained quadratic problem. The variational inequality implies that

$$\left\langle \boldsymbol{p}_{k}-\bar{\boldsymbol{p}}_{k},-\boldsymbol{B}_{k}^{-1}\bar{\boldsymbol{g}}_{k}-\bar{\boldsymbol{p}}_{k}\right\rangle _{\boldsymbol{B}_{k}}\leq0.$$

Since p_k is the solution of the relaxed SQP subproblem at x_k with exact gradient $\nabla f(x_k)$, we similarly have

$$\left\langle \bar{\boldsymbol{p}}_{k} - \boldsymbol{p}_{k}, -\boldsymbol{B}_{k}^{-1} \nabla f(\boldsymbol{x}_{k}) - \boldsymbol{p}_{k} \right\rangle_{\boldsymbol{B}_{k}} \leq 0.$$

Summing up the above two inequalities, we have

(B.3)

$$0 \ge \langle \boldsymbol{p}_{k} - \bar{\boldsymbol{p}}_{k}, -\boldsymbol{B}_{k}^{-1}\bar{\boldsymbol{g}}_{k} - \bar{\boldsymbol{p}}_{k} + \boldsymbol{B}_{k}^{-1}\nabla f(\boldsymbol{x}_{k}) + \boldsymbol{p}_{k} \rangle_{\boldsymbol{B}_{k}}$$

$$= \|\boldsymbol{p}_{k} - \bar{\boldsymbol{p}}_{k}\|_{\boldsymbol{B}_{k}}^{2} + \langle \bar{\boldsymbol{p}}_{k} - \boldsymbol{p}_{k}, \bar{\boldsymbol{g}}_{k} - \nabla f(\boldsymbol{x}_{k}) \rangle$$

$$\ge \|\boldsymbol{p}_{k} - \bar{\boldsymbol{p}}_{k}\|_{\boldsymbol{B}_{k}}^{2} - \|\boldsymbol{p}_{k} - \bar{\boldsymbol{p}}_{k}\|_{2} \cdot \|\bar{\boldsymbol{g}}_{k} - \nabla f(\boldsymbol{x}_{k})\|_{2}.$$

Note that $\|\boldsymbol{p}_k - \bar{\boldsymbol{p}}_k\|_{\boldsymbol{B}_k}^2 \ge \kappa_1 \|\boldsymbol{p}_k - \bar{\boldsymbol{p}}_k\|_2^2$, combining with Assumption 3, we complete the proof.

LEMMA 18. Suppose that Assumptions 2 and 3 hold, then

(B.4)
$$\mathbb{E}_{k}\left[\left|\left(\nabla f(\boldsymbol{x}_{k})-\bar{\boldsymbol{g}}_{k}\right)^{\top}\bar{\boldsymbol{p}}_{k}\right|\right] \leq M_{\sigma}M_{\boldsymbol{\ell},\boldsymbol{u}},$$

(B.5)
$$\mathbb{E}_{k}\left[\left|\nabla f(\boldsymbol{x}_{k})^{\top}\boldsymbol{p}_{k}-\bar{\boldsymbol{g}}_{k}^{\top}\bar{\boldsymbol{p}}_{k}\right|\right] \leq M_{\sigma}M_{\boldsymbol{\ell},\boldsymbol{u}}+2\left(M_{\nabla f}+M_{\sigma}\right)M_{\boldsymbol{\ell},\boldsymbol{u}},$$

and

(B.6)
$$\mathbb{E}_{k}\left[\left|\boldsymbol{p}_{k}^{\top}\boldsymbol{B}_{k}\boldsymbol{p}_{k}-\bar{\boldsymbol{p}}_{k}^{\top}\boldsymbol{B}_{k}\bar{\boldsymbol{p}}_{k}\right|\right] \leq 2\kappa_{1}^{-1}\kappa_{2}M_{\boldsymbol{\ell},\boldsymbol{u}}M_{\sigma}.$$

PROOF. The first relation is straightforward since $\mathbb{E}_k[\|\bar{\boldsymbol{g}}_k - \nabla f(\boldsymbol{x}_k)\|_2] \leq M_{\sigma}$ and $\|\bar{\boldsymbol{p}}_k\|_2 \leq M_{\ell,\boldsymbol{u}}$. By triangle inequalities, we have

$$\mathbb{E}_{k}\left[\left|\nabla f(\boldsymbol{x}_{k})^{\top}\boldsymbol{p}_{k}-\bar{\boldsymbol{g}}_{k}^{\top}\bar{\boldsymbol{p}}_{k}\right|\right]$$
$$=\mathbb{E}_{k}\left[\left|\left(\nabla f(\boldsymbol{x}_{k})-\bar{\boldsymbol{g}}_{k}\right)^{\top}\boldsymbol{p}_{k}+\bar{\boldsymbol{g}}_{k}^{\top}\left(\boldsymbol{p}_{k}-\bar{\boldsymbol{p}}_{k}\right)\right|\right]$$
$$\leq M_{\sigma}M_{\boldsymbol{\ell},\boldsymbol{u}}+2\left(M_{\nabla f}+M_{\sigma}\right)M_{\boldsymbol{\ell},\boldsymbol{u}},$$

and

$$\mathbb{E}_{k}\left[\left|\boldsymbol{p}_{k}^{\top}\boldsymbol{B}_{k}\boldsymbol{p}_{k}-\bar{\boldsymbol{p}}_{k}^{\top}\boldsymbol{B}_{k}\bar{\boldsymbol{p}}_{k}\right|\right]$$
$$=\mathbb{E}_{k}\left[\left|\left(\boldsymbol{p}_{k}-\bar{\boldsymbol{p}}_{k}\right)^{\top}\boldsymbol{B}_{k}\left(\boldsymbol{p}_{k}+\bar{\boldsymbol{p}}_{k}\right)\right|\right]$$
$$\leq 2\kappa_{2}M_{\boldsymbol{\ell},\boldsymbol{u}}\mathbb{E}_{k}\left[\left\|\boldsymbol{p}_{k}-\bar{\boldsymbol{p}}_{k}\right\|\right] \leq 2\kappa_{1}^{-1}\kappa_{2}M_{\boldsymbol{\ell},\boldsymbol{u}}M_{\sigma}.$$

LEMMA 19. Suppose that Assumptions 1, 2 and 3 hold, then

(B.7)
$$\mathbb{E}_{k}\left[\left|\Delta q(\boldsymbol{x}_{k},\boldsymbol{p}_{k},\nabla f(\boldsymbol{x}_{k}),\boldsymbol{B}_{k};\rho_{k})-\Delta q(\boldsymbol{x}_{k},\bar{\boldsymbol{p}}_{k},\bar{\boldsymbol{g}}_{k},\boldsymbol{B}_{k};\rho_{k})\right|\right]\\\leq M_{\sigma}M_{\boldsymbol{\ell},\boldsymbol{u}}+2\left(M_{\nabla f}+M_{\sigma}\right)M_{\boldsymbol{\ell},\boldsymbol{u}}+\kappa_{1}^{-1}\kappa_{2}M_{\boldsymbol{\ell},\boldsymbol{u}}M_{\sigma}.$$

PROOF. Note that $\Delta q(\boldsymbol{x}_k, \boldsymbol{p}_k, \nabla f(\boldsymbol{x}_k), \boldsymbol{B}_k; \rho_k) = -\nabla f(\boldsymbol{x}_k)^\top \boldsymbol{p}_k - \frac{1}{2} \boldsymbol{p}_k^\top \boldsymbol{B}_k \boldsymbol{p}_k + \rho_k \theta_k \|\boldsymbol{c}(\boldsymbol{x}_k)\|_2$, where the last term $\rho_k \theta_k \|\boldsymbol{c}(\boldsymbol{x}_k)\|_2$ is independent of $\nabla f(\boldsymbol{x}_k)$ and \boldsymbol{p}_k . Using results in Lemma 18, we have

$$\mathbb{E}_{k} \left[\left| \Delta q(\boldsymbol{x}_{k}, \boldsymbol{p}_{k}, \nabla f(\boldsymbol{x}_{k}), \boldsymbol{B}_{k}; \rho_{k}) - \Delta q(\boldsymbol{x}_{k}, \bar{\boldsymbol{p}}_{k}, \bar{\boldsymbol{g}}_{k}, \boldsymbol{B}_{k}; \rho_{k}) \right| \right] \\ = \mathbb{E}_{k} \left[\left| -\nabla f(\boldsymbol{x}_{k})^{\top} \boldsymbol{p}_{k} + \bar{\boldsymbol{g}}_{k}^{\top} \bar{\boldsymbol{p}}_{k} - \frac{1}{2} \boldsymbol{p}_{k}^{\top} \boldsymbol{B}_{k} \boldsymbol{p}_{k} + \frac{1}{2} \bar{\boldsymbol{p}}_{k}^{\top} \boldsymbol{B}_{k} \bar{\boldsymbol{p}}_{k} \right| \right] \\ \leq \mathbb{E}_{k} \left[\left| -\nabla f(\boldsymbol{x}_{k})^{\top} \boldsymbol{p}_{k} + \bar{\boldsymbol{g}}_{k}^{\top} \bar{\boldsymbol{p}}_{k} \right| \right] + \frac{1}{2} \mathbb{E}_{k} \left[\left| \boldsymbol{p}_{k}^{\top} \boldsymbol{B}_{k} \boldsymbol{p}_{k} - \bar{\boldsymbol{p}}_{k}^{\top} \boldsymbol{B}_{k} \bar{\boldsymbol{p}}_{k} \right| \right] \\ \leq M_{\sigma} M_{\boldsymbol{\ell}, \boldsymbol{u}} + 2 \left(M_{\nabla f} + M_{\sigma} \right) M_{\boldsymbol{\ell}, \boldsymbol{u}} + \kappa_{1}^{-1} \kappa_{2} M_{\boldsymbol{\ell}, \boldsymbol{u}} M_{\sigma}.$$

LEMMA 20. In line 10 of Algorithm 2, the step size α_k is selected from the interval $[\alpha_k^{\min}, \alpha_k^{\max}] := \left[\frac{\xi_k \gamma_k}{\kappa_{\nabla f} + \rho_k \kappa_{\nabla c}}, \frac{\xi_k \gamma_k}{\kappa_{\nabla f} + \rho_k \kappa_{\nabla c}} + \varrho \gamma_k^2\right]$, then

(B.8)
$$\mathbb{E}_{k}\left[\alpha_{k}\nabla f(\boldsymbol{x}_{k})^{\top}(\bar{\boldsymbol{p}}_{k}-\boldsymbol{p}_{k})\right] \leq \varrho \gamma_{k}^{2} M_{\nabla f} \kappa_{1}^{-1} M_{\sigma} + \alpha_{k}^{min} \mathbb{E}_{k}\left[\nabla f(\boldsymbol{x}_{k})^{\top}(\bar{\boldsymbol{p}}_{k}-\boldsymbol{p}_{k})\right],$$

under Assumptions 2 and 3.

PROOF. Denote the event
$$C_k = \{\nabla f(\boldsymbol{x}_k)^\top (\bar{\boldsymbol{p}}_k - \boldsymbol{p}_k) \ge 0\}$$
, then

$$\mathbb{E}_k \left[\alpha_k \nabla f(\boldsymbol{x}_k)^\top (\bar{\boldsymbol{p}}_k - \boldsymbol{p}_k) \right]$$

$$= \mathbb{E}_k \left[\alpha_k \nabla f(\boldsymbol{x}_k)^\top (\bar{\boldsymbol{p}}_k - \boldsymbol{p}_k) |C_k \right] + \mathbb{E}_k \left[\alpha_k \nabla f(\boldsymbol{x}_k)^\top (\bar{\boldsymbol{p}}_k - \boldsymbol{p}_k) |C_k^c \right]$$

$$\leq \alpha_k^{\max} \mathbb{E}_k \left[\nabla f(\boldsymbol{x}_k)^\top (\bar{\boldsymbol{p}}_k - \boldsymbol{p}_k) |C_k \right] + \alpha_k^{\min} \mathbb{E}_k \left[\nabla f(\boldsymbol{x}_k)^\top (\bar{\boldsymbol{p}}_k - \boldsymbol{p}_k) |C_k^c \right]$$

$$= \alpha_k^{\min} \mathbb{E}_k \left[\nabla f(\boldsymbol{x}_k)^\top (\bar{\boldsymbol{p}}_k - \boldsymbol{p}_k) \right] + \left(\alpha_k^{\max} - \alpha_k^{\min} \right) \mathbb{E}_k \left[\nabla f(\boldsymbol{x}_k)^\top (\bar{\boldsymbol{p}}_k - \boldsymbol{p}_k) |C_k \right]$$

LEMMA 21. Under Assumptions 1, 2 and 3, if $\sum_{k=0}^{\infty} \gamma_k = \infty$, $\sum_{k=0}^{\infty} \gamma_k^2 < \infty$ and $\sum_{k=\bar{K}}^{\infty} \alpha_k^{\min} \mathbb{E}[\|\bar{p}_k - p_k\|_2] < \infty$, then ∞

(B.9)
$$\sum_{k=\bar{K}}^{\infty} \alpha_k^{\min} \Delta q(\boldsymbol{x}_k, \boldsymbol{p}_k, \nabla f(\boldsymbol{x}_k), \boldsymbol{B}_k; \bar{\rho}) < \infty, \text{ almost surely.}$$

It further implies that

(B.10)
$$\liminf_{k \to \infty} \Delta q(\boldsymbol{x}_k, \boldsymbol{p}_k, \nabla f(\boldsymbol{x}_k), \boldsymbol{B}_k; \bar{\rho}) = 0, \text{ almost surely.}$$

PROOF. We only consider the case when ρ_k and ξ_k becomes stable, i.e., $\rho_k = \bar{\rho}$ and $\xi_k = \bar{\xi}$ for all $k \ge \bar{K}$. It follows from Assumptions 1, 2 and 3 that

$$\begin{aligned} \text{(B.11)} \\ & \mathbb{E}_{k} \left[\phi(\boldsymbol{x}_{k+1}, \bar{\rho}) - \phi(\boldsymbol{x}_{k}, \bar{\rho}) \right] \\ = \mathbb{E}_{k} \left[f(\boldsymbol{x}_{k} + \alpha_{k} \bar{\boldsymbol{p}}_{k}) - f(\boldsymbol{x}_{k}) + \bar{\rho} \left(\|\boldsymbol{c}(\boldsymbol{x}_{k} + \alpha_{k} \bar{\boldsymbol{p}}_{k}) \|_{2} - \|\boldsymbol{c}(\boldsymbol{x}_{k}) \|_{2} \right) \right] \\ \leq \mathbb{E}_{k} \left[\alpha_{k} \nabla f(\boldsymbol{x}_{k})^{\top} \bar{\boldsymbol{p}}_{k} + \frac{\kappa_{\nabla f}}{2} \alpha_{k}^{2} \| \bar{\boldsymbol{p}}_{k} \|_{2}^{2} + \bar{\rho} \left(\left\| \boldsymbol{c}(\boldsymbol{x}_{k}) + \alpha_{k} \nabla \boldsymbol{c}(\boldsymbol{x}_{k})^{\top} \bar{\boldsymbol{p}}_{k} \right) \right\|_{2}^{2} - \|\boldsymbol{c}(\boldsymbol{x}_{k}) \|_{2}^{2} + \frac{\kappa_{\nabla c}}{2} \alpha_{k}^{2} \| \bar{\boldsymbol{p}}_{k} \|_{2}^{2} \right) \right] \\ = \mathbb{E}_{k} \left[\alpha_{k} \nabla f(\boldsymbol{x}_{k})^{\top} \bar{\boldsymbol{p}}_{k} - \alpha_{k} \bar{\rho} \theta_{k} \| \boldsymbol{c}(\boldsymbol{x}_{k}) \|_{2}^{2} + \frac{\kappa_{\nabla f} + \bar{\rho} \kappa_{\nabla c}}{2} \alpha_{k}^{2} \| \bar{\boldsymbol{p}}_{k} \|_{2}^{2} \right] \\ = \mathbb{E}_{k} \left[\alpha_{k} \nabla f(\boldsymbol{x}_{k})^{\top} \boldsymbol{p}_{k} + \alpha_{k} \nabla f(\boldsymbol{x}_{k})^{\top} \left(\bar{\boldsymbol{p}}_{k} - \boldsymbol{p}_{k} \right) - \alpha_{k} \bar{\rho} \theta_{k} \| \boldsymbol{c}(\boldsymbol{x}_{k}) \|_{2}^{2} + \frac{\kappa_{\nabla f} + \bar{\rho} \kappa_{\nabla c}}{2} \alpha_{k}^{2} \| \bar{\boldsymbol{p}}_{k} \|_{2}^{2} \right] \\ = \mathbb{E}_{k} \left[\alpha_{k} \nabla f(\boldsymbol{x}_{k}) \bar{\boldsymbol{p}}_{k} + \alpha_{k} \nabla f(\boldsymbol{x}_{k})^{\top} \left(\bar{\boldsymbol{p}}_{k} - \boldsymbol{p}_{k} \right) + \frac{\kappa_{\nabla f} + \bar{\rho} \kappa_{\nabla c}}{2} \alpha_{k}^{2} \| \bar{\boldsymbol{p}}_{k} \|_{2}^{2} \right] \\ = \mathbb{E}_{k} \left[\alpha_{k} \Delta q(\boldsymbol{x}_{k}, \boldsymbol{p}_{k}, \nabla f(\boldsymbol{x}_{k}), \boldsymbol{B}_{k}, \bar{\rho}) - \frac{\alpha_{k}}{2} \boldsymbol{p}_{k}^{\top} \boldsymbol{B}_{k} \boldsymbol{p}_{k} + \alpha_{k} \nabla f(\boldsymbol{x}_{k})^{\top} \left(\bar{\boldsymbol{p}}_{k} - \boldsymbol{p}_{k} \right) + \frac{\kappa_{\nabla f} + \bar{\rho} \kappa_{\nabla c}}{2} \alpha_{k}^{2} \| \bar{\boldsymbol{p}}_{k} \|_{2}^{2} \right] \\ \leq \mathbb{E}_{k} \left[-\alpha_{k} \Delta q(\boldsymbol{x}_{k}, \boldsymbol{p}_{k}, \nabla f(\boldsymbol{x}_{k}), \boldsymbol{B}_{k}, \bar{\rho}) - \frac{\alpha_{k}}{2} \boldsymbol{p}_{k}^{\top} \boldsymbol{B}_{k} \boldsymbol{p}_{k} + \alpha_{k} \nabla f(\boldsymbol{x}_{k})^{\top} \left(\bar{\boldsymbol{p}}_{k} - \boldsymbol{p}_{k} \right) \\ + \frac{1}{2} \alpha_{k} \gamma_{k} \Delta q(\boldsymbol{x}_{k}, \bar{\boldsymbol{p}}_{k}, \nabla \bar{f}(\boldsymbol{x}_{k}), \boldsymbol{B}_{k}, \bar{\rho}) \right], \\ = \mathbb{E}_{k} \left[\left(-\alpha_{k} + \frac{1}{2} \alpha_{k} \gamma_{k} \left(\Delta q(\boldsymbol{x}_{k}, \bar{\boldsymbol{p}}_{k}, \nabla \bar{f}(\boldsymbol{x}_{k}), \boldsymbol{B}_{k}, \bar{\rho}) - \Delta q(\boldsymbol{x}_{k}, \boldsymbol{p}_{k}, \nabla f(\boldsymbol{x}_{k}), \boldsymbol{B}_{k}, \bar{\rho}) \right) \right], \end{aligned}$$

where the last inequality comes from the selection of ξ_k and α_k in lines 10 and 11, respectively. Without the loss of generality, we assume that $\gamma_k \leq 1$ and continue from Equation (B.11),

$$\mathbb{E}_{k} \left[\phi(\boldsymbol{x}_{k+1}, \bar{\rho}) - \phi(\boldsymbol{x}_{k}, \bar{\rho}) \right] \\
\leq \mathbb{E}_{k} \left[-\frac{1}{2} \alpha_{k}^{\min} \Delta q(\boldsymbol{x}_{k}, \boldsymbol{p}_{k}, \nabla f(\boldsymbol{x}_{k}), \boldsymbol{B}_{k}, \bar{\rho}) \right] + M_{\nabla f} \alpha_{k}^{\min} \mathbb{E}_{k} \left[\| \bar{\boldsymbol{p}}_{k} - \boldsymbol{p}_{k} \|_{2} \right] + \varrho \gamma_{k}^{2} M_{\nabla f} \kappa_{1}^{-1} M_{\sigma} \\
+ \frac{1}{2} \alpha_{k}^{\max} \gamma_{k} \left(M_{\sigma} M_{\boldsymbol{\ell}, \boldsymbol{u}} + 2 \left(M_{\nabla f} + M_{\sigma} \right) M_{\boldsymbol{\ell}, \boldsymbol{u}} + \kappa_{1}^{-1} \kappa_{2} M_{\boldsymbol{\ell}, \boldsymbol{u}} M_{\sigma} \right),$$

where the inequality is due to Lemmas 19 and 20. It further implies that

$$\mathbb{E}_{k}\left[\phi(\boldsymbol{x}_{k+1},\bar{\rho}) - \min_{\boldsymbol{\ell} \leq \boldsymbol{x} \leq \boldsymbol{u}} \phi(\boldsymbol{x},\bar{\rho}) | \mathcal{F}_{k-1}\right]$$

$$\leq \phi(\boldsymbol{x}_{k},\bar{\rho}) - \min_{\boldsymbol{\ell} \leq \boldsymbol{x} \leq \boldsymbol{u}} \phi(\boldsymbol{x},\bar{\rho}) - \frac{1}{2} \sum_{k=\bar{K}}^{\bar{K}+K} \alpha_{k}^{\min} \mathbb{E}_{k} \left[\Delta q(\boldsymbol{x}_{k},\boldsymbol{p}_{k},\nabla f(\boldsymbol{x}_{k}),\boldsymbol{B}_{k},\bar{\rho}) | \mathcal{F}_{k-1}\right]$$

$$(B.12) + M_{\nabla f} \sum_{k=\bar{K}}^{\bar{K}+K} \alpha_{k}^{\min} \mathbb{E}_{k} \left[\| \bar{\boldsymbol{p}}_{k} - \boldsymbol{p}_{k} \|_{2} | \mathcal{F}_{k-1} \right] + \varrho M_{\nabla f} \kappa_{1}^{-1} M_{\sigma} \sum_{k=\bar{K}}^{\bar{K}+K} \gamma_{k}^{2}$$

$$+ \frac{1}{2} \left(M_{\sigma} M_{\boldsymbol{\ell},\boldsymbol{u}} + 2 \left(M_{\nabla f} + M_{\sigma} \right) M_{\boldsymbol{\ell},\boldsymbol{u}} + \kappa_{1}^{-1} \kappa_{2} M_{\boldsymbol{\ell},\boldsymbol{u}} M_{\sigma} \right) \sum_{k=\bar{K}}^{\bar{K}+K} \alpha_{k}^{\max} \gamma_{k}.$$

Since $\sum_{k=0}^{\infty} \gamma_k^2 < \infty$, $\alpha_k^{\min} = \mathcal{O}(\gamma_k)$ and $\alpha_k^{\max} = \mathcal{O}(\gamma_k + \gamma_k^2)$, we have $\sum_{k=0}^{\infty} \alpha_k^{\max} \gamma_k < \infty$. Note that $\mathbb{E}\left[\sum_{k=\bar{K}}^{\infty} \alpha_k^{\min} \mathbb{E}_k [\|\bar{p}_k - p_k\|_2 |\mathcal{F}_{k-1}]\right] = \sum_{k=\bar{K}}^{\infty} \alpha_k^{\min} \mathbb{E}[\|\bar{p}_k - p_k\|_2] < \infty$ shows that $\sum_{k=\bar{K}}^{\infty} \alpha_k^{\min} \mathbb{E}_k [\|\bar{p}_k - p_k\|_2 |\mathcal{F}_{k-1}] < \infty$. We conclude from Equation (B.12), the step size $\alpha_k^{\min} = \mathcal{O}(\gamma_k)$, the assumption $\sum_{k=\bar{K}}^{\infty} \alpha_k^{\min} \mathbb{E}_k [\|\bar{p}_k - p_k\|_2 |\mathcal{F}_{k-1}] < \infty$ and Robbins-Siegmund theorem [59] that Equation (B.9) holds. Moreover, since $\sum_{k=0}^{\infty} \gamma_k = \infty$, together with Equation (B.9), we obtain Equation (B.10).

B.2. Proof for Theorem 2. Note that although $\{\gamma_k\}$ is the pre-defined sequence in the algorithm, the only difference between α_k and γ_k is a constant. Therefore, $\alpha_k^{\min} = \iota_1(k + 1)^{-1}$ implies that $\gamma_k = \iota_3(k + 1)^{-1}$ for some $\iota_3 > 0$. For simplicity, we directly discuss the behavior of the sequence related to α_k^{\min} rather than γ_k . We need the techniques and notations in the proof of Lemma 16, where all conditions are satisfied and $\mathbb{E}\left[\|\bar{\boldsymbol{g}}_k - \nabla f(\boldsymbol{x}_k)\|_2^2\right] \le 2\mathbb{E}\left[\|\mathcal{W}_1\|_2^2\right] + 2\mathbb{E}\left[\|\mathcal{W}_2\|_2^2\right]$. In details,

$$\mathbb{E}\left[\left\|\mathcal{W}_{1}\right\|_{2}^{2}\right] = \sum_{i=0}^{k} \left(\prod_{j=i+1}^{k} (1-\beta_{k})\right)^{2} \beta_{i}^{2} \mathbb{E}\left[\left\|\boldsymbol{g}_{i}-\nabla f(\boldsymbol{x}_{i})\right\|_{2}^{2}\right]$$
$$\leq \sigma_{g}^{2} \sum_{i=0}^{k} \left(\prod_{j=i+1}^{k} (1-\beta_{k})\right)^{2} \beta_{i}^{2} = \mathcal{O}\left(\beta_{k}\right),$$

by utilizing Lemma 15. Similarly, for $\|\mathcal{W}_2\|_2$, we have

$$\|\mathcal{W}_2\|_2 \le M_{\ell,\boldsymbol{u}} \sum_{i=1}^k \left(\prod_{j=i}^k (1-\beta_j) \right) \alpha_{i-1}^{\min} = \mathcal{O}\left(\alpha_k^{\min} / \beta_k \right).$$

Therefore, we conclude that $\mathbb{E}[\|\bar{\boldsymbol{g}}_k - \nabla f(\boldsymbol{x}_k)\|_2] \leq \mathcal{O}(\beta_k + \alpha_k^{\min}/\beta_k)$. Since $\alpha_k^{\min} = \iota_1(k+1)^{-b_1}$, it is not difficult to verify that $\sum_{k=\bar{K}}^{\infty} \alpha_k^{\min} \mathbb{E}[\|\bar{\boldsymbol{p}}_k - \boldsymbol{p}_k\|_2] < \infty$, if $b_1 + \frac{b_2}{2} > 1$ and $2b_1 - b_2 > 1$. We equivalently require that $b_1 \in (\frac{3}{4}, 1]$ and $b_2 \in (2 - 2b_1, 2b_1 - 1)$.

B.3. Proof for Theorem 3. We first show that the Problem (3.6) is convex and then the corresponding solution $(\lambda_k^*, \mu_{1,k}^*, \mu_{2,k}^*)$ is well-defined. The stability of quadratic problems in Lemma 23 is an generalization of Lemma 17.

LEMMA 22. Problem (3.6) is convex, i.e., the Hessian matrix $\nabla^2 F(\lambda, \mu_1, \mu_2; x)$ is positive semi-definite for any x, λ, μ_1 and μ_2 .

PROOF. The direct computation of the Hessian matrix for $F(\lambda, \mu_1, \mu_2; x)$ is (B.13)

$$\nabla^2 F(\boldsymbol{\lambda}, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2; \boldsymbol{x}) = \begin{pmatrix} 2\nabla \boldsymbol{c}(\boldsymbol{x})^\top \nabla \boldsymbol{c}(\boldsymbol{x}) & -2\nabla \boldsymbol{c}(\boldsymbol{x}) \\ -2\nabla \boldsymbol{c}(\boldsymbol{x})^\top & 2\boldsymbol{I} + 2\text{diag}\left((\boldsymbol{x}-\boldsymbol{\ell})^2\right) & -2\boldsymbol{I} \\ 2\nabla \boldsymbol{c}(\boldsymbol{x})^\top & -2\boldsymbol{I} & 2\boldsymbol{I} + 2\text{diag}\left((\boldsymbol{x}-\boldsymbol{u})^2\right) \end{pmatrix}.$$

For any vector $\boldsymbol{w} = (\boldsymbol{w}_1, \boldsymbol{w}_2, \boldsymbol{w}_3) \in \mathbb{R}^r \times \mathbb{R}^d \times \mathbb{R}^d$, $\boldsymbol{w}^\top \nabla^2 F(\boldsymbol{\lambda}, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2; \boldsymbol{x}) \boldsymbol{w} = \|\nabla \boldsymbol{c}(\boldsymbol{x}) \boldsymbol{w}_1 - \boldsymbol{w}_2 + \boldsymbol{w}_3\|_2^2 + \|(\boldsymbol{x} - \boldsymbol{\ell}) \odot \boldsymbol{w}_2\|_2^2 + \|(\boldsymbol{x} - \boldsymbol{u}) \odot \boldsymbol{w}_3\|_2^2 \ge 0.$ Therefore, Problem (3.6) is convex. LEMMA 23 (Stability of Quadratic Programs, Theorem 2.1 in [21]). For two constrained strongly convex quadratic problems

$$oldsymbol{y}^* \in \min_{oldsymbol{y} \in \Lambda} oldsymbol{g}^ op oldsymbol{y} + rac{1}{2}oldsymbol{y}^ op oldsymbol{Q} oldsymbol{y},$$

and

$$oldsymbol{y}^{**} \in \min_{oldsymbol{y} \in \Lambda} oldsymbol{g}^{\prime op} oldsymbol{y} + rac{1}{2} oldsymbol{y}^{ op} oldsymbol{Q}^{\prime} oldsymbol{y},$$

where the feasible region Λ is convex. Suppose that $\max\{\|\boldsymbol{y}^*\|_2, \|\boldsymbol{y}^{**}\|_2\} \leq M_y$ for some $M_y > 0$. If $\epsilon = \max\{\|\boldsymbol{g} - \boldsymbol{g}'\|_2, \|\boldsymbol{Q} - \boldsymbol{Q}'\|_2\}$ and both two matrices $\boldsymbol{Q}, \boldsymbol{Q}'$ are positive definite with $v_1 \boldsymbol{I} \leq \boldsymbol{Q}, \boldsymbol{Q}' \leq v_2 \boldsymbol{I}$, for some $0 < v_1 \leq v_2$. Then, the following holds

$$\|\boldsymbol{y}^* - \boldsymbol{y}^{**}\|_2 \le \epsilon v_1^{-1} (1 + M_y)$$

LEMMA 24. Under assumptions in Theorem 3, we have

$$\lim_{k\to\infty} F(\boldsymbol{\lambda}_k^*, \boldsymbol{\mu}_{1,k}^*, \boldsymbol{\mu}_{2,k}^*; \boldsymbol{x}_k) = 0, \text{ almost surely.}$$

PROOF. Denote $(\lambda_k^{\text{sub}}, \mu_{1,k}^{\text{sub}}, \mu_{2,k}^{\text{sub}})$ as Lagrangian multipliers of the relaxed SQP subproblem at x_k with full gradient $\nabla f(x_k)$. It follows that

$$\begin{split} F(\boldsymbol{\lambda}_{k}^{*}, \boldsymbol{\mu}_{1,k}^{*}, \boldsymbol{\mu}_{2,k}^{*}; \boldsymbol{x}_{k}) &\leq F(\boldsymbol{\lambda}_{k}^{\text{sub}}, \boldsymbol{\mu}_{1,k}^{\text{sub}}, \boldsymbol{\mu}_{2,k}^{\text{sub}}; \boldsymbol{x}_{k}) \leq \|\boldsymbol{B}_{k}\boldsymbol{p}_{k}\|_{2}^{2} + \left\|\boldsymbol{\mu}_{1,k}^{\text{sub}} \odot \boldsymbol{p}_{k}\right\|_{2}^{2} + \left\|\boldsymbol{\mu}_{2,k}^{\text{sub}} \odot \boldsymbol{p}_{k}\right\|_{2}^{2} \\ &\leq (\kappa_{2} + 2M_{\text{Lag}}) \left\|\boldsymbol{p}_{k}\right\|_{2}^{2}, \end{split}$$

then

$$\liminf_{k\to\infty} F(\boldsymbol{\lambda}_k^*, \boldsymbol{\mu}_{1,k}^*, \boldsymbol{\mu}_{2,k}^*; \boldsymbol{x}_k) = 0,$$

by $\liminf_{k \to \infty} \|\boldsymbol{p}_k\|_2 = 0$. Suppose that $\limsup_{k \to \infty} F(\boldsymbol{\lambda}_k^*, \boldsymbol{\mu}_{1,k}^*, \boldsymbol{\mu}_{2,k}^*; \boldsymbol{x}_k) > 0$, we can find a sufficiently small number $\varepsilon > 0$ and two infinite sequences $\{m_i\}$ and $\{n_i\}$ with $\bar{K} \le m_i < n_i$, such that

$$F(\boldsymbol{\lambda}_{m_i}^*, \boldsymbol{\mu}_{1,m_i}^*, \boldsymbol{\mu}_{2,m_i}^*; \boldsymbol{x}_{m_i}) > 2\varepsilon, \quad \|\boldsymbol{p}_{n_i}\|_2 \le \sqrt{\frac{\varepsilon}{\kappa_2 + M_{\text{Lag}}}},$$

and

$$\|\boldsymbol{p}_k\|_2 \ge \sqrt{rac{arepsilon}{\kappa_2 + M_{\text{Lag}}}}, ext{ for } m_i \le k < n_i.$$

Note that we can always achieve it due to the following derivation

(B.14)

$$F(\boldsymbol{\lambda}_{k}^{*}, \boldsymbol{\mu}_{1,k}^{*}, \boldsymbol{\mu}_{2,k}^{*}; \boldsymbol{x}_{k}) = \min_{\boldsymbol{\lambda}, \boldsymbol{\mu}_{1} \ge \boldsymbol{0}, \boldsymbol{\mu}_{2} \ge \boldsymbol{0}} F(\boldsymbol{\lambda}, \boldsymbol{\mu}_{1}, \boldsymbol{\mu}_{2}; \boldsymbol{x}_{k})$$

$$\leq \min_{\boldsymbol{\lambda}, \boldsymbol{\mu}_{1} \ge \boldsymbol{0}, \boldsymbol{\mu}_{2} \ge \boldsymbol{0}} \left\{ F(\boldsymbol{\lambda}, \boldsymbol{\mu}_{1}, \boldsymbol{\mu}_{2}; \boldsymbol{x}_{k}) + \frac{\varepsilon}{6M_{\text{Lag}}^{2}} \| (\boldsymbol{\lambda}, \boldsymbol{\mu}_{1}, \boldsymbol{\mu}_{2}) \|_{2}^{2} \right\}$$

$$\leq F(\boldsymbol{\lambda}_{k}^{\text{sub}}, \boldsymbol{\mu}_{1,k}^{\text{sub}}, \boldsymbol{\mu}_{2,k}^{\text{sub}}; \boldsymbol{x}_{k}) + \frac{\varepsilon}{6M_{\text{Lag}}^{2}} \| (\boldsymbol{\lambda}_{k}^{\text{sub}}, \boldsymbol{\mu}_{1,k}^{\text{sub}}, \boldsymbol{\mu}_{2,k}^{\text{sub}}) \|_{2}^{2}$$

$$\leq \| \boldsymbol{B}_{k} \boldsymbol{p}_{k} \|_{2}^{2} + \| \boldsymbol{\mu}_{1,k}^{\text{sub}} \odot \boldsymbol{p}_{k} \|_{2}^{2} + \| \boldsymbol{\mu}_{2,k}^{\text{sub}} \odot \boldsymbol{p}_{k} \|_{2}^{2} + \frac{\varepsilon}{2}$$

$$\leq (\kappa_{2} + 2M_{\text{Lag}}) \| \boldsymbol{p}_{k} \|_{2}^{2} + \frac{\varepsilon}{2}.$$

Here, $F(\boldsymbol{\lambda}_{m_i}^*, \boldsymbol{\mu}_{1,m_i}^*, \boldsymbol{\mu}_{2,m_i}^*; \boldsymbol{x}_{m_i}) > 2\varepsilon$ automatically implies that $\|\boldsymbol{p}_{m_i}\|_2 \ge \sqrt{\frac{3\varepsilon}{2(\kappa_2 + M_{\text{Lag}})}}$. Since $\liminf_{k \to \infty} \|\boldsymbol{p}_k\|_2 = 0$, there must exists $n_i > m_i$ such that $\|\boldsymbol{p}_{n_i}\|_2 \le \sqrt{\frac{\varepsilon}{\kappa_2 + M_{\text{Lag}}}}$. Let

$$\widetilde{F}(\boldsymbol{\lambda}, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2; \boldsymbol{x}) = F(\boldsymbol{\lambda}, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2; \boldsymbol{x}_k) + \frac{\varepsilon}{6M_{\text{Lag}}^2} \|(\boldsymbol{\lambda}, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2)\|_2^2$$

It also implies that $\widetilde{F}(\boldsymbol{\lambda}_{m_i}^{**}, \boldsymbol{\mu}_{1,m_i}^{**}, \boldsymbol{\mu}_{2,m_i}^{**}; \boldsymbol{x}_{m_i}) \geq 2\varepsilon$ and $\widetilde{F}(\boldsymbol{\lambda}_{n_i}^{**}, \boldsymbol{\mu}_{1,n_i}^{**}, \boldsymbol{\mu}_{2,n_i}^{**}; \boldsymbol{x}_{n_i}) \leq (\kappa_2 + 2M_{\text{Lag}}) \|\boldsymbol{p}_{n_i}\|_2^2 + \frac{\varepsilon}{2} \leq \frac{3}{2}\varepsilon$, where $(\boldsymbol{\lambda}_{m_i}^{**}, \boldsymbol{\mu}_{1,m_i}^{**}, \boldsymbol{\mu}_{2,m_i}^{**}) \in \min_{\boldsymbol{\lambda}, \boldsymbol{\mu}_1 \geq 0, \boldsymbol{\mu}_2 \geq 0} \widetilde{F}(\boldsymbol{\lambda}, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2; \boldsymbol{x}_{m_i})$ and $(\boldsymbol{\lambda}_{n_i}^{**}, \boldsymbol{\mu}_{1,n_i}^{**}, \boldsymbol{\mu}_{2,n_i}^{**}) \in \min_{\boldsymbol{\lambda}, \boldsymbol{\mu}_1 \geq 0, \boldsymbol{\mu}_2 \geq 0} \widetilde{F}(\boldsymbol{\lambda}, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2; \boldsymbol{x}_{n_i})$. Note that the function $\widetilde{F}(\boldsymbol{\lambda}, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2; \boldsymbol{x})$ is strictly positive-definite with

$$\frac{\varepsilon}{6M_{\text{Lag}}^2} \le \left\|\nabla^2 \widetilde{F}(\boldsymbol{\lambda}, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2; \boldsymbol{x})\right\|_2 \le 2\left(M_{\nabla c}^2 + 4M_{\nabla c}^2 + 2M_{\boldsymbol{\ell}, \boldsymbol{u}}^2 + 4\right).$$

For simplicity, we denote $oldsymbol{w}_k=(oldsymbol{\lambda}_k^{**},oldsymbol{\mu}_{2,k}^{**},oldsymbol{\mu}_{2,k}^{**}),$ then

$$\frac{\varepsilon}{6M_{\text{Lag}}^2} \|\boldsymbol{w}_k\|_2^2 \leq \min_{\boldsymbol{\lambda}, \boldsymbol{\mu}_1 \geq \boldsymbol{0}, \boldsymbol{\mu}_2 \geq \boldsymbol{0}} \widetilde{F}(\boldsymbol{\lambda}, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2; \boldsymbol{x}) \leq (\kappa_2 + 2M_{\text{Lag}}) \|\boldsymbol{p}_k\|_2^2 + \frac{\varepsilon}{2},$$

and thus

(B.15)
$$\|\boldsymbol{w}_{k}\|_{2} \leq \sqrt{\frac{6M_{\text{Lag}}^{2}\left(\kappa_{2}+2M_{\text{Lag}}\right)M_{\boldsymbol{\ell},\boldsymbol{u}}^{2}}{\varepsilon}} + 3M_{\text{Lag}},$$

for all $k \in \mathbb{N}$.

We first write $\widetilde{F}(\boldsymbol{\lambda}, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2; \boldsymbol{x}_k)$ into the general quadratic form that

$$\widetilde{F}(\boldsymbol{\lambda}, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2; \boldsymbol{x}_k) = \|\nabla f(\boldsymbol{x}_k)\|_2^2 + \boldsymbol{q}_k^\top \boldsymbol{w} + \frac{1}{2} \boldsymbol{w}^\top \boldsymbol{Q}_k \boldsymbol{w},$$

where

$$\boldsymbol{q}_{k} = \begin{pmatrix} 2\nabla \boldsymbol{c}(\boldsymbol{x}_{k})^{\top} \nabla f(\boldsymbol{x}_{k}) \\ -2\nabla f(\boldsymbol{x}_{k}) \\ 2\nabla f(\boldsymbol{x}_{k}) \end{pmatrix}$$

and

$$\boldsymbol{Q}_{k} = \begin{pmatrix} 2\nabla \boldsymbol{c}(\boldsymbol{x}_{k})^{\top} \nabla \boldsymbol{c}(\boldsymbol{x}_{k}) & -2\nabla \boldsymbol{c}(\boldsymbol{x}_{k}) \\ -2\nabla \boldsymbol{c}(\boldsymbol{x}_{k})^{\top} & 2\boldsymbol{I} + 2\text{diag}\left((\boldsymbol{x}_{k} - \boldsymbol{\ell})^{2}\right) & -2\boldsymbol{I} \\ 2\nabla \boldsymbol{c}(\boldsymbol{x}_{k})^{\top} & -2\boldsymbol{I} & 2\boldsymbol{I} + 2\text{diag}\left((\boldsymbol{x}_{k} - \boldsymbol{u})^{2}\right) \end{pmatrix} + \frac{\varepsilon}{6M_{\text{Lag}}^{2}}\boldsymbol{I}.$$

The smoothness of the objective f(x) and the constraints c(x) show that

(B.16)
$$\begin{aligned} \|\boldsymbol{q}_{k+1} - \boldsymbol{q}_k\|_2 &\leq 2\left(\kappa_{\nabla c}M_{\nabla f} + M_{\nabla c}\kappa_{\nabla f} + 2\kappa_{\nabla f}\right)\|\boldsymbol{x}_{k+1} - \boldsymbol{x}_k\|_2 \\ &\leq 2\left(\kappa_{\nabla c}M_{\nabla f} + M_{\nabla c}\kappa_{\nabla f} + 2\kappa_{\nabla f}\right)M_{\boldsymbol{\ell},\boldsymbol{u}}\alpha_k, \end{aligned}$$

and

(B.17)
$$\begin{aligned} \|\boldsymbol{Q}_{k+1} - \boldsymbol{Q}_k\|_2 &\leq 4 \left(M_{\nabla c} \kappa_{\nabla c} + 2\kappa_{\nabla c} + 2M_{\boldsymbol{\ell}, \boldsymbol{u}} \right) \|\boldsymbol{x}_{k+1} - \boldsymbol{x}_k\|_2 \\ &\leq 4 \left(M_{\nabla c} \kappa_{\nabla c} + 2\kappa_{\nabla c} + 2M_{\boldsymbol{\ell}, \boldsymbol{u}} \right) M_{\boldsymbol{\ell}, \boldsymbol{u}} \alpha_k. \end{aligned}$$

Then

$$\begin{split} & \left| \widetilde{F}(\boldsymbol{\lambda}_{k+1}^{**}, \boldsymbol{\mu}_{1,k+1}^{**}, \boldsymbol{\mu}_{2,k+1}^{**}; \boldsymbol{x}_{k+1}) - \widetilde{F}(\boldsymbol{\lambda}_{k}^{**}, \boldsymbol{\mu}_{1,k}^{**}, \boldsymbol{\mu}_{2,k}^{**}; \boldsymbol{x}_{k}) \right| \\ & \leq \left| \boldsymbol{q}_{k+1}^{\top} \boldsymbol{w}_{k+1} + \frac{1}{2} \boldsymbol{w}_{k+1}^{\top} \boldsymbol{Q}_{k+1} \boldsymbol{w}_{k+1} - \boldsymbol{q}_{k}^{\top} \boldsymbol{w}_{k} - \frac{1}{2} \boldsymbol{w}_{k}^{\top} \boldsymbol{Q}_{k} \boldsymbol{w}_{k} \right| + \left| \|\nabla f(\boldsymbol{x}_{k+1})\|_{2}^{2} - \|\nabla f(\boldsymbol{x}_{k})\|_{2}^{2} \right| \\ & \leq \left| \boldsymbol{q}_{k+1}^{\top} \boldsymbol{w}_{k+1} + \frac{1}{2} \boldsymbol{w}_{k+1}^{\top} \boldsymbol{Q}_{k+1} \boldsymbol{w}_{k+1} - \boldsymbol{q}_{k}^{\top} \boldsymbol{w}_{k+1} - \frac{1}{2} \boldsymbol{w}_{k+1}^{\top} \boldsymbol{Q}_{k} \boldsymbol{w}_{k+1} \right| \\ & - \left| \boldsymbol{q}_{k}^{\top} \boldsymbol{w}_{k+1} + \frac{1}{2} \boldsymbol{w}_{k+1}^{\top} \boldsymbol{Q}_{k} \boldsymbol{w}_{k+1} - \boldsymbol{q}_{k}^{\top} \boldsymbol{w}_{k} - \frac{1}{2} \boldsymbol{w}_{k}^{\top} \boldsymbol{Q}_{k} \boldsymbol{w}_{k} \right| + \left| \|\nabla f(\boldsymbol{x}_{k+1})\|_{2}^{2} - \|\nabla f(\boldsymbol{x}_{k})\|_{2}^{2} \right| \\ & \leq \left\| \boldsymbol{w}_{k+1} \right\|_{2} \left\| \boldsymbol{q}_{k+1} - \boldsymbol{q}_{k} \right\|_{2} + \frac{1}{2} \left\| \boldsymbol{w}_{k+1} \right\|_{2}^{2} \left\| \boldsymbol{Q}_{k+1} - \boldsymbol{Q}_{k} \right\|_{2} \\ & + \left\| \boldsymbol{q}_{k} \right\|_{2} \left\| \boldsymbol{w}_{k+1} - \boldsymbol{w}_{k} \right\|_{2} + \frac{1}{2} \left\| \boldsymbol{w}_{k+1} \right\|_{2} \left\| \boldsymbol{Q}_{k} \right\|_{2} \left\| \boldsymbol{w}_{k+1} - \boldsymbol{w}_{k} \right\|_{2} \\ & + \frac{1}{2} \left\| \boldsymbol{w}_{k} \right\|_{2} \left\| \boldsymbol{Q}_{k} \right\|_{2} \left\| \boldsymbol{w}_{k+1} - \boldsymbol{w}_{k} \right\|_{2} + \left\| \nabla f(\boldsymbol{x}_{k+1}) - \nabla f(\boldsymbol{x}_{k}) \right\|_{2} (\left\| \nabla f(\boldsymbol{x}_{k+1}) \right\|_{2} + \left\| \nabla f(\boldsymbol{x}_{k}) \right\|_{2}) \right| \\ & \text{Using Lemma 23, Equations (B.15), (B.16) and (B.17), and \boldsymbol{Q}_{k} \succeq \frac{\varepsilon}{\varepsilon} \mathbf{M}^{2}. \end{aligned}$$

Using Lemma 23, Equations (B.15), (B.16) and (B.17), and $Q_k \simeq \frac{1}{6M_{Lag}^2} \mathbf{1}$, we have $\|\boldsymbol{w}_{k+1} - \boldsymbol{w}_k\|_2 = \mathcal{O}\left(\frac{\alpha_k}{\varepsilon^{3/2}}\right)$, where we omit some universal and uncritical constants. Combining it with Equations (B.15), (B.16) and (B.17), we have

$$\left|\widetilde{F}(\boldsymbol{\lambda}_{k+1}^{**}, \boldsymbol{\mu}_{1,k+1}^{**}, \boldsymbol{\mu}_{2,k+1}^{**}; \boldsymbol{x}_{k+1}) - \widetilde{F}(\boldsymbol{\lambda}_{k}^{**}, \boldsymbol{\mu}_{1,k}^{**}, \boldsymbol{\mu}_{2,k}^{**}; \boldsymbol{x}_{k})\right| \leq M_{F} \frac{\alpha_{k}}{\varepsilon^{2}},$$

for a universal constant $M_F > 0$, where the constant is independent of α_k , k and ε .

Therefore, it follows from the above inequalities and our construction of the sequences $\{m_i\}$ and $\{n_i\}$ that

(B.18)
$$\frac{1}{2}\varepsilon \leq \widetilde{F}(\boldsymbol{\lambda}_{m_{i}}^{*}, \boldsymbol{\mu}_{1,m_{i}}^{*}, \boldsymbol{\mu}_{2,m_{i}}^{*}; \boldsymbol{x}_{m_{i}}) - \widetilde{F}(\boldsymbol{\lambda}_{n_{i}}^{*}, \boldsymbol{\mu}_{1,n_{i}}^{*}, \boldsymbol{\mu}_{2,n_{i}}^{*}; \boldsymbol{x}_{n_{i}}) \\
\leq \sum_{k=m_{i}}^{n_{i}-1} \left| \widetilde{F}(\boldsymbol{\lambda}_{k}^{*}, \boldsymbol{\mu}_{1,k}^{*}, \boldsymbol{\mu}_{2,k}^{*}; \boldsymbol{x}_{k}) - \widetilde{F}(\boldsymbol{\lambda}_{k+1}^{*}, \boldsymbol{\mu}_{1,k}^{*}, \boldsymbol{\mu}_{2,k}^{*}; \boldsymbol{x}_{k}) \right| \\
\leq \sum_{k=m_{i}}^{n_{i}-1} M_{F} \frac{\alpha_{k}}{\varepsilon^{2}}.$$

Summing up both two side from i = 1 to ∞ , we have

$$\infty = \sum_{i=1}^{\infty} \frac{1}{2M_F} \varepsilon^3 \le \sum_{i=1}^{\infty} \sum_{k=m_i}^{n_i-1} \alpha_k.$$

However, $\|\boldsymbol{p}_k\|_2 \ge \sqrt{\frac{\varepsilon}{\kappa_2 + M_{\text{Lag}}}}$ for $m_i \le k \le n_i - 1$, which further implies that

$$\sum_{i=1}^{\infty}\sum_{k=m_i}^{n_i-1}\alpha_k \leq \frac{\kappa_2 + M_{\text{Lag}}}{\varepsilon}\sum_{i=1}^{\infty}\sum_{k=m_i}^{n_i-1}\alpha_k \|\boldsymbol{p}_k\|_2^2 \leq \frac{\kappa_2 + M_{\text{Lag}}}{\varepsilon}\sum_{k=\bar{K}}^{\infty}\alpha_k \|\boldsymbol{p}_k\|_2^2 < \infty.$$

It is a contradiction. Therefore, we complete the proof that $\lim_{k\to\infty} F(\lambda_k^*, \mu_{1,k}^*, \mu_{2,k}^*; x_k) = 0.$

LEMMA 25. Under assumptions in Theorem 3, we have

$$\lim_{k\to\infty} \boldsymbol{c}(\boldsymbol{x}_k) = \boldsymbol{0}, \text{ almost surely.}$$

PROOF. The proof scheme is similar. For completeness, we provide the details here. Suppose that $\limsup_{k\to\infty} \|\boldsymbol{c}(\boldsymbol{x}_k)\|_2 = 0$ but $\liminf_{k\to\infty} \|\boldsymbol{c}(\boldsymbol{x}_k)\|_2 = 0$. Then we can find a sufficiently small number $\varepsilon > 0$ and two infinite sequences $\{m_i\}$ and $\{n_i\}$ with $\bar{K} \leq m_i < n_i$, such that

$$\|\boldsymbol{c}(\boldsymbol{x}_{m_i})\|_2 > 2\varepsilon, \quad \|\boldsymbol{c}(\boldsymbol{x}_{n_i})\|_2 < \varepsilon,$$

and

$$\|\boldsymbol{c}(\boldsymbol{x}_k)\|_2 \ge \varepsilon$$
, for $m_i \le k < n_i$.

It follows from the definition of the sequence that

$$egin{aligned} arepsilon &\leq \|oldsymbol{c}(oldsymbol{x}_{m_i})\|_2 - \|oldsymbol{c}(oldsymbol{x}_{n_i})\|_2 \ &\leq \sum_{k=m_i}^{n_i-1} \|oldsymbol{c}(oldsymbol{x}_k)\|_2 - \|oldsymbol{c}(oldsymbol{x}_{k+1})\|_2 \ &\leq \sum_{k=m_i}^{n_i-1} \|oldsymbol{c}(oldsymbol{x}_k) - oldsymbol{c}(oldsymbol{x}_{k+1})\|_2 \ &\leq \kappa_c M_{oldsymbol{\ell},oldsymbol{u}} \sum_{k=m_i}^{n_i-1} lpha_k, \quad ext{for all } i \in \mathbb{N}. \end{aligned}$$

Multiplying both two sides by ε and by the fact that $\|\boldsymbol{c}(\boldsymbol{x}_k)\|_2 \ge \varepsilon$, for $m_i \le k < n_i$, we have

$$\varepsilon^2 \leq \kappa_c M_{\ell, \boldsymbol{u}} \sum_{k=m_i}^{n_i-1} \alpha_k \| \boldsymbol{c}(\boldsymbol{x}_k) \|_2, \quad \text{for all } i \in \mathbb{N},$$

which implies that $\infty \leq \sum_{i=1}^{\infty} \sum_{k=m_i}^{n_i-1} \alpha_k \| \boldsymbol{c}(\boldsymbol{x}_k) \|_2 \leq \sum_{k=\bar{K}}^{\infty} \alpha_k \| \boldsymbol{c}(\boldsymbol{x}_k) \|_2 < \infty$. It is a contradiction.

Combining with Lemmas 24 and 25, we finish the proof for Theorem 3.

APPENDIX C: PROOF FOR THEOREM 4

C.1. Proof for Lemma 3. We proceed by prove each of the four conclusions in turn.

C.1.1. *Proof for Conclusion 1*. Note that

$$oldsymbol{p}_k \in rgmin_{oldsymbol{p} \in \Omega_k}
abla f(oldsymbol{x}_k)^{ op} oldsymbol{p} + rac{1}{2}oldsymbol{p}^{ op} oldsymbol{B}_k oldsymbol{p},$$

where
$$\Omega_k = \{ \boldsymbol{p} : \boldsymbol{c}(\boldsymbol{x}_k) + \nabla \boldsymbol{c}(\boldsymbol{x}_k)^\top \boldsymbol{p} = \boldsymbol{0} \} \cap \{ \boldsymbol{p} : \boldsymbol{\ell} \leq \boldsymbol{x}_k + \boldsymbol{p} \leq \boldsymbol{u} \}$$
, and let
 $\boldsymbol{p}^* \in \operatorname*{arg\,min}_{\boldsymbol{p} \in \Omega^*} \nabla f(\boldsymbol{x}^*)^\top \boldsymbol{p} + \frac{1}{2} \boldsymbol{p}^\top \boldsymbol{B}_k \boldsymbol{p},$

and

$$oldsymbol{p}_k^* \in rgmin_{oldsymbol{p} \in \Omega^*}
abla f(oldsymbol{x}_k)^{ op} oldsymbol{p} + rac{1}{2} oldsymbol{p}^{ op} oldsymbol{B}_k oldsymbol{p},$$

where $\Omega^* = \{ \boldsymbol{p} : \boldsymbol{c}(\boldsymbol{x}^*) + \nabla \boldsymbol{c}(\boldsymbol{x}^*)^\top \boldsymbol{p} = \boldsymbol{0} \} \cap \{ \boldsymbol{p} : \boldsymbol{\ell} \leq \boldsymbol{x}^* + \boldsymbol{p} \leq \boldsymbol{u} \}$. Lemma 23 shows that $\| \boldsymbol{p}_k^* - \boldsymbol{p}^* \|_2 \leq C \| \boldsymbol{x}_k - \boldsymbol{x}^* \|_2$ for some C > 0, since both \boldsymbol{p}_k^* and \boldsymbol{p}^* are bounded by $M_{\boldsymbol{\ell},\boldsymbol{u}}$. In the next part, for simplicity, we use the same notation C to denote some universal constants. We slightly rewrite the formulation for \boldsymbol{p}_k and \boldsymbol{p}_k^* that

$$\hat{\boldsymbol{p}}_k \in \operatorname*{arg\,min}_{\boldsymbol{p}\in\widehat{\Omega}_k} \frac{1}{2} \left\| \boldsymbol{p} \right\|_{\boldsymbol{B}_k}^2,$$

and

$$\hat{\boldsymbol{p}}_k^* \in \operatorname*{arg\,min}_{\boldsymbol{p}\in\widehat{\Omega}^*} \frac{1}{2} \|\boldsymbol{p}\|_{\boldsymbol{B}_k}^2,$$

where $\widehat{\Omega}_k = \{ \boldsymbol{p} + \boldsymbol{B}_k^{-1} \nabla f(\boldsymbol{x}_k) : \boldsymbol{p} \in \Omega_k \}$ and $\widehat{\Omega}^* = \{ \boldsymbol{p} + \boldsymbol{B}_k^{-1} \nabla f(\boldsymbol{x}_k) : \boldsymbol{p} \in \Omega^* \}$. Then $\|\boldsymbol{p}_k - \boldsymbol{p}_k^*\|_2 = \|\hat{\boldsymbol{p}}_k - \hat{\boldsymbol{p}}_k^*\|_2$. By Proposition 3.1 in [21] and results in [37], there exists $\hat{\boldsymbol{p}}^{*\prime} \in \widehat{\Omega}^*$ and $\hat{\boldsymbol{p}}_k' \in \widehat{\Omega}_k$, such that $\|\hat{\boldsymbol{p}}^{*\prime} - \hat{\boldsymbol{p}}_k\|_{\boldsymbol{B}_k} \leq C \|\boldsymbol{x}_k - \boldsymbol{x}^*\|_2$ and $\|\hat{\boldsymbol{p}}_k' - \hat{\boldsymbol{p}}_k^*\|_{\boldsymbol{B}_k} \leq C \|\boldsymbol{x}_k - \boldsymbol{x}^*\|_2$ for some C > 0, then

$$\|\hat{\boldsymbol{p}}_{k}^{*}\|_{\boldsymbol{B}_{k}} \leq \|\hat{\boldsymbol{p}}^{*\prime}\|_{\boldsymbol{B}_{k}} \leq \|\hat{\boldsymbol{p}}_{k}\|_{\boldsymbol{B}_{k}} + C \|\boldsymbol{x}_{k} - \boldsymbol{x}^{*}\|_{2},$$

and

$$\|\hat{p}_{k}\|_{B_{k}} \leq \|\hat{p}_{k}'\|_{B_{k}} \leq \|\hat{p}_{k}^{*}\|_{B_{k}} + C \|x_{k} - x^{*}\|_{2}$$

Equipped with the above inequalities and the optimality condition that $\langle \hat{p}_k^*, \hat{p}^{*\prime} - \hat{p}_k^* \rangle \ge 0$, we have

$$\begin{split} \left\| \hat{p}^{*\prime} \right\|_{B_{k}}^{2} &= \left\| \hat{p}_{k}^{*} + \hat{p}^{*\prime} - \hat{p}_{k}^{*} \right\|_{B_{k}}^{2} \\ &= \left\| \hat{p}_{k}^{*} \right\|_{B_{k}}^{2} + 2\left\langle \hat{p}_{k}^{*}, \hat{p}^{*\prime} - \hat{p}_{k}^{*} \right\rangle + \left\| \hat{p}^{*\prime} - \hat{p}_{k}^{*} \right\|_{B_{k}}^{2} \\ &\geq \left\| \hat{p}_{k}^{*} \right\|_{B_{k}}^{2} + \left\| \hat{p}^{*\prime} - \hat{p}_{k}^{*} \right\|_{B_{k}}^{2}. \end{split}$$

Therefore,

$$\begin{split} \left\| \hat{p}^{*'} - \hat{p}_{k}^{*} \right\|_{B_{k}}^{2} &\leq \left\| \hat{p}^{*'} \right\|_{B_{k}}^{2} - \left\| \hat{p}_{k}^{*} \right\|_{B_{k}}^{2} \\ &= \left\| \hat{p}^{*'} - \hat{p}_{k} + \hat{p}_{k} \right\|_{B_{k}}^{2} - \left\| \hat{p}_{k}^{*} \right\|_{B_{k}}^{2} \\ &\leq \left\| \hat{p}^{*'} - \hat{p}_{k} \right\|_{B_{k}}^{2} + 2 \left\| \hat{p}^{*'} - \hat{p}_{k} \right\|_{B_{k}} \left\| \hat{p}_{k} \right\|_{B_{k}} + \left\| \hat{p}_{k} \right\|_{B_{k}}^{2} - \left\| \hat{p}_{k}^{*} \right\|_{B_{k}}^{2} \\ &\leq \left\| \hat{p}^{*'} - \hat{p}_{k} \right\|_{B_{k}}^{2} + 2 \left\| \hat{p}^{*'} - \hat{p}_{k} \right\|_{B_{k}} \left\| \hat{p}_{k} \right\|_{B_{k}} + \left\| \left\| \hat{p}_{k} \right\|_{B_{k}} - \left\| \hat{p}_{k}^{*} \right\|_{B_{k}} \right\| \left(\left\| \hat{p}_{k} \right\|_{B_{k}} + \left\| \hat{p}_{k}^{*} \right\|_{B_{k}} \right) \\ &\leq C \left\| \boldsymbol{x}_{k} - \boldsymbol{x}^{*} \right\|_{2}, \end{split}$$

for some C > 0. Thus

(C.1)
$$\begin{aligned} \|\hat{\boldsymbol{p}}_{k} - \hat{\boldsymbol{p}}_{k}^{*}\|_{\boldsymbol{B}_{k}}^{2} &\leq \|\hat{\boldsymbol{p}}_{k} - \hat{\boldsymbol{p}}^{*\prime} + \hat{\boldsymbol{p}}^{*\prime} - \hat{\boldsymbol{p}}_{k}^{*}\|_{\boldsymbol{B}_{k}}^{2} \\ &\leq 2 \|\hat{\boldsymbol{p}}_{k} - \hat{\boldsymbol{p}}^{*\prime}\|_{\boldsymbol{B}_{k}}^{2} + 2 \|\hat{\boldsymbol{p}}^{*\prime} - \hat{\boldsymbol{p}}_{k}^{*}\|_{\boldsymbol{B}_{k}}^{2} \\ &\leq C \|\boldsymbol{x}_{k} - \boldsymbol{x}^{*}\|_{2}, \end{aligned}$$

for some C > 0. The facts that $\|\boldsymbol{p}_k - \boldsymbol{p}_k^*\|_2 = \|\hat{\boldsymbol{p}}_k - \hat{\boldsymbol{p}}_k^*\|_2 \leq C\sqrt{\|\boldsymbol{x}_k - \boldsymbol{x}^*\|_2}$ and $\|\boldsymbol{p}_k^* - \boldsymbol{p}^*\|_2 \leq C\|\boldsymbol{x}_k - \boldsymbol{x}^*\|_2$ for some C > 0. Note that $\boldsymbol{p}^* = \boldsymbol{0}$ for any positive-definite matrix \boldsymbol{B}_k since \boldsymbol{x}^* is a local solution of Problem (1.1). We complete the proof as $\boldsymbol{x}_k \to \boldsymbol{x}^*$ in Assumption 5.

C.1.2. Proof for Conclusion 2. We revisit the definition of
$$\bar{g}_k$$
 and have that
 $\bar{g}_k - \nabla f(x_k) = \beta_k (g_k - \nabla f(x_k)) + (1 - \beta_k) (\bar{g}_{k-1} - \nabla f(x_{k-1})) + (1 - \beta_k) (\nabla f(x_{k-1}) - \nabla f(x_k)) + (1 - \beta_k) (\nabla f(x_{k-1}) - \nabla f(x_{k-2})) + (1 - \beta_{k-1}) (\nabla f(x_{k-2}) - \nabla f(x_{k-1})) + (1 - \beta_k) (\nabla f(x_{k-1}) - \nabla f(x_k)) + (1 - \beta_{k-1}) (\nabla f(x_{k-2}) - \nabla f(x_{k-1})) + (1 - \beta_k) (\nabla f(x_{k-1}) - \nabla f(x_k)) = \cdots$

$$= \sum_{i=0}^k \left(\prod_{j=i+1}^k (1 - \beta_j) \right) \beta_i (g_i - \nabla f(x_i)) + \sum_{i=1}^k \left(\prod_{j=i}^k (1 - \beta_j) \right) (\nabla f(x_{i-1}) - \nabla f(x_i)) + \sum_{i=1}^k \left(\prod_{j=i}^k (1 - \beta_j) \right) (\nabla f(x_{i-1}) - \nabla f(x_i)) + \sum_{i=1}^k \left(\prod_{j=i}^k (1 - \beta_j) \right) (\nabla f(x_{i-1}) - \nabla f(x_i)) + \sum_{i=1}^k \left(\prod_{j=i}^k (1 - \beta_j) \right) (\nabla f(x_{i-1}) - \nabla f(x_i)) + \sum_{i=1}^k \left(\prod_{j=i}^k (1 - \beta_j) \right) (\nabla f(x_{i-1}) - \nabla f(x_i)) + \sum_{i=1}^k \left(\prod_{j=i}^k (1 - \beta_j) \right) (\nabla f(x_{i-1}) - \nabla f(x_i)) + \sum_{i=1}^k \left(\prod_{j=i}^k (1 - \beta_j) \right) (\nabla f(x_{i-1}) - \nabla f(x_i)) + \sum_{i=1}^k \left(\prod_{j=i}^k (1 - \beta_j) \right) (\nabla f(x_{i-1}) - \nabla f(x_i)) + \sum_{i=1}^k \left(\prod_{j=i}^k (1 - \beta_j) \right) (\nabla f(x_{i-1}) - \nabla f(x_i)) + \sum_{i=1}^k \left(\prod_{j=i}^k (1 - \beta_j) \right) (\nabla f(x_{i-1}) - \nabla f(x_i)) + \sum_{i=1}^k \left(\prod_{j=i}^k (1 - \beta_j) \right) (\nabla f(x_{i-1}) - \nabla f(x_i)) + \sum_{i=1}^k \left(\prod_{j=i}^k (1 - \beta_j) \right) (\nabla f(x_{i-1}) - \nabla f(x_i)) + \sum_{i=1}^k \left(\prod_{j=i}^k (1 - \beta_j) \right) (\nabla f(x_{i-1}) - \nabla f(x_i)) + \sum_{i=1}^k \left(\prod_{j=i}^k (1 - \beta_j) \right) (\nabla f(x_{i-1}) - \nabla f(x_i)) + \sum_{i=1}^k \left(\prod_{j=i}^k (1 - \beta_j) \right) (\nabla f(x_{i-1}) - \nabla f(x_i)) + \sum_{i=1}^k \left(\prod_{j=i}^k (1 - \beta_j) \right) (\nabla f(x_{i-1}) - \nabla f(x_i)) + \sum_{i=1}^k \left(\prod_{j=i}^k (1 - \beta_j) \right) (\nabla f(x_{i-1}) - \nabla f(x_i)) + \sum_{i=1}^k \left(\prod_{j=i}^k (1 - \beta_j) \right) (\nabla f(x_{i-1}) - \nabla f(x_i)) + \sum_{i=1}^k \left(\prod_{j=i}^k (1 - \beta_j) \right) (\nabla f(x_{i-1}) - \nabla f(x_i)) + \sum_{i=1}^k \left(\prod_{j=i}^k (1 - \beta_j) \right) (\nabla f(x_i) - \nabla f(x_i)) + \sum_{i=i}^k \left(\prod_{j=i}^k (1 - \beta_j) \right) (\nabla f(x_i) - \nabla f(x_i)) + \sum_{i=i}^k \left(\prod_{j=i}^k (1 - \beta_j) \right) (\nabla f(x_i) - \nabla f(x_i)) + \sum_{i=i}^k \left(\prod_{j=i}^k (1 - \beta_j) \right) (\nabla f(x_i) - \nabla f(x_i)) + \sum_{i=i}^k \left(\prod_{j=i}^k (1 - \beta_j) \right) (\nabla f(x_i) - \nabla f(x_i)) + \sum_{i=i}^k \left(\prod_{j=i}^k (1 - \beta_j) \right) (\nabla f(x_i)$$

Here,

$$\left\|\mathcal{W}_{2,k}\right\|_{2} \leq \sum_{i=1}^{k} \left(\prod_{j=i}^{k} (1-\beta_{j})\right) \alpha_{i-1} M_{\ell,\boldsymbol{u}},$$

then $\mathcal{W}_{2,k} \to 0$ as $k \to \infty$ since $\lim_{k\to\infty} \alpha_{i-1}/\beta_i = 0$. We apply Theorem 4.4 in [34] with $\gamma = 2, \phi(k) = 1, \alpha = -1/2, \text{ and } X_{k,h} = \left(\prod_{h'=h+1}^k (1-\beta_{h'})\right) \beta_h/(\sqrt{\beta_k}k^{1/2+\epsilon}) (\boldsymbol{g}_h - \nabla f(\boldsymbol{x}_h))$ for any sufficiently small $\varepsilon > 0$, as well as the Borel-Cantelli Lemma that the martingale difference array satisfies $\|\mathcal{W}_{1,k}\|_2 \to 0$, almost surely.

C.1.3. Proof for Condition 3. We will show that there exists a sufficiently small $\varepsilon^* > 0$ such that if $\|\bar{g}_k - \nabla f(x_k)\|_2 \le \varepsilon^*$, $\mathcal{I}(x_k + \bar{p}_k) = \mathcal{I}(x_k + p_k)$ and $\mathcal{J}(x_k + \bar{p}_k) = \mathcal{J}(x_k + p_k)$ hold. The almost sure convergence of $\bar{g}_k - \nabla f(x_k)$ implies that $\|\bar{g}_k - \nabla f(x_k)\|_2 \le \varepsilon^*$ holds for some sufficiently large $k \ge K^*$. Let $(\lambda_k^{\text{sub}}, \mu_{1,k}^{\text{sub}}, \mu_{2,k}^{\text{sub}})$ and $(\bar{\lambda}_k^{\text{sub}}, \bar{\mu}_{1,k}^{\text{sub}}, \bar{\mu}_{2,k}^{\text{sub}})$ be the Lagrangian multipliers of the relaxed SQP subproblem with the full gradient $\nabla f(x_k)$ and the stochastic averaged gradient \bar{g}_k , respectively. Let $(\lambda^*, \mu_1^*, \mu_2^*)$ be the Lagrangian multiplier for Problem (1.1) at x^* and denote $\epsilon = \min\{\{(\mu_1^*)_i : i \in \mathcal{I}(x^*)\} \cup \{(\mu_2^*)_i : i \in \mathcal{J}(x^*)\}\} > 0$ due to the strictly complementary slackness condition. Since p_k is the optimal solution of the strongly convex quadratic SQP subproblem, the KKT condition shows that $\nabla f(x_k) + B_k p_k + \nabla c(x_k) \lambda_k^{\text{sub}} - \mu_{1,k}^{\text{sub}} + \mu_{2,k}^{\text{sub}} = 0$. Taking $k \to \infty$ $(x_k \to x^*$ and $p_k \to \nabla f(x_k) + B_k p_k + \nabla c(x_k) \lambda_k^{\text{sub}} - \mu_{1,k}^{\text{sub}} + \mu_{2,k}^{\text{sub}} = 0$.

NA ET AL.

0), it follows from the LICQ at \boldsymbol{x}^* that $(\boldsymbol{\lambda}_k^{\text{sub}}, \boldsymbol{\mu}_{1,k}^{\text{sub}}, \boldsymbol{\mu}_{2,k}^{\text{sub}}) \rightarrow (\boldsymbol{\lambda}^*, \boldsymbol{\mu}_1^*, \boldsymbol{\mu}_2^*)$. So there exists sufficiently large $K^* > 0$ such that $(\boldsymbol{\mu}_{1,k}^{\text{sub}})_i > \frac{3}{4}\epsilon$ for all $i \in \mathcal{I}(\boldsymbol{x}^*)$ and $(\boldsymbol{\mu}_{2,k}^{\text{sub}})_i > \frac{3}{4}\epsilon$ for all $i \in \mathcal{J}(\boldsymbol{x}^*)$. Therefore, $\boldsymbol{x}_k + \boldsymbol{p}_k$ has the same active and inactive set as \boldsymbol{x}^* , i.e., $\mathcal{I}(\boldsymbol{x}_k + \boldsymbol{p}_k) = \mathcal{I}(\boldsymbol{x}^*)$ and $\mathcal{J}(\boldsymbol{x}_k + \boldsymbol{p}_k) = \mathcal{J}(\boldsymbol{x}^*)$. Denote $\epsilon' = \max\{(\boldsymbol{x}^* - \boldsymbol{\ell})_i, (\boldsymbol{u} - \boldsymbol{x}^*)_i : i \notin \mathcal{I}(\boldsymbol{x}^*)$ and $i \notin \mathcal{J}(\boldsymbol{x}^*)\}$. When K^* is sufficiently large, we have $\max\{(\boldsymbol{x}_k + \boldsymbol{p}_k - \boldsymbol{\ell})_i, (\boldsymbol{u} - \boldsymbol{x}_k - \boldsymbol{p}_k)_i : i \notin \mathcal{I}(\boldsymbol{x}^*)\}$ and $i \notin \mathcal{J}(\boldsymbol{x}^*)\} \geq \frac{3}{4}\epsilon'$. Lemma 23 shows that $\|\bar{\boldsymbol{p}}_k - \boldsymbol{p}_k\|_2 \leq (1 + M_{\boldsymbol{\ell},\boldsymbol{u}})\kappa_1^{-1}\varepsilon^*$ under the assumption that $\|\bar{\boldsymbol{g}}_k - \nabla f(\boldsymbol{x}_k)\|_2 \leq \varepsilon^*$ when $k \geq K^*$. If ε^* is sufficiently small such that $(1 + M_{\boldsymbol{\ell},\boldsymbol{u}})\kappa_1^{-1}\varepsilon^* \leq \frac{1}{4}\epsilon'$, then $\max\{(\boldsymbol{x}_k + \bar{\boldsymbol{p}}_k - \boldsymbol{\ell})_i, (\boldsymbol{u} - \boldsymbol{x}_k - \bar{\boldsymbol{p}}_k)_i : i \notin \mathcal{I}(\boldsymbol{x}^*)\}\} \geq \frac{1}{2}\epsilon'$.

The LICQ condition implies that columns of $[\nabla c(\boldsymbol{x}^*), [-\boldsymbol{I}]_{\mathcal{I}(\boldsymbol{x}^*)}, [\boldsymbol{I}]_{\mathcal{J}(\boldsymbol{x}^*)}]$ are linearly independent and $[\nabla c(\boldsymbol{x}^*), [-\boldsymbol{I}]_{\mathcal{I}(\boldsymbol{x}^*)}, [\boldsymbol{I}]_{\mathcal{J}(\boldsymbol{x}^*)}]^{\top} [\nabla c(\boldsymbol{x}^*), [-\boldsymbol{I}]_{\mathcal{I}(\boldsymbol{x}^*)}, [\boldsymbol{I}]_{\mathcal{J}(\boldsymbol{x}^*)}] \succeq \kappa_0 \boldsymbol{I}$ for some $\kappa_0 > 0$. By the smoothness of $\boldsymbol{c}(\boldsymbol{x})$, there exists sufficiently large K^* such that $[\nabla c(\boldsymbol{x}_k), [-\boldsymbol{I}]_{\mathcal{I}(\boldsymbol{x}^*)}, [\boldsymbol{I}]_{\mathcal{J}(\boldsymbol{x}^*)}]^{\top} [\nabla c(\boldsymbol{x}_k), [-\boldsymbol{I}]_{\mathcal{I}(\boldsymbol{x}^*)}] \succeq \frac{1}{2}\kappa_0 \boldsymbol{I}$ for all $k \geq K^*$. The KKT condition of the SQP subproblem at \boldsymbol{x}_k with $\bar{\boldsymbol{g}}_k$ shows that $\bar{\boldsymbol{g}}_k + \boldsymbol{B}_k \bar{\boldsymbol{p}}_k + \nabla c(\boldsymbol{x}_k) \bar{\boldsymbol{\lambda}}_k^{\text{sub}} - \bar{\boldsymbol{\mu}}_{1,k}^{\text{sub}} + \bar{\boldsymbol{\mu}}_{2,k}^{\text{sub}} = \boldsymbol{0}$. Since $\max\{(\boldsymbol{x}_k + \bar{\boldsymbol{p}}_k - \ell)_i, (\boldsymbol{u} - \boldsymbol{x}_k - \bar{\boldsymbol{p}}_k)_i : i \notin \mathcal{I}(\boldsymbol{x}^*) \text{ and } i \notin \mathcal{J}(\boldsymbol{x}^*)\} \geq \frac{1}{2}\epsilon', (\bar{\boldsymbol{\mu}}_{1,k}^{\text{sub}})_i = 0 \text{ for } i \notin \mathcal{I}(\boldsymbol{x}^*) \text{ and } i \notin \mathcal{J}(\boldsymbol{x}^*)$. Therefore,

$$\begin{split} & \left\| \nabla \boldsymbol{c}(\boldsymbol{x}_{k}) \left(\bar{\boldsymbol{\lambda}}_{k}^{\text{sub}} - \boldsymbol{\lambda}_{k}^{\text{sub}} \right) - \left(\bar{\boldsymbol{\mu}}_{1,k}^{\text{sub}} - \boldsymbol{\mu}_{1,k}^{\text{sub}} \right) + \left(\bar{\boldsymbol{\mu}}_{2,k}^{\text{sub}} - \boldsymbol{\mu}_{2,k}^{\text{sub}} \right) \right\|_{2} \\ & \leq \left\| \bar{\boldsymbol{g}}_{k} - \nabla f(\boldsymbol{x}_{k}) \right\|_{2} + \left\| \boldsymbol{B}_{k} \bar{\boldsymbol{p}}_{k} - \boldsymbol{B}_{k} \boldsymbol{p}_{k} \right\|_{2} \\ & \leq \left(1 + (1 + M_{\boldsymbol{\ell},\boldsymbol{u}}) \kappa_{1}^{-1} \kappa_{2} \right) \varepsilon^{*}, \end{split}$$

and

$$\left\| \begin{pmatrix} \bar{\boldsymbol{\lambda}}_{k}^{\text{sub}} - \boldsymbol{\lambda}_{k}^{\text{sub}} \\ \begin{bmatrix} \bar{\boldsymbol{\mu}}_{1,k}^{\text{sub}} - \boldsymbol{\mu}_{1,k}^{\text{sub}} \end{bmatrix}_{\mathcal{I}(\boldsymbol{x}^{*})} \\ \begin{bmatrix} \bar{\boldsymbol{\mu}}_{2,k}^{\text{sub}} - \boldsymbol{\mu}_{2,k}^{\text{sub}} \end{bmatrix}_{\mathcal{J}(\boldsymbol{x}^{*})} \end{pmatrix} \right\|_{2} \leq 2\kappa_{0}^{-1} \left(M_{\nabla c} + 2 \right) \left(1 + (1 + M_{\boldsymbol{\ell},\boldsymbol{u}})\kappa_{1}^{-1}\kappa_{2} \right) \varepsilon^{*}$$

We let ε^* to be small enough such that the right-hand side of the above inequality is less than $\frac{1}{4}\epsilon$, i.e., $2\kappa_0^{-1}(M_{\nabla c}+2)\left(1+(1+M_{\ell,u})\kappa_1^{-1}\kappa_2\right)\varepsilon^* \leq \frac{1}{4}\epsilon$. Then, together with $(\boldsymbol{\mu}_{1,k}^{\mathrm{sub}})_i > \frac{3}{4}\epsilon$ for $i \in \mathcal{I}(\boldsymbol{x}^*)$ and $(\boldsymbol{\mu}_{2,k}^{\mathrm{sub}})_i > \frac{3}{4}\epsilon$ for $i \in \mathcal{J}(\boldsymbol{x}^*)$, we have $(\bar{\boldsymbol{\mu}}_{1,k}^{\mathrm{sub}})_i > \frac{1}{2}\epsilon$ for $i \in \mathcal{I}(\boldsymbol{x}^*)$ and $(\bar{\boldsymbol{\mu}}_{2,k}^{\mathrm{sub}})_i > \frac{1}{2}\epsilon$ for $i \in \mathcal{J}(\boldsymbol{x}^*)$. It implies that both $\boldsymbol{x}_k + \boldsymbol{p}_k$ and $\boldsymbol{x}_k + \bar{\boldsymbol{p}}_k$ can correctly identify the active and inactive sets of constraints at \boldsymbol{x}^* . Therefore, $\mathcal{I}(\boldsymbol{x}_k + \bar{\boldsymbol{p}}_k) = \mathcal{I}(\boldsymbol{x}_k + \boldsymbol{p}_k) = \mathcal{I}(\boldsymbol{x}^*)$ and $\mathcal{J}(\boldsymbol{x}_k + \bar{\boldsymbol{p}}_k) = \mathcal{J}(\boldsymbol{x}_k + \boldsymbol{p}_k) = \mathcal{J}(\boldsymbol{x}^*)$.

C.1.4. Proof for Conclusion 4. Equipped with the fact that $\bar{g}_k - \nabla f(x_k) \to 0$ almost surely, the condition $\|\bar{g}_k - \nabla f(x_k)\|_2 \leq \varepsilon^*$ always holds when k is sufficiently large. By the proof in the previous section, we know that $(\bar{\lambda}_k^{\text{sub}}, [\bar{\mu}_{1,k}^{\text{sub}}]_{\mathcal{I}(x^*)}, [\bar{\mu}_{2,k}^{\text{sub}}]_{\mathcal{J}(x^*)}) \to (\lambda^*, [\mu_1^*]_{\mathcal{I}(x^*)}, [\mu_2^*]_{\mathcal{J}(x^*)})$. The update scheme for dual variables in Step 6 shows the following recursion

$$\boldsymbol{\lambda}_{k+1} = \prod_{j=K^*}^k \left(1 - \alpha_j\right) \boldsymbol{\lambda}_{K^*} + \sum_{i=K^*}^k \prod_{j=i+1}^k \left(1 - \alpha_j\right) \alpha_i \bar{\boldsymbol{\lambda}}_i^{\text{sub}},$$

then $\lambda_k \to \lambda^*$ almost surely. Similar convergence results hold for $([\mu_{1,k}]_{\mathcal{I}(\boldsymbol{x}^*)}, [\mu_{2,k}]_{\mathcal{J}(\boldsymbol{x}^*)}) \to ([\mu_1^*]_{\mathcal{I}(\boldsymbol{x}^*)}, [\mu_2^*]_{\mathcal{J}(\boldsymbol{x}^*)})$ and for dual variables indexed on inactive sets $\mathcal{I}^-(\boldsymbol{x}^*)$ and $\mathcal{J}^-(\boldsymbol{x}^*)$.

C.2. Proof for Lemma 4. The definition of the averaged Hessian matrix B_k shows that (C.2)

$$\begin{split} \|\boldsymbol{B}_{k}-\boldsymbol{B}^{*}\|_{2} &\leq \left\|\frac{1}{k}\sum_{i=1}^{k}\nabla^{2}f(\boldsymbol{x}_{i};\zeta_{i})-\nabla^{2}f(\boldsymbol{x}_{i})\right\|_{2} \\ &+\frac{1}{k}\sum_{i=1}^{k}\left\|\nabla^{2}f(\boldsymbol{x}_{i})+\sum_{j=1}^{r}\left(\boldsymbol{\lambda}_{i}\right)_{j}\nabla^{2}c_{j}(\boldsymbol{x}_{i})-\nabla^{2}f(\boldsymbol{x}^{*})-\sum_{j=1}^{r}\left(\boldsymbol{\lambda}^{*}\right)_{j}\nabla^{2}c_{j}(\boldsymbol{x}^{*})\right\|_{2} \\ &+\|\boldsymbol{\Delta}_{k}\|_{2} \\ &\leq \left\|\frac{1}{k}\sum_{i=1}^{k}\nabla^{2}f(\boldsymbol{x}_{i};\zeta_{i})-\nabla^{2}f(\boldsymbol{x}_{i})\right\|_{2}+\frac{\Upsilon_{\nabla^{2}\mathcal{L}}}{k}\sum_{i=1}^{k}\left\|\left(\frac{\boldsymbol{x}_{i}-\boldsymbol{x}^{*}}{\boldsymbol{\lambda}_{i}-\boldsymbol{\lambda}^{*}}\right)\right\|_{2}+\|\boldsymbol{\Delta}_{k}\|_{2}, \end{split}$$

for some $\Upsilon_{\nabla^2 \mathcal{L}} > 0$ due to the compactness of iterates and smoothness of $\nabla^2 f(x)$ and $\nabla^2 c(x)$. The first term converges to 0 almost surely by the strong law of large number, while the second term converges to 0 almost surely by the Stolz-Cesaro theorem. Since Δ_k acts as a regularization term for the positive definiteness of B_k and B^* is positive definite, we deduce that $\Delta_k = 0$ when k is sufficiently large. Moreover,

$$\|\boldsymbol{H}_{k} - \boldsymbol{H}^{*}\|_{2} \leq \|\boldsymbol{B}_{k} - \boldsymbol{B}^{*}\|_{2} + \kappa_{\nabla c} \|\boldsymbol{x}_{k} - \boldsymbol{x}^{*}\|_{2}$$

implies that $H_k \rightarrow H^*$ almost surely.

C.3. Proof for Theorem 4.

LEMMA 26. When H_k is sufficiently close to H^* , there exists a constant $\Upsilon_L > 0$, such that

(C.3)
$$\left\| \boldsymbol{H}_{k}^{-1} - (\boldsymbol{H}^{*})^{-1} \right\|_{2} \leq \Upsilon_{L} \left\| \boldsymbol{H}_{k} - \boldsymbol{H}^{*} \right\|_{2}.$$

Then

$$\left\|\boldsymbol{H}_{k}^{-1}\right\|_{2}, \left\|\left(\boldsymbol{H}^{*}\right)^{-1}\right\|_{2} \leq \Upsilon_{H},$$

for some $\Upsilon_H > 0$. Moreover, $H_k^{-1} \to (H^*)^{-1}$ almost surely.

PROOF. First, we build the relationship between $H_k^{-1} - (H^*)^{-1}$ and $H_k - H^*$. Note that

(C.4)

$$0 = (H^*)^{-1} H^* - H_k^{-1} H_k$$

$$= (H^*)^{-1} H^* - (H^*)^{-1} H_k + (H^*)^{-1} H_k - H_k^{-1} H_k$$

$$= (H^*)^{-1} (H^* - H_k) + ((H^*)^{-1} - H_k^{-1}) H_k,$$

then

(C.5)
$$\left\| \boldsymbol{H}_{k}^{-1} - (\boldsymbol{H}^{*})^{-1} \right\|_{2} \leq \frac{\left\| (\boldsymbol{H}^{*})^{-1} \right\|_{2} \left\| \boldsymbol{H}_{k} - \boldsymbol{H}^{*} \right\|_{2}}{\lambda_{\min} (\boldsymbol{H}_{k})} \leq \Upsilon_{L} \left\| \boldsymbol{H}_{k} - \boldsymbol{H}^{*} \right\|_{2},$$

for some $\Upsilon_L > 0$, since we can assume that $\lambda_{\min}(H_k) > \frac{1}{2}\lambda_{\min}(H^*)$ without the loss of generality, when H_k is sufficiently close to H^* . The boundedness of $||H_k^{-1}||_2$ is a direct result of Equation (C.3). The almost sure convergence of $H_k^{-1} \to (H^*)^{-1}$ is straightforward from Equation (C.3) and Lemma 4.

LEMMA 27. Algorithm 3 generates a sequence $\{(x_k, \lambda_k, \mu_{1,k}, \mu_{2,k})\}$ satisfying

$$egin{pmatrix} oldsymbol{x}_{k+1}-oldsymbol{x}^*\ oldsymbol{\lambda}_{k+1}-oldsymbol{\lambda}^*\ [oldsymbol{\mu}_{1,k+1}-oldsymbol{\mu}_1^*]_{\mathcal{I}(oldsymbol{x}^*)}\ [oldsymbol{\mu}_{2,k+1}-oldsymbol{\mu}_2^*]_{\mathcal{J}(oldsymbol{x}^*)} \end{pmatrix}=\mathcal{Q}_{1,k}+\mathcal{Q}_{2,k}+\mathcal{Q}_{3,k},$$

and

$$\begin{pmatrix} [\boldsymbol{\mu}_{1,k+1} - \boldsymbol{\mu}_1^*]_{\mathcal{I}^-(\boldsymbol{x}^*)} \\ [\boldsymbol{\mu}_{2,k+1} - \boldsymbol{\mu}_2^*]_{\mathcal{J}^-(\boldsymbol{x}^*)} \end{pmatrix} = \prod_{i=K^*}^k (1 - \alpha_i^{\min}) \begin{pmatrix} [\boldsymbol{\mu}_{1,K^*} - \boldsymbol{\mu}_1^*]_{\mathcal{I}^-(\boldsymbol{x}^*)} \\ [\boldsymbol{\mu}_{2,K^*} - \boldsymbol{\mu}_2^*]_{\mathcal{J}^-(\boldsymbol{x}^*)} \end{pmatrix},$$

where

$$\begin{aligned} \mathcal{Q}_{1,k} &= \sum_{i=K^*}^k \prod_{j=i+1}^k (1 - \alpha_j^{\min}) \alpha_i^{\min} \phi_i, \\ \mathcal{Q}_{2,k} &= \sum_{i=K^*}^k \left(\prod_{j=i+1}^k (1 - \alpha_j^{\min}) \right) (\alpha_i - \alpha_i^{\min}) \begin{pmatrix} \bar{p}_k \\ \Delta \lambda_k \\ [\Delta \mu_{1,k}]_{\mathcal{I}(\boldsymbol{x}^*)} \\ [\Delta \mu_{2,k}]_{\mathcal{J}(\boldsymbol{x}^*)} \end{pmatrix}, \\ \mathcal{Q}_{3,k} &= \prod_{i=K^*}^k (1 - \alpha_i^{\min}) \begin{pmatrix} \boldsymbol{x}_{K^*} - \boldsymbol{x}^* \\ \boldsymbol{\lambda}_{K^*} - \boldsymbol{\lambda}^* \\ [\boldsymbol{\mu}_{1,K^*} - \boldsymbol{\mu}_1^*]_{\mathcal{I}(\boldsymbol{x}^*)} \\ [\boldsymbol{\mu}_{2,K^*} - \boldsymbol{\mu}_2^*]_{\mathcal{J}(\boldsymbol{x}^*)} \end{pmatrix} + \sum_{i=K^*}^k \prod_{j=i+1}^k (1 - \alpha_j^{\min}) \alpha_i^{\min} \boldsymbol{\delta}_i, \end{aligned}$$

and

$$\begin{split} \phi_i &= -\boldsymbol{H}_i^{-1} \begin{pmatrix} \bar{\boldsymbol{g}}_i - \nabla f(\boldsymbol{x}_i) \\ \boldsymbol{0} \\ \boldsymbol{0} \end{pmatrix}, \\ \boldsymbol{\delta}_i &= -(\boldsymbol{H}^*)^{-1} \, \psi_i - \left(\boldsymbol{H}_i^{-1} - (\boldsymbol{H}^*)^{-1}\right) \begin{pmatrix} \nabla f(\boldsymbol{x}_i) + \nabla \boldsymbol{c}(\boldsymbol{x}_i) \boldsymbol{\lambda}_i - \boldsymbol{\mu}_{1,i} + \boldsymbol{\mu}_{2,i} \\ \boldsymbol{c}(\boldsymbol{x}_i) \\ [\boldsymbol{\ell} - \boldsymbol{x}_i]_{\mathcal{I}(\boldsymbol{x}^*)} \\ [\boldsymbol{x}_i - \boldsymbol{u}]_{\mathcal{J}(\boldsymbol{x}^*)} \end{pmatrix}, \\ \psi_i &= \begin{pmatrix} \nabla f(\boldsymbol{x}_i) + \nabla \boldsymbol{c}(\boldsymbol{x}_i) \boldsymbol{\lambda}_i - \boldsymbol{\mu}_{1,i} + \boldsymbol{\mu}_{2,i} \\ \boldsymbol{c}(\boldsymbol{x}_i) \\ [\boldsymbol{\ell} - \boldsymbol{x}_i]_{\mathcal{I}(\boldsymbol{x}^*)} \\ [\boldsymbol{x}_i - \boldsymbol{u}]_{\mathcal{J}(\boldsymbol{x}^*)} \end{pmatrix} - \boldsymbol{H}^* \begin{pmatrix} \boldsymbol{x}_i - \boldsymbol{x}^* \\ \boldsymbol{\lambda}_i - \boldsymbol{\lambda}^* \\ [\boldsymbol{\mu}_{1,i} - \boldsymbol{\mu}_1^*]_{\mathcal{I}(\boldsymbol{x}^*)} \\ [\boldsymbol{\mu}_{2,i} - \boldsymbol{\mu}_2^*]_{\mathcal{J}(\boldsymbol{x}^*)} \end{pmatrix}. \end{split}$$

PROOF. By the update scheme of Algorithm 3, we have

$$\begin{pmatrix} \boldsymbol{x}_{k+1} - \boldsymbol{x}^{*} \\ \boldsymbol{\lambda}_{k+1} - \boldsymbol{\lambda}^{*} \\ [\boldsymbol{\mu}_{1,k+1} - \boldsymbol{\mu}_{1}^{*}]_{\mathcal{I}(\boldsymbol{x}^{*})} \\ [\boldsymbol{\mu}_{2,k+1} - \boldsymbol{\mu}_{2}^{*}]_{\mathcal{J}(\boldsymbol{x}^{*})} \end{pmatrix} = \begin{pmatrix} \boldsymbol{x}_{k} - \boldsymbol{x}^{*} \\ \boldsymbol{\lambda}_{k} - \boldsymbol{\lambda}^{*} \\ [\boldsymbol{\mu}_{1,k} - \boldsymbol{\mu}_{2}^{*}]_{\mathcal{J}(\boldsymbol{x}^{*})} \\ [\boldsymbol{\mu}_{2,k} - \boldsymbol{\mu}_{2}^{*}]_{\mathcal{J}(\boldsymbol{x}^{*})} \end{pmatrix} + (\alpha_{k} - \alpha_{k}^{\min} + \alpha_{k}^{\min}) \begin{pmatrix} \bar{\boldsymbol{p}}_{k} \\ [\boldsymbol{\Delta}\boldsymbol{\mu}_{1,k}]_{\mathcal{I}(\boldsymbol{x}^{*})} \\ [\boldsymbol{\Delta}\boldsymbol{\mu}_{2,k}]_{\mathcal{J}(\boldsymbol{x}^{*})} \end{pmatrix}$$
$$= \begin{pmatrix} \boldsymbol{x}_{k} - \boldsymbol{x}^{*} \\ \boldsymbol{\lambda}_{k} - \boldsymbol{\lambda}^{*} \\ [\boldsymbol{\mu}_{1,k} - \boldsymbol{\mu}_{1}^{*}]_{\mathcal{I}(\boldsymbol{x}^{*})} \\ [\boldsymbol{\mu}_{2,k} - \boldsymbol{\mu}_{2}^{*}]_{\mathcal{J}(\boldsymbol{x}^{*})} \end{pmatrix} + (\alpha_{k} - \alpha_{k}^{\min}) \begin{pmatrix} \bar{\boldsymbol{p}}_{k} \\ [\boldsymbol{\Delta}\boldsymbol{\lambda}_{k} \\ [\boldsymbol{\Delta}\boldsymbol{\mu}_{1,k}]_{\mathcal{I}(\boldsymbol{x}^{*})} \\ [\boldsymbol{\Delta}\boldsymbol{\mu}_{2,k}]_{\mathcal{J}(\boldsymbol{x}^{*})} \end{pmatrix} - \alpha_{k}^{\min} \boldsymbol{H}_{k}^{-1} \begin{pmatrix} \nabla f(\boldsymbol{x}_{k}) + \nabla \boldsymbol{c}(\boldsymbol{x}_{k})\boldsymbol{\lambda}_{k} - \boldsymbol{\mu}_{1,k} + \boldsymbol{\mu}_{2,k} \\ \boldsymbol{c}(\boldsymbol{x}_{k}) \\ [\boldsymbol{\ell} - \boldsymbol{x}_{k}]_{\mathcal{I}(\boldsymbol{x}^{*})} \\ [\boldsymbol{\ell} - \boldsymbol{x}_{k}]_{\mathcal{I}(\boldsymbol{x}^{*})} \end{pmatrix}$$

 $+ \alpha_k^{\min} \phi_k$

$$= \begin{pmatrix} \boldsymbol{x}_{k} - \boldsymbol{x}^{*} \\ \boldsymbol{\lambda}_{k} - \boldsymbol{\lambda}^{*} \\ [\boldsymbol{\mu}_{1,k} - \boldsymbol{\mu}_{1}^{*}]_{\mathcal{I}(\boldsymbol{x}^{*})} \\ [\boldsymbol{\mu}_{2,k} - \boldsymbol{\mu}_{2}^{*}]_{\mathcal{J}(\boldsymbol{x}^{*})} \end{pmatrix} + (\alpha_{k} - \alpha_{k}^{\min}) \begin{pmatrix} \bar{\boldsymbol{p}}_{k} \\ [\Delta \boldsymbol{\mu}_{1,k}]_{\mathcal{I}(\boldsymbol{x}^{*})} \\ [\Delta \boldsymbol{\mu}_{2,k}]_{\mathcal{J}(\boldsymbol{x}^{*})} \end{pmatrix} + \alpha_{k}^{\min} \boldsymbol{\phi}_{k} \\ - \alpha_{k}^{\min} \left(\boldsymbol{H}_{k}^{-1} - (\boldsymbol{H}^{*})^{-1} \right) \begin{pmatrix} \nabla f(\boldsymbol{x}_{k}) + \nabla \boldsymbol{c}(\boldsymbol{x}_{k})\boldsymbol{\lambda}_{k} - \boldsymbol{\mu}_{1,k} + \boldsymbol{\mu}_{2,k} \\ \boldsymbol{c}(\boldsymbol{x}_{k}) \\ [\boldsymbol{\ell} - \boldsymbol{x}_{k}]_{\mathcal{I}(\boldsymbol{x}^{*})} \\ [\boldsymbol{x}_{k} - \boldsymbol{u}]_{\mathcal{J}(\boldsymbol{x}^{*})} \end{pmatrix} \\ - \alpha_{k}^{\min} \left(\boldsymbol{H}^{*} \right)^{-1} \begin{pmatrix} \nabla f(\boldsymbol{x}_{k}) + \nabla \boldsymbol{c}(\boldsymbol{x}_{k})\boldsymbol{\lambda}_{k} - \boldsymbol{\mu}_{1,k} + \boldsymbol{\mu}_{2,k} \\ \boldsymbol{c}(\boldsymbol{x}_{k}) \\ [\boldsymbol{\ell} - \boldsymbol{x}_{k}]_{\mathcal{I}(\boldsymbol{x}^{*})} \\ [\boldsymbol{k}_{k} - \boldsymbol{u}]_{\mathcal{J}(\boldsymbol{x}^{*})} \end{pmatrix} \\ = (1 - \alpha_{k}^{\min}) \begin{pmatrix} \boldsymbol{x}_{k} - \boldsymbol{x}^{*} \\ \boldsymbol{\lambda}_{k} - \boldsymbol{\lambda}^{*} \\ [\boldsymbol{\mu}_{1,k} - \boldsymbol{\mu}_{1}^{*}]_{\mathcal{I}(\boldsymbol{x}^{*})} \end{pmatrix} + (\alpha_{k} - \alpha_{k}^{\min}) \begin{pmatrix} \boldsymbol{\bar{p}}_{k} \\ [\Delta \boldsymbol{\lambda}_{k} \\ [\Delta \boldsymbol{\lambda}_{k}] \\ [\Delta \boldsymbol{\mu}_{1,k}]_{\mathcal{I}(\boldsymbol{x}^{*})} \end{pmatrix} + \alpha_{k}^{\min} \boldsymbol{\phi}_{k} \end{cases}$$

$$\left(\begin{bmatrix} \boldsymbol{\mu}_{2,k} - \boldsymbol{\mu}_{2}^{*} \end{bmatrix}_{\mathcal{J}(\boldsymbol{x}^{*})} \right) \left(\begin{bmatrix} \Delta \boldsymbol{\mu}_{2,k} \end{bmatrix}_{\mathcal{J}(\boldsymbol{x}^{*})} \right)$$

$$- \alpha_{k}^{\min} \left(\boldsymbol{H}_{k}^{-1} - (\boldsymbol{H}^{*})^{-1} \right) \begin{pmatrix} \nabla f(\boldsymbol{x}_{k}) + \nabla \boldsymbol{c}(\boldsymbol{x}_{k}) \boldsymbol{\lambda}_{k} - \boldsymbol{\mu}_{1,k} + \boldsymbol{\mu}_{2,k} \\ \boldsymbol{c}(\boldsymbol{x}_{k}) \\ [\boldsymbol{\ell} - \boldsymbol{x}_{k}]_{\mathcal{I}(\boldsymbol{x}^{*})} \\ [\boldsymbol{k}_{k} - \boldsymbol{u}]_{\mathcal{J}(\boldsymbol{x}^{*})} \end{pmatrix}$$

$$- \alpha_{k}^{\min} \left(\boldsymbol{H}^{*} \right)^{-1} \boldsymbol{\psi}_{k}.$$

We then obtain the result by applying the above equation recursively.

$$\left\|\mathcal{Q}_{2,k}\right\|_{2} = \mathcal{O}\left(\alpha_{k}^{\min}\right).$$

PROOF. Under the boundedness of the generated iterates and the almost sure convergence that $\|\bar{g}_k - \nabla f(x_k)\|_2 \le \varepsilon^*$, we have that the iterates in Equation (4.3) are bounded for all $k \ge K^*$, i.e.,

$$\left\| \begin{pmatrix} \bar{\boldsymbol{p}}_k \\ \Delta \boldsymbol{\lambda}_k \\ [\Delta \boldsymbol{\mu}_{1,k}]_{\mathcal{I}(\boldsymbol{x}^*)} \\ [\Delta \boldsymbol{\mu}_{2,k}]_{\mathcal{J}(\boldsymbol{x}^*)} \end{pmatrix} \right\|_2 \le M_\Delta,$$

for some $M_{\Delta}>0,$ due to the LICQ condition. Recall that

$$\mathcal{Q}_{2,k} = \sum_{i=K^*}^k \left(\prod_{j=i+1}^k (1-\alpha_j^{\min}) \right) (\alpha_i - \alpha_i^{\min}) \begin{pmatrix} \bar{p}_k \\ \Delta \lambda_k \\ [\Delta \mu_{1,k}]_{\mathcal{I}(\boldsymbol{x}^*)} \\ [\Delta \mu_{2,k}]_{\mathcal{J}(\boldsymbol{x}^*)} \end{pmatrix},$$

which shows that

 $b_1 = 1$, using Lemma 15.

$$\left\|\mathcal{Q}_{2,k}\right\|_{2} = \mathcal{O}\left(\alpha_{k}^{\min}\right),$$

since $|\alpha_i - \alpha_i^{\min}| \le (\iota_0/\iota_1^2)(\alpha_k^{\min})^2$ and Lemma 15.

LEMMA 29. Under Assumptions 5 and 6, and suppose that $\iota_1 > b_2$ if $b_1 = 1$, then (C, 6) $\mathbb{E}\left[\|\mathcal{O}_{1,t}\|^2\right] = \mathcal{O}\left(\beta_{t_1} + (\alpha_t^{\min})^2/\beta_t^2\right).$

(C.6)
$$\mathbb{E}\left[\left\|\mathcal{Q}_{1,k}\right\|_{2}^{2}\right] = \mathcal{O}\left(\beta_{k} + (\alpha_{k}^{\min})^{2}/\beta_{k}^{2}\right)$$

PROOF. By the definition of $\bar{g}_i - \nabla f(x_i)$, we have

$$\begin{split} \begin{pmatrix} \bar{g}_i - \nabla f(\boldsymbol{x}_i) \\ \boldsymbol{0} \\ \boldsymbol{0} \\ \boldsymbol{0} \end{pmatrix} &= \sum_{h=K^*}^{i} \begin{pmatrix} \prod_{h'=h+1}^{i} (1-\beta_{h'}) \end{pmatrix} \beta_h \begin{pmatrix} \boldsymbol{g}_h - \nabla f(\boldsymbol{x}_h) \\ \boldsymbol{0} \\ \boldsymbol{0} \end{pmatrix} \\ &+ \sum_{h=K^*}^{i} \begin{pmatrix} \prod_{h'=h}^{i} (1-\beta_{h'}) \end{pmatrix} \begin{pmatrix} \nabla f(\boldsymbol{x}_{h-1}) - \nabla f(\boldsymbol{x}_h) \\ \boldsymbol{0} \\ \boldsymbol{0} \end{pmatrix} \\ &:= \mathcal{W}_{1,i} + \mathcal{W}_{2,i}. \end{split}$$

$$\begin{split} &(\mathbf{C}.7)\\ &\mathbb{E}\left[\|\mathcal{Q}_{1,k}\|_{2}^{2}\right] \\ \leq & \Upsilon_{H}^{2}\mathbb{E}\left[\left(\sum_{i=K^{*}}^{k}\prod_{j=i+1}^{k}\left(1-\alpha_{j}^{\min}\right)\alpha_{i}^{\min}\|\mathcal{W}_{1,i}+\mathcal{W}_{2,i}\|_{2}\right)^{2}\right] \\ \leq & \Upsilon_{H}^{2}\sum_{i=K^{*}}^{k}\prod_{j=i+1}^{k}\left(1-\alpha_{j}^{\min}\right)\alpha_{i}^{\min}\sum_{i'=K^{*}}^{k}\prod_{j'=i'+1}^{k}\left(1-\alpha_{j'}^{\min}\right)\alpha_{i'}^{\min}\mathbb{E}\left[\|\mathcal{W}_{1,i}+\mathcal{W}_{2,i}\|_{2}\|\mathcal{W}_{1,i'}+\mathcal{W}_{2,i'}\|_{2}\right] \\ \leq & \Upsilon_{H}^{2}\sum_{i=K^{*}}^{k}\prod_{j=i+1}^{k}\left(1-\alpha_{j}^{\min}\right)\alpha_{i}^{\min}\sum_{i'=K^{*}}^{k}\prod_{j'=i'+1}^{k}\left(1-\alpha_{j'}^{\min}\right)\alpha_{i'}^{\min}\sqrt{\mathbb{E}\left[\|\mathcal{W}_{1,i}+\mathcal{W}_{2,i}\|_{2}^{2}\right]}\sqrt{\mathbb{E}\left[\|\mathcal{W}_{1,i'}+\mathcal{W}_{2,i'}\|_{2}^{2}\right]} \\ \leq & \Upsilon_{H}^{2}\left(\sum_{i=K^{*}}^{k}\prod_{j=i+1}^{k}\left(1-\alpha_{j}^{\min}\right)\alpha_{i}^{\min}\sqrt{\mathbb{E}\left[\|\mathcal{W}_{1,i}+\mathcal{W}_{2,i}\|_{2}^{2}\right]}\right)^{2}. \end{split}$$
Note that $\mathbb{E}\left[\|\mathcal{W}_{1,i}\|_{2}^{2}\right] \leq M_{\sigma}\sum_{h=K^{*}}^{i}\prod_{h'=h+1}^{i}\left(1-\beta_{h'}\right)^{2}\beta_{h}^{2} \leq 2M_{\sigma}\beta_{i} \text{ and } \|\mathcal{W}_{2,i}\|_{2}^{2} \leq M_{\ell_{e}u}\left(\sum_{h=K^{*}}^{i}\prod_{h'=h+1}^{i}\left(1-\beta_{h'}\right)\alpha_{h}\right)^{2} \leq 2M_{\ell_{e}u}^{2}\alpha_{i}^{2}\beta_{i}^{2} = \mathcal{O}\left(\beta_{i}+\left(\alpha_{i}^{\min}\right)^{2}/\beta_{i}^{2}\right), \text{ for } i \text{ suf-ficiently large, then }\mathbb{E}\left[\|\mathcal{Q}_{1,k}\|_{2}^{2}\right] = \mathcal{O}\left(\beta_{k}+\left(\alpha_{k}^{\min}\right)^{2}/\beta_{k}^{2}\right). \end{aligned}$

LEMMA 30. Under Assumptions 5 and 6, and suppose that $\iota_1 > b_2$ if $b_1 = 1$, then

$$\mathbb{E}\left[\left\|\mathcal{Q}_{3,k}\right\|_{2}^{2}\right] = \mathcal{O}\left(\beta_{k} + (\alpha_{k}^{\min})^{2}/\beta_{k}^{2}\right),$$

and

$$\mathbb{E}\left[\left\|\begin{pmatrix}\boldsymbol{x}_{k}-\boldsymbol{x}^{*}\\\boldsymbol{\lambda}_{k}-\boldsymbol{\lambda}^{*}\\[\boldsymbol{\mu}_{1,k}-\boldsymbol{\mu}_{1}^{*}]_{\mathcal{I}(\boldsymbol{x}^{*})}\\[\boldsymbol{\mu}_{2,k}-\boldsymbol{\mu}_{2}^{*}]_{\mathcal{J}(\boldsymbol{x}^{*})}\end{pmatrix}\right\|_{2}^{2}\right]=\mathcal{O}\left(\beta_{k}+(\alpha_{k}^{\min})^{2}/\beta_{k}^{2}\right).$$

PROOF. Recall the definition of $\mathcal{Q}_{3,k}$ that

$$\mathcal{Q}_{3,k} = \prod_{i=K^*}^k (1-\alpha_i^{\min}) \begin{pmatrix} \boldsymbol{x}_{K^*} - \boldsymbol{x}^* \\ \boldsymbol{\lambda}_{K^*} - \boldsymbol{\lambda}^* \\ [\boldsymbol{\mu}_{1,K^*} - \boldsymbol{\mu}_1^*]_{\mathcal{I}(\boldsymbol{x}^*)} \\ [\boldsymbol{\mu}_{2,K^*} - \boldsymbol{\mu}_2^*]_{\mathcal{J}(\boldsymbol{x}^*)} \end{pmatrix} + \sum_{i=K^*}^k \prod_{j=i+1}^k (1-\alpha_j^{\min}) \alpha_i^{\min} \boldsymbol{\delta}_i,$$

we have the following recursion

(C.8)
$$\mathcal{Q}_{3,k+1} = \left(1 - \alpha_{k+1}^{\min}\right)\mathcal{Q}_{3,k} + \alpha_{k+1}^{\min}\boldsymbol{\delta}_{k+1}.$$

Here,

$$\begin{split} \|\boldsymbol{\delta}_{k+1}\|_{2} &\leq \left\| (\boldsymbol{H}^{*})^{-1} \right\|_{2} \|\boldsymbol{\psi}_{k+1}\|_{2} + \left\| \boldsymbol{H}_{k+1}^{-1} - (\boldsymbol{H}^{*})^{-1} \right\|_{2} \|\nabla \mathcal{L}_{k+1}\|_{2} \\ &\leq \kappa_{\nabla \mathcal{L}} \Upsilon_{H} \left\| \begin{pmatrix} \boldsymbol{x}_{k+1} - \boldsymbol{x}^{*} \\ \boldsymbol{\lambda}_{k+1} - \boldsymbol{\lambda}^{*} \\ [\boldsymbol{\mu}_{1,k+1} - \boldsymbol{\mu}_{1}^{*}]_{\mathcal{I}(\boldsymbol{x}^{*})} \\ [\boldsymbol{\mu}_{2,k+1} - \boldsymbol{\mu}_{2}^{*}]_{\mathcal{J}(\boldsymbol{x}^{*})} \end{pmatrix} \right\|_{2}^{2} + \kappa_{\nabla \mathcal{L}} \Upsilon_{L} \|\boldsymbol{H}_{k+1} - \boldsymbol{H}^{*}\|_{2} \left\| \begin{pmatrix} \boldsymbol{x}_{k+1} - \boldsymbol{x}^{*} \\ \boldsymbol{\lambda}_{k+1} - \boldsymbol{\lambda}^{*} \\ [\boldsymbol{\mu}_{1,k+1} - \boldsymbol{\mu}_{1}^{*}]_{\mathcal{I}(\boldsymbol{x}^{*})} \\ [\boldsymbol{\mu}_{2,k+1} - \boldsymbol{\mu}_{2}^{*}]_{\mathcal{J}(\boldsymbol{x}^{*})} \end{pmatrix} \right\|_{2} \\ &:= \varepsilon_{k+1} \left\| \begin{pmatrix} \boldsymbol{x}_{k+1} - \boldsymbol{x}^{*} \\ \boldsymbol{\lambda}_{k+1} - \boldsymbol{\lambda}^{*} \\ [\boldsymbol{\mu}_{1,k+1} - \boldsymbol{\mu}_{1}^{*}]_{\mathcal{I}(\boldsymbol{x}^{*})} \\ [\boldsymbol{\mu}_{2,k+1} - \boldsymbol{\mu}_{2}^{*}]_{\mathcal{J}(\boldsymbol{x}^{*})} \end{pmatrix} \right\|_{2} \\ &\leq \varepsilon_{k+1} \left(\| \mathcal{Q}_{1,k} \|_{2} + \| \mathcal{Q}_{2,k} \|_{2} + \| \mathcal{Q}_{3,k} \|_{2} \right), \end{split}$$

where we define

$$\varepsilon_k := \kappa_{\nabla \mathcal{H}} \Upsilon_L \left\| \begin{pmatrix} \boldsymbol{x}_k - \boldsymbol{x}^* \\ \boldsymbol{\lambda}_k - \boldsymbol{\lambda}^* \\ [\boldsymbol{\mu}_{1,k} - \boldsymbol{\mu}_1^*]_{\mathcal{I}(\boldsymbol{x}^*)} \\ [\boldsymbol{\mu}_{2,k} - \boldsymbol{\mu}_2^*]_{\mathcal{J}(\boldsymbol{x}^*)} \end{pmatrix} \right\|_2 + \kappa_{\nabla \mathcal{L}} \Upsilon_L \| \boldsymbol{H}_k - \boldsymbol{H}^* \|_2.$$

Then, for any $a \in (0,1)$ and there exists the corresponding threshold $K_a \ge K^*$ such that $\varepsilon_{k+1} \le a$ and

$$\|\mathcal{Q}_{3,k+1}\|_{2} \leq \left(1 - (1 - a)\alpha_{k+1}^{\min}\right) \|\mathcal{Q}_{3,k}\|_{2} + a\alpha_{k+1}^{\min} \cdot \left(\|\mathcal{Q}_{1,k}\|_{2} + \|\mathcal{Q}_{2,k}\|_{2}\right),$$

for all $k \ge K_a$, as $\varepsilon_k \to 0$ almost surely. We then develop the recursion for $\|\mathcal{Q}_{3,k+1}\|_2$ that

$$\begin{aligned} &(\mathbf{C}.9) \\ &\|\mathcal{Q}_{3,k+1}\|_{2} \\ &\leq \left(1 - (1 - a)\alpha_{k+1}^{\min}\right) \|\mathcal{Q}_{3,k}\|_{2} + a\alpha_{k+1}^{\min} \cdot \left(\|\mathcal{Q}_{1,k}\|_{2} + \|\mathcal{Q}_{2,k}\|_{2}\right) \\ &\leq \left(1 - (1 - a)\alpha_{k+1}^{\min}\right) \left(1 - (1 - a)\alpha_{k}^{\min}\right) \|\mathcal{Q}_{3,k-1}\|_{2} \\ &+ \left(1 - (1 - a)\alpha_{k+1}^{\min}\right) a\alpha_{k}^{\min} \cdot \left(\|\mathcal{Q}_{1,k-1}\|_{2} + \|\mathcal{Q}_{2,k-1}\|_{2}\right) \\ &+ a\alpha_{k+1}^{\min} \cdot \left(\|\mathcal{Q}_{1,k}\|_{2} + \|\mathcal{Q}_{2,k}\|_{2}\right) \\ &\leq \cdots \\ &\leq \prod_{j=K_{a}+1}^{k+1} \left(1 - (1 - a)\alpha_{j}^{\min}\right) \|\mathcal{Q}_{3,K_{a}}\|_{2} + \sum_{i=K_{a}+1}^{k+1} \left(\prod_{j=i+1}^{k+1} \left(1 - (1 - a)\alpha_{j}^{\min}\right)\right) a\alpha_{i}^{\min} \cdot \left(\|\mathcal{Q}_{1,i-1}\|_{2} + \|\mathcal{Q}_{2,i-1}\|_{2}\right), \end{aligned}$$

and thus

$$\begin{split} & \mathbb{E}\left[\|\mathcal{Q}_{3,k}\|_{2}^{2}\right] \\ &\leq 2\left(\prod_{j=K_{a}+1}^{k}\left(1-(1-a)\alpha_{j}^{\min}\right)\|\mathcal{Q}_{3,K_{a}}\|_{2}\right)^{2}+2\sum_{i=K_{a}+1}^{k}\left(\prod_{j=i+1}^{k}\left(1-(1-a)\alpha_{j}^{\min}\right)\right)a\alpha_{i}^{\min} \\ &\cdot\sum_{i'=K_{a}+1}^{k}\left(\prod_{j'=i'+1}^{k}\left(1-(1-a)\alpha_{j'}^{\min}\right)\right)a\alpha_{i'}^{\min}\cdot\mathbb{E}\left[\left(\|\mathcal{Q}_{1,i-1}\|_{2}+\|\mathcal{Q}_{2,i-1}\|_{2}\right)\left(\|\mathcal{Q}_{1,i'-1}\|_{2}+\|\mathcal{Q}_{2,i'-1}\|_{2}\right)\right] \\ &\leq 2\left(\prod_{j=K_{a}+1}^{k}\left(1-(1-a)\alpha_{j}^{\min}\right)\|\mathcal{Q}_{3,K_{a}}\|_{2}\right)^{2}+2\sum_{i=K_{a}+1}^{k}\left(\prod_{j=i+1}^{k}\left(1-(1-a)\alpha_{j}^{\min}\right)\right)a\alpha_{i}^{\min} \\ &\cdot\sum_{i'=K_{a}+1}^{k}\left(\prod_{j'=i'+1}^{k}\left(1-(1-a)\alpha_{j'}^{\min}\right)\right)a\alpha_{i'}^{\min}\cdot\sqrt{\mathbb{E}\left[\left(\|\mathcal{Q}_{1,i-1}\|_{2}+\|\mathcal{Q}_{2,i-1}\|_{2}\right)^{2}\right]}\sqrt{\mathbb{E}\left[\left(\|\mathcal{Q}_{1,i'-1}\|_{2}+\|\mathcal{Q}_{2,i'-1}\|_{2}\right)^{2}\right]} \\ &\leq 2\left(\prod_{j=K_{a}+1}^{k}\left(1-(1-a)\alpha_{j}^{\min}\right)\|\mathcal{Q}_{3,K_{a}}\|_{2}\right)^{2} \\ &+2\left(\sum_{i=K_{a}+1}^{k}\left(1-(1-a)\alpha_{j}^{\min}\right)\|\mathcal{Q}_{3,K_{a}}\|_{2}\right)^{2} \\ &+2\left(\sum_{i=K_{a}+1}^{k}\left(1-(1-a)\alpha_{j}^{\min}\right)\right)a\alpha_{i'}^{\min}\sqrt{\mathbb{E}\left[\left(\|\mathcal{Q}_{1,i-1}\|_{2}+\|\mathcal{Q}_{2,i-1}\|_{2}\right)^{2}\right]}\right)^{2}. \\ \\ &\text{Here, the fact that }\mathbb{E}\left[\left(\|\mathcal{Q}_{1,i}\|_{2}+\|\mathcal{Q}_{2,i}\|_{2}\right)^{2}\right] \leq 2\mathbb{E}\left[\|\mathcal{Q}_{1,i}\|_{2}^{2}+\|\mathcal{Q}_{2,i}\|_{2}^{2}\right] = \mathcal{O}\left(\beta_{i}+(\alpha_{i}^{\min})^{2}/\beta_{i}^{2}\right) \\ &\text{implies }\mathbb{E}\left[\|\mathcal{Q}_{3,k}\|_{2}^{2}\right] = \mathcal{O}\left(\beta_{k}+(\alpha_{k}^{\min})^{2}/\beta_{k}^{2}\right). \\ &\text{Here, we require that }\iota_{1}>\frac{b_{2}}{2\left(1-a\right)} \\ &\mathbb{E}\left[\|\mathcal{Q}_{2,k}\|_{2}^{2}\right] \text{ and }\mathbb{E}\left[\|\mathcal{Q}_{3,k}\|_{2}^{2}\right] \text{ are at least of the order }\mathcal{O}\left(\beta_{k}+(\alpha_{k}^{\min})^{2}/\beta_{k}^{2}\right). \\ &\text{Here, we require that }\iota_{1}>\frac{b_{2}}{2\left(1-a\right)} \\ &\text{ if } b_{1}=1. \\ &\text{Sinc } a \in (0,1) \\ &\text{ can be arbitrarily close to 0, we \\ &\text{ know } \frac{b_{2}}{2\left(1-a\right)} \leq b_{2} < 1 \\ &\text{ if } a \leq \frac{1}{2}, \\ &\text{ and } u_{1}>b_{2}. \end{bmatrix}$$

LEMMA 31. Under Assumptions 5 and 6, and suppose that $\iota_1 > b_2$ if $b_1 = 1$, then we have

(C.10)
$$\mathbb{E}\left[\|\boldsymbol{H}_{k}-\boldsymbol{H}^{*}\|_{2}^{2}\right] = \mathcal{O}\left(\beta_{k}+(\alpha_{k}^{\min})^{2}/\beta_{k}^{2}\right)$$

and

(C.11)
$$\mathbb{E}\left[\left\|\boldsymbol{H}_{k}^{-1}-(\boldsymbol{H}^{*})^{-1}\right\|_{2}^{2}\right] = \mathcal{O}\left(\beta_{k}+(\alpha_{k}^{\min})^{2}/\beta_{k}^{2}\right).$$

PROOF. We revisit the result in Lemma 26 that $\left\| \boldsymbol{H}_{k}^{-1} - (\boldsymbol{H}^{*})^{-1} \right\|_{2} \leq \Upsilon_{L} \| \boldsymbol{H}_{k} - \boldsymbol{H}^{*} \|_{2}$. Then, in the left part of the proof, we mainly show the first equality.

$$\begin{aligned} \|\boldsymbol{H}_{k}-\boldsymbol{H}^{*}\|_{2} &\leq \left\|\frac{1}{k+1}\sum_{i=0}^{k}\left(\nabla^{2}f(\boldsymbol{x}_{i};\zeta_{i})-\nabla^{2}f(\boldsymbol{x}_{i})\right)\right\|_{2}+\frac{\kappa_{\nabla^{2}f}}{k+1}\sum_{i=0}^{k}\left\|\begin{pmatrix}\boldsymbol{x}_{i}-\boldsymbol{x}^{*}\\\boldsymbol{\lambda}_{i}-\boldsymbol{\lambda}^{*}\\[\boldsymbol{\mu}_{1,i}-\boldsymbol{\mu}_{1}^{*}]_{\mathcal{I}(\boldsymbol{x}^{*})}\\[\boldsymbol{\mu}_{2,i}-\boldsymbol{\mu}_{2}^{*}]_{\mathcal{J}(\boldsymbol{x}^{*})}\end{pmatrix}\right\|_{2} \\ &+\kappa_{\nabla c}\left\|\begin{pmatrix}\boldsymbol{x}_{k}-\boldsymbol{x}^{*}\\\boldsymbol{\lambda}_{k}-\boldsymbol{\lambda}^{*}\\[\boldsymbol{\mu}_{1,k}-\boldsymbol{\mu}_{1}^{*}]_{\mathcal{I}(\boldsymbol{x}^{*})}\\[\boldsymbol{\mu}_{2,k}-\boldsymbol{\mu}_{2}^{*}]_{\mathcal{J}(\boldsymbol{x}^{*})}\end{pmatrix}\right\|_{2}.\end{aligned}$$

Note that Δ_k is the modification to the positive-definiteness of B_k . If Δ_k is the matrix with the smallest ℓ_2 -norm such that B_k is positive definite, then $\|\Delta_k\|_2 \leq \|B_k - \nabla^2 f(x^*) - \sum_{i=1}^r (\lambda^*)_i \nabla^2 c_i(x^*)\|_2$. Here, the strong law of large number shows that (C.13)

$$\left\|\frac{1}{k+1}\sum_{i=0}^{k} \left(\nabla^2 f(\boldsymbol{x}_i; \zeta_i) - \nabla^2 f(\boldsymbol{x}_i)\right)\right\|_2 = o\left(\sqrt{\frac{(\log k)^{1+\nu}}{k}}\right) = \mathcal{O}\left(\sqrt{\beta_k} + \alpha_k^{\min}/\beta_k\right), \text{ almost surely,}$$

for any $\nu > 0$. It further shows that H_k (resp. B_k) converges to H^* (resp. B^*) almost surely. Then

$$\mathbb{E}\left[\left\|\boldsymbol{H}_{k}-\boldsymbol{H}^{*}\right\|_{2}^{2}\right] \leq 3\mathbb{E}\left[\left\|\frac{1}{k+1}\sum_{i=0}^{k}\left(\nabla^{2}f(\boldsymbol{x}_{i};\zeta_{i})-\nabla^{2}f(\boldsymbol{x}_{i})\right)\right\|_{2}^{2}\right] \\ +\frac{3\kappa_{\nabla^{2}f}^{2}}{k+1}\sum_{i=0}^{k}\mathbb{E}\left[\left\|\begin{pmatrix}\boldsymbol{x}_{i}-\boldsymbol{x}^{*}\\\boldsymbol{\lambda}_{i}-\boldsymbol{\lambda}^{*}\\ [\boldsymbol{\mu}_{1,i}-\boldsymbol{\mu}_{1}^{*}]_{\mathcal{I}(\boldsymbol{x}^{*})}\\ [\boldsymbol{\mu}_{2,i}-\boldsymbol{\mu}_{2}^{*}]_{\mathcal{J}(\boldsymbol{x}^{*})}\end{pmatrix}\right\|_{2}^{2}\right] + 3\kappa_{\nabla^{2}c}^{2}\mathbb{E}\left[\left\|\begin{pmatrix}\boldsymbol{x}_{k}-\boldsymbol{x}^{*}\\\boldsymbol{\lambda}_{k}-\boldsymbol{\lambda}^{*}\\ [\boldsymbol{\mu}_{1,k}-\boldsymbol{\mu}_{1}^{*}]_{\mathcal{I}(\boldsymbol{x}^{*})}\\ [\boldsymbol{\mu}_{2,k}-\boldsymbol{\mu}_{2}^{*}]_{\mathcal{J}(\boldsymbol{x}^{*})}\end{pmatrix}\right\|_{2}^{2}\right] \\ = \mathcal{O}\left(\beta_{k}+(\alpha_{k}^{\min})^{2}/\beta_{k}^{2}\right).$$

LEMMA 32. Under Assumptions 5 and 6, and suppose that $\iota_1 > b_2$ if $b_1 = 1$, then

$$\mathbb{E}\left[\left\|\mathcal{W}_{2,k}\right\|_{2}^{2}\right] = o\left(\alpha_{k}^{\min}\right).$$

PROOF. For simplicity, we denote

$$\boldsymbol{v}_{k} = \begin{pmatrix} -\bar{\boldsymbol{g}}_{k} - \boldsymbol{\lambda}_{k} \nabla \boldsymbol{c}(\boldsymbol{x}_{k}) + \boldsymbol{\mu}_{1,k} - \boldsymbol{\mu}_{2,k} \\ -\boldsymbol{c}(\boldsymbol{x}_{k}) \\ [\boldsymbol{x}_{k} - \boldsymbol{\ell}]_{\mathcal{I}(\boldsymbol{x}^{*})} \\ [\boldsymbol{u} - \boldsymbol{x}_{k}]_{\mathcal{J}(\boldsymbol{x}^{*})} \end{pmatrix}, \text{ and } \boldsymbol{v}^{*} = \begin{pmatrix} -\nabla f(\boldsymbol{x}^{*}) - \boldsymbol{\lambda}^{*} \nabla \boldsymbol{c}(\boldsymbol{x}^{*}) + \boldsymbol{\mu}_{1}^{*} - \boldsymbol{\mu}_{2}^{*} \\ -\boldsymbol{c}(\boldsymbol{x}^{*}) \\ [\boldsymbol{x}^{*} - \boldsymbol{\ell}]_{\mathcal{I}(\boldsymbol{x}^{*})} \\ [\boldsymbol{u} - \boldsymbol{x}^{*}]_{\mathcal{J}(\boldsymbol{x}^{*})} \end{pmatrix} = \boldsymbol{0}.$$

Then, there exist some $\kappa_v>0$ such that

$$egin{aligned} \|oldsymbol{v}_k-oldsymbol{v}^*\|_2 &\leq \|oldsymbol{ar{g}}_k-
abla f(oldsymbol{x}_k)\|_2+\kappa_v \left\|egin{pmatrix}oldsymbol{x}_k-oldsymbol{x}^*\oldsymbol{\lambda}_k-oldsymbol{\lambda}^*\ [oldsymbol{\mu}_{1,k}-oldsymbol{\mu}_1^*]_{\mathcal{I}(oldsymbol{x}^*)}\ [oldsymbol{\mu}_{2,k}-oldsymbol{\mu}_2^*]_{\mathcal{J}(oldsymbol{x}^*)}\end{pmatrix}
ight\|_2. \end{aligned}$$

We further have

$$\mathbb{E}\left[\left\|\bar{\boldsymbol{p}}_{k}\right\|_{2}^{2}\right] = \mathbb{E}\left[\left\|\boldsymbol{H}_{k}^{-1}\boldsymbol{v}_{k}\right\|_{2}^{2}\right] = \mathbb{E}\left[\left\|\boldsymbol{H}_{k}^{-1}\boldsymbol{v}_{k}-\boldsymbol{H}_{k}^{-1}\boldsymbol{v}^{*}\right\|_{2}^{2}\right]$$

$$\leq \Upsilon_{H}^{2}\mathbb{E}\left[\left\|\boldsymbol{v}_{k}-\boldsymbol{v}^{*}\right\|_{2}^{2}\right]$$

$$(C.14) \qquad \leq 2\Upsilon_{H}^{2}\left(\mathbb{E}\left[\left\|\bar{\boldsymbol{g}}_{k}-\nabla f(\boldsymbol{x}_{k})\right\|_{2}^{2}\right]+\kappa_{v}^{2}\mathbb{E}\left[\left\|\left(\begin{pmatrix}\boldsymbol{x}_{k}-\boldsymbol{x}^{*}\\\boldsymbol{\lambda}_{k}-\boldsymbol{\lambda}^{*}\\ [\boldsymbol{\mu}_{1,k}-\boldsymbol{\mu}_{1}^{*}]_{\mathcal{I}(\boldsymbol{x}^{*})}\\ [\boldsymbol{\mu}_{2,k}-\boldsymbol{\mu}_{2}^{*}]_{\mathcal{J}(\boldsymbol{x}^{*})}\end{pmatrix}\right\|_{2}^{2}\right]\right)$$

$$= \mathcal{O}\left(\beta_{k}\right).$$

Then

$$\mathbb{E}\left[\left\|\mathcal{W}_{2,k}\right\|_{2}^{2}\right] \leq \left(\sum_{h=K^{*}}^{k} \prod_{h'=h+1}^{k} \left(1-\beta_{h'}\right) \alpha_{h-1} \sqrt{\mathbb{E}\left[\left\|\bar{\boldsymbol{p}}_{h-1}\right\|_{2}^{2}\right]}\right)^{2} = \mathcal{O}\left(\left(\alpha_{k}^{\min}\right)^{2}/\beta_{k}\right) = o\left(\alpha_{k}^{\min}\right)^{2}/\beta_{k}\right)$$

After putting back N times, we have

$$\mathbb{E}\left[\left\|\mathcal{W}_{1,k}\right\|_{2}^{2}\right] = \mathcal{O}\left(\beta_{k}\right),$$

and

$$\mathbb{E}\left[\left\|\mathcal{W}_{2,k}\right\|_{2}^{2}\right] = \mathcal{O}\left(\beta_{k}\sum_{i=1}^{N}\left(\alpha_{k}/\beta_{k}\right)^{2i} + \left(\alpha_{k}/\beta_{k}\right)^{2(N+1)}\right).$$

Observe that $\beta_k \sum_{i=1}^N (\alpha_k^{\min}/\beta_k)^{2i} = o(\alpha_k^{\min})$ and N can be any arbitrarily large integer, then

$$\mathbb{E}\left[\left\|\mathcal{W}_{2,k}\right\|_{2}^{2}\right] = o\left(\alpha_{k}^{\min}\right),$$

under the condition that $b_2 < b_1$.

LEMMA 33. Denote

$$\mathcal{E}_{1,k}^{*} = \sum_{i=K^{*}}^{k} \prod_{j=i+1}^{k} (1 - \alpha_{j}^{\min}) \alpha_{i}^{\min} \left((\boldsymbol{H}^{*})^{-1} - \boldsymbol{H}_{i}^{-1} \right) \begin{pmatrix} \bar{\boldsymbol{g}}_{i} - \nabla f(\boldsymbol{x}_{i}) \\ \boldsymbol{0} \\ \boldsymbol{0} \\ \boldsymbol{0} \end{pmatrix}.$$

Under Assumptions 5 and 6, and suppose that $\iota_1 > b_2$ if $b_1 = 1$, then

(C.15)
$$\mathbb{E}\left[\left\|\mathcal{E}_{1,k}^{*}\right\|_{2}\right] = \mathcal{O}\left(\beta_{k}\right),$$

and

(C.16)
$$\mathbb{E}\left[\left\|\mathcal{Q}_{3,k}\right\|_{2}\right] = \mathcal{O}\left(\beta_{k}\right).$$

Proof.

$$\mathbb{E}\left[\left\|\mathcal{E}_{1,k}^{*}\right\|_{2}\right] \leq \sum_{i=K^{*}}^{k} \prod_{j=i+1}^{k} \left(1-\alpha_{j}^{\min}\right) \alpha_{i}^{\min} \sqrt{\mathbb{E}\left[\left\|\boldsymbol{H}_{k}^{-1}-(\boldsymbol{H}^{*})^{-1}\right\|_{2}^{2}\right]} \sqrt{\mathbb{E}\left[\left\|\boldsymbol{\mathcal{W}}_{1,i}+\boldsymbol{\mathcal{W}}_{2,i}\right\|_{2}^{2}\right]} \\ = \mathcal{O}\left(\beta_{k}\right).$$

$$\mathbb{E}\left[\|\mathcal{Q}_{3,k}\|_{2}\right] = \prod_{j=K_{a}+1}^{k} \left(1 - (1-a)\alpha_{j}^{\min}\right) \|\mathcal{Q}_{3,K_{a}}\|_{2} + \sum_{i=K_{a}+1}^{k} \left(\prod_{j=i+1}^{k} \left(1 - (1-a)\alpha_{j}^{\min}\right)\right) \alpha_{i}^{\min} \cdot \sqrt{\mathbb{E}\left[\varepsilon_{i}^{2}\right]} \sqrt{\mathbb{E}\left[\left(\|\mathcal{Q}_{1,i-1}\|_{2} + \|\mathcal{Q}_{2,i-1}\|_{2}\right)^{2}\right]}.$$

Here,

$$\mathbb{E}\left[\varepsilon_{i}^{2}\right] = \mathcal{O}\left(\mathbb{E}\left[\|\mathcal{Q}_{1,i}\|_{2}^{2} + \|\mathcal{Q}_{2,i}\|_{2}^{2} + \|\mathcal{Q}_{3,i}\|_{2}^{2} + \|\mathbf{H}_{i+1} - \mathbf{H}^{*}\|_{2}^{2}\right]\right) = \mathcal{O}\left(\beta_{i}\right)$$

and

$$\mathbb{E}\left[\left(\left\|\mathcal{Q}_{1,i}\right\|_{2}+\left\|\mathcal{Q}_{2,i}\right\|_{2}\right)^{2}\right]=\mathcal{O}\left(\beta_{i}\right)$$

complete the first part of the proof.

LEMMA 34. Let

$$\mathcal{Q}_{1,k} = \mathcal{Q}_{1,k}^* + \mathcal{E}_{1,k}^*,$$

,

where

$$\begin{aligned} \mathcal{Q}_{1,k}^{*} &= \sum_{i=K^{*}}^{k} \prod_{j=i+1}^{k} (1-\alpha_{j}^{\min}) \alpha_{i}^{\min} (-H^{*})^{-1} \begin{pmatrix} \bar{g}_{i} - \nabla f(\boldsymbol{x}_{i}) \\ \boldsymbol{0} \\ \boldsymbol{0} \\ \boldsymbol{0} \end{pmatrix} \\ &:= \sum_{i=K^{*}}^{k} \prod_{j=i+1}^{k} (1-\alpha_{j}^{\min}) \alpha_{i}^{\min} (-H^{*})^{-1} \left(\mathcal{W}_{1,i} + \mathcal{W}_{2,i} \right). \end{aligned}$$

and

$$\mathcal{E}_{1,k}^{*} = \sum_{i=K^{*}}^{k} \prod_{j=i+1}^{k} (1-\alpha_{j}^{\min}) \alpha_{i}^{\min} \left((-H_{i})^{-1} - (-H^{*})^{-1} \right) \begin{pmatrix} \bar{g}_{i} - \nabla f(x_{i}) \\ 0 \\ 0 \\ 0 \end{pmatrix}.$$

Under Assumptions 5 and 6, and suppose that $\iota_1 > b_2$ if $b_1 = 1$, then

(C.17)
$$\frac{1}{\sqrt{\alpha_k^{min}}} \mathcal{Q}_{1,k}^* \to \mathcal{N}(\mathbf{0}, \Theta \mathbf{\Omega}^*).$$

PROOF. Let
(C.18)

$$\mathcal{Q}_{1,k}^{**} := \sum_{i=K^*}^k \prod_{j=i+1}^k (1-\alpha_j^{\min}) \alpha_i^{\min} (-H^*)^{-1} \mathcal{W}_{1,i}$$

$$= \sum_{i=K^*}^k \prod_{j=i+1}^k (1-\alpha_j^{\min}) \alpha_i^{\min} \sum_{h=K^*}^i \left(\prod_{h'=h+1}^i (1-\beta_{h'})\right) \beta_h (-H^*)^{-1} \begin{pmatrix} g_h - \nabla f(x_h) \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

$$= \sum_{h=K^*}^k \sum_{i=h}^k \prod_{j=i+1}^k (1-\alpha_j^{\min}) \alpha_i^{\min} \prod_{h'=h+1}^i (1-\beta_{h'}) \beta_h (-H^*)^{-1} \begin{pmatrix} g_h - \nabla f(x_h) \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

$$:= \sum_{h=K^*}^k a_{h,k} (-H^*)^{-1} \begin{pmatrix} g_h - \nabla f(x_h) \\ 0 \\ 0 \\ 0 \end{pmatrix} := \sum_{h=K^*}^k s_{h,k},$$

where $a_{h,k} = \sum_{i=h}^{k} \prod_{j=i+1}^{k} \left(1 - \alpha_{j}^{\min}\right) \alpha_{i}^{\min} \prod_{h'=h+1}^{i} (1 - \beta_{h'}) \beta_{h}$ and $s_{h,k}$ are independent for different *h*. The asymptotic normality can be implied by the central limit theorem for the martingale difference array. Before that, we first verify the corresponding conditions. We first denote

$$m{\phi}_h^* = (-m{H}^*)^{-1} egin{pmatrix} m{g}_h -
abla f(m{x}_h) \\ m{0} \\ m{0} \end{pmatrix},$$

and note that $\mathbb{E}\left[\phi_h^*\phi_h^{*\top}|\mathcal{F}_{h-1}\right] \to \Omega^*$ as $h \to \infty$ almost surely, since the smoothness of $f(\boldsymbol{x},\xi)$ shows that

$$\begin{split} \mathbf{\Lambda}_{i} &:= \mathbb{E}\left[\mathbf{g}_{i}\mathbf{g}_{i}^{\top} - \nabla f(\mathbf{x}_{i})\nabla f(\mathbf{x}_{i})^{\top}|\mathcal{F}_{i-1}\right] - \mathbb{E}\left[\nabla f(\mathbf{x}^{*};\xi)\nabla f(\mathbf{x}^{*};\xi)^{\top} - \nabla f(\mathbf{x}^{*})\nabla f(\mathbf{x}^{*})^{\top}\right] \\ &= \mathbb{E}\left[\nabla f(\mathbf{x}_{i};\xi)\nabla f(\mathbf{x}_{i};\xi)^{\top} - \nabla f(\mathbf{x}^{*};\xi)\nabla f(\mathbf{x}_{i};\xi)^{\top}|\mathcal{F}_{i-1}\right] + \mathbb{E}\left[\nabla f(\mathbf{x}^{*};\xi)\nabla f(\mathbf{x}_{i};\xi)^{\top} - \nabla f(\mathbf{x}^{*};\xi)\nabla f(\mathbf{x}^{*};\xi)^{\top}|\mathcal{F}_{i-1}\right] \\ &+ \nabla f(\mathbf{x}_{i})\nabla f(\mathbf{x}_{i})^{\top} - \nabla f(\mathbf{x}^{*})\nabla f(\mathbf{x}_{i})^{\top} + \nabla f(\mathbf{x}^{*})\nabla f(\mathbf{x}_{i})^{\top} - \nabla f(\mathbf{x}^{*})\nabla f(\mathbf{x}^{*})^{\top} \\ &\leq \kappa_{\nabla f} \|\mathbf{x}_{i} - \mathbf{x}^{*}\|_{2} \left(\sqrt{\mathbb{E}\left[\|\nabla f(\mathbf{x}_{i};\xi)\|_{2}^{2}|\mathcal{F}_{i-1}\right]} + \sqrt{\mathbb{E}\left[\|\nabla f(\mathbf{x}^{*};\xi)\|_{2}^{2}\right]} + \|\nabla f(\mathbf{x}_{i})\|_{2} + \|\nabla f(\mathbf{x}^{*})\|_{2}\right) \\ &\leq 4\kappa_{\nabla f}M_{\nabla f}\|\mathbf{x}_{i} - \mathbf{x}^{*}\|_{2} \rightarrow 0, \text{ as } i \rightarrow \infty. \end{split}$$

Then,

$$\begin{split} &\sum_{h=K^{*}}^{k} \mathbb{E}\left[a_{h,k}^{2}\phi_{h}^{*}\phi_{h}^{*\top}|\mathcal{F}_{h-1}\right] \\ &= \sum_{i=K^{*}}^{k} \prod_{j=i+1}^{k} (1-\alpha_{j}^{\min})\alpha_{i}^{\min} \cdot \sum_{i'=K^{*}}^{k} \prod_{j'=i'+1}^{k} (1-\alpha_{j'}^{\min})\alpha_{i'}^{\min} \sum_{h=K^{*}}^{\min\{i,i'\}} \left(\prod_{h'=h+1}^{i} (1-\beta_{h'})\right) \left(\prod_{h'=h+1}^{i'} (1-\beta_{h'})\right) \beta_{h}^{2} \\ &\cdot \mathbb{E}\left[\phi_{h}^{*}\phi_{h}^{*\top}|\mathcal{F}_{h-1}\right] \\ &= 2\sum_{i=K^{*}}^{k} \prod_{j=i+1}^{k} (1-\alpha_{j}^{\min})\alpha_{i}^{\min} \cdot \sum_{i'=K^{*}}^{i} \prod_{j'=i'+1}^{k} (1-\alpha_{j'}^{\min})\alpha_{i'}^{\min} \sum_{h=K^{*}}^{i'} \left(\prod_{h'=h+1}^{i} (1-\beta_{h'})\right) \left(\prod_{h'=h+1}^{i'} (1-\beta_{h'})\right) \beta_{h}^{2} \\ &\cdot \mathbb{E}\left[\phi_{h}^{*}\phi_{h}^{*\top}|\mathcal{F}_{h-1}\right] - \sum_{i=K^{*}}^{k} \prod_{j=i+1}^{k} (1-\alpha_{j}^{\min})^{2} (\alpha_{i}^{\min})^{2} \sum_{h=K^{*}}^{i} \left(\prod_{h'=h+1}^{i} (1-\beta_{h'})\right)^{2} \beta_{h}^{2} \mathbb{E}\left[\phi_{h}^{*}\phi_{h}^{*\top}|\mathcal{F}_{h-1}\right] \\ &= 2\sum_{i=K^{*}}^{k} \prod_{j=i+1}^{k} (1-\alpha_{j}^{\min})^{2} \alpha_{i}^{\min} \cdot \sum_{i'=K^{*}}^{i} \prod_{j'=i'+1}^{i} (1-\alpha_{j'}^{\min})^{2} (\alpha_{i}^{\min})^{2} \sum_{h=K^{*}}^{i} \left(\prod_{h'=h+1}^{i} (1-\beta_{h'})\right)^{2} \beta_{h}^{2} \mathbb{E}\left[\phi_{h}^{*}\phi_{h}^{*\top}|\mathcal{F}_{h-1}\right] - \sum_{i=K^{*}}^{k} \prod_{j'=i'+1}^{k} (1-\alpha_{j'}^{\min})^{2} (\alpha_{i}^{\min})^{2} \sum_{h=K^{*}}^{i} \left(\prod_{h'=h+1}^{i} (1-\beta_{h'})\right)^{2} \beta_{h}^{2} \mathbb{E}\left[\phi_{h}^{*}\phi_{h}^{*\top}|\mathcal{F}_{h-1}\right] - \sum_{i=K^{*}}^{k} \prod_{j'=i'+1}^{k} (1-\alpha_{j'}^{\min})^{2} (\alpha_{i}^{\min})^{2} \sum_{h=K^{*}}^{i} \left(\prod_{h'=h+1}^{i} (1-\beta_{h'})\right)^{2} \beta_{h}^{2} \mathbb{E}\left[\phi_{h}^{*}\phi_{h}^{*\top}|\mathcal{F}_{h-1}\right]. \end{split}$$

Note that

$$\lim_{i' \to \infty} \beta_i^{-1} \sum_{h=K^*}^{i'} \left(\prod_{h'=h+1}^{i'} (1-\beta_{h'}) \right)^2 \beta_h^2 \mathbb{E} \left[\phi_h^* \phi_h^{*\top} | \mathcal{F}_{h-1} \right] = \frac{1}{2} \Omega^*,$$
(C.19)

$$\lim_{i \to \infty} (\alpha_i^{\min})^{-1} \sum_{i'=K^*}^{i} \prod_{j'=i'+1}^{i} (1-\alpha_{j'}^{\min}) (1-\beta_{j'}) \alpha_{i'}^{\min} \beta_{i'} = 1,$$
(C.19)

$$\lim_{k \to \infty} (\alpha_k^{\min})^{-1} \sum_{i=K^*}^{k} \prod_{j=i+1}^{k} (1-\alpha_j^{\min})^2 (\alpha_i^{\min})^2 = \Theta := \begin{cases} 1/2, & \text{if } b_1 < 1, \\ 1/\left(2-\frac{1}{\iota_1}\right), & \text{if } b_1 = 1, \end{cases}$$

$$\lim_{i \to \infty} (\alpha_k^{\min})^{-1} \sum_{i=K^*}^{k} \prod_{j=i+1}^{k} (1-\alpha_j^{\min})^2 (\alpha_i^{\min})^2 \beta_i = 0,$$

where we require that $\iota_1 > \frac{1}{2}$ if $b_1 = 1$. Therefore,

$$\lim_{k \to \infty} (\alpha_k^{\min})^{-1} \sum_{h=K^*}^k \mathbb{E} \left[a_{h,k}^2 \phi_h^* \phi_h^{*\top} | \mathcal{F}_{h-1} \right] = \Theta \Omega^*.$$

We then verify the Lindeberg condition. It is equivalent to showing that

(C.20)
$$\lim_{k \to \infty} \frac{1}{\alpha_k^{\min}} \sum_{h=K^*}^k a_{h,k}^2 \mathbb{E} \left[\| \boldsymbol{\phi}_h^* \|_2^2 \cdot \mathbf{1}_{\| a_{h,k} \boldsymbol{\phi}_h^* \|_2 \ge \epsilon(\alpha_k^{\min})^{1/2}} | \mathcal{F}_{h-1} \right]$$
$$\leq \lim_{k \to \infty} \frac{1}{\epsilon(\alpha_k^{\min})^{3/2}} \sum_{h=K^*}^k a_{h,k}^3 \mathbb{E} \left[\| \boldsymbol{\phi}_h^* \|_2^3 | \mathcal{F}_{h-1} \right] \le \lim_{k \to \infty} \frac{\Upsilon_{\boldsymbol{\phi}}}{\epsilon(\alpha_k^{\min})^{3/2}} \sum_{h=K^*}^k a_{h,k}^3 = 0.$$

Suppose that $X_1, X_2, \dots, X_k, \dots$ are i.i.d. 1-dimensional random variables with zero mean and unit 3-moment, i.e., $\mathbb{E}\left[X_i^3\right] = 1$ for all $i \in \mathbb{N}$, then $\sum_{h=K^*}^k a_{h,k}^3 = \mathbb{E}\left[\left(\sum_{h=K^*}^k a_{h,k}X_h\right)^3\right]$. The equivalent form

$$\sum_{h=K^*}^k a_{h,k} X_h = \sum_{i=K^*}^k \prod_{j=i+1}^k (1-\alpha_j^{\min}) \alpha_i^{\min} \sum_{h=K^*}^i \left(\prod_{h'=h+1}^i (1-\beta_{h'})\right) \beta_h X_h$$

further shows

$$\begin{split} &\sum_{h=K^*}^k a_{h,k}^3 = \mathbb{E}\left[\left(\sum_{h=K^*}^k a_{h,k} X_h\right)^3\right] \\ = \mathbb{E}\left[\left(\sum_{i=K^*}^k \prod_{j=i+1}^k (1-\alpha_j^{\min})\alpha_i^{\min} \sum_{h=K^*}^i \left(\prod_{h'=h+1}^i (1-\beta_{h'})\right)\beta_h X_h\right)^3\right] \\ &\leq 6\sum_{i=K^*}^k \prod_{j=i+1}^k (1-\alpha_j^{\min})\alpha_i^{\min} \cdot \sum_{i'=K^*}^i \prod_{j'=i'+1}^k (1-\alpha_{j'}^{\min})\alpha_{i''}^{\min} \cdot \sum_{i''=K^*}^{i'} \prod_{j''=i''+1}^k (1-\alpha_{j''}^{\min})\alpha_{i''}^{\min} \\ &\cdot \sum_{h=K^*}^{i''} \left(\prod_{h'=h+1}^i (1-\beta_{h'})\right) \left(\prod_{h'=h+1}^{i'} (1-\beta_{h'})\right) \left(\prod_{h'=h+1}^{i''} (1-\beta_{h'})\right) \beta_h^3 \mathbb{E}\left[X_h^3\right] \\ &= 6\sum_{i=K^*}^k \prod_{j=i+1}^k (1-\alpha_{j^{\min}}^{\min})^3 \alpha_i^{\min} \cdot \sum_{i'=K^*}^i \prod_{j'=i'+1}^i (1-\alpha_{j''}^{\min})^2 (1-\beta_{j'}) \alpha_{i''}^{\min} \cdot \sum_{i''=K^*}^{i'} \prod_{j''=i''+1}^{i''} (1-\alpha_{j''}^{\min}) (1-\beta_{j''})^2 \alpha_{i''}^{\min} \\ &\cdot \sum_{h=K^*}^{i''} \left(\prod_{h'=h+1}^{i''} (1-\beta_{h'})\right)^3 \beta_h^3. \end{split}$$

Similarly, note that

(C.21)

$$\sum_{h=K^{*}}^{i''} \left(\prod_{h'=h+1}^{i''} (1-\beta_{h'}) \right)^{3} \beta_{h}^{3} = \frac{1}{3} = \mathcal{O}\left(\beta_{i''}^{2}\right) \\
\sum_{i''=K^{*}}^{i'} \prod_{j''=i''+1}^{i'} (1-\alpha_{j''}^{\min})(1-\beta_{j''})^{2} \alpha_{i''}^{\min} \beta_{i''}^{2} = \mathcal{O}\left(\alpha_{i'}^{\min} \beta_{i'}\right), \\
\sum_{i'=K^{*}}^{i} \prod_{j'=i'+1}^{i} (1-\alpha_{j'}^{\min})^{2} (1-\beta_{j'})(\alpha_{i'}^{\min})^{2} \beta_{i'} = \mathcal{O}\left((\alpha_{i}^{\min})^{2}\right), \\
\sum_{i=K^{*}}^{k} \prod_{j=i+1}^{k} (1-\alpha_{j}^{\min})^{3} (\alpha_{i}^{\min})^{3} = \mathcal{O}\left((\alpha_{k}^{\min})^{2}\right),$$

where we require that $\iota_1 > b_2$ if $b_1 = 1$. The above results imply that $\sum_{h=K^*}^k a_{h,k}^3 = \mathcal{O}\left((\alpha_k^{\min})^2\right)$, thus the Lindeberg condition is satisfied. By the central limit theorem for martingale difference array (also called Lévy's theorem), we deduce that

$$\frac{1}{\sqrt{\alpha_k^{\min}}} \mathcal{Q}_{1,k}^{**} \xrightarrow{d} \mathcal{N}(\mathbf{0}, \Theta \mathbf{\Omega}^*).$$

Denote

$$\mathcal{E}_{1,k}^{**} = \mathcal{Q}_{1,k}^{*} - \mathcal{Q}_{1,k}^{**} := \sum_{i=K^*}^{k} \prod_{j=i+1}^{k} (1 - \alpha_j^{\min}) \alpha_i^{\min} (-\boldsymbol{H}^*)^{-1} \mathcal{W}_{2,i},$$

according to Lemma, we have

$$\mathbb{E}\left[\left\|\mathcal{E}_{1,k}^{**}\right\|_{2}\right] \leq \Upsilon_{H} \sum_{i=K^{*}}^{k} \prod_{j=i+1}^{k} (1-\alpha_{j}^{\min})\alpha_{i}^{\min} \sqrt{\mathbb{E}\left[\left\|\mathcal{W}_{2,i}\right\|_{2}^{2}\right]} = o\left(\sqrt{\alpha_{k}^{\min}}\right),$$

where we require that $\iota_1 > b_2$ if $b_1 = 1$. By Slutsky's theorem,

$$\frac{1}{\sqrt{\alpha_k^{\min}}} \mathcal{Q}_{1,k}^* \xrightarrow{d} \mathcal{N}(\mathbf{0}, \Theta \mathbf{\Omega}^*) \,.$$

Proof for Theorem 4: it is a direct result from Lemmas 28, 33 and 34.

C.4. Proof for Theorem 5. The second relation is implied by the first one because the proof in Lemma 31 and the almost sure convergence of primal-dual iterates jointly show that $\left\| \boldsymbol{H}_{k}^{-1} - (\boldsymbol{H}^{*})^{-1} \right\|_{2} \rightarrow 0$, almost surely. We are left to show the first relation. Note that

$$\begin{split} \|\boldsymbol{\Sigma}_{k} - \boldsymbol{\Sigma}^{*}\|_{2} &= \left\| \frac{1}{k+1} \sum_{i=0}^{k} \boldsymbol{g}_{i} \boldsymbol{g}_{i}^{\top} - \mathbb{E} \left[\nabla f(\boldsymbol{x}^{*}; \zeta) \nabla f(\boldsymbol{x}^{*}; \zeta)^{\top} \right] \right\|_{2} \\ &+ \left\| \left(\frac{1}{k+1} \sum_{i=0}^{k} \boldsymbol{g}_{i} \right) \left(\frac{1}{k+1} \sum_{i=0}^{k} \boldsymbol{g}_{i} \right)^{\top} - \nabla f(\boldsymbol{x}^{*}) \nabla f(\boldsymbol{x}^{*})^{\top} \right\|_{2} \right. \\ &\left\| \frac{1}{k+1} \sum_{i=0}^{k} \boldsymbol{g}_{i} \boldsymbol{g}_{i}^{\top} - \mathbb{E} \left[\nabla f(\boldsymbol{x}^{*}; \zeta) \nabla f(\boldsymbol{x}^{*}; \zeta)^{\top} \right] \right\|_{2} \\ &= \left\| \frac{1}{k+1} \sum_{i=0}^{k} \boldsymbol{g}_{i} \boldsymbol{g}_{i}^{\top} - \mathbb{E} \left[\nabla f(\boldsymbol{x}_{i}; \zeta) \nabla f(\boldsymbol{x}_{i}; \zeta)^{\top} |\mathcal{F}_{i-1} \right] \right\|_{2} \\ &+ \left\| \frac{1}{k+1} \sum_{i=0}^{k} \mathbb{E} \left[\nabla f(\boldsymbol{x}_{i}; \zeta) \nabla f(\boldsymbol{x}_{i}; \zeta)^{\top} |\mathcal{F}_{i-1} - \nabla f(\boldsymbol{x}^{*}; \zeta) \nabla f(\boldsymbol{x}^{*}; \zeta)^{\top} \right] \right\|_{2} \end{split}$$

The strong law of large number shows that

$$\left\|\frac{1}{k+1}\sum_{i=0}^{k}\boldsymbol{g}_{i}\boldsymbol{g}_{i}^{\top}-\mathbb{E}\left[\nabla f(\boldsymbol{x}_{i};\boldsymbol{\zeta})\nabla f(\boldsymbol{x}_{i};\boldsymbol{\zeta})^{\top}|\mathcal{F}_{i-1}\right]\right\|_{2}=o\left(\sqrt{\frac{\left(\log k\right)^{1+\nu}}{k}}\right),$$

for any $\nu > 0$, almost surely. The almost sure convergence of iterates (i.e., $x_k \to x^*$) implies that

$$\left\|\frac{1}{k+1}\sum_{i=0}^{k} \mathbb{E}\left[\nabla f(\boldsymbol{x}_{i};\boldsymbol{\zeta})\nabla f(\boldsymbol{x}_{i};\boldsymbol{\zeta})^{\top} - \nabla f(\boldsymbol{x}^{*};\boldsymbol{\zeta})\nabla f(\boldsymbol{x}^{*};\boldsymbol{\zeta})^{\top}|\mathcal{F}_{i-1}\right]\right\|_{2} \to 0,$$

almost surely. Similarly, for the second term

$$\left\|\frac{1}{k+1}\sum_{i=0}^{k}\boldsymbol{g}_{i}-\nabla f(\boldsymbol{x}^{*})\right\|_{2} \leq \left\|\frac{1}{k+1}\sum_{i=0}^{k}\left(\boldsymbol{g}_{i}-\nabla f(\boldsymbol{x}_{k})\right)\right\|_{2}+\left\|\frac{1}{k+1}\sum_{i=0}^{k}\left(\nabla f(\boldsymbol{x}_{k})-\nabla f(\boldsymbol{x}^{*})\right)\right\|_{2},$$

the strong law of large number also shows

$$\left\|\frac{1}{k+1}\sum_{i=0}^{k}\left(\boldsymbol{g}_{i}-\nabla f(\boldsymbol{x}_{k})\right)\right\|_{2}=o\left(\sqrt{\frac{\left(\log k\right)^{1+\nu}}{k}}\right),$$

for any $\nu > 0$ almost surely, and

$$\left\|\frac{1}{k+1}\sum_{i=0}^{k}\left(\nabla f(\boldsymbol{x}_{k})-\nabla f(\boldsymbol{x}^{*})\right)\right\|_{2}\to 0,$$

almost surely. Therefore, we complete the proof.