

Online Inference of Constrained Optimization: Primal-Dual Optimality and Sequential Quadratic Programming

Yihang Gao

Department of Mathematics, National University of Singapore

Michael K. Ng

Department of Mathematics, Hong Kong Baptist University

Michael W. Mahoney

ICSI and Department of Statistics, University of California, Berkeley

Sen Na

School of Industrial and Systems Engineering, Georgia Institute of Technology

Abstract

We study online statistical inference for the solutions of stochastic optimization problems with equality and inequality constraints. Such problems are prevalent in statistics and machine learning, encompassing constrained M -estimation, physics-informed models, safe reinforcement learning, and algorithmic fairness. We develop a stochastic sequential quadratic programming (SSQP) method to solve these problems, where the step direction is computed by sequentially 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 stochastic problems lies in dealing with the bias in the step direction. As such, we apply a momentum-style gradient moving-average technique within SSQP to debias the step. We show that our method achieves global almost-sure convergence and exhibits local asymptotic normality with an *optimal* primal-dual limiting covariance matrix in the sense of Hájek and Le Cam. In addition, we provide a plug-in covariance matrix estimator for practical inference. To our knowledge, the proposed SSQP method is the *first fully online method* that attains primal-dual asymptotic minimax optimality without relying on projection operators onto the constraint set, which are generally intractable for nonlinear problems. Through extensive experiments on benchmark nonlinear problems, as well as on constrained generalized linear models and portfolio allocation problems using both synthetic and real data, we demonstrate superior performance of our method, showing that the method and its asymptotic behavior not only solve constrained stochastic problems efficiently but also provide valid and practical online inference in real-world applications.

Keywords: constrained model inference; primal-dual minimax optimality; stochastic SQP; gradient momentum

1 Introduction

We consider stochastic optimization problems with equality and box inequality constraints:

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) = \mathbb{E}_{\zeta \sim \mathcal{P}} [F(\mathbf{x}; \zeta)] \quad \text{s.t.} \quad \mathbf{c}(\mathbf{x}) = \mathbf{0}, \quad \boldsymbol{\ell} \leq \mathbf{x} \leq \mathbf{u}. \quad (1)$$

Here, $\boldsymbol{\ell}, \mathbf{u} \in \mathbb{R}^d$ denote the lower and upper bounds, respectively, with " \leq " representing element-wise comparison; $F(\cdot; \zeta) : \mathbb{R}^d \rightarrow \mathbb{R}$ denotes a realization of the stochastic objective f ; and $\mathbf{c} : \mathbb{R}^d \rightarrow \mathbb{R}^m$ encodes deterministic equality constraints. Throughout the paper, we assume $f, \mathbf{c}, F(\cdot; \zeta)$ are twice continuously differentiable while potentially nonconvex and nonlinear. Problem (1) readily accommodates nonlinear inequality constraints $\mathbf{c}_{\mathcal{I}}(\mathbf{x}) \leq \mathbf{0}$ by introducing slack variables, i.e., reformulating them as $\mathbf{c}_{\mathcal{I}}(\mathbf{x}) + \mathbf{y} = \mathbf{0}$ and $\mathbf{y} \geq \mathbf{0}$.

Problem (1) is ubiquitous with many statistical and machine learning applications, where constraints can encode prior domain knowledge, ensure models' identifiability, and reduce models' intrinsic dimensionality. We introduce concrete motivating examples in Appendix A. Given the ubiquity of Problem (1), statisticians aim to estimate its (local) solution \mathbf{x}^* and perform statistical inference. Arguably, the most primitive estimator is constrained M -estimator, where we draw n samples $\zeta_1, \dots, \zeta_n \stackrel{\text{iid}}{\sim} \mathcal{P}$ and estimate the population loss f by the empirical loss \hat{f}_n :

$$\min_{\mathbf{x} \in \mathbb{R}^d} \hat{f}_n(\mathbf{x}) := \frac{1}{n} \sum_{i=1}^n F(\mathbf{x}; \zeta_i) \quad \text{s.t.} \quad \mathbf{c}(\mathbf{x}) = \mathbf{0}, \quad \boldsymbol{\ell} \leq \mathbf{x} \leq \mathbf{u}. \quad (2)$$

In fact, it is known that the above constrained M -estimation attains *asymptotic optimality* in Hájek and Le Cam's sense (Hájek, 1972; Le Cam, 1972). Roughly speaking, under certain regularity conditions, the (local) minimizer $\hat{\mathbf{x}}_n$ of Problem (2) exhibits asymptotic normality with the smallest covariance matrix, given by

$$\sqrt{n} (\hat{\mathbf{x}}_n - \mathbf{x}^*) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \mathbf{L}^\dagger \text{Cov}(\nabla F(\mathbf{x}^*; \zeta)) \mathbf{L}^\dagger), \quad (3)$$

where $\mathbf{L} = \mathbf{P}_{\text{Null}(\mathbf{J}^*)} \nabla_{\mathbf{x}}^2 \mathcal{L}^* \mathbf{P}_{\text{Null}(\mathbf{J}^*)}$ and \dagger denotes the Moore-Penrose pseudoinverse. Here, we let

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) := f(\mathbf{x}) + \boldsymbol{\lambda}^\top \mathbf{c}(\mathbf{x}) + \boldsymbol{\mu}_1^\top (\boldsymbol{\ell} - \mathbf{x}) + \boldsymbol{\mu}_2^\top (\mathbf{x} - \mathbf{u}) \quad (4)$$

be the Lagrangian function of (1) with $\boldsymbol{\lambda} \in \mathbb{R}^m$ and $\boldsymbol{\mu} = (\boldsymbol{\mu}_1, \boldsymbol{\mu}_2) \in \mathbb{R}^{2d}$ being the dual variables

associated with the equality and inequality constraints; $\nabla_{\mathbf{x}}^2 \mathcal{L}^* = \nabla_{\mathbf{x}}^2 \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$ be the Lagrangian Hessian at the primal-dual solution $(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$ (with respect to \mathbf{x}); and $\mathbf{P}_{\text{Null}(\mathbf{J}^*)}$ denote the projection matrix onto the null space of $\mathbf{J}^* = \mathbf{J}(\mathbf{x}^*)$, where $\mathbf{J}(\mathbf{x})$ is the Jacobian of the active constraints at \mathbf{x} (see (9)). We refer to Duchi and Ruan, 2021; Davis et al., 2024 for the rigorous statement of the result (3); and see Theorem 2.5 for primal-dual generalization.

Although $\widehat{\mathbf{x}}_n$ enjoys nice properties, offline M -estimation requires dealing with a full batch of samples, leading to significant computational costs. Over the past decades, online methods that can handle online streaming data have been preferred. One of the most fundamental online methods is Stochastic Gradient Descent (SGD) (Robbins and Monro, 1951; Kiefer and Wolfowitz, 1952). Recent work has significantly advanced local asymptotic analysis of SGD (and its variates) to enable online inference of \mathbf{x}^* . See, e.g., Polyak and Juditsky, 1992; Toulis and Airolidi, 2017; Fang et al., 2018; Toulis et al., 2021; Chen et al., 2020; 2024; Zhu et al., 2021; Lee et al., 2022.

The above literature all studied unconstrained stochastic optimization, while to accommodate constraints, Duchi and Ruan, 2021 designed a Projected Riemannian Stochastic Gradient method, which to our knowledge is the *first fully online* method capable of achieving asymptotic optimality for constrained model parameters \mathbf{x}^* . The authors also identified a gap between problems with linear and nonlinear constraints, illustrating the reason why vanilla Projected SGD,

$$\mathbf{x}_{k+1} = \text{Proj}_{\Omega}(\mathbf{x}_k - \alpha_k \nabla F(\mathbf{x}_k; \zeta_k)), \quad \Omega = \{\mathbf{x} \in \mathbb{R}^d : \mathbf{c}(\mathbf{x}) = \mathbf{0}, \ell \leq \mathbf{x} \leq \mathbf{u}\}, \quad (5)$$

can be sub-optimal for nonlinear constraints. Notably, the gap has recently been closed by Davis et al., 2024, showing that the averaged Projected SGD, $\bar{\mathbf{x}}_n = \frac{1}{n} \sum_{k=0}^{n-1} \mathbf{x}_k$, actually attains the same asymptotic optimality for nonlinear constraints as displayed in (3). A consistent limiting covariance estimator has also been proposed in the follow-up work Jiang et al., 2025.

Despite these advances in online constrained inference, existing methods rely on a projection operator onto the feasible set $\text{Proj}_{\Omega}(\cdot)$, making them, as stated in Duchi and Ruan, 2021, *more intellectually intriguing than practically applicable*. In particular, projections require solving non-

linear equations that is often intractable and needs global knowledge of the feasible set Ω . However, many problems (cf. CUTEst in Section 5.1) only have access to black-box oracles that return the local constraint value and Jacobian at a given point, making the projection inapplicable. Even in the ideal case where we have global information of the feasible set, the projection at each step may only be approximately computed by solving an intensive optimization subproblem; yet the resulting approximation error is typically neglected in existing analyses, which rely heavily on exact computation of $\text{Proj}_\Omega(\cdot)$. Furthermore, a key distinction between constrained and unconstrained problems lies in the presence of dual multipliers. The dual component is not involved in existing methods, yet it plays a crucial role in stationarity certification and active-set identification (ensured by complementary slackness theorem; see (7)). Thus, performing primal-dual joint inference is of fundamental interest. This motivation leads to the following question: *What is the asymptotic minimax optimality for the primal-dual solution $(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$ of Problem (1), and how to obtain the optimal joint online estimator?*

1.1 Backbone: Sequential Quadratic Programming

In this paper, we first establish primal-dual asymptotic minimax optimality of $(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$, and then design a *Stochastic Sequential Quadratic Programming* (SSQP) method to attain the lower bound. The idea is to *localize* the model along with constraints, and utilize dual information of constrained local model when updating the estimate. In particular, at the k -th step, we perform a quadratic approximation of the loss and a linear approximation of the constraints in (1):

$$\min_{\bar{\Delta}\mathbf{x}_k \in \mathbb{R}^d} \nabla F(\mathbf{x}_k; \zeta_k)^\top \bar{\Delta}\mathbf{x}_k + \frac{1}{2} \bar{\Delta}\mathbf{x}_k^\top \bar{\mathbf{B}}_k \bar{\Delta}\mathbf{x}_k, \quad \text{s.t.} \quad \mathbf{c}_k + \nabla \mathbf{c}_k \bar{\Delta}\mathbf{x}_k = \mathbf{0}, \quad \boldsymbol{\ell} \leq \mathbf{x}_k + \bar{\Delta}\mathbf{x}_k \leq \mathbf{u}. \quad (6)$$

Here, we denote the constraint value and Jacobian by $\mathbf{c}_k = \mathbf{c}(\mathbf{x}_k)$ and $\nabla \mathbf{c}_k = \nabla \mathbf{c}(\mathbf{x}_k)$; and $\bar{\mathbf{B}}_k$, rather than estimating $\nabla^2 f_k$, estimates the Lagrangian Hessian $\nabla_{\mathbf{x}}^2 \mathcal{L}_k$ to leverage the curvature information of nonlinear constraints. After solving (6) and obtaining its primal-dual solution $(\bar{\Delta}\mathbf{x}_k, \bar{\boldsymbol{\lambda}}_k^{\text{sub}}, \bar{\boldsymbol{\mu}}_k^{\text{sub}})$, we update the iterate as

$$(\mathbf{x}_{k+1}, \boldsymbol{\lambda}_{k+1}, \boldsymbol{\mu}_{k+1}) = (\mathbf{x}_k, \boldsymbol{\lambda}_k, \boldsymbol{\mu}_k) + \alpha_k (\bar{\Delta}\mathbf{x}_k, \bar{\boldsymbol{\lambda}}_k^{\text{sub}} - \boldsymbol{\lambda}_k, \bar{\boldsymbol{\mu}}_k^{\text{sub}} - \boldsymbol{\mu}_k).$$

Although the above scheme streamlines the estimation procedure by solving a linear–quadratic program at each step, two fundamental challenges remain. First, even if the original nonlinear constraints $\mathbf{c}(\mathbf{x}) = \mathbf{0}$, $\boldsymbol{\ell} \leq \mathbf{x} \leq \mathbf{u}$ admit a nonempty feasible set Ω (cf. (5)), their linearization in (6) may become infeasible. Second, the inequality constraints truncate the noise in a way such that, even if $\nabla F(\mathbf{x}_k; \zeta_k)$ is an unbiased estimate of the true gradient ∇f_k , the resulting step direction $\bar{\Delta}\mathbf{x}_k$ can still be biased (cf. Remark 3.3), thereby preventing convergence of the scheme. To address the two challenges, we respectively introduce a constraint relaxation technique and a gradient (and Hessian) moving average technique. Although momentum-based gradient method has been widely studied in optimization context since the introduction in Polyak, 1964, our method is the *first* attempt to incorporate both gradient and Hessian momentum into inequality-constrained problems and show promising asymptotic properties to facilitate inference. Furthermore, our method allows an adaptive random stepsize α_k , which is often preferred in practice.

Specifically, with all above designs, we show that the *last* SSQP iterate exhibits asymptotic normality, with an optimal primal-dual covariance matrix matching the theoretical lower bound. In particular, the marginal primal covariance of the SSQP estimator also matches that of the offline constrained M -estimator and online averaged Projected SGD estimator, as shown in (3). Moreover, we propose a plug-in limiting covariance matrix estimator to facilitate practical inference. We demonstrate the superiority of our method through extensive experiments on benchmark nonlinear problems, constrained generalized linear models, and portfolio allocation problems with both synthetic data and real **Fama-French Portfolios** and **Chicago Air Pollution** data. Note that projection-based methods may not be applicable or competitive in the experiments (e.g., Duchi and Ruan, 2021 considered a linear regression problem with $d = 2$).

1.2 Related literature and contribution

Inspired by the profound success of deterministic SQP in numerical optimization (Nocedal and Wright, 2006), recent research has extended SQP framework to stochastic settings (Na et al.,

2022; 2023; Na, 2025; Curtis et al., 2023; 2024; Berahas et al., 2021; 2023; Fang et al., 2024a; b; 2025). However, existing work has primarily focused on equality-constrained problems. Only Na et al., 2023; Curtis et al., 2024 considered inequality constraints, yet both required strong constraint qualification conditions to circumvent infeasibility issue and a gradually increasing batch size for convergence, making them inapplicable in online settings. Moreover, all above numerical work studied only convergence of SSQP but fell short on uncertainty quantification. In contrast, we not only inherit the strengths of the above numerical methods by allowing adaptive stepsizes, but also significantly refine their designs to fit in *fully online* settings. Furthermore, we bridge the two worlds of optimization and statistics, leveraging the SSQP methods to enable online constrained statistical inference, thereby addressing practical limitations of projection-based methods as stated in Duchi and Ruan, 2021; Davis et al., 2024; Jiang et al., 2025.

This paper also relates to the growing literature on online inference via second-order methods. Bercu et al., 2020; Boyer and Godichon-Baggioni, 2022; Leluc and Portier, 2023; Cénac et al., 2025 studied online Newton (or conditioned SGD) methods for (non)linear regression problems, establishing asymptotic normality for the resulting Newton estimators. Kuang et al., 2025 proposed a consistent covariance estimator for more general sketched Newton methods. Na and Mahoney, 2025 established asymptotic normality for a sketched SSQP method under equality constraints, while Du et al., 2025 leveraged a random scaling technique to construct pivotal statistics. Due to the lack of inequality constraints, none of these works address the challenges of infeasibility and biased step directions discussed in Section 1.1. In fact, to our knowledge, establishing asymptotic distribution of online methods with gradient (and Hessian) momentum is open even in unconstrained settings, while this work establishes such results for inequality-constrained problems. We clearly demonstrate how the (adaptive) stepsize and momentum weight are inter-related: gradient averaging reduces direction bias to facilitate convergence, while past gradient information is also forgotten fast enough to attain optimal asymptotic normality.

1.3 Organization and notation

Section 2 introduces the preliminaries of Problem (1), including the infeasibility issue, constraint qualification, and constraint relaxation. We also establish the primal–dual asymptotic minimax optimality of Problem (1). Section 3 introduces the proposed SSQP method utilizing gradient moving average, and establishes its global almost-sure convergence guarantee. Section 4 presents the asymptotic normality result and proposes a consistent plug-in covariance matrix estimator. Section 5 reports extensive experimental results on both synthetic and real data, followed by conclusions and future work discussions in Section 6.

Throughout the paper, we use $\|\cdot\|$ to denote the ℓ_2 norm for vectors and the spectral norm for matrices. We use boldface lowercase and capital letters (\mathbf{a} and \mathbf{A}) to denote vectors and matrices, respectively. For a positive integer m , $[m] := \{1, \dots, m\}$. For an index set $\mathcal{A} \subseteq [m]$, $\mathcal{A}^- := [m] \setminus \mathcal{A}$ and $|\mathcal{A}|$ denotes its cardinality. For a vector $\mathbf{a} \in \mathbb{R}^m$, $\mathbf{a}_{\mathcal{A}} \in \mathbb{R}^{|\mathcal{A}|}$ denotes the subvector containing only entries with indices in \mathcal{A} . Similarly, for a matrix $\mathbf{A} \in \mathbb{R}^{m \times d}$, $\mathbf{A}_{\mathcal{A}} \in \mathbb{R}^{|\mathcal{A}| \times d}$ denotes the submatrix containing only rows with indices in \mathcal{A} . Without ambiguity, the i -th entry of a vector \mathbf{a} is written as either $[\mathbf{a}]_i$ or a_i . We use $\mathcal{O}(\cdot)$ to denote the standard big- \mathcal{O} notation: $a_k = \mathcal{O}(b_k)$ if $|a_k/b_k| \leq C$ for some constant $C > 0$ and large enough k . We use \odot to denote the Hadamard (element-wise) product. For notational brevity, we write $\mathbf{c}_k := \mathbf{c}(\mathbf{x}_k)$ and $f_k := f(\mathbf{x}_k)$ (similarly for $\nabla \mathbf{c}_k$, etc.) as abbreviations.

2 Preliminary and Asymptotic Optimality

Recall the Lagrangian function $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$ of Problem (1) is defined in (4). Under proper constraint qualifications, a necessary condition for \mathbf{x}^* being a local solution of (1) is the KKT conditions: there exist dual multipliers $\boldsymbol{\lambda}^* \in \mathbb{R}^d$ and $\boldsymbol{\mu}^* = (\boldsymbol{\mu}_1^*, \boldsymbol{\mu}_2^*) \in \mathbb{R}^{2d}$ such that

$$\begin{aligned} \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) &= \mathbf{0}, & \mathbf{c}(\mathbf{x}^*) &= \mathbf{0}, & \boldsymbol{\ell} &\leq \mathbf{x}^* \leq \mathbf{u}, \\ \boldsymbol{\mu}^* &\geq \mathbf{0}, & (\boldsymbol{\mu}_1^*)^\top (\boldsymbol{\ell} - \mathbf{x}^*) &= (\boldsymbol{\mu}_2^*)^\top (\mathbf{x}^* - \mathbf{u}) = 0. \end{aligned} \tag{7}$$

At any \mathbf{x} , we let $\nabla \mathbf{c}(\mathbf{x}) \in \mathbb{R}^{m \times d}$ be the Jacobian of $\mathbf{c}(\mathbf{x})$. We denote active sets for inequalities as

$$\mathcal{A}_\ell(\mathbf{x}) = \{i \in [d] : [\mathbf{x}]_i = [\ell]_i\} \quad \text{and} \quad \mathcal{A}_\mathbf{u}(\mathbf{x}) = \{i \in [d] : [\mathbf{x}]_i = [\mathbf{u}]_i\}. \quad (8)$$

The active constraints include equality and active inequality constraints, $(\mathbf{c}(\mathbf{x}); [\ell - \mathbf{x}]_{\mathcal{A}_\ell(\mathbf{x}); [\mathbf{x} - \mathbf{u}]_{\mathcal{A}_\mathbf{u}(\mathbf{x})})$; and their full Jacobian is given by

$$\mathbf{J}(\mathbf{x}) := \begin{pmatrix} \nabla \mathbf{c}(\mathbf{x}) \\ -\mathbf{I}_{\mathcal{A}_\ell(\mathbf{x})} \\ \mathbf{I}_{\mathcal{A}_\mathbf{u}(\mathbf{x})} \end{pmatrix} \in \mathbb{R}^{(m+|\mathcal{A}_\ell(\mathbf{x})|+|\mathcal{A}_\mathbf{u}(\mathbf{x})|) \times d}, \quad (9)$$

where $\mathbf{I} \in \mathbb{R}^{d \times d}$ is the d -dimensional identity matrix.

2.1 Infeasibility and constraint relaxation

Compared to the projection operator $\text{Proj}_\Omega(\cdot)$ in (5), the constraint linearization in (6) requires only local information and leads to a much simpler subproblem. However, an intrinsic difficulty of linearization is the issue of *infeasibility*: even if the original constraint set $\Omega \neq \emptyset$, the linearized constraints may end up with an empty set.

Consider an illustrative example with $d = m = 1$, where we have $c(x) = x^2 - 4$ and constraints $0 \leq x \leq 3$. We know the constraint set $\Omega = \{x = \pm 2\} \cap [0, 3] = \{2\}$ is nonempty. However, for any $x_k \in (0, 3]$, the linearized equality constraint is $x_k^2 - 4 + 2x_k \bar{\Delta}x_k = 0 \Leftrightarrow \bar{\Delta}x_k = (4 - x_k^2)/(2x_k)$, depicted by the red curve in Figure 1; while the linearized inequality constraint is $-x_k \leq \bar{\Delta}x_k \leq 3 - x_k$, depicted by the shaded region. From Figure 1, we see that the linearized equality and inequality constraints intersect only if $x_k \geq 3 - \sqrt{5}$, suggesting an infeasibility issue for Problem (6) when $x_k \in (0, 3 - \sqrt{5})$. Thus, we need constraint relaxation in this regard.

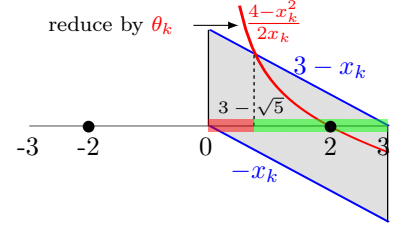


Figure 1: *Insights of our constraint relaxation.*

The constraint relaxation introduced in this paper stems from the observation that the nonlinear curve in Figure 1 can be flattened by a proper scaling factor $\theta_k \in (0, 1]$, so that the curve has more intersections with the shaded region. In particular, we consider the relaxed version of (6):

$$\theta_k \cdot \mathbf{c}_k + \nabla \mathbf{c}_k \bar{\Delta} \mathbf{x}_k = \mathbf{0}, \quad \ell \leq \mathbf{x}_k + \bar{\Delta} \mathbf{x}_k \leq \mathbf{u}. \quad (10)$$

Naturally, we should investigate the conditions under which a relaxation parameter θ_k exists or fails to exist. We find that this is closely tied to the *extended generalized Mangasarian-Fromowitz constraint qualification* (EGMFCQ, [Xu et al., 2015](#), Definition 2.4), which is weaker than *linear independence constraint qualification* (LICQ) that needs to be imposed for asymptotic minimax optimality analysis nevertheless; see [Duchi and Ruan, 2021](#), Assumption B and [Davis et al., 2024](#), Example 2.1 for references. Thus, our relaxation does not comes with stronger conditions.

Definition 2.1 (EGMFCQ v.s. LICQ). *The extended generalized Mangasarian-Fromowitz constraint qualification is said to hold at $\mathbf{x} \in \mathbb{R}^d$ satisfying $\boldsymbol{\ell} \leq \mathbf{x} \leq \mathbf{u}$ if the two conditions are met:*

- (a) *The Jacobian $\nabla \mathbf{c}(\mathbf{x}) \in \mathbb{R}^{m \times d}$ has full row rank.*
- (b) *There exists a vector $\mathbf{z} \in \mathbb{R}^d$ such that*

$$\mathbf{c}(\mathbf{x}) + \nabla \mathbf{c}(\mathbf{x})\mathbf{z} = \mathbf{0}, \quad [\mathbf{z}]_i > 0 \text{ for } i \in \mathcal{A}_\ell(\mathbf{x}), \quad [\mathbf{z}]_i < 0 \text{ for } i \in \mathcal{A}_\mathbf{u}(\mathbf{x}). \quad (11)$$

The linear independence constraint qualification is said to hold at $\mathbf{x} \in \mathbb{R}^d$ if the Jacobian of active constraints $\mathbf{J}(\mathbf{x})$ defined in (9) has full row rank.

Setting $\mathbf{p} = -(\mathbf{c}(\mathbf{x}); \mathbf{1}_{\mathcal{A}_\ell(\mathbf{x})}; \mathbf{1}_{\mathcal{A}_\mathbf{u}(\mathbf{x})})$, we note that LICQ ensures the linear system $\mathbf{J}(\mathbf{x})\mathbf{z} = \mathbf{p}$ has a solution, and such a solution satisfies (11). Hence, LICQ implies EGMFCQ. We further establish the connection between the constraint relaxation in (10) and EGMFCQ.

Theorem 2.2. *For any point $\boldsymbol{\ell} \leq \mathbf{x} \leq \mathbf{u}$, if EGMFCQ holds at \mathbf{x} , there exists a threshold $\bar{\theta} \in (0, 1]$ such that for any $\theta \in [0, \bar{\theta}]$,*

$$\Omega(\mathbf{x}; \theta) := \{\mathbf{z} \in \mathbb{R}^d : \theta \cdot \mathbf{c}(\mathbf{x}) + \nabla \mathbf{c}(\mathbf{x})\mathbf{z} = \mathbf{0}, \boldsymbol{\ell} \leq \mathbf{x} + \mathbf{z} \leq \mathbf{u}\} \neq \emptyset.$$

Conversely, suppose a sequence of points $\boldsymbol{\ell} \leq \mathbf{x}_k \leq \mathbf{u}, \forall k \geq 0$, admits a sequence of "sharp" $\theta_k \in (0, 1]$ in the sense that $\Omega(\mathbf{x}_k, \theta_k) \neq \emptyset$ but $\Omega(\mathbf{x}_k, \min\{2\theta_k, 1\}) = \emptyset$. If there exists a subsequence $\{k_l\}_l$ such that $\lim_{l \rightarrow \infty} \theta_{k_l} = 0$, then EGMFCQ fails to hold at any accumulation point of $\{\mathbf{x}_{k_l}\}_l$.

Theorem 2.2 suggests that the condition of EGMFCQ ensures the existence of $\theta_k > 0$ that makes the linearized constraints (10) feasible, i.e., $\Omega(\mathbf{x}_k; \theta_k) \neq \emptyset$. Note that EGMFCQ or even LICQ is indeed widely imposed in constrained optimization problems, which essentially charac-

terizes the quality of linear approximation to nonlinear constraints (Nocedal and Wright, 2006). Based on Theorem 2.2, we can design the algorithm with constraint relaxation in the following way: at each step, we initialize $\theta_k = 1$ and backtrack reducing θ_k until $\Omega(\mathbf{x}_k; \theta_k) \neq \emptyset$. Note that checking $\Omega(\mathbf{x}_k; \theta_k) \neq \emptyset$ can be efficiently done by solving a convex quadratic program:

$$\min_{\mathbf{z} \in \mathbb{R}^d} \|\theta_k \mathbf{c}_k + \nabla \mathbf{c}_k \mathbf{z}\|^2 \quad \text{s.t.} \quad \boldsymbol{\ell} \leq \mathbf{x}_k + \mathbf{z} \leq \mathbf{u}. \quad (12)$$

We have $\Omega(\mathbf{x}_k; \theta_k) \neq \emptyset$ if and only if the optimal value of the above program is 0. Furthermore, if the selected sequence $\{\theta_k\}$ is not bounded away from zero, Theorem 2.2 suggests that a subsequence of $\{\mathbf{x}_k\}$ will converge to a point that fails EGMFCQ, and certainly fails LICQ. In that case, we should re-initiate the method to explore other local solutions, since any estimation procedures with asymptotic normality, either offline M -estimator or online projection-based estimator, target only solutions that satisfy LICQ (Duchi and Ruan, 2021; Davis et al., 2024).

2.2 Primal-dual asymptotic optimality

In this section, we explore the primal-dual asymptotic minimax optimality of Problem (1). In particular, we target a KKT triplet $(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$ satisfying (7) and investigate what is the best estimator one can expect in the sense of Hájek and Le Cam. To that end, we introduce some assumptions and notations beforehand.

Assumption 2.3 (Strict Complementarity). *We assume $[\boldsymbol{\mu}_1^*]_{\mathcal{A}_\ell(\mathbf{x}^*)} > \mathbf{0}$ and $[\boldsymbol{\mu}_2^*]_{\mathcal{A}_u(\mathbf{x}^*)} > \mathbf{0}$.*

Assumption 2.4 (Second-Order Sufficient Condition). *The Lagrangian Hessian with respect to \mathbf{x} , $\nabla_{\mathbf{x}}^2 \mathcal{L}^* = \nabla_{\mathbf{x}}^2 \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$, is positive definite in the null space of active constraint Jacobian $\{\mathbf{z} \in \mathbb{R}^d : \mathbf{J}(\mathbf{x}^*)\mathbf{z} = \mathbf{0}\}$. In other words, there exists $\omega > 0$ such that*

$$\mathbf{z}^\top \nabla_{\mathbf{x}}^2 \mathcal{L}^* \mathbf{z} \geq \omega \|\mathbf{z}\|^2 \quad \text{for all } \mathbf{z} \in \{\mathbf{z} \in \mathbb{R}^d : \mathbf{J}(\mathbf{x}^*)\mathbf{z} = \mathbf{0}\}.$$

Assumptions 2.3 and 2.4, together with LICQ at \mathbf{x}^* , are standard in the sensitivity analysis of constrained nonlinear optimization (Shapiro, 1990; Bonnans and Shapiro, 2000), and are also imposed for minimax optimality analysis of the primal solution \mathbf{x}^* in the literature (Duchi and

Ruan, 2021; Davis et al., 2024). These assumptions essentially require that the KKT triplet $(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$ is a strict local solution of (1), admitting a local convexity structure in the sense that the Lagrangian is strongly convex in the tangent space of the constraints manifold at \mathbf{x}^* .

For any $\mathbf{x} \in \mathbb{R}^d$ and $\boldsymbol{\mu} = (\boldsymbol{\mu}_1, \boldsymbol{\mu}_2) \in \mathbb{R}^{2d}$, let us denote active dual components by $\boldsymbol{\mu}_{\mathcal{A}(\mathbf{x})} = ([\boldsymbol{\mu}_1]_{\mathcal{A}_\ell(\mathbf{x})}, [\boldsymbol{\mu}_2]_{\mathcal{A}_u(\mathbf{x})}) \in \mathbb{R}^{|\mathcal{A}_\ell(\mathbf{x})|+|\mathcal{A}_u(\mathbf{x})|}$ and inactive dual components by $\boldsymbol{\mu}_{\mathcal{A}(\mathbf{x})^-} = ([\boldsymbol{\mu}_1]_{\mathcal{A}_\ell(\mathbf{x})^-}, [\boldsymbol{\mu}_2]_{\mathcal{A}_u(\mathbf{x})^-}) \in \mathbb{R}^{2d-|\mathcal{A}_\ell(\mathbf{x})|-|\mathcal{A}_u(\mathbf{x})|}$. For notational simplicity, we write $\boldsymbol{\mu}_{\mathcal{A}^*}^* = \boldsymbol{\mu}_{\mathcal{A}(\mathbf{x}^*)}^*$ and note that $\boldsymbol{\mu}_{(\mathcal{A}^*)^-}^* = \boldsymbol{\mu}_{\mathcal{A}(\mathbf{x}^*)^-}^* = \mathbf{0}$ by the KKT conditions in (7). We also define the Lagrangian Hessian matrix with respect to $(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}_{\mathcal{A}(\mathbf{x})})$, a covariance matrix, and a sandwich matrix as

$$\mathbf{H}^* = \begin{pmatrix} \nabla_{\mathbf{x}}^2 \mathcal{L}^* & (\mathbf{J}^*)^\top \\ \mathbf{J}^* & \mathbf{0} \end{pmatrix}, \quad \boldsymbol{\Sigma}^* = \begin{pmatrix} \text{Cov}(\nabla F(\mathbf{x}^*; \zeta)) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}, \quad \boldsymbol{\Omega}^* = (\mathbf{H}^*)^{-1} \boldsymbol{\Sigma}^* (\mathbf{H}^*)^{-1}. \quad (13)$$

The non-singularity of \mathbf{H}^* is ensured by LICQ (Definition 2.1) and SOSC (Assumption 2.4).

The lower-bound analysis closely relates to the sensitivity analysis of Problem (1) with respect to the distribution \mathcal{P} . We will explore how the primal-dual solution as a function of \mathcal{P} , $(\mathbf{x}^*(\mathcal{P}), \boldsymbol{\lambda}^*(\mathcal{P}), \boldsymbol{\mu}^*(\mathcal{P}))$, varies when \mathcal{P} is perturbed slightly. As such, we measure the perturbation size via the ϕ -divergence, defined as

$$D_\phi(\mathcal{P}', \mathcal{P}) = \int \phi \left(\frac{d\mathcal{P}'}{d\mathcal{P}} \right) d\mathcal{P},$$

where $\phi : (0, \infty) \rightarrow \mathbb{R}$ is a C^3 -smooth convex function satisfying $\phi(1) = 0$. We define an admissible neighborhood $\mathcal{B}(\mathcal{P}; \epsilon)$ of \mathcal{P} to consist of all probability distributions \mathcal{P}' with $D_\phi(\mathcal{P}', \mathcal{P}) \leq \epsilon$, such that there exists a primal-dual solution $(\mathbf{x}^*(\mathcal{P}'), \boldsymbol{\lambda}^*(\mathcal{P}'), \boldsymbol{\mu}^*(\mathcal{P}'))$ to the perturbed Problem (1). It is well known that the inactive constraint set $\mathcal{A}(\mathbf{x}^*(\mathcal{P}))^-$ is preserved under small perturbations (cf. Shapiro, 1990 Theorem 1). That is, $\mathcal{A}_\ell(\mathbf{x}^*(\mathcal{P}'))^- = \mathcal{A}_\ell(\mathbf{x}^*(\mathcal{P}))^-$ and $\mathcal{A}_u(\mathbf{x}^*(\mathcal{P}'))^- = \mathcal{A}_u(\mathbf{x}^*(\mathcal{P}))^-$ for ϵ small enough; and we trivially have $\boldsymbol{\mu}_{(\mathcal{A}^*)^-}^*(\mathcal{P}) = [\boldsymbol{\mu}^*(\mathcal{P})]_{\mathcal{A}(\mathbf{x}^*(\mathcal{P}))^-} = \mathbf{0}$ in the neighborhood of \mathcal{P} . Below, we focus on the nontrivial part $\mathbf{w}^*(\mathcal{P}) = (\mathbf{x}^*(\mathcal{P}), \boldsymbol{\lambda}^*(\mathcal{P}), \boldsymbol{\mu}_{\mathcal{A}^*}^*(\mathcal{P}))$.

The next theorem asymptotically lower bounds the performance of any estimation procedure for finding the solution to (1), by considering adversarially chosen small perturbations of the unperturbed distribution \mathcal{P} .

Theorem 2.5 (Local primal-dual minimax optimality). *Suppose Assumptions 2.3, 2.4, and LICQ hold at $(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$. Let $\ell : \mathbb{R}^{d+m+|\mathcal{A}_\ell^*|+|\mathcal{A}_u^*|} \rightarrow [0, \infty)$ be any symmetric, quasiconvex, and lower semicontinuous function, and let $\mathbf{w}_k(\mathcal{P}') : \zeta_{1:k} \rightarrow \mathbb{R}^{d+m+|\mathcal{A}_\ell^*|+|\mathcal{A}_u^*|}$ be a sequence of estimators based on k samples $\zeta_{1:k} \stackrel{iid}{\sim} \mathcal{P}'$. Then, the following inequality holds:*

$$\lim_{c \rightarrow +\infty} \liminf_{k \rightarrow +\infty} \sup_{\mathcal{P}' \in \mathcal{B}(\mathcal{P}; c/k)} \mathbb{E}_{\mathcal{P}'_k} \left[\ell \left(\sqrt{k} (\mathbf{w}_k(\mathcal{P}') - \mathbf{w}^*(\mathcal{P}')) \right) \right] \geq \mathbb{E} [\ell(Z)],$$

where $\mathbb{E}_{\mathcal{P}'_k}[\cdot]$ is taken over k iid samples $\zeta_{1:k}$ and $Z \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Omega}^*)$.

Compared with (Davis et al., 2024, Theorem 3.2) and (Duchi and Ruan, 2021, Theorem 1), our result jointly captures both the primal solution and the active dual solution (since inactive dual solution remains zero). Notably, following a similar calculation as in (Na and Mahoney, 2025, Remark 5.9), we observe that the marginal covariance of $\boldsymbol{\Omega}^*$ with respect to the primal solution matches those in the above works (i.e., (3)). Thus, Theorem 2.5 provides a complete and sharp characterization of local asymptotic optimality for constrained estimation problems.

3 SSQP with Polyak’s Momentum

In this section, we introduce the proposed online SSQP estimation procedure for Problem (1), which leverages the constraint relaxation idea discussed in Section 2.1 and employs gradient (and Hessian) moving average schemes based on Polyak’s momentum to debias the step direction.

3.1 Estimation procedure

The estimation procedure involves tuning parameters $\tau \in (0, 1)$, $\psi \geq 0$, $p \geq 1$, and one stepsize and two averaging weight sequences $\alpha_k, \beta_k, \gamma_k > 0$. Given the current k -th primal-dual iterate $(\mathbf{x}_k, \boldsymbol{\lambda}_k, \boldsymbol{\mu}_k)$ with $\boldsymbol{\mu}_k = (\boldsymbol{\mu}_{1,k}, \boldsymbol{\mu}_{2,k})$, we initialize $\theta_k = 1$ and perform

$$\theta_k \leftarrow \tau \theta_k \quad \text{until} \quad \text{convex problem (12) has an optimal value 0.} \quad (14)$$

The above scheme is supposed to terminate in finite time by Theorem 2.2. In particular, if θ_k falls below an infinitesimal threshold, our theory implies that \mathbf{x}_k is approaching an undesirable

point at which even EGMFCQ (let alone LICQ) does not hold.

With the selected θ_k from (14), we then draw a single sample $\zeta_k \sim \mathcal{P}$ and estimate the objective gradient and Hessian as $\nabla F(\mathbf{x}_k; \zeta_k)$ and $\nabla^2 F(\mathbf{x}_k; \zeta_k)$, respectively. Instead of using these noisy estimates directly, we perform gradient and Hessian moving averages:

$$\bar{\mathbf{g}}_k = (1 - \beta_k)\bar{\mathbf{g}}_{k-1} + \beta_k \nabla F(\mathbf{x}_k; \zeta_k) \quad \text{and} \quad \bar{\mathbf{Q}}_k = (1 - \gamma_k)\bar{\mathbf{Q}}_{k-1} + \gamma_k \nabla^2 F(\mathbf{x}_k; \zeta_k). \quad (15)$$

The Lagrangian Hessian with respect to \mathbf{x} depends on the constraint curvature information, which we estimate as

$$\bar{\mathbf{B}}_k = \bar{\mathbf{Q}}_k + \sum_{i=1}^m [\boldsymbol{\lambda}_k]_i \nabla^2 \mathbf{c}_i(\mathbf{x}_k) + \boldsymbol{\Delta}_k, \quad (16)$$

where $\boldsymbol{\Delta}_k \in \mathbb{R}^{d \times d}$ is a Hessian modification term. In fact, if only the global convergence¹ of the estimation procedure is of interest, we do not require any second-order information and can simply set $\bar{\mathbf{B}}_k = \mathbf{I}$. As such, the importance of the two averaging weight sequences, β_k and γ_k , is different: the choice of β_k is more delicate than that of γ_k , since first-order information plays a more crucial role in debiasing the step and facilitating asymptotic normality (cf. Theorem 4.6).

With the above estimators, we construct an SSQP subproblem similar to (6), but using averaged gradient (and Hessian) estimators along with the constraint relaxation:

$$\min_{\bar{\Delta}\mathbf{x}_k \in \mathbb{R}^d} \bar{\mathbf{g}}_k^\top \bar{\Delta}\mathbf{x}_k + \frac{1}{2} \bar{\Delta}\mathbf{x}_k^\top \bar{\mathbf{B}}_k \bar{\Delta}\mathbf{x}_k, \quad \text{s.t.} \quad \theta_k \mathbf{c}_k + \nabla \mathbf{c}_k \bar{\Delta}\mathbf{x}_k = \mathbf{0}, \quad \boldsymbol{\ell} \leq \mathbf{x}_k + \bar{\Delta}\mathbf{x}_k \leq \mathbf{u}. \quad (17)$$

Since solving nonconvex inequality-constrained quadratic program without any constraint qualifications can be NP-hard (Burer and Letchford, 2009), we require the Hessian modification $\boldsymbol{\Delta}_k$ to convexify $\bar{\mathbf{Q}}_k + \sum_{i=1}^m [\boldsymbol{\lambda}_k]_i \nabla^2 \mathbf{c}_i(\mathbf{x}_k)$ such that $\bar{\mathbf{B}}_k$ is positive definite (cf. Assumption 3.5). We slightly abuse notation by letting $\bar{\Delta}\mathbf{x}_k$ denote the primal solution of (17), and $(\bar{\boldsymbol{\lambda}}_k^{\text{sub}}, \bar{\boldsymbol{\mu}}_k^{\text{sub}})$ with $\bar{\boldsymbol{\mu}}_k^{\text{sub}} = (\bar{\boldsymbol{\mu}}_{1,k}^{\text{sub}}, \bar{\boldsymbol{\mu}}_{2,k}^{\text{sub}})$ denote any (one of) dual solutions. We then update the iterate as follows:

$$(\mathbf{x}_{k+1}, \boldsymbol{\lambda}_{k+1}, \boldsymbol{\mu}_{k+1}) = (\mathbf{x}_k, \boldsymbol{\lambda}_k, \boldsymbol{\mu}_k) + \bar{\alpha}_k \cdot (\bar{\Delta}\mathbf{x}_k, \bar{\boldsymbol{\lambda}}_k^{\text{sub}} - \boldsymbol{\lambda}_k, \bar{\boldsymbol{\mu}}_k^{\text{sub}} - \boldsymbol{\mu}_k), \quad (18)$$

where $\bar{\alpha}_k > 0$ is an adaptive, potentially random stepsize that satisfies a safeguard condition:

¹ In our context of nonlinear problems, global convergence refers to convergence to a stationary point from *any* initialization, as opposed to convergence to a global solution, which is generally unachievable without specific problem structures. That said, these two notions coincide in the case of convex problems.

$\alpha_k \leq \bar{\alpha}_k \leq \alpha_k + \psi \alpha_k^p$ for the tuning parameters $\psi > 0$ and $p \geq 1$.

Remark 3.1 (Uniqueness of subproblem dual solution). *We should emphasize that the dual solution $(\bar{\lambda}_k^{sub}, \bar{\mu}_k^{sub})$ of the strongly convex SSQP subproblem (17) is not necessarily unique, as we do not impose strong constraint qualifications. In fact, in global convergence analysis, when evaluating the KKT residual at the iterate \mathbf{x}_k , we do not rely directly on the dual iterates (λ_k, μ_k) or $(\bar{\lambda}_k^{sub}, \bar{\mu}_k^{sub})$. Note that the KKT condition is defined solely in terms of the primal variable (see (7)); therefore, we only desire the existence of dual multipliers, possibly not unique, under which the KKT residual converges to zero. In other words, even though the dual variables of individual SSQP subproblems may not be unique, the global convergence of the KKT residual with respect to the primal iterates can still be established.*

On the other hand, our local analysis in Section 4 targets a local solution \mathbf{x}^* at which LICQ holds, as required by minimax optimality guarantee in Section 2. In that case, we can examine the local convergence behavior of the dual variables (λ_k, μ_k) , since the SSQP subproblem admits a unique dual solution $(\bar{\lambda}_k^{sub}, \bar{\mu}_k^{sub})$ under LICQ. In particular, as \mathbf{x}_k approaches \mathbf{x}^* satisfying LICQ, the inactive constraints at \mathbf{x}^* for the original problem (1) remain inactive at $\mathbf{x}_k + \bar{\Delta} \mathbf{x}_k$ for the SSQP subproblem (17), implying (17) has a unique dual solution $(\bar{\lambda}_k^{sub}, \bar{\mu}_k^{sub})$. Therefore, our local convergence analysis on the dual variable in Section 4 (cf. Lemma 4.4) is well justified.

Remark 3.2 (Adaptive stepsize). *Compared to online projection-based methods (Duchi and Ruan, 2021; Davis et al., 2024), we allow using a random adaptive stepsize $\bar{\alpha}_k$, as long as $\bar{\alpha}_k$ that may depend on the step direction is well controlled within the interval $[\alpha_k, \alpha_k + \psi \alpha_k^p]$. Here, the interval length $\psi \alpha_k^p$ measures the degree of adaptivity. Setting $\psi = 0$ corresponds to deterministic stepsize. The adaptive designs have been proposed for SSQP schemes in numerical optimization literature (Curtis et al., 2024; Berahas et al., 2021; Fang et al., 2024a; Na and Mahoney, 2025), where the stepsize is chosen via either line-search or trust-region strategies. All the proposed stepsize rules in those works satisfy the safeguard condition. Motivated by promising empirical*

performance of adaptive methods, we adopt a similar design of stepsize in our estimation procedure. This additional flexibility, however, makes the statistical inference analysis more involved, as it requires controlling the effect of adaptivity on the limiting distribution.

Remark 3.3 (Bias in step direction). Compared to SSQP designs for equality-constrained problems, a crucial challenge posed by inequality constraints is the bias in the step direction estimation.

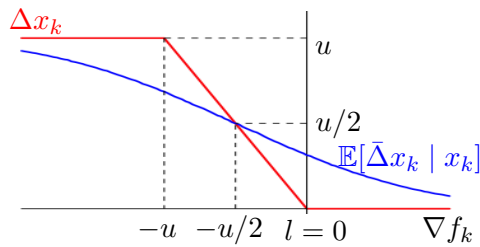


Figure 2: Illustration of bias in step direction estimation.

Consider the following example on \mathbb{R} , where we let $\bar{B}_k = 1$ and suppress equality constraints in (6): $\min_{\bar{\Delta}x_k \in \mathbb{R}} 0.5\bar{\Delta}x_k^2 + \nabla F(x_k; \zeta_k)\bar{\Delta}x_k$ subject to $0 = l \leq \bar{\Delta}x_k \leq u$. Without inequality constraints, this reduces to SGD with $\Delta x_k = -\nabla f_k = \mathbb{E}[-\nabla F(x_k; \zeta_k) | x_k] = \mathbb{E}[\bar{\Delta}x_k | x_k]$, as long as $\nabla F(x_k; \zeta_k)$ is unbiased. However, with inequality constraints, we have $\Delta x_k = -\text{Project}(\nabla f_k, [-u, -l])$ shown as the red line in Figure 2, differing from

$$\mathbb{E}[\bar{\Delta}x_k | x_k] = -\mathbb{E}[\text{Project}(\nabla F(x_k; \zeta_k), [-u, -l]) | x_k],$$

shown as the blue line in Figure 2 with Gaussian noise $\nabla F(x_k; \zeta_k) \sim \mathcal{N}(\nabla f_k, \sigma^2)$. In other words, inequality constraints truncate the gradient noise in a nonlinear manner, such that an unbiased gradient estimate does not necessarily yield an unbiased step direction estimate. To address this issue, we leverage momentum-based averaging estimators in (15). We show that $\bar{\Delta}x_k$ with averaged gradient estimators satisfies $\mathbb{E}[\|\bar{\Delta}x_k - \Delta x_k\|^2] \rightarrow 0$ as $k \rightarrow \infty$, i.e., it is asymptotically unbiased (cf. Lemmas D.1 and D.2). One significant contribution of this paper is to further show that gradient momentum simultaneously preserves asymptotic optimality while debiasing step directions, as long as the weight sequence β_k is chosen properly according to the stepsize α_k .

3.2 Global almost-sure convergence

In this section, we establish global almost-sure convergence for the proposed SSQP estimation procedure in Section 3.1. We show that the KKT residual of \mathbf{x}_k converges to zero from any initialization. Our analysis generalizes existing SSQP literature on numerical optimization (Na

et al., 2023; Curtis et al., 2024) and relaxes restrictive constraint qualification conditions.

We define an ℓ_2 -penalized Lyapunov function to quantify the progress toward stationarity:

$$\varphi_\rho(\mathbf{x}) = f(\mathbf{x}) + \rho \|\mathbf{c}(\mathbf{x})\| \quad \text{for } \rho > 0. \quad (19)$$

The penalty term biases the feasibility error of equality constraints, while inequality constraints are not penalized since SSQP iterates always stay within the box constraints (cf. Lemma C.3). In contrast to unconstrained problems, where methods aim to reduce the objective function f , reducing f alone is not sufficient to justify the iteration progress toward stationarity in constrained problems, since a small objective value may be achieved by severely violating the constraints. We will show that our method decreases the Lyapunov function (19) for sufficiently large $\rho > 0$, thereby vanishing the KKT residual, which is defined at the primal-dual triplet by

$$\mathbf{R}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \begin{pmatrix} \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \\ \mathbf{c}(\mathbf{x}) \\ \boldsymbol{\mu}_1 \odot (\boldsymbol{\ell} - \mathbf{x}) \\ \boldsymbol{\mu}_2 \odot (\mathbf{x} - \mathbf{u}) \end{pmatrix}. \quad (20)$$

Here, \odot denotes the Hadamard elementwise product and each entry of $\mathbf{R}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$ corresponds to one KKT condition in (7). Again, we do not measure violation of conditions $\boldsymbol{\mu} \geq 0$ and $\boldsymbol{\ell} \leq \mathbf{x} \leq \mathbf{u}$ since they hold at every step. Thus, $\mathbf{R}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \mathbf{0}$ suggests we arrive at stationarity.

3.2.1 Assumptions

Now, we state the assumptions for global analysis. Assumption 3.4 is tied to our constraint relaxation technique, which has been justified by Theorem 2.2 under EGMFCQ (weaker than LICQ). Assumptions 3.5 and 3.6 are standard regularity conditions commonly imposed in the literature.

Assumption 3.4. *There exists a deterministic constant $\bar{\theta} \in (0, 1]$ such that the linearized constraints (10) are feasible for any relaxation parameter $\theta_k \in [0, \bar{\theta}]$. That is, $\Omega(\mathbf{x}_k; \theta_k) \neq \emptyset$ for $\theta_k \in [0, \bar{\theta}]$.*

Assumption 3.4 requires the relaxation parameter θ_k to be bounded away from zero to ensure subproblem feasibility. By Theorem 2.2, if θ_k is not bounded below, any accumulation point of the iterates \mathbf{x}_k fails to satisfy EGMFCQ and, certainly, LICQ. In that case, it becomes intractable

to perform statistical inference for such an accumulation point even using offline M -estimation or online projected methods, since asymptotic normality result in (3) of both methods, as well as primal-dual minimax optimality in Theorem 2.5, crucially relies on LICQ.

Assumption 3.5. *We assume $\nabla f, \nabla \mathbf{c}$ are Lipschitz continuous; that is, for all $\mathbf{x}, \mathbf{y} \leq \mathbf{u}$,*

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq \kappa_{\nabla f} \|\mathbf{x} - \mathbf{y}\|, \quad \|\nabla \mathbf{c}(\mathbf{x}) - \nabla \mathbf{c}(\mathbf{y})\| \leq \kappa_{\nabla \mathbf{c}} \|\mathbf{x} - \mathbf{y}\|$$

for some constants $\kappa_{\nabla f}, \kappa_{\nabla \mathbf{c}} > 0$. Furthermore, the Hessian modification term $\Delta_k \in \mathbb{R}^{d \times d}$ in (16) ensures $\bar{\mathbf{B}}_k$ satisfy $\kappa_1 \mathbf{I} \preceq \bar{\mathbf{B}}_k \preceq \kappa_2 \mathbf{I}$ for some $0 < \kappa_1 \leq \kappa_2$. In addition, we assume the dual multipliers of the SQP subproblem (17) with the exact true gradient ∇f_k , denoted by $(\boldsymbol{\lambda}_k^{\text{sub}}, \boldsymbol{\mu}_k^{\text{sub}})$, are bounded; that is, there exists $M_{\text{dual}} > 0$ such that $\max\{\|\boldsymbol{\lambda}_k^{\text{sub}}\|, \|\boldsymbol{\mu}_k^{\text{sub}}\|\} \leq M_{\text{dual}}$.

The Lipschitz condition is standard and naturally holds within a closed bounded feasible set (Curtis et al., 2024; Berahas et al., 2023). The condition on $\bar{\mathbf{B}}_k$ ensures the subproblem (17) is strongly convex and admits a unique solution; this requirement can be readily relaxed to the SOSC condition, as stated in Assumption 2.4, when a strong constraint qualification condition, e.g., LICQ, is imposed (see discussion after Assumption 4.1 in Section 4). The boundedness of dual multipliers in (deterministic) SQP subproblems is a conventional assumption in constrained optimization. We show in Theorem E.1 (Appendix E) that if the accumulation point of the iteration sequence \mathbf{x}_k is feasible and satisfies EGMFCQ, then the multipliers indeed remain bounded. This result provides a theoretical justification for the boundedness assumption.

Assumption 3.6. *We assume, for any $k \geq 0$, $\nabla F(\mathbf{x}_k; \zeta_k)$ is an unbiased estimate of ∇f_k with bounded variance. In particular, $\mathbb{E}_k[\nabla F(\mathbf{x}_k; \zeta_k)] = \nabla f_k$ and $\mathbb{E}_k[\|\nabla F(\mathbf{x}_k; \zeta_k) - \nabla f_k\|^2] \leq \sigma_g^2$ for some $\sigma_g^2 > 0$. Here, $\mathbb{E}_k[\cdot]$ denotes the conditional expectation given \mathbf{x}_k .*

The unbiasedness and bounded variance of stochastic gradients are standard assumptions in both unconstrained and constrained stochastic optimization (Polyak and Juditsky, 1992; Chen et al., 2020; Curtis et al., 2024; Na and Mahoney, 2025).

3.2.2 Almost-sure convergence

The Lyapunov function φ_ρ in (19) serves as a surrogate of the original constrained problem, providing a unified measure of objective value and constraint violation. Drawing intuition from unconstrained optimization, a decrease in the Lyapunov function can be interpreted as progress toward stationarity of constrained problem. Thus, we analyze the decrease of the Lyapunov function $\varphi_\rho(\mathbf{x})$ along the sequence of iterates. As such, we consider a local quadratic model of $\varphi_\rho(\mathbf{x})$:

$$\varphi_\rho^{\text{loc}}(\mathbf{x}, \Delta\mathbf{x}, \mathbf{B}) := f(\mathbf{x}) + \nabla f(\mathbf{x})^\top \Delta\mathbf{x} + \frac{1}{2} \Delta\mathbf{x}^\top \mathbf{B} \Delta\mathbf{x} + \rho \|\mathbf{c}(\mathbf{x}) + \nabla \mathbf{c}(\mathbf{x}) \Delta\mathbf{x}\|, \quad (21)$$

and its difference measures the reduction of φ_ρ at \mathbf{x} along the step $\Delta\mathbf{x}$:

$$\Delta\varphi_\rho^{\text{loc}}(\mathbf{x}, \Delta\mathbf{x}, \mathbf{B}) := \varphi_\rho^{\text{loc}}(\mathbf{x}, \mathbf{0}, \mathbf{B}) - \varphi_\rho^{\text{loc}}(\mathbf{x}, \Delta\mathbf{x}, \mathbf{B}). \quad (22)$$

The next lemma establishes a reduction argument for the local model difference $\Delta\varphi_\rho^{\text{loc}}(\mathbf{x}, \Delta\mathbf{x}, \mathbf{B})$, showing that the reduction at each step is proportional to both the step magnitude $\|\Delta\mathbf{x}_k\|^2$ and the constraint violation $\|\mathbf{c}_k\|$. Recall from Section 3.1 that $\bar{\Delta}\mathbf{x}_k$ and $(\bar{\boldsymbol{\lambda}}_k^{\text{sub}}, \bar{\boldsymbol{\mu}}_k^{\text{sub}})$ denote the primal and dual solutions of the subproblem (17), while $\Delta\mathbf{x}_k$ and $(\boldsymbol{\lambda}_k^{\text{sub}}, \boldsymbol{\mu}_k^{\text{sub}})$ denote their counterparts when replacing gradient estimate $\bar{\mathbf{g}}_k$ by the true gradient ∇f_k in (17).

Lemma 3.7. *Under Assumptions 3.4 and 3.5, for any $\nu \in (0, 1)$, as long as $\rho \geq \bar{\rho} := \frac{M_{\text{dual}}}{(1-\nu)\tau\theta}$,*

$$\Delta\varphi_\rho^{\text{loc}}(\mathbf{x}_k, \Delta\mathbf{x}_k, \bar{\mathbf{B}}_k) \geq \frac{1}{2} \Delta\mathbf{x}_k^\top \bar{\mathbf{B}}_k \Delta\mathbf{x}_k + \nu\rho\theta_k \|\mathbf{c}_k\| \quad \text{for all } k \geq 0.$$

Since $\varphi_\rho^{\text{loc}}(\mathbf{x}, \Delta\mathbf{x}, \mathbf{B})$ closely approximates the Lyapunov function $\varphi_\rho(\mathbf{x})$, the above result implies a decrease in φ_ρ itself. Consequently, the convergence of the iterates toward a stationary point is closely linked to the vanishing of $\|\Delta\mathbf{x}_k\|$ and $\|\mathbf{c}_k\|$.

Theorem 3.8. *Under Assumptions 3.4, 3.5, and 3.6, we specify the stepsize and averaging weight sequences as $\alpha_k = \iota_1(k+1)^{-b_1}$ and $\beta_k = \iota_2(k+1)^{-b_2}$, for constants $\iota_1, \iota_2 > 0$ and exponents b_1, b_2 satisfying $b_1 \in (0.75, 1]$ and $b_2 \in (2 - 2b_1, 2b_1 - 1)$. Then,*

$$\liminf_{k \rightarrow \infty} \{\|\Delta\mathbf{x}_k\| + \|\mathbf{c}_k\|\} = 0 \quad \text{almost surely.}$$

Finally, we establish the convergence of the KKT residual. Note that the KKT condition (7)

is defined in terms of the primal variable only; without LICQ, the associated dual multipliers of Problem (1) are not unique (see Remark 3.1). Thus, it suffices to establish the existence of (possibly nonunique) dual multipliers such that, together with \mathbf{x}_k , the KKT residual converges to zero. As such, given \mathbf{x}_k of SSQP, we consider the multipliers defined by a least-squares problem:

$$\begin{aligned} \min_{\boldsymbol{\lambda}, \boldsymbol{\mu}} \quad & \|\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}_k, \boldsymbol{\lambda}, \boldsymbol{\mu})\|^2 + \|\boldsymbol{\mu}_1 \odot (\boldsymbol{\ell} - \mathbf{x}_k)\|^2 + \|\boldsymbol{\mu}_2 \odot (\mathbf{x}_k - \mathbf{u})\|^2, \\ \text{s.t.} \quad & \boldsymbol{\mu} \geq \mathbf{0}. \end{aligned} \tag{23}$$

We denote (one of) the optimal solution to (23) by $(\boldsymbol{\lambda}_k^*, \boldsymbol{\mu}_k^*)$. The next theorem shows that the SSQP iterate \mathbf{x}_k , coupled with $(\boldsymbol{\lambda}_k^*, \boldsymbol{\mu}_k^*)$, yields a vanishing KKT residual.

Theorem 3.9. *Under the same conditions as in Theorem 3.8, we have*

$$\lim_{k \rightarrow \infty} \mathbf{R}(\mathbf{x}_k, \boldsymbol{\lambda}_k^*, \boldsymbol{\mu}_k^*) = \mathbf{0} \quad \text{almost surely.}$$

We note that in global analysis, two algorithmic components, the moving averaged second-order information of $\bar{\mathbf{B}}_k$ in (16) and the dual update in (18), are not essential; one may simply set $\bar{\mathbf{B}}_k = \mathbf{I}$ and suppress the dual update if only the consistency of the estimation is of interest. However, these components are crucial for local inference analysis, as they ensure the minimax optimality of the SSQP method. On the other hand, our method involves two key sequences: α_k for the primal updates and β_k for gradient averaging, which exhibit *competing* effects. In particular, a larger α_k accelerates the movement of the iterates but may impede the stabilization of the averaged gradients, while a smaller β_k improves gradient accuracy but slows the adaptation of both the gradient sequence and the primal iterates. These coupled dynamics highlight the importance of carefully balancing the decay rates of α_k and β_k , as specified in Theorem 3.8, to ensure simultaneous convergence of the averaged gradients and the primal iterates.

4 Asymptotic Normality and Inference

We now set the stage to present the local asymptotic analysis of our SSQP method in Section 3.1.

We establish that the joint primal-dual variables, with proper stepsize control and averaging

weight sequences $\{\alpha_k, \beta_k, \gamma_k\}$, converge in distribution to a normal variable with limiting covariance $\mathbf{\Omega}^*$ defined in (13). The resulting limiting distribution attains the minimax lower bound established in Theorem 2.5, thereby demonstrating that the proposed method is asymptotically optimal in the sense of Hájek and Le Cam. Furthermore, we provide a consistent plug-in estimator for the limiting covariance matrix, which enables practical online statistical inference, including hypothesis testing and confidence interval construction.

4.1 Asymptotic normality of the last iterate

We first introduce a local assumption on the limiting point \mathbf{x}^* to which the sequence $\{\mathbf{x}_k\}$ converges. This assumption is standard in both the lower-bound analysis (cf. Theorem 2.5) and the upper-bound analyses of offline M -estimation and online projection-based methods (Duchi and Ruan, 2021; Davis et al., 2024).

Assumption 4.1. *The iteration sequence \mathbf{x}_k converges almost surely to a strict local solution \mathbf{x}^* , which satisfies: (1) LICQ in Definition 2.1 (implying the uniqueness of the dual solution $(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$); (2) strict complementarity in Assumption 2.3; and (3) second-order sufficient condition (SOSC) in Assumption 2.4.*

With Assumption 4.1, we now elaborate on the role of the Hessian matrix $\bar{\mathbf{B}}_k$. Note that the global analysis does not rely on any second-order information; hence, one may simply set $\bar{\mathbf{B}}_k = \mathbf{I}$ in the SSQP subproblem (17) and all results in Section 3.2 still hold. In contrast, for the local analysis, where \mathbf{x}_k is sufficiently close to a strict local solution \mathbf{x}^* that satisfies LICQ and SOSC, we expect $\bar{\mathbf{B}}_k$ to serve as a good approximation to the Lagrangian Hessian $\nabla_{\mathbf{x}}^2 \mathcal{L}_k$, which motivates the Hessian averaging update in (15). In fact, the conditions of LICQ and SOSC at \mathbf{x}^* extend to a local neighborhood. Thus, as long as \mathbf{x}_k is close to \mathbf{x}^* (which is the case) and $\bar{\mathbf{B}}_k$ is close to $\nabla_{\mathbf{x}}^2 \mathcal{L}_k$, these conditions apply to the SSQP subproblem (17) and further imply the *local* uniqueness of its primal-dual solution. Since no constraint qualification is imposed at \mathbf{x}_k globally, the Hessian information in the subproblem is not helpful for determining its solutions in the global

analysis. Note that Hessian approximation in the local analysis is necessary even for deterministic (quasi-)second-order methods, so called the Dennis–Moré condition; see [Dennis and Moré, 1974](#), ([Nocedal and Wright, 2006](#), Chapter 18), and ([Liu and Yuan, 2000](#), Section 5) for details.

To precisely state the above argument, let us consider any $\epsilon > 0$ such that

$$0 < \epsilon \leq 0.5 \min \{ [\mathbf{u} - \boldsymbol{\ell}]_i, [\mathbf{x}^* - \boldsymbol{\ell}]_j, [\mathbf{u} - \mathbf{x}^*]_j : i \in \mathcal{A}^*, j \in (\mathcal{A}^*)^- \}, \quad (24)$$

where we recall from Section 2 that $\mathcal{A}^* = \mathcal{A}_{\boldsymbol{\ell}}^* \cup \mathcal{A}_{\mathbf{u}}^*$ denotes the index set of active inequality constraints and $(\mathcal{A}^*)^- = [d] \setminus \mathcal{A}^*$ denotes its inactive complement. We define two sets at \mathbf{x}_k ,

$$\mathcal{A}_{\boldsymbol{\ell},k}(\epsilon) = \{i : [\mathbf{x}_k - \boldsymbol{\ell}]_i \leq \epsilon\} \quad \text{and} \quad \mathcal{A}_{\mathbf{u},k}(\epsilon) = \{i : [\mathbf{u} - \mathbf{x}_k]_i \leq \epsilon\}, \quad (25)$$

as our guess of true active sets $\mathcal{A}_{\boldsymbol{\ell}}^*$ and $\mathcal{A}_{\mathbf{u}}^*$, and let $\mathcal{A}_k(\epsilon) := \mathcal{A}_{\boldsymbol{\ell},k}(\epsilon) \cup \mathcal{A}_{\mathbf{u},k}(\epsilon)$. We will show next that $\mathcal{A}_{\boldsymbol{\ell},k}(\epsilon) = \mathcal{A}_{\boldsymbol{\ell}}^*$, $\mathcal{A}_{\mathbf{u},k}(\epsilon) = \mathcal{A}_{\mathbf{u}}^*$, and both also coincide with the active set of SSQP subproblem (17) when k is large enough, indicating that our method successfully identifies the true active set. Now, since SOSC holds at \mathbf{x}^* (cf. Assumption 2.4), given $\mathcal{A}_{\boldsymbol{\ell},k}(\epsilon)$ and $\mathcal{A}_{\mathbf{u},k}(\epsilon)$ in (25), we require the regularizer $\boldsymbol{\Delta}_k$ to convexify the estimate of $\nabla_{\mathbf{x}}^2 \mathcal{L}_k$, i.e., $\bar{\mathbf{Q}}_k + \sum_{i=1}^m [\boldsymbol{\lambda}_k]_i \nabla^2 \mathbf{c}_i(\mathbf{x}_k)$ in (16), to ensure $\mathbf{z}^\top \bar{\mathbf{B}}_k \mathbf{z} \geq \omega \|\mathbf{z}\|^2, \forall \mathbf{z} \in \{\mathbf{z} : \mathbf{J}_k(\epsilon) \mathbf{z} = \mathbf{0}\}$, where $\mathbf{J}_k(\epsilon) = (\nabla \mathbf{c}_k; -\mathbf{I}_{\mathcal{A}_{\boldsymbol{\ell},k}(\epsilon)}; \mathbf{I}_{\mathcal{A}_{\mathbf{u},k}(\epsilon)})$ is active constraints Jacobian. The following lemma shows that, under this construction of $\bar{\mathbf{B}}_k$, the SSQP subproblem admits a *unique* solution in a neighborhood of $\mathbf{0}$ when \mathbf{x}_k is close to \mathbf{x}^* . In particular, if we fail to solve the subproblem (17) with the above $\bar{\mathbf{B}}_k$, we should turn to safely regularize $\bar{\mathbf{B}}_k$ to be positive definite to ensure the subproblem well-defined.

Lemma 4.2. *Under Assumptions 3.6 and 4.1, we let ϵ satisfy (24) and $\beta_k = \iota_2(k+1)^{-b_2}$ satisfy $\iota_2 > 0, b_2 \in (0.5, 1]$. For any run of the method, there exists a (potential random) threshold $K^* < \infty$ such that for all $k \geq K^*$,*

- (a) *The relaxation parameter satisfies $\theta_k = 1$.*
- (b) *The estimated active sets satisfy $\mathcal{A}_{\boldsymbol{\ell},k}(\epsilon) = \mathcal{A}_{\boldsymbol{\ell}}^*$ and $\mathcal{A}_{\mathbf{u},k}(\epsilon) = \mathcal{A}_{\mathbf{u}}^*$.*
- (c) *The subproblems (17) with both the averaged gradient $\bar{\mathbf{g}}_k$ and true gradient ∇f_k admit unique local solutions in a neighborhood of $\mathbf{0}$; furthermore, $\bar{\boldsymbol{\Delta}} \mathbf{x}_k \rightarrow \mathbf{0}$ and $\boldsymbol{\Delta} \mathbf{x}_k \rightarrow \mathbf{0}$ as $k \rightarrow \infty$.*

(d) The active sets identified by the subproblems satisfy $\mathcal{A}_\ell(\mathbf{x}_k + \bar{\Delta}\mathbf{x}_k) = \mathcal{A}_\ell(\mathbf{x}_k + \Delta\mathbf{x}_k) = \mathcal{A}_\ell^*$ and $\mathcal{A}_u(\mathbf{x}_k + \bar{\Delta}\mathbf{x}_k) = \mathcal{A}_u(\mathbf{x}_k + \Delta\mathbf{x}_k) = \mathcal{A}_u^*$.

We additionally emphasize that, by Lemma 4.2(a) and under LICQ at \mathbf{x}^* , the relaxation parameter θ_k eventually equals one so that the linearized constraints become locally exact. We next explore almost-sure convergence of the dual iterates $(\boldsymbol{\lambda}_k, \boldsymbol{\mu}_k)$ and the Hessian approximation $\bar{\mathbf{B}}_k$ of the SSQP method. For the latter purpose, we require the stochastic Hessian estimates to be unbiased and have bounded variance.

Assumption 4.3. We assume, for any $k \geq 0$, $\nabla^2 F(\mathbf{x}_k; \zeta_k)$ is an unbiased estimate of $\nabla^2 f_k$ with bounded variance. In particular, $\mathbb{E}_k[\nabla^2 F(\mathbf{x}_k; \zeta_k)] = \nabla^2 f_k$ and $\mathbb{E}_k[\|\nabla^2 F(\mathbf{x}_k; \zeta_k) - \nabla^2 f_k\|^2] \leq \sigma_H^2$ for some $\sigma_H^2 > 0$. Here, $\mathbb{E}_k[\cdot]$ denotes the conditional expectation given \mathbf{x}_k .

Lemma 4.4. Under Assumptions 3.6 and 4.1, let $\alpha_k = \iota_1(k+1)^{-b_1}$ and $\beta_k = \iota_2(k+1)^{-b_2}$ satisfy $\iota_1, \iota_2 > 0, b_1 \in (0, 1], b_2 \in (0.5, 1]$. Then, we have $(\mathbf{x}_k, \boldsymbol{\lambda}_k, \boldsymbol{\mu}_k) \rightarrow (\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$ as $k \rightarrow \infty$ almost surely. Furthermore, suppose Assumption 4.3 holds and let the averaging weight sequence of the Hessian in (15) satisfy $\gamma_k = \iota_3(k+1)^{-b_3}$ with $\iota_3 > 0, b_3 \in (0.5, 1]$. Then, we also have $\bar{\mathbf{B}}_k \rightarrow \nabla_{\mathbf{x}}^2 \mathcal{L}^*$ as $k \rightarrow \infty$ almost surely.

Before stating normality result, we strengthen the bounded variance condition on stochastic gradients in Assumption 3.6 to a bounded $(2+\delta)$ moment. For notational simplicity, we reuse the symbol σ_g^2 from Assumption 3.6 to denote the upper bound of $(2+\delta)$ moment. This moment condition is mild and standard in existing literature on both unconstrained and constrained stochastic approximation methods (Duchi and Ruan, 2021; Polyak and Juditsky, 1992; Toulis et al., 2021; Chen et al., 2020; 2024; Zhu et al., 2021; Lee et al., 2022; Na and Mahoney, 2025). Some methods may require even higher-order moments; for example, projected SGD in Davis et al., 2024 imposes a bounded fourth moment to establish normality.

Assumption 4.5. We assume, for any $k \geq 0$, $\nabla F(\mathbf{x}_k; \zeta_k)$ is unbiased and has bounded $(2+\delta)$ moment for some $\delta > 0$. In particular, $\mathbb{E}_k[\|\nabla F(\mathbf{x}_k; \zeta_k) - \nabla f_k\|^{2+\delta}] \leq \sigma_g^2$ for some $\sigma_g^2 > 0$. Fur-

thermore, we assume $\mathbb{E}_k[\nabla F(\mathbf{x}_k; \zeta_k) \nabla F(\mathbf{x}_k; \zeta_k)^\top] \rightarrow \mathbb{E}[\nabla F(\mathbf{x}^*; \zeta) \nabla F(\mathbf{x}^*; \zeta)^\top]$ as $\mathbf{x}_k \rightarrow \mathbf{x}^*$.

We are now ready to present the asymptotic normality result for the SSQP iterates. Following the notation in the lower-bound analysis of Theorem 2.5, we denote $\mathbf{w}_k = (\mathbf{x}_k, \boldsymbol{\lambda}_k, [\boldsymbol{\mu}_k]_{\mathcal{A}^*})$ and $\mathbf{w}^* = (\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}_{\mathcal{A}^*}^*)$ as the nontrivial part of the primal-dual solution, noting that $\boldsymbol{\mu}_{(\mathcal{A}^*)^-}^* = \mathbf{0}$.

Theorem 4.6. *Consider the SSQP method in Section 3.1 and suppose Assumptions 4.1, 4.3, 4.5 hold. We specify the stepsize and averaging weight sequences $\alpha_k = \iota_1(k+1)^{-b_1}$, $\beta_k = \iota_2(k+1)^{-b_2}$, $\gamma_k = \iota_3(k+1)^{-b_3}$, and the stepsize adaptivity gap parameter p in (18) to satisfy*

$$\max \left\{ 0.5, \frac{2-2\delta}{2+\delta} \right\} < b_1 \leq 1, \quad 0.5 < b_2 < b_1, \quad 0.5 < b_3 \leq 1, \quad p > 1.5 - \frac{b_2}{2b_1},$$

and $\iota_1 > 2/3$ if $b_1 = 1$ and $\iota_3 > 0.25b_1$ if $b_3 = 1$. Then, we have (recall $\boldsymbol{\Omega}^*$ defined in (13))

$$1/\sqrt{\bar{\alpha}_k} \cdot (\mathbf{w}_k - \mathbf{w}^*) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \eta \cdot \boldsymbol{\Omega}^*) \quad \text{with} \quad \eta = \begin{cases} 0.5 & \text{if } b_1 < 1, \\ \iota_1/(2\iota_1 - 1) & \text{if } b_1 = 1. \end{cases}$$

For the inactive dual components, $[\boldsymbol{\mu}_k]_{(\mathcal{A}^*)^-}$ vanishes almost surely with a rate $\|[\boldsymbol{\mu}_k]_{(\mathcal{A}^*)^-}\| = o(k^{-b})$ for any $b > 0$ if $b_1 < 1$ and $O(k^{-\iota_1})$ if $b_1 = 1$.

We consider a specific setup where $\alpha_k = (k+1)^{-1}$. The following corollary shows that SSQP with this stepsize attains optimal \sqrt{k} -consistency with the limiting covariance $\boldsymbol{\Omega}^*$ matching the minimax lower bound in the sense of Hájek and Le Cam, as established in Theorem 2.5.

Corollary 4.7. *Under Assumptions 4.1, 4.3, 4.5, we let $\iota_1 = b_1 = 1$, $b_2, b_3 \in (0.5, 1)$, $\iota_2, \iota_3 > 0$, $p \geq 1.5$. Then, $\sqrt{k}(\mathbf{w}_k - \mathbf{w}^*) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \boldsymbol{\Omega}^*)$.*

Theorem 4.6 establishes asymptotic normality for primal and active dual variables of SSQP under general setups of parameters. To our knowledge, this result illustrates that SSQP with gradient momentum is the *first* method to attain joint primal-dual asymptotic optimality, extending prior works Duchi and Ruan, 2021; Davis et al., 2024 that explore only primal asymptotic optimality and rely on expensive projection operators. The main idea of SSQP is to perform sequential linear-quadratic approximations to original nonlinearly constrained problem.

The primary technical challenge in our analysis arises from handling the gradient momentum.

In particular, existing literature on asymptotic normality of first- and second-order stochastic approximation methods relies crucially on the conditional independence of stochastic gradients (Duchi and Ruan, 2021; Davis et al., 2024; Polyak and Juditsky, 1992; Toulis et al., 2021; Chen et al., 2020; 2024; Zhu et al., 2021; Lee et al., 2022; Leluc and Portier, 2023; Na and Mahoney, 2025). In contrast, the averaged gradients in our setting are inherently dependent across iterations. As a result, our analysis must carefully examine the interplay between the stepsize sequence α_k and the averaging weight sequence β_k . We show that achieving asymptotic normality requires the weights β_k to decay more slowly than the stepsizes α_k , i.e., $b_2 < b_1$. Although gradient momentum helps reduce stochastic noise and allows the method to mimic the behavior of its deterministic counterparts, the slowly decaying weights ensure that the method still effectively uses a gradient estimate $\bar{\mathbf{g}}_k$ that is close to the stochastic gradient $\nabla F(\mathbf{x}_k; \zeta_k)$, thereby preserving asymptotically optimal behavior of the iterates. We also note that the choice of b_1 is always feasible for any $\delta > 0$, with larger values of δ permitting larger feasible ranges. When $\delta = 1$ (i.e., $\nabla F(\mathbf{x}_k; \zeta_k)$ has a bounded third moment), we can take $b_1 \in (0.5, 1]$. In addition, the adaptivity gap condition is satisfied as long as $p \geq 1.5$, and we remark that numerical literature often sets $p = 2$ (Berahas et al., 2021). Finally, the conditions on ι_1, ι_3 address the corner cases when b_1 or $b_3 = 1$, and setting $\iota_1 = \iota_3 = 1$ is always sufficient to satisfy conditions.

4.2 Practical inference: plug-in covariance estimation

To perform online statistical inference based on Theorem 4.6, we must also estimate the limiting covariance matrix $\mathbf{\Omega}^*$. Existing literature considers either plug-in estimators (Chen et al., 2020; Na and Mahoney, 2025) or batch-means estimators (Zhu et al., 2021; Kuang et al., 2025). In this paper, we provide a simple plug-in estimator since the dominant computational cost of SSQP lies in solving the subproblem instead of covariance estimation. We also note that the plug-in estimator does not require any additional gradient evaluations. Define

$$\mathbf{H}_k = \begin{pmatrix} \bar{\mathbf{B}}_k & \mathbf{J}_k^\top \\ \mathbf{J}_k & \mathbf{0} \end{pmatrix}, \quad \mathbf{\Sigma}_k = \begin{pmatrix} \text{Cov}(\{\nabla F(\mathbf{x}_i; \zeta_i)\}_{i=0}^k) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}, \quad \mathbf{\Omega}_k = \mathbf{H}_k^{-1} \mathbf{\Sigma}_k \mathbf{H}_k^{-1},$$

where $\bar{\mathbf{B}}_k$ is from (16); $\mathbf{J}_k = (\nabla \mathbf{c}_k; -\mathbf{I}_{\mathcal{A}_\ell(\mathbf{x}_k + \bar{\Delta} \mathbf{x}_k)}; \mathbf{I}_{\mathcal{A}_u(\mathbf{x}_k + \bar{\Delta} \mathbf{x}_k)})$ is the Jacobian of the active constraints with the active set identified by the SSQP subproblem (cf. Lemma 4.2); and the sample covariance is defined as

$$\begin{aligned} \text{Cov}(\{\nabla F(\mathbf{x}_i; \zeta_i)\}_{i=0}^k) &= \frac{1}{k+1} \sum_{i=0}^k \nabla F(\mathbf{x}_i; \zeta_i) \nabla F(\mathbf{x}_i; \zeta_i)^\top \\ &\quad - \left(\frac{1}{k+1} \sum_{i=0}^k \nabla F(\mathbf{x}_i; \zeta_i) \right) \left(\frac{1}{k+1} \sum_{i=0}^k \nabla F(\mathbf{x}_i; \zeta_i) \right)^\top. \end{aligned}$$

Note that $\mathbf{\Omega}_k$ can be updated recursively during the iterations, as it relies on the same stochastic gradient evaluations as the averaged gradient $\bar{\mathbf{g}}_k$, namely $\{\nabla F(\mathbf{x}_i; \zeta_i)\}_{i=0}^k$. Consequently, the plug-in estimator $\mathbf{\Omega}_k$ incurs negligible additional computational cost.

The following theorem establishes the almost-sure convergence of $\mathbf{\Omega}_k$. We strengthen the bounded $(2+\delta)$ moment condition of gradient estimates to bounded 4-th moment, which is also standard for other plug-in covariance estimators (Chen et al., 2020; Na and Mahoney, 2025).

Theorem 4.8. *Under the conditions of Theorem 4.6 and further assuming $\delta \geq 2$ in Assumption 4.5, we have $\mathbf{\Omega}_k \rightarrow \mathbf{\Omega}^*$ as $k \rightarrow \infty$ almost surely.*

With the above theorem, we are now able to construct confidence intervals or regions for $\mathbf{w}^* = (\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}_{\mathcal{A}^*}^*)$. For example, fixing a desired coverage probability $1-q$ with $q \in (0, 1)$, the $100(1-q)\%$ confidence region of \mathbf{w}^* is given by

$$P(\mathbf{w}^* \in \mathcal{C}_{k,q}) \rightarrow 1 - q \quad \text{as} \quad k \rightarrow \infty, \quad (26)$$

where $\mathcal{C}_{k,q} = \{\mathbf{w} : (\mathbf{w} - \mathbf{w}_k)^\top \mathbf{\Omega}_k^\dagger (\mathbf{w} - \mathbf{w}_k) / \bar{\alpha}_k \leq \chi_{d,1-q}^2\}$. Here, $\chi_{d,1-q}^2$ is the $(1-q)$ -quantile of χ_d^2 distribution and $\mathbf{\Omega}_k^\dagger$ is the pseudo-inverse of $\mathbf{\Omega}_k$. The result (26) follows from the distribution of quadratic forms and the fact that $\text{rank}(\mathbf{\Omega}^*) = d$ (Christensen, 2020, Corollary 1.3.6a).

5 Experimental Studies

In this section, we perform comprehensive experiments to demonstrate the effectiveness of the SSQP estimation procedure in Section 3.1, validate asymptotic normality results, and illustrate

its applicability to practical statistical inference tasks, such as performing hypothesis testing and constructing confidence intervals for constrained model parameters. We evaluate the empirical performance of our method on a diverse set of tasks, including benchmark nonlinear constrained optimization problems from the CUTEst library (Gould et al., 2014; Fowkes et al., 2022), constrained regression problems involving linear, logistic, and Poisson models, as well as portfolio allocation problems. We study both synthetic data and real **Fama-French Portfolios** and **Chicago Air Pollution** data.

We use the same set of hyper-parameters for all experiments. Specifically, we set $\tau = 0.5$, $\psi = 1$, and $p = 2$. The stepsize control sequence is set to $\alpha_k = (k+1)^{-0.751}$, while the weight sequences for the gradient and Hessian are set to $\beta_k = (k+1)^{-0.501}$ and $\gamma_k = (k+1)^{-1}$, respectively, which satisfy the conditions in Theorems 3.9 and 4.6. Our implementation code is provided in the public repository: <https://github.com/yihang-gao/SSQP>.

5.1 CUTEst benchmark problems

The CUTEst library collects numerous constrained problems that are widely used for benchmarking optimization methods. We apply our method to a subset of CUTEst problems and introduce stochastic perturbations to gradients and Hessians in order to simulate noisy environments. In particular, at each step, we let $\nabla F(\mathbf{x}_k; \zeta_k) = \nabla f_k + \mathbf{E}_{k,\nabla f}$ and $\nabla^2 F(\mathbf{x}_k; \zeta_k) = \nabla^2 f_k + \mathbf{E}_{k,\nabla^2 f}$, where the deterministic quantities ∇f_k and $\nabla^2 f_k$ are provided by the CUTEst package, and $\mathbf{E}_{k,\nabla f}$ and $\mathbf{E}_{k,\nabla^2 f}$ denote the noise vector and matrix, respectively. We study two types of noise: light-tailed Gaussian noise, where each entry $\mathbf{E}_{k,\nabla f}, \mathbf{E}_{k,\nabla^2 f}$ follows $\mathcal{N}(0, \sigma^2)$, and heavy-tailed t -distribution noise, where each entry $\mathbf{E}_{k,\nabla f}, \mathbf{E}_{k,\nabla^2 f}$ follows t_{df} . We vary $\sigma^2 \in \{1, 10^{-1}, 10^{-2}, 10^{-4}\}$ and $df \in \{3, 4, 5, 9\}$. Note that the noise of t_9 closely approximates that of $\mathcal{N}(0, 1)$.

We compare our SSQP method with two state-of-the-art baselines: ActiveSet-SSQP (Na et al., 2023) and Biased-SSQP (Curtis et al., 2024). The ActiveSet-SSQP method adaptively increases the batch size to reduce the noise in the step direction, and is hence not an online method

and is more sample-intensive. In contrast, the Biased-SSQP method is an online method but lacks bias reduction through gradient momentum.

• **Global convergence.** For each problem instance, we run $K = 10^5$ iterations; and for each method, we evaluate it using the KKT residual $\|\mathbf{R}_K\|$ defined in (20) and the feasibility error $\|\mathbf{c}_K\|$. The performance comparison is shown in Figure 3. From the figure, we observe that Biased-SSQP exhibits noticeable deviation from the optimal solution, primarily due to the bias in its step direction induced by the inequality constraints. ActiveSet-SSQP, on the other hand, performs robustly even under high noise levels but heavily depends on increasing batch size to ensure a sufficiently accurate step direction. In contrast, our proposed method requires only a single sample per iteration to estimate both the gradient and Hessian, while effectively reducing stochastic noise and debiasing the step through gradient averaging. We see that when the noise level is moderate (Gaussian noise with $\sigma^2 \in \{10^{-1}, 10^{-2}, 10^{-4}\}$), our method achieves performance comparable to, or even better than, ActiveSet-SSQP.

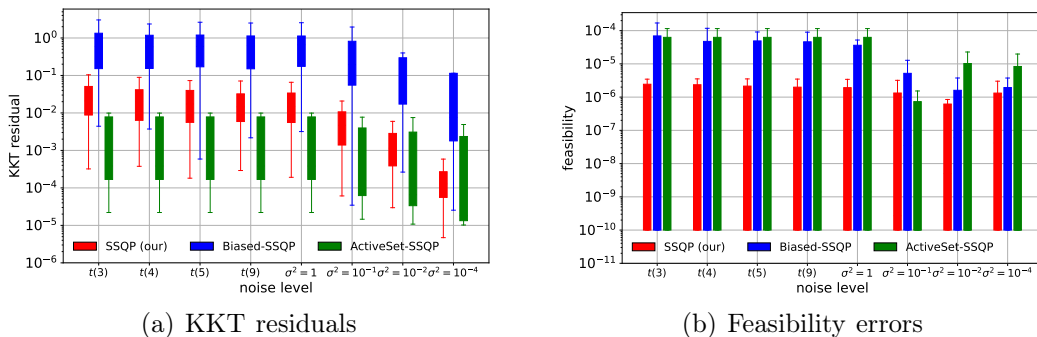


Figure 3: *Boxplots of KKT residuals and feasibility errors for CUTEst problems. For each noise setting, three boxplots are shown, corresponding to the SSQP (ours), Biased-SSQP, and ActiveSet-SSQP methods.*

• **Local normality.** We next examine local asymptotic normality of the SSQP iterates. For each problem, we estimate the averaged model parameter $\mathbf{1}^\top \mathbf{x}^*/d$ and construct its nominal 95% confidence interval. The performance of the constructed intervals is evaluated in terms of their empirical coverage rate (CovRate) and average length (AvgLen), each averaged over 200 independent runs. Results of 5 HS-type CUTEst problems are reported in Table 1. We find that the confidence intervals produced by SSQP achieve empirical coverage probabilities that closely

match the nominal 95% level, thereby providing strong empirical support for the theoretical asymptotic normality guarantee established in Theorem 4.6. Furthermore, the average length of the confidence intervals increases as the noise level grows, a behavior also consistent with our theoretical expectations, since the asymptotic covariance matrix $\mathbf{\Omega}^*$ is proportional to $\text{Cov}(\nabla F(\mathbf{x}^*; \zeta))$ as defined in (13).

Table 1: *The coverage rate (CovRate) and length of confidence intervals (AvgLen) for 5 HS-type CUTEst constrained problems. The standard deviation of the interval length is also reported.*

Problem	Noise Level (σ^2)	Gaussian		Noise Level (df)	Student t	
		CovRate(%)	AvgLen		CovRate(%)	AvgLen
HS41	1	97.0	2.50E-2 (7.73E-4)	3	86.0	3.77E-2 (1.90E-3)
	10^{-1}	97.5	7.59E-3 (7.03E-5)	4	93.0	3.06E-2 (1.28E-3)
	10^{-2}	97.0	2.40E-3 (8.69E-6)	5	94.0	2.79E-2 (9.25E-4)
	10^{-4}	97.5	2.40E-4 (5.95E-7)	9	97.0	2.45E-2 (7.10E-4)
HS65	1	94.5	1.87E-3 (6.82E-6)	3	96.5	3.18E-3 (1.66E-4)
	10^{-1}	94.5	5.92E-4 (1.59E-6)	4	95.0	2.59E-3 (1.72E-5)
	10^{-2}	95.0	1.87E-4 (4.96E-7)	5	95.0	2.37E-3 (1.11E-5)
	10^{-4}	94.5	1.87E-5 (4.97E-8)	9	94.5	2.08E-3 (8.36E-6)
HS68	1	97.0	2.31E-1 (4.85E-2)	3	95.5	3.00E-1 (1.32E-1)
	10^{-1}	98.0	5.09E-2 (2.33E-3)	4	94.5	2.08E-1 (6.09E-2)
	10^{-2}	98.5	1.58E-2 (2.23E-4)	5	95.0	1.81E-1 (5.07E-2)
	10^{-4}	95.5	1.58E-3 (4.56E-6)	9	94.5	1.48E-1 (3.52E-2)
HS71	1	97.0	1.95E-3 (1.44E-5)	3	94.0	3.34E-3 (1.23E-4)
	10^{-1}	96.5	6.17E-4 (1.93E-6)	4	96.0	2.74E-3 (6.79E-3)
	10^{-2}	96.5	1.95E-4 (5.20E-7)	5	96.5	2.49E-3 (2.51E-5)
	10^{-4}	98.5	1.95E-5 (5.08E-8)	9	95.0	2.19E-3 (2.12E-5)
HS81	1	94.5	3.49E-2 (3.17E-3)	3	91.0	5.04E-2 (7.56E-3)
	10^{-1}	97.0	1.13E-2 (4.77E-5)	4	94.0	4.21E-2 (3.51E-3)
	10^{-2}	98.0	3.58E-3 (9.63E-6)	5	94.5	3.88E-2 (2.42E-3)
	10^{-4}	98.0	3.59E-4 (9.22E-7)	9	95.0	3.43E-2 (2.10E-3)

• **Effectiveness of gradient averaging.** Empirically, our online method consistently outperforms the other online Biased-SSQP method. This improvement is primarily attributed to the momentum-style gradient moving-average scheme: as iterations proceed, the averaged gradient progressively approximates the true population gradient, enabling the method to emulate the behavior of deterministic SQP while maintaining full stochasticity. In contrast, Biased-SSQP updates its iterates using noisy gradients, leading to inherent biases in step estimates that accumulate over time and result in deviations from the optimal solution. As shown in Figure 4, we visualize the difference between the averaged gradients and the true gradients across iterations. The results align with our theory that averaged gradients (solid lines) progressively approach

the true gradients; while non-averaged gradients (dashed lines) exhibit constant variance.

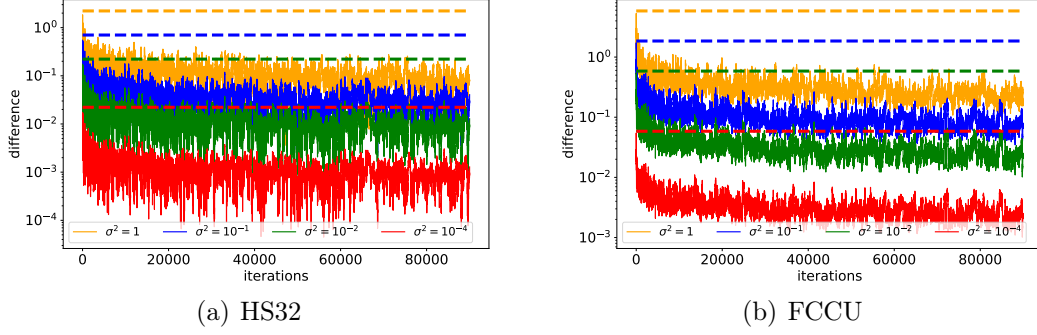


Figure 4: *Differences between the averaged gradients and the true gradients on HS32 and FCCU problems. Solid lines: trajectories of gradient difference between the averaged gradients and the exact gradients during iterations, i.e., $\|\bar{\mathbf{g}}_k - \nabla f_k\|$. Dashed lines: expected error without averaging, i.e., $\mathbb{E}[\|\nabla F(\mathbf{x}_k; \zeta_k) - \nabla f_k\| \mid \mathbf{x}_k]$.*

5.2 Constrained regression problems

We further implement our method on constrained regression problems, including both linear and logistic regression models. In this study, each sample is denoted by the covariate-response pair $\zeta_k = (\zeta_k^a, \zeta_k^b)$, and the response is generated based on different regression models. In particular, for linear regression, we let $\zeta_k^b = \zeta_k^{a\top} \mathbf{x}^* + \varepsilon_k$ with $\varepsilon_k \stackrel{iid}{\sim} \mathcal{N}(0, 1)$ and set $F(\mathbf{x}; \zeta_k) = \frac{1}{2}(\zeta_k^b - \zeta_k^{a\top} \mathbf{x})^2$. For logistic regression, we let $P(\zeta_k^b \mid \zeta_k^a) = 1 / \{1 + \exp(-\zeta_k^b \cdot \zeta_k^{a\top} \mathbf{x}^*)\}$ with $\zeta_k^b \in \{-1, 1\}$ and set $F(\mathbf{x}; \zeta_k) = \log(1 + \exp(-\zeta_k^b \cdot \zeta_k^{a\top} \mathbf{x}))$. For both models, we set the true model parameters as $\mathbf{x}^* = (\frac{3}{2d}, \dots, \frac{3}{2d}, \frac{1}{2d}, \dots, \frac{1}{2d})$, where the first half of the components are $\frac{3}{2d}$ and the second half are $\frac{1}{2d}$. We then impose the probability simplex constraints $\Omega := \{\mathbf{x} \in \mathbb{R}^d : \mathbf{1}^\top \mathbf{x} = 1, \mathbf{x} \geq \mathbf{0}\}$ in the estimation procedure. We generate covariate vector ζ_k^a from $\zeta_k^a \sim \mathcal{N}(\boldsymbol{\mu}^a, \boldsymbol{\Sigma}^a)$, where $\boldsymbol{\mu}^a = (1, \dots, 1, -1, \dots, -1)$. Following unconstrained regression settings (Chen et al., 2020; Zhu et al., 2021), we consider three structures for $\boldsymbol{\Sigma}^a$: (i) Identity: $\boldsymbol{\Sigma}^a = \mathbf{I}$; (ii) Toeplitz: $[\boldsymbol{\Sigma}^a]_{ij} = r^{|i-j|}$ with $r \in \{0.4, 0.5, 0.6\}$; and (iii) Equi-correlation: $[\boldsymbol{\Sigma}^a]_{ij} = r$ for $i \neq j$ and $[\boldsymbol{\Sigma}^a]_{ii} = 1$ with $r \in \{0.1, 0.2, 0.3\}$. For each setting of $\boldsymbol{\Sigma}^a$, we vary the problem dimension $d \in \{5, 10, 20, 30\}$ and construct 95% confidence interval for $\boldsymbol{\mu}^{a\top} \mathbf{x}^*$, i.e., the parameters difference of two groups.

We measure the performance by repeating each experiment 200 times with varying random seeds and reporting the empirical coverage rate (CovRate) and average length (AvgLen) of the

constructed confidence intervals. The results are summarized in Tables 2 and 3. From the tables, we observe that the empirical coverage probabilities are close to the nominal 95% level across all settings, thereby providing strong empirical support for our local asymptotic normality guarantees and demonstrating the practical effectiveness of our method. Furthermore, the average length of the confidence intervals remains on the order of 10^{-2} , consistent with the findings reported by Chen et al., 2020; Zhu et al., 2021; Na and Mahoney, 2025. The low standard deviation of the interval lengths relative to their means further indicates the robustness of the proposed inference procedure across different covariance structures.

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	Dim d	CovRate(%)	AvgLen	Dim d	CovRate(%)	AvgLen
Identity	5	93.5	3.73E-2 (1.74E-4)	20	92.5	4.00E-2 (1.33E-4)
	10	96.5	3.91E-2 (1.47E-4)	30	92.5	4.03E-2 (1.53E-4)
Toeplitz ($r = 0.4$)	5	94.0	3.71E-2 (1.68E-4)	20	96.0	3.93E-2 (1.38E-4)
	10	94.5	3.82E-2 (1.62E-4)	30	93.0	3.98E-2 (1.52E-4)
Toeplitz ($r = 0.5$)	5	94.0	3.74E-2 (1.67E-4)	20	96.0	3.91E-2 (1.38E-4)
	10	95.5	3.82E-2 (1.60E-4)	30	93.0	3.95E-2 (1.61E-4)
Toeplitz ($r = 0.6$)	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 ($r = 0.1$)	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)
EquiCorr ($r = 0.2$)	5	92.5	3.79E-2 (1.59E-4)	20	93.5	4.02E-2 (1.26E-4)
	10	95.0	3.94E-2 (1.50E-4)	30	96.0	4.05E-2 (1.44E-4)
EquiCorr ($r = 0.3$)	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)

5.3 Portfolio allocation: Fama-French dataset

We implement our method on portfolio allocation problems, using 30 portfolios selected from the Fama–French Portfolios dataset, subject to the gross-exposure constraint (Fan, 2007; Fan et al., 2012; Du et al., 2022): $\Omega := \{\mathbf{x} \in \mathbb{R}^d : \mathbf{1}^\top \mathbf{x} = 1, \|\mathbf{x}\|_1 \leq c\}$, where we set $c = 3$ and \mathbf{x} denotes the portfolio weight vector; a negative weight signifies shorting an asset. Let $\boldsymbol{\mu}^a, \boldsymbol{\Sigma}^a$ denote the mean vector and covariance matrix of the asset returns $\zeta^a \in \mathbb{R}^{30}$. We consider four widely studied portfolio models: (i) Global Minimum Variance (**GMV**): $\min_{\mathbf{x} \in \Omega} \mathbf{x}^\top \boldsymbol{\Sigma}^a \mathbf{x}$; (ii) Mean-Variance (**MV**): $\min_{\mathbf{x} \in \Omega} -\mathbf{x}^\top \boldsymbol{\mu}^a + \mathbf{x}^\top \boldsymbol{\Sigma}^a \mathbf{x}$; (iii) Exponential Utility (**EXP**): $\min_{\mathbf{x} \in \Omega} \mathbb{E}[\exp(-\eta_1 \mathbf{x}^\top \zeta^a)]$, where

Table 3: *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	Dim d	CovRate(%)	AvgLen	Dim d	CovRate(%)	AvgLen
Identity	5	96.5	4.46E-2 (7.97E-5)	20	94.5	5.87E-2 (7.13E-5)
	10	94.5	5.87E-2 (7.13E-5)	30	93.0	7.34E-2 (7.90E-5)
Toeplitz ($r = 0.4$)	5	94.5	4.46E-2 (9.06E-5)	20	92.5	6.86E-2 (1.01E-4)
	10	95.5	5.83E-2 (8.59E-5)	30	93.5	7.30E-2 (1.13E-4)
Toeplitz ($r = 0.5$)	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)
Toeplitz ($r = 0.6$)	5	94.5	4.47E-2 (9.63E-5)	20	92.5	6.82E-2 (1.19E-4)
	10	94.0	5.83E-2 (8.77E-5)	30	94.5	7.26E-2 (1.32E-4)
EquiCorr ($r = 0.1$)	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 ($r = 0.2$)	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 ($r = 0.3$)	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)

$\eta_1 > 0$ is the risk-aversion parameter and set to $\eta_1 = 0.1$; and (iv) Logarithmic Utility (**LOG**):

$\min_{\mathbf{x} \in \Omega} -\mathbb{E}[\log(\mathbf{x}^\top \boldsymbol{\zeta}^a + \eta_2)]$, where $\eta_2 > 0$ is a regularization parameter ensuring the validity of the logarithm and set to $\eta_2 = 15$.

For each month, we use historical daily data from the preceding year as training samples and apply both the SSQP method (dealing with one sample at a time) and the M -estimation method (using full samples to estimate $\boldsymbol{\mu}^a, \boldsymbol{\Sigma}^a$) to solve the resulting constrained stochastic problem and obtain a portfolio weight vector. This weight vector is then held fixed throughout the following month for evaluation. Specifically, the out-of-sample performance of the estimated portfolio weights is evaluated using four standard metrics computed over 30 months during the period 2021–2023: accumulative return, maximum drawdown, Sharpe ratio, and Sortino ratio. The accumulative return quantifies the overall gain or loss of a portfolio strategy, while the other three metrics balance return with risk aspect: the maximum drawdown measures the largest observed decline from a peak to a trough; the Sharpe ratio relates the portfolio’s mean return to its total risk measured by the standard deviation of returns; and the Sortino ratio refines this measure by considering only downside volatility, i.e., the standard deviation of negative returns. The results are summarized in Table 4.

Table 4: *The comparison of SSQP and M -estimation on Fama-French Portfolios Dataset from 2021-2023. For each metric, the best-performing model is shown in bold.*

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

From Table 4, we see that SSQP achieves comparable performance to offline M -estimation across all metrics. In particular, the performance differences across all four models and four metrics are within 0.15, with the largest observed in the global minimum variance model under the Sortino ratio metric, where SSQP outperforms M -estimation by 0.14. Among different portfolio models, we observe that the portfolio strategy based on the logarithmic utility model achieves the highest accumulative return, which is consistent with the empirical findings of [Du et al., 2022](#). The mean-variance model yields higher Sharpe and Sortino ratios, while the simple equal-weight strategy yields the smallest maximum drawdown.

Figure 5 further visualizes the estimated weights of two representative stocks under two portfolio models, exponential utility and logarithmic utility. From the figure, we observe a strong correlation between the temporal adjustments of the weights and the dynamics of the corresponding stock returns. In particular, as the return decreases, the weight assigned to the stock also decreases. Moreover, an abrupt drop in the stock return is promptly followed by a widening blue band, indicating an increase in the estimated standard deviation of the associated weight. This behavior aligns with economic intuition and demonstrates the interpretability and validity of the estimated weights under constrained model estimation.

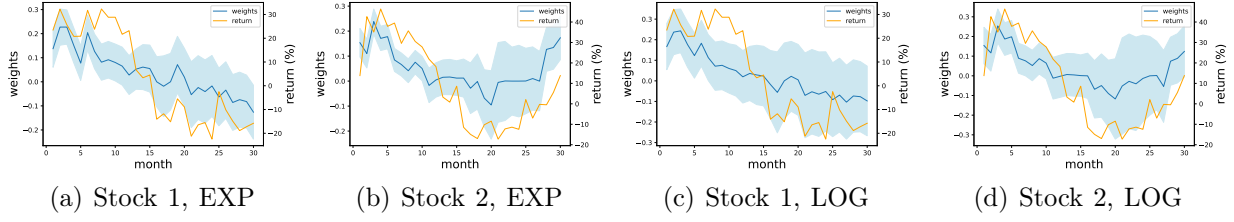


Figure 5: *Trajectories of portfolio weights and corresponding stock returns. We show two stocks under two portfolio models. The blue lines depict the predicted weights for the stock, while the shaded blue bands indicate the estimated standard deviations of these weights, computed based on the derived asymptotic normality results. The yellow lines depict the returns of the same stock.*

5.4 Poisson regression: Chicago air pollution data

In this section, we compare unconstrained and constrained methods for generalized linear regression problems. We use daily air pollution and death rate data for Chicago (Wood, 2017) to analyze the relationship between death rates and air pollution levels. Following Toulis and Airoldi, 2017, we fit a Poisson regression model that regresses the death counts (ζ^b) on six variables: Intercept, Time, PM_{10} , $\text{PM}_{2.5}$, SO_2 , and O_3 (ζ^a). In particular, we assume $\zeta^b | \zeta^a \sim \text{Poisson}(\lambda(\zeta^a))$ with $\log \lambda(\zeta^a) = \zeta^{a\top} \mathbf{x}^*$. In this study, we may want to incorporate the prior domain knowledge that higher pollutant concentrations are expected to increase mortality, and thus impose non-negativity constraints on the coefficients of the pollution covariates. Accordingly, we compare the following two methods:

$$\begin{aligned} \min_{\mathbf{x}} \quad & \mathbb{E}_{\zeta}[\zeta^b \zeta^{a\top} \mathbf{x} - \exp(\zeta^{a\top} \mathbf{x})], \\ \text{s.t.} \quad & \mathbf{x}_{3:6} \geq \mathbf{0}. \end{aligned} \quad (27)$$

The above unconstrained Poisson model is fitted using the `statsmodels` package in Python (Seabold and Perktold, 2010), while the constrained model is fitted using the proposed SSQP method. We summarize the estimated coefficients along with their 95% confidence intervals and p -values in Table 5 (left). From this table, we observe that the unconstrained method estimates the coefficient of O_3 to be negative and significant, which contradicts prior domain knowledge and can be difficult to interpret. In contrast, the constrained method correctly identifies the coefficient of O_3 as active. Among other variables, the two methods produce comparable coefficient

estimates. For example, the coefficients estimated by the unconstrained method lie within the 95% confidence intervals constructed by the constrained method, and vice versa. Furthermore, both methods identify the Intercept, Time, and SO₂ as significant, while PM₁₀ as insignificant. The only difference occurs with PM_{2.5}, which is insignificant under the unconstrained method but significant under the constrained method.

Table 5: *Summary of two Poisson regression models (27) (left) and (28) (right) applied to the Chicago air pollution data.*

Var	Method	Coeff (10 ⁻²)	95% CI (10 ⁻²)	p-Value
Intercept	Uncons	4.6968	[4.690, 4.704]	<0.001
	Cons	4.6974	[4.692, 4.703]	<0.001
Time	Uncons	0.95	[0.17, 1.74]	0.008
	Cons	1.13	[0.64, 1.63]	<0.001
PM₁₀	Uncons	0.42	[-0.57, 1.42]	0.396
	Cons	0.13	[-0.56, 0.79]	0.371
PM_{2.5}	Uncons	0.72	[-0.08, 1.52]	0.103
	Cons	0.65	[0.02, 1.28]	0.023
SO₂	Uncons	1.38	[0.58, 2.20]	0.001
	Cons	2.08	[1.43, 2.73]	<0.001
O₃	Uncons	-2.97	[-3.70, -2.24]	<0.001
	Cons	0.00	active	

Var	Method	Coeff (10 ⁻²)	95% CI (10 ⁻²)	p-Value
Intercept	Uncons	4.6972	[4.690, 4.704]	<0.001
	Cons	4.6973	[4.692, 4.703]	<0.001
Time	Uncons	1.21	[0.53, 1.89]	0.001
	Cons	1.13	[0.64, 1.63]	<0.001
PM₁₀	Uncons	-0.86	[-1.82, 0.10]	0.062
	Cons	0.11	[-0.52, 0.74]	0.362
PM_{2.5}	Uncons	1.37	[0.51, 2.23]	0.001
	Cons	0.65	[0.01, 1.28]	0.022
SO₂	Uncons	2.06	[1.28, 2.84]	<0.001
	Cons	2.08	[1.42, 2.73]	<0.001

When we encounter a negative coefficient of O₃ after fitting the model, we may next perform model modification by simply removing this variable for sake of interpretability. As such, we next consider a sub-model comparison where we remove O₃ variable:

$$\begin{aligned}
\min_{\mathbf{x}} \mathbb{E}_{\zeta}[\zeta^b \zeta^{a^\top} \mathbf{x} - \exp(\zeta^{a^\top} \mathbf{x})], & \quad \min_{\mathbf{x}} \mathbb{E}_{\zeta}[\zeta^b \zeta^{a^\top} \mathbf{x} - \exp(\zeta^{a^\top} \mathbf{x})], \\
\text{s.t. } \mathbf{x}_{3:5} \geq \mathbf{0}. & \quad \text{s.t. } \mathbf{x}_{3:5} \geq \mathbf{0}.
\end{aligned} \tag{28}$$

The results are summarized in Table 5 (right). From this table, we observe that the two methods output largely similar results. In particular, both methods identify the Intercept, Time, and SO₂ as significant, and PM₁₀ as insignificant, consistent with the conclusions drawn from the full model. Notably, both methods also identify PM_{2.5} as significant, which aligns with the constrained method fitted on the full model. To summarize, the unconstrained method may yield inconsistency between fitting the full and sub models, while our constrained method preserves consistency when active variables on the boundary are removed. This illustrates the benefits of applying our constrained methods when incorporating prior domain knowledge into model fitting.

6 Conclusion and Future Work

In this paper, we studied online statistical inference for the solutions of stochastic optimization problems with equality and inequality constraints. We developed a stochastic sequential quadratic programming method that incorporates a moving-average gradient scheme to correct the bias in the stochastic step direction induced by inequality constraints. We established global almost-sure convergence and proved that the proposed method achieves local asymptotic normality with a minimax-optimal primal–dual limiting covariance matrix in the sense of Hájek and Le Cam. Furthermore, we proposed a plug-in covariance matrix estimator for practical inference. Extensive experiments on benchmark nonlinear problems, as well as on constrained generalized linear models and portfolio allocation tasks using both synthetic and real data, demonstrated the superior empirical performance of our method and confirmed its effectiveness of online constrained inference in practice.

As for future directions, it would be significant to provide a non-asymptotic analysis that quantifies how quickly the stochastic iterates approach the limiting distribution. For example, [Anastasiou et al., 2019](#) derived a non-asymptotic convergence rate of averaged SGD to a normal distribution by applying a non-asymptotic martingale central limit theorem. Establishing comparable bounds for SSQP methods on constrained problems would further highlight the effectiveness of this type of methods. Furthermore, recent work has shown promising advances in applying SGD methods to high-dimensional settings ([Li et al., 2025](#)). Developing a high-dimensional theory for constrained model estimation would require incorporating regularization techniques such as sparsity-inducing penalties, manifold constraints, or low-rank constraints to enable valid analysis. Integrating these techniques with the proposed SSQP method can open new pathways toward scalable, structure-aware online inference for high-dimensional constrained problems.

References

- Hájek, J. (1972). Local asymptotic minimax and admissibility in estimation. *Proceedings of the sixth Berkeley symposium on mathematical statistics and probability*. Vol. 1, p. 175–194.
- Le Cam, L. (1972). Limits of experiments. *Proceedings of the Sixth Berkeley Symposium on Mathematical Statistics and Probability*. Vol. 1. University of California Press Berkeley-Los Angeles, p. 245–261.
- Duchi, J. C. and F. Ruan (2021). Asymptotic optimality in stochastic optimization. *The Annals of Statistics*, 49(1).
- Davis, D., D. Drusvyatskiy, and L. Jiang (2024). Asymptotic normality and optimality in non-smooth stochastic approximation. *The Annals of Statistics*, 52(4).
- Robbins, H. and S. Monro (1951). A Stochastic Approximation Method. *The Annals of Mathematical Statistics*, 22(3): p. 400–407.
- Kiefer, J. and J. Wolfowitz (1952). Stochastic Estimation of the Maximum of a Regression Function. *The Annals of Mathematical Statistics*, 23(3): p. 462–466.
- Polyak, B. T. and A. B. Juditsky (1992). Acceleration of Stochastic Approximation by Averaging. *SIAM Journal on Control and Optimization*, 30(4): p. 838–855.
- Toulis, P. and E. M. Airoldi (2017). Asymptotic and finite-sample properties of estimators based on stochastic gradients. *The Annals of Statistics*, 45(4).
- Fang, Y., J. Xu, and L. Yang (2018). Online bootstrap confidence intervals for the stochastic gradient descent estimator. *Journal of Machine Learning Research*, 19(78): p. 1–21.
- Toulis, P., T. Horel, and E. M. Airoldi (2021). The Proximal Robbins–Monro Method. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 83(1): p. 188–212.
- Chen, X., J. D. Lee, X. T. Tong, and Y. Zhang (2020). Statistical inference for model parameters in stochastic gradient descent. *The Annals of Statistics*, 48(1).

- Chen, X., Z. Lai, H. Li, and Y. Zhang (2024). Online Statistical Inference for Stochastic Optimization via Kiefer-Wolfowitz Methods. *Journal of the American Statistical Association*, 119(548): p. 2972–2982.
- Zhu, W., X. Chen, and W. B. Wu (2021). Online Covariance Matrix Estimation in Stochastic Gradient Descent. *Journal of the American Statistical Association*, 118(541): p. 393–404.
- Lee, S., Y. Liao, M. H. Seo, and Y. Shin (2022). Fast and Robust Online Inference with Stochastic Gradient Descent via Random Scaling. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(7): p. 7381–7389.
- Jiang, L. et al. (2025). Online covariance estimation in nonsmooth stochastic approximation. *The 38th Annual Conference on Learning Theory (COLT)*.
- Polyak, B. (1964). Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics*, 4(5): p. 1–17.
- Nocedal, J. and S. J. Wright (2006). Numerical optimization. Springer New York.
- Na, S., M. Anitescu, and M. Kolar (2022). An adaptive stochastic sequential quadratic programming with differentiable exact augmented lagrangians. *Mathematical Programming*, 199(1–2): p. 721–791.
- (2023). Inequality constrained stochastic nonlinear optimization via active-set sequential quadratic programming. *Mathematical Programming*, 202(1–2): p. 279–353.
- Na, S. (2025). Derivative-Free Sequential Quadratic Programming for Equality-Constrained Stochastic Optimization. *arXiv preprint arXiv:2510.22458*.
- Curtis, F. E., M. J. O’Neill, and D. P. Robinson (2023). Worst-case complexity of an SQP method for nonlinear equality constrained stochastic optimization. *Mathematical Programming*.
- Curtis, F. E., D. P. Robinson, and B. Zhou (2024). Sequential Quadratic Optimization for Stochastic Optimization with Deterministic Nonlinear Inequality and Equality Constraints. *SIAM Journal on Optimization*, 34(4): p. 3592–3622.

- Berahas, A. S., F. E. Curtis, D. Robinson, and B. Zhou (2021). Sequential Quadratic Optimization for Nonlinear Equality Constrained Stochastic Optimization. *SIAM Journal on Optimization*, 31(2): p. 1352–1379.
- Berahas, A. S., F. E. Curtis, M. J. O’Neill, and D. P. Robinson (2023). A Stochastic Sequential Quadratic Optimization Algorithm for Nonlinear-Equality-Constrained Optimization with Rank-Deficient Jacobians. *Mathematics of Operations Research*.
- Fang, Y., S. Na, M. W. Mahoney, and M. Kolar (2024a). Fully Stochastic Trust-Region Sequential Quadratic Programming for Equality-Constrained Optimization Problems. *SIAM Journal on Optimization*, 34(2): p. 2007–2037.
- Fang, Y., S. Na, M. W. Mahoney, and M. Kolar (2024b). Trust-Region Sequential Quadratic Programming for Stochastic Optimization with Random Models. *arXiv preprint arXiv:2409.15734*.
- Fang, Y., J. Lavaei, and S. Na (2025). High Probability Complexity Bounds of Trust-Region Stochastic Sequential Quadratic Programming with Heavy-Tailed Noise. *arXiv preprint arXiv:2503.19091*.
- Bercu, B., A. Godichon, and B. Portier (2020). An Efficient Stochastic Newton Algorithm for Parameter Estimation in Logistic Regressions. *SIAM Journal on Control and Optimization*, 58(1): p. 348–367.
- Boyer, C. and A. Godichon-Baggioni (2022). On the asymptotic rate of convergence of Stochastic Newton algorithms and their Weighted Averaged versions. *Computational Optimization and Applications*, 84(3): p. 921–972.
- Leluc, R. and F. Portier (2023). Asymptotic Analysis of Conditioned Stochastic Gradient Descent. *Transactions on Machine Learning Research*.
- Cénac, P., A. Godichon-Baggioni, and B. Portier (2025). An efficient averaged stochastic Gauss-Newton algorithm for estimating parameters of nonlinear regressions models. *Bernoulli*, 31(1).

- Kuang, W., M. Anitescu, and S. Na (2025). Online covariance matrix estimation in sketched newton methods. *arXiv preprint arXiv:2502.07114*.
- Na, S. and M. Mahoney (2025). Statistical inference of constrained stochastic optimization via sketched sequential quadratic programming. *Journal of Machine Learning Research*, 26(33): p. 1–75.
- Du, X., W. Zhu, W. B. Wu, and S. Na (2025). Online Statistical Inference of Constrained Stochastic Optimization via Random Scaling. *arXiv preprint arXiv:2505.18327*.
- Xu, M., J. J. Ye, and L. Zhang (2015). Smoothing SQP Methods for Solving Degenerate Nonsmooth Constrained Optimization Problems with Applications to Bilevel Programs. *SIAM Journal on Optimization*, 25(3): p. 1388–1410.
- Shapiro, A. (1990). On differential stability in stochastic programming. *Mathematical Programming*, 47(1–3): p. 107–116.
- Bonnans, J. F. and A. Shapiro (2000). Perturbation Analysis of Optimization Problems. Springer New York.
- Burer, S. and A. N. Letchford (2009). On Nonconvex Quadratic Programming with Box Constraints. *SIAM Journal on Optimization*, 20(2): p. 1073–1089.
- Dennis, J. E. and J. J. Moré (1974). A characterization of superlinear convergence and its application to quasi-Newton methods. *Mathematics of Computation*, 28(126): p. 549–560.
- Liu, X.-w. and Y.-x. Yuan (2000). A Robust Algorithm for Optimization with General Equality and Inequality Constraints. *SIAM Journal on Scientific Computing*, 22(2): p. 517–534.
- Christensen, R. (2020). Plane Answers to Complex Questions: The Theory of Linear Models. Springer International Publishing.
- Gould, N. I. M., D. Orban, and P. L. Toint (2014). CUTEst: a Constrained and Unconstrained Testing Environment with safe threads for mathematical optimization. *Computational Optimization and Applications*, 60(3): p. 545–557.

- Fowkes, J., L. Roberts, and Á. Bűrmen (2022). PyCUTEst: an open source Python package of optimization test problems. *Journal of Open Source Software*, 7(78): p. 4377.
- Fan, J. (2007). Variable screening in high-dimensional feature space. *Proceedings of the 4th International Congress of Chinese Mathematicians*. Vol. 2. Citeseer, p. 735–747.
- Fan, J., J. Zhang, and K. Yu (2012). Vast Portfolio Selection With Gross-Exposure Constraints. *Journal of the American Statistical Association*, 107(498): p. 592–606.
- Du, J.-H., Y. Guo, and X. Wang (2022). High-Dimensional Portfolio Selection with Cardinality Constraints. *Journal of the American Statistical Association*, 118(542): p. 779–791.
- Wood, S. N. (2017). Generalized Additive Models: An Introduction with R. Chapman and Hall/CRC.
- Seabold, S. and J. Perktold (2010). Statsmodels: Econometric and statistical modeling with python. *Proceedings of the 9th Python in Science Conference*. Vol. 57. (61).
- Anastasiou, A., K. Balasubramanian, and M. A. Erdogdu (2019). Normal approximation for stochastic gradient descent via non-asymptotic rates of martingale CLT. *Conference on Learning Theory (COLT)*. PMLR, p. 115–137.
- Li, J., Z. Lou, J. Schmidt-Hieber, and W. B. Wu (2025). Statistical Guarantees for High-Dimensional Stochastic Gradient Descent. *Conference on Neural Information Processing Systems (NeurIPS)*.
- Na, S., Z. Yang, Z. Wang, and M. Kolar (2019). High-dimensional Varying Index Coefficient Models via Stein’s Identity. *Journal of Machine Learning Research*, 20: p. 152–1.
- Na, S. and M. Kolar (2021). High-dimensional index volatility models via Stein’s identity. *Bernoulli*, 27(2).
- Shapiro, A. (1985). Asymptotic distribution of test statistics in the analysis of moment structures under inequality constraints. *Biometrika*, 72(1): p. 133–144.

- Zafar, M. B., I. Valera, M. Gomez-Rodriguez, and K. P. Gummadi (2019). Fairness constraints: A flexible approach for fair classification. *The Journal of Machine Learning Research*, 20(1): p. 2737–2778.
- Sen, P. K. (1979). Asymptotic Properties of Maximum Likelihood Estimators Based on Conditional Specification. *The Annals of Statistics*, 7(5).
- Dupacova, J. and R. Wets (1988). Asymptotic Behavior of Statistical Estimators and of Optimal Solutions of Stochastic Optimization Problems. *The Annals of Statistics*, 16(4).
- Nagaraj, N. K. and W. A. Fuller (1991). Estimation of the Parameters of Linear Time Series Models Subject to Nonlinear Restrictions. *The Annals of Statistics*, 19(3).
- Shapiro, A. (2000). On the asymptotics of constrained local M -estimators. *The Annals of Statistics*, 28(3).
- Cuomo, S. et al. (2022). Scientific Machine Learning Through Physics-Informed Neural Networks: Where we are and What’s Next. *Journal of Scientific Computing*, 92(3).
- Chen, R. T., Y. Rubanova, J. Bettencourt, and D. K. Duvenaud (2018). Neural ordinary differential equations. *Advances in Neural Information Processing Systems (NeurIPS)*, 31.
- Wang, S., H. Wang, and P. Perdikaris (2021). Learning the solution operator of parametric partial differential equations with physics-informed DeepONets. *Science Advances*, 7(40).
- Raissi, M., P. Perdikaris, and G. Karniadakis (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: p. 686–707.
- Krishnapriyan, A. et al. (2021). Characterizing possible failure modes in physics-informed neural networks. *Advances in Neural Information Processing Systems (NeurIPS)*, 34.
- Négiar, G., M. W. Mahoney, and A. Krishnapriyan (2023). Learning differentiable solvers for systems with hard constraints. *International Conference on Learning Representations*.
- Robinson, S. M. (1976). Stability Theory for Systems of Inequalities, Part II: Differentiable Nonlinear Systems. *SIAM Journal on Numerical Analysis*, 13(4): p. 497–513.

- Vaart, A. W. v. d. (1998). Asymptotic Statistics. Cambridge University Press.
- Daniel, J. W. (1973). Stability of the solution of definite quadratic programs. *Mathematical Programming*, 5(1): p. 41–53.
- Robbins, H. and D. Siegmund (1971). A convergence theorem for non negative almost supermartingales and some applications. *Optimizing Methods in Statistics*. Elsevier, p. 233–257.
- Duflo, M. (1997). Random iterative models. Ed. by S. S. Wilson. Applications of mathematics 34. Aus dem Franz. übers. Berlin: Springer. ISBN: 3540571000.
- Dontchev, A. L. and R. T. Rockafellar (2009). Implicit Functions and Solution Mappings: A View from Variational Analysis. Springer New York. ISBN: 9780387878218.
- Hall, P. and C. C. Heyde (2014). Martingale limit theory and its application. Academic press.

Appendix: Online Inference of Constrained Optimization: Primal-Dual Optimality and Sequential Quadratic Programming

A Motivating Examples

The question of when constraints are needed in statistical machine learning problems may be divided into two scenarios. The first scenario refers to problems that are otherwise ill-posed and impossible to solve. For example, in principal component analysis, we estimate the leading eigenvector of the population covariance matrix Σ . However, maximizing $\mathbf{x}^\top \Sigma \mathbf{x}$ is not well-defined without the constraint $\|\mathbf{x}\|^2 = 1$. The second scenario refers to problems where constraints incorporate prior model information and enable data to be explained by simpler hypotheses. We introduce several examples that can be cast into the form of Problem (1).

A.1 Constrained regression

Given a set of feature-response data $\{(\zeta_1^a, \zeta_1^b), \dots, (\zeta_n^a, \zeta_n^b)\}$, we consider fitting a smooth monotone function $z : \zeta^a \rightarrow \zeta^b$. Monotonic regression has many applications, such as estimating cumulative distribution functions, and monotonic relationships between two variables are prevalent in various studies. For example, demand is typically a non-increasing function of price, and the rate of chemical processes is an increasing function of temperature. See our study on environmental data in Section 5.4. A simple approach to fit a monotone function is to select some basis functions $\{z_1(\cdot), \dots, z_d(\cdot)\}$ and define

$$z(\zeta^a; \mathbf{x}) = \sum_{i=1}^d x_i z_i(\zeta^a),$$

where $\mathbf{x} = (x_1, \dots, x_d)$ are regression coefficients. Then, we minimize the mean-squared error and impose the monotonicity, e.g., non-increasing, via constraints on the coefficients \mathbf{x} :

$$\min_{\mathbf{x}} \frac{1}{n} \sum_{i=1}^n F(\mathbf{x}; \zeta_i) := \frac{1}{n} \sum_{i=1}^n (\zeta_i^b - z(\zeta_i^a; \mathbf{x}))^2 \quad \text{s.t.} \quad z'(\omega_j; \mathbf{x}) := \sum_{i=1}^d x_i z'_i(\omega_j) \leq 0, \quad j \in [m],$$

where $\omega_1, \dots, \omega_m$ are a fine grid of points over the range of ζ^a .

Constraints are also imposed in other statistical problems. In portfolio optimization (cf.

Section 5.3), \mathbf{x} refers to the asset weights, and it is common to constrain the estimation within the probability simplex $\{\mathbf{x} \in \mathbb{R}^d : \mathbf{1}^\top \mathbf{x} = 1, \mathbf{x} \geq \mathbf{0}\}$ (Du et al., 2022). Alternatively, one may control the risk exposures to a known threshold by imposing box constraints $\|\mathbf{x}\|_\infty \leq c$, or control allocations among sectors by imposing affine constraints $A\mathbf{x} = \mathbf{b}$, where a negative weight signifies shorting the asset (Fan, 2007; Fan et al., 2012). In semiparametric index models, we impose $\{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|^2 = 1, x_1 > 0\}$ to ensure models’ identifiability (Na et al., 2019; Na and Kolar, 2021). In factor analysis, constraints can prevent Heywood cases that cause negative estimates for the variance (Shapiro, 1985). In algorithmic fairness, constraints are applied to prevent classifiers from yielding disparate outcomes based on sensitive features such as gender and ethnicity (Zafar et al., 2019). For more examples of constrained regression models, we refer to Sen, 1979; Dupacova and Wets, 1988; Nagaraj and Fuller, 1991; Shapiro, 2000.

A.2 Physics-informed machine learning

A recent trending topic in machine learning is its applications on scientific problems, where models must adhere to domain knowledge (such as physical laws) that is often described by partial differential equation (PDE) constraints (Cuomo et al., 2022). Consider a smooth neural network model $z(\cdot; \mathbf{x}) : \zeta^a \rightarrow \zeta^b$, where $\zeta^a = (s, t)$ is the spatial-temporal pair, ζ^b is the measurement of some quantity, and \mathbf{x} denotes the network parameters (e.g., weight matrices). If $z(\cdot; \mathbf{x})$ models the transport of the quantity, then it should satisfy the transport equation (γ is given):

$$\partial_t z(\zeta^a; \mathbf{x}) + \gamma \cdot \partial_s z(\zeta^a; \mathbf{x}) = \mathbf{0}. \quad (\text{A.1})$$

Applying PDE constraints to the model $z(\cdot; \mathbf{x})$ leads to many important network architectures, including Neural ODEs (Chen et al., 2018), DeepONets (Wang et al., 2021), and physics-informed neural networks (Raissi et al., 2019). In these problems, the desired model is solved from the following constrained optimization:

$$\min_{\mathbf{x}} \mathbb{E}_{\zeta \sim \mathcal{P}}[F(\mathbf{x}; \zeta)] \quad \text{s.t.} \quad \mathcal{C}_{\text{PDE}}(z(\omega_j; \mathbf{x})) = \mathbf{0}, \quad j \in [m],$$

where $F(\mathbf{x}; \zeta)$ is the data fitting loss (e.g., mean-squared error in Section A.1), $(\omega_1, \dots, \omega_m)$ are sampled points in the space of $\zeta^{\mathbf{a}}$ including both interior and boundary points, and $\mathcal{C}_{\text{PDE}}(\cdot)$ is a PDE operator with respect to $\zeta^{\mathbf{a}}$ that encodes domain knowledge (e.g., transport equation as in (A.1)). Numerous research has shown that failure to enforce these constraints can lead to serious interpretability and instability issues in neural network models (Krishnapriyan et al., 2021; Négier et al., 2023).

B Preparation Definitions and Lemmas

Lemma B.1 (Theorem 3 in Robinson, 1976). *If $\ell \leq \bar{\mathbf{x}} \leq \mathbf{u}$ satisfies EGMFCQ, then there exists a neighborhood $\mathcal{B}(\bar{\mathbf{x}}; \bar{r}) := \{\ell \leq \mathbf{x} \leq \mathbf{u} : \|\mathbf{x} - \bar{\mathbf{x}}\| \leq \bar{r}\}$ for some sufficiently small radius $\bar{r} > 0$, such that all points in the neighborhood also satisfy EGMFCQ.*

Definition B.2 (Regular mapping sequence, Vaart, 1998). *Let $\mathcal{W} \subseteq \mathbb{R}^r$ be a neighborhood of the origin. A sequence of mappings $\Gamma_k : \mathcal{W} \rightarrow \mathbb{R}^d$ is said to be regular with derivative $\mathbf{D} \in \mathbb{R}^{d \times r}$ if*

$$\lim_{k \rightarrow \infty} \sqrt{k} (\Gamma_k(\mathbf{w}) - \Gamma_k(\mathbf{0})) = \mathbf{D}\mathbf{w}, \quad \text{for all } \mathbf{w} \in \mathcal{W}.$$

Lemma B.3 (Hájek and Le Cam’s local minimax theorem, Vaart, 1998). *Let $\{\mathcal{T}_k, \mathcal{F}_k, \mathcal{P}_{k,\mathbf{v}}\}_{\mathbf{v} \in \mathcal{W}}$ be a locally asymptotically normal family with precision $\mathbf{Q} \succeq \mathbf{0}$. Let $\Gamma_k : \mathcal{W} \rightarrow \mathbb{R}^d$ be a regular mapping sequence with derivative \mathbf{D} , and let $\ell : \mathbb{R}^d \rightarrow [0, \infty)$ be any symmetric, quasiconvex, and lower semicontinuous function. Then, for any sequence of estimators $T_k : \mathcal{T}_k \rightarrow \mathbb{R}^d$, the following inequality holds:*

$$\sup_{\mathcal{W}_0 \subseteq \mathcal{W}, |\mathcal{W}_0| < \infty} \liminf_{k \rightarrow \infty} \max_{\mathbf{v} \in \mathcal{W}_0} \mathbb{E}_{\mathcal{P}_{k,\mathbf{v}}} [\ell(\sqrt{k}(T_k - \Gamma_k(\mathbf{v})))] \geq \mathbb{E} [\ell(Z)],$$

where $Z \sim \mathcal{N}(\mathbf{0}, \mathbf{D}\mathbf{Q}^{-1}\mathbf{D}^\top)$ when \mathbf{Q} is invertible. If \mathbf{Q} is singular, then the above inequality holds with $Z \sim \mathcal{N}(\mathbf{0}, \mathbf{D}(\mathbf{Q} + \gamma\mathbf{I})^{-1}\mathbf{D}^\top)$ for any $\gamma > 0$. Here, $\mathbb{E}_{\mathcal{P}_{k,\mathbf{v}}}[\cdot]$ is taken over the randomness of samples \mathcal{T}_k .

Lemma B.4 (Lemma B.3 in Na and Mahoney, 2025). *Let $\alpha_k = \iota_1(k+1)^{-b_1}$ and $\beta_k = \iota_2(k+1)^{-b_2}$ for some $\iota_1, \iota_2 > 0$ and $b_1, b_2 > 0$. Then the following statements hold:*

(a) Define $\chi = 0$ if $0 < b_2 < 1$ and $\chi = -b_1/\iota_2$ if $b_2 = 1$. If $\sum_{t=1}^l a_t + \chi > 0$, then

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

Moreover, for any $b \in \mathbb{R}$ and any sequence e_k satisfying $\lim_{k \rightarrow \infty} e_k \rightarrow 0$, we have

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

(b) Suppose $b_2 = 1$ and let $\chi = -b_1/\iota_2$. If $\sum_{t=1}^l a_t + p\chi > 0$ for some $p \in (0, 1)$, then

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

(c) If $0 < b_2 < b_1 \leq 1$, then

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

Lemma B.5 (Stability of Quadratic Programs, Theorem 2.1 in [Daniel, 1973](#)). *Consider two constrained strongly convex quadratic problems*

$$\mathbf{y}^* \in \min_{\mathbf{y} \in \Omega} \mathbf{g}^\top \mathbf{y} + \frac{1}{2} \mathbf{y}^\top \mathbf{Q} \mathbf{y} \quad \text{and} \quad \mathbf{y}^{**} \in \min_{\mathbf{y} \in \Omega} \mathbf{g}'^\top \mathbf{y} + \frac{1}{2} \mathbf{y}^\top \mathbf{Q}' \mathbf{y},$$

where the feasible region Ω is convex. Suppose $\mathbf{Q}, \mathbf{Q}' \succeq \nu \mathbf{I}$ and $\max\{\|\mathbf{y}^*\|, \|\mathbf{y}^{**}\|\} \leq M_y$ for some $\nu, M_y > 0$. Then,

$$\|\mathbf{y}^* - \mathbf{y}^{**}\| \leq \nu^{-1} (1 + M_y) \max\{\|\mathbf{g} - \mathbf{g}'\|, \|\mathbf{Q} - \mathbf{Q}'\|\}.$$

C Proofs of Main Results

C.1 Proof of Theorem 2.2

Suppose EGMFCQ holds at $\boldsymbol{\ell} \leq \mathbf{x} \leq \mathbf{u}$, we denote $\mathcal{A}(\mathbf{x}) = \mathcal{A}_{\boldsymbol{\ell}}(\mathbf{x}) \cup \mathcal{A}_{\mathbf{u}}(\mathbf{x})$, $\mathcal{A}(\mathbf{x})^- = [d] \setminus \mathcal{A}(\mathbf{x})$, and let \mathbf{z} satisfy (11). If $\mathbf{z} = \mathbf{0}$, then we define $\bar{\theta} = 1$. Otherwise, we define

$$\bar{\theta} = \frac{\epsilon}{\|\mathbf{z}\|} \quad \text{where} \quad \epsilon = \min \{ [\mathbf{u} - \boldsymbol{\ell}]_i, [\mathbf{x} - \boldsymbol{\ell}]_j, [\mathbf{u} - \mathbf{x}]_j : i \in \mathcal{A}(\mathbf{x}), j \in \mathcal{A}(\mathbf{x})^- \}. \quad (\text{C.1})$$

Note that $\bar{\theta} > 0$ since $\epsilon > 0$, although we do not restrict $\bar{\theta} \leq 1$. We claim $\Omega(\mathbf{x}; \bar{\theta}) \neq \emptyset$. In fact,

$\bar{\theta}\mathbf{z} \in \Omega(\mathbf{x}; \bar{\theta})$ since $\bar{\theta}\mathbf{c}(\mathbf{x}) + \nabla\mathbf{c}(\mathbf{x})(\bar{\theta}\mathbf{z}) \stackrel{(11)}{=} \mathbf{0}$ and

$$\begin{aligned}\ell_{\mathcal{A}_\ell(\mathbf{x})} &\stackrel{(8)}{=} \mathbf{x}_{\mathcal{A}_\ell(\mathbf{x})} \stackrel{(11)}{\leq} [\mathbf{x} + \bar{\theta}\mathbf{z}]_{\mathcal{A}_\ell(\mathbf{x})} \stackrel{(8)}{=} [\ell + \bar{\theta}\mathbf{z}]_{\mathcal{A}_\ell(\mathbf{x})} \leq [\ell + \epsilon\mathbf{1}]_{\mathcal{A}_\ell(\mathbf{x})} \stackrel{(C.1)}{\leq} \mathbf{u}_{\mathcal{A}_\ell(\mathbf{x})}, \\ \ell_{\mathcal{A}_u(\mathbf{x})} &\stackrel{(C.1)}{\leq} [\mathbf{u} - \epsilon\mathbf{1}]_{\mathcal{A}_u(\mathbf{x})} \leq [\mathbf{u} + \bar{\theta}\mathbf{z}]_{\mathcal{A}_u(\mathbf{x})} \stackrel{(8)}{=} [\mathbf{x} + \bar{\theta}\mathbf{z}]_{\mathcal{A}_u(\mathbf{x})} \stackrel{(11)}{\leq} \mathbf{x}_{\mathcal{A}_u(\mathbf{x})} \stackrel{(8)}{=} \mathbf{u}_{\mathcal{A}_u(\mathbf{x})}, \\ \ell_{\mathcal{A}(\mathbf{x})^-} &\stackrel{(C.1)}{\leq} [\mathbf{x} - \epsilon\mathbf{1}]_{\mathcal{A}(\mathbf{x})^-} \leq [\mathbf{x} + \bar{\theta}\mathbf{z}]_{\mathcal{A}(\mathbf{x})^-} \leq [\mathbf{x} + \epsilon\mathbf{1}]_{\mathcal{A}(\mathbf{x})^-} \stackrel{(C.1)}{\leq} \mathbf{u}_{\mathcal{A}(\mathbf{x})^-},\end{aligned}$$

where $\mathbf{1} = (1, \dots, 1) \in \mathbb{R}^d$ denotes the all-one vector and the above display also uses $-\mathbf{1} \leq \mathbf{z}/\|\mathbf{z}\| \leq \mathbf{1}$. From the above display, we know $\ell \leq \mathbf{x} + \bar{\theta}\mathbf{z} \leq \mathbf{u}$; thus, $\bar{\theta}\mathbf{z} \in \Omega(\mathbf{x}; \bar{\theta})$. Now, we show that $\Omega(\mathbf{x}; \theta) \neq \emptyset$ for any $\theta \in [0, \bar{\theta}]$. In fact, for any $\mathbf{v} \in \Omega(\mathbf{x}; \bar{\theta})$, we consider $\theta/\bar{\theta} \cdot \mathbf{v}$. First, we have $\theta\mathbf{c}(\mathbf{x}) + \nabla\mathbf{c}(\mathbf{x})(\theta/\bar{\theta}\mathbf{v}) = \mathbf{0}$ since $\bar{\theta}\mathbf{c}(\mathbf{x}) + \nabla\mathbf{c}(\mathbf{x})\mathbf{v} = \mathbf{0}$. Second, we bound $\mathbf{x} + \theta/\bar{\theta}\mathbf{v}$ by considering two cases based on the sign of \mathbf{v} . We apply $\ell \leq \mathbf{x} + \mathbf{v} \leq \mathbf{u}$ and have

$$\ell_{\{i:v_i \geq 0\}} \leq \mathbf{x}_{\{i:v_i \geq 0\}} \leq [\mathbf{x} + \theta/\bar{\theta}\mathbf{v}]_{\{i:v_i \geq 0\}} \leq [\mathbf{x} + \mathbf{v}]_{\{i:v_i \geq 0\}} \leq \mathbf{u}_{\{i:v_i \geq 0\}},$$

$$\ell_{\{i:v_i < 0\}} \leq [\mathbf{x} + \mathbf{v}]_{\{i:v_i < 0\}} \leq [\mathbf{x} + \theta/\bar{\theta}\mathbf{v}]_{\{i:v_i < 0\}} \leq \mathbf{x}_{\{i:v_i < 0\}} \leq \mathbf{u}_{\{i:v_i < 0\}}.$$

Thus, $\ell \leq \mathbf{x} + \theta/\bar{\theta}\mathbf{v} \leq \mathbf{u}$. This completes the proof of the first part of the theorem.

For the second part of the theorem, without loss of generality, we suppose $\lim_{k \rightarrow \infty} \theta_k = 0$ and $\lim_{k \rightarrow \infty} \mathbf{x}_k = \mathbf{x}^*$. Suppose \mathbf{x}^* satisfies EGMFCQ, then there exists \mathbf{z}^* such that (11) holds. On the other hand, by Lemma B.1 we know \mathbf{x}_k satisfies EGMFCQ for k large enough, implying that $\nabla\mathbf{c}_k$ has full row rank. Noting that

$$\mathbf{z}^* = -(\nabla\mathbf{c}^*)^\top (\nabla\mathbf{c}^* (\nabla\mathbf{c}^*)^\top)^{-1} \mathbf{c}^* + (\mathbf{I} - (\nabla\mathbf{c}^*)^\top (\nabla\mathbf{c}^* (\nabla\mathbf{c}^*)^\top)^{-1} \nabla\mathbf{c}^*) \mathbf{z}^*,$$

we define

$$\mathbf{z}_k = -\nabla\mathbf{c}_k^\top (\nabla\mathbf{c}_k \nabla\mathbf{c}_k^\top)^{-1} \mathbf{c}_k + (\mathbf{I} - \nabla\mathbf{c}_k^\top (\nabla\mathbf{c}_k \nabla\mathbf{c}_k^\top)^{-1} \nabla\mathbf{c}_k) \mathbf{z}^*$$

and have $\mathbf{c}_k + \nabla\mathbf{c}_k \mathbf{z}_k = \mathbf{0}$ and $\mathbf{z}_k \rightarrow \mathbf{z}^*$ as $k \rightarrow \infty$. Since $\mathcal{A}(\mathbf{x}^*)^- \subseteq \mathcal{A}(\mathbf{x}_k)^-$, we know $\mathcal{A}(\mathbf{x}_k) \subseteq \mathcal{A}(\mathbf{x}^*)$. Since each entry of $\mathbf{z}_{\mathcal{A}(\mathbf{x}^*)}^*$ is nonzero by (11), we know each entry of $[\mathbf{z}_k]_{\mathcal{A}(\mathbf{x}_k)}$ is also nonzero and has the same sign as $\mathbf{z}_{\mathcal{A}(\mathbf{x}_k)}^*$. This shows that the constructed \mathbf{z}_k satisfies (11) for \mathbf{x}_k . Let $l_k := \inf\{\|\mathbf{z}\| : \mathbf{z} \text{ satisfies (11) at } \mathbf{x}_k\}$. For any \mathbf{z} satisfies (11) at \mathbf{x}_k with k large

enough, we can follow the analysis of the first part, define

$$\epsilon_k = \min\{0.5[\mathbf{u} - \boldsymbol{\ell}]_i, [\mathbf{x}_k - \boldsymbol{\ell}]_j, [\mathbf{u} - \mathbf{x}_k]_j : i \in \mathcal{A}(\mathbf{x}^*), j \in \mathcal{A}(\mathbf{x}^*)^-\},$$

and show $\Omega(\mathbf{x}_k; \epsilon_k/\|\mathbf{z}\|) \neq \emptyset$. Since $\epsilon_k \geq \epsilon$ for some $\epsilon > 0$ independent of k , we know $\theta_k \rightarrow 0$ implies $\limsup_{k \rightarrow \infty} l_k = \infty$. This leads to a contradiction since $\infty = \limsup_{k \rightarrow \infty} l_k \leq \|\mathbf{z}^*\| =: l^* < \infty$. Thus, EGMFCQ does not hold at \mathbf{x}^* . This completes the proof of the second part of the theorem.

C.2 Proof of Theorem 2.5

First, by Assumption 2.3, we know $\boldsymbol{\mu}_{\mathcal{A}^*}^* > \mathbf{0}$. For any KKT triplet $(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$ that is close to $(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$, the corresponding active and inactive sets coincide with those at $(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$. In particular, for any $i \in \mathcal{A}^* = \mathcal{A}_\ell^* \cup \mathcal{A}_u^*$, we have $\boldsymbol{\mu}_{\mathcal{A}^*}^* > \mathbf{0}$ since $\boldsymbol{\mu}_{\mathcal{A}^*}^* > \mathbf{0}$, which implies $\mathcal{A}^* \subseteq \mathcal{A}(\mathbf{x})$ due to the KKT conditions (7) at $(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$. On the other hand, since $\boldsymbol{\ell}_{(\mathcal{A}^*)^-} < \mathbf{x}_{(\mathcal{A}^*)^-}^* < \mathbf{u}_{(\mathcal{A}^*)^-}$, the above inequalities also hold for $\mathbf{x}_{(\mathcal{A}^*)^-}$, leading to $(\mathcal{A}^*)^- \subseteq \mathcal{A}(\mathbf{x})^-$ and hence $\mathcal{A}(\mathbf{x}) \subseteq \mathcal{A}^*$. This shows $\mathcal{A}(\mathbf{x}) = \mathcal{A}^*$.

With this observation, we define the Lagrangian gradient with respect to $\mathbf{w} = (\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}_{\mathcal{A}^*}) \in \mathbb{R}^{d+m+|\mathcal{A}_\ell^*|+|\mathcal{A}_u^*|}$ as

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = \begin{pmatrix} \nabla f(\mathbf{x}) + \nabla \mathbf{c}(\mathbf{x})^\top \boldsymbol{\lambda} - \mathbf{I}_{\mathcal{A}_\ell^*}^\top [\boldsymbol{\mu}_1]_{\mathcal{A}_\ell^*} + \mathbf{I}_{\mathcal{A}_u^*}^\top [\boldsymbol{\mu}_2]_{\mathcal{A}_u^*} \\ \mathbf{c}(\mathbf{x}) \\ [\boldsymbol{\ell} - \mathbf{x}]_{\mathcal{A}_\ell^*} \\ [\mathbf{x} - \mathbf{u}]_{\mathcal{A}_u^*} \end{pmatrix} \in \mathbb{R}^{d+m+|\mathcal{A}_\ell^*|+|\mathcal{A}_u^*|}, \quad (\text{C.2})$$

and define the mapping $\sigma(\boldsymbol{\delta}) = \{\mathbf{w} : \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = \boldsymbol{\delta}\}$. This mapping characterizes the set of primal-dual points satisfying the perturbed optimality conditions up to a residual vector $\boldsymbol{\delta}$. In fact, $\sigma(\boldsymbol{\delta})$ is a single-valued mapping in a neighborhood of $(\boldsymbol{\delta}, \mathbf{w}) = (\mathbf{0}, \mathbf{w}^*)$ with $\nabla \sigma(\boldsymbol{\delta})|_{\boldsymbol{\delta}=\mathbf{0}} = (\mathbf{H}^*)^{-1}$. To see this, let $\hat{\mathbf{w}}, \bar{\mathbf{w}} \in \sigma(\boldsymbol{\delta})$ for sufficiently small $\boldsymbol{\delta}$. By SOSC and LICQ in Assumption 2.4 and Definition 2.1, we know $\nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{w}^*) = \mathbf{H}^*$ is nonsingular. Thus, there exists a neighborhood of \mathbf{w}^* within which $\nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{w})$ remains nonsingular. We then have

$$\mathbf{0} = \nabla_{\mathbf{w}} \mathcal{L}(\hat{\mathbf{w}}) - \nabla_{\mathbf{w}} \mathcal{L}(\bar{\mathbf{w}}) = \left(\int_0^1 \nabla_{\mathbf{w}}^2 \mathcal{L}(\bar{\mathbf{w}} + t(\hat{\mathbf{w}} - \bar{\mathbf{w}})) dt \right) \cdot (\hat{\mathbf{w}} - \bar{\mathbf{w}}). \quad (\text{C.3})$$

Since both $\hat{\mathbf{w}}$ and $\bar{\mathbf{w}}$ lie sufficiently close to \mathbf{w}^* , the integral matrix is nonsingular, which implies

$\widehat{\mathbf{w}} - \bar{\mathbf{w}} = \mathbf{0}$. Hence, $\sigma(\cdot)$ is locally single-valued. Furthermore, by (C.3) we have

$$\boldsymbol{\delta} = \nabla_{\mathbf{w}} \mathcal{L}(\widehat{\mathbf{w}}) - \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}^*) = \left(\int_0^1 \nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{w}^* + t(\widehat{\mathbf{w}} - \mathbf{w}^*)) dt \right) \cdot (\widehat{\mathbf{w}} - \mathbf{w}^*),$$

leading to

$$\sigma(\boldsymbol{\delta}) = \widehat{\mathbf{w}} = \mathbf{w}^* + \left(\int_0^1 \nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{w}^* + t(\widehat{\mathbf{w}} - \mathbf{w}^*)) dt \right)^{-1} \boldsymbol{\delta}.$$

Differentiating $\sigma(\boldsymbol{\delta})$ with respect to $\boldsymbol{\delta}$ at $\boldsymbol{\delta} = \mathbf{0}$ yields

$$\nabla \sigma(\boldsymbol{\delta})|_{\boldsymbol{\delta}=\mathbf{0}} = \lim_{\boldsymbol{\delta} \rightarrow \mathbf{0}} \left(\int_0^1 \nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{w}^* + t(\widehat{\mathbf{w}} - \mathbf{w}^*)) dt \right)^{-1} = (\nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{w}^*))^{-1} = (\mathbf{H}^*)^{-1}.$$

Let us use \mathcal{T} to denote the corresponding sample space of the probability measure \mathcal{P} . We define the function class

$$\mathcal{G} = \left\{ \mathbf{g} : \mathcal{T} \rightarrow \mathbb{R}^{d+m+|\mathcal{A}_{\ell}^*|+|\mathcal{A}_{\mathbf{u}}^*|} : \mathbb{E}_{\zeta \sim \mathcal{P}} [\mathbf{g}(\zeta)] = \mathbf{0}, \mathbb{E}_{\zeta \sim \mathcal{P}} \|\mathbf{g}(\zeta)\|^2 < \infty \right\}.$$

Fix an arbitrary function $\mathbf{g} \in \mathcal{G}$, and let $h : \mathbb{R} \rightarrow [-1, 1]$ be any C^3 -smooth function whose first three derivatives are globally bounded and that satisfies $h(t) = t$ for $t \in [-1/2, 1/2]$. For each $\mathbf{v} \in \mathbb{R}^{d+m+|\mathcal{A}_{\ell}^*|+|\mathcal{A}_{\mathbf{u}}^*|}$, we define a perturbed distribution whose density with respect to \mathcal{P} is

$$d\mathcal{P}_{\mathbf{v}}(\zeta) = \frac{1 + h(\mathbf{v}^\top \mathbf{g}(\zeta))}{C(\mathbf{v})} d\mathcal{P}(\zeta) \quad \text{where} \quad C(\mathbf{v}) = 1 + \int h(\mathbf{v}^\top \mathbf{g}(\zeta)) d\mathcal{P}(\zeta).$$

Each vector \mathbf{v} specifies a perturbed optimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^d} f_{\mathbf{v}}(\mathbf{x}) = \mathbb{E}_{\zeta \sim \mathcal{P}_{\mathbf{v}}} [F(\mathbf{x}; \zeta)] \quad \text{s.t.} \quad \mathbf{c}(\mathbf{x}) = \mathbf{0}, \quad \boldsymbol{\ell} \leq \mathbf{x} \leq \mathbf{u}. \quad (\text{C.4})$$

We further define the associated perturbed optimality mapping $\nabla_{\mathbf{w}} \mathcal{L}_{\mathbf{v}}(\mathbf{w})$ analogously to (C.2), except that $\nabla f(\mathbf{x})$ is replaced by $\nabla f_{\mathbf{v}}(\mathbf{x})$. We introduce the following two lemmas regarding the mapping $\nabla_{\mathbf{w}} \mathcal{L}_{\mathbf{v}}(\mathbf{w})$, which are proved in Appendices D.1 and D.2, respectively.

Lemma C.1. *The mapping $\nabla_{\mathbf{w}} \mathcal{L}_{\mathbf{v}}(\mathbf{w})$ is continuously differentiable in a neighborhood of $(\mathbf{w}, \mathbf{v}) = (\mathbf{w}^*, \mathbf{0})$, with derivatives given by*

$$\nabla_{\mathbf{w}}^2 \mathcal{L}_{\mathbf{v}}(\mathbf{w}^*)|_{\mathbf{v}=\mathbf{0}} = \nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{w}^*) = \mathbf{H}^*, \quad \nabla_{\mathbf{v}\mathbf{w}} \mathcal{L}_{\mathbf{v}}(\mathbf{w})|_{\mathbf{v}=\mathbf{0}} = \mathbb{E}_{\zeta \sim \mathcal{P}} [\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}; \zeta) \mathbf{g}(\zeta)^\top], \quad (\text{C.5})$$

where $\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}; \zeta)$ has the same form as $\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w})$ in (C.2) but replaces $\nabla f(\mathbf{x})$ by $\nabla F(\mathbf{x}; \zeta)$.

Lemma C.2. *The solution mapping $S(\mathbf{v}) = \{\mathbf{w} : \nabla_{\mathbf{w}} \mathcal{L}_{\mathbf{v}}(\mathbf{w}) = \mathbf{0}\}$ admits a single-valued lo-*

calization $s(\mathbf{v})$ in a neighborhood of $(\mathbf{w}, \mathbf{v}) = (\mathbf{w}^*, \mathbf{0})$. Moreover, $s(\cdot)$ is differentiable at $\mathbf{v} = \mathbf{0}$ with Jacobian $\nabla s(\mathbf{0}) = -(\mathbf{H}^*)^{-1} \mathbb{E}_{\zeta \sim \mathcal{P}} [\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}^*; \zeta) \mathbf{g}(\zeta)^\top]$.

With the above two lemmas, we then consider the mapping sequence $\Gamma_k(\mathbf{v}) = s(\mathbf{v}/\sqrt{k})$, where $s(\cdot)$ is the single-valued localization of the solution map from Lemma C.2. In fact, $\Gamma_k(\mathbf{v})$ is regular with derivative $\nabla s(\mathbf{0})$ (cf. Definition B.2). To see this, we just note that by the differentiability of $s(\mathbf{v})$ at $\mathbf{v} = \mathbf{0}$, we have for all small enough \mathbf{v} ,

$$\lim_{k \rightarrow \infty} \sqrt{k} (\Gamma_k(\mathbf{v}) - \Gamma_k(\mathbf{0})) = \lim_{k \rightarrow \infty} \sqrt{k} (s(\mathbf{v}/\sqrt{k}) - s(\mathbf{0})) = \nabla s(\mathbf{0}) \mathbf{v}.$$

Now, we are ready to combine all above pieces to finalize the proof of Theorem 2.5. For each $\mathbf{v} \in \mathbb{R}^{d+m+|\mathcal{A}_{\mathbf{t}}^*|+|\mathcal{A}_{\mathbf{u}}^*|}$, we define the product probability space $(\mathcal{T}_k, \mathcal{F}_k, \mathcal{P}_{k,\mathbf{v}}) := (\mathcal{T}, \mathcal{F}, \mathcal{P}_{\mathbf{v}/\sqrt{k}})^{\otimes k}$, that is, the k -fold product of $(\mathcal{T}, \mathcal{F}, \mathcal{P}_{\mathbf{v}/\sqrt{k}})$. It has been shown in (Duchi and Ruan, 2021, Lemma 8.3) that the sequence $(\mathcal{T}_k, \mathcal{F}_k, \mathcal{P}_{k,\mathbf{v}})$ is locally asymptotically normal with the precision matrix $\mathbf{Q} = \mathbb{E}_{\zeta \sim \mathcal{P}} [\mathbf{g}(\zeta) \mathbf{g}(\zeta)^\top]$. Let us specify

$$\mathbf{g}(\zeta) = \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}^*; \zeta) - \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}^*) = \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}^*; \zeta),$$

then we have $\Sigma^* = \mathbb{E}_{\zeta \sim \mathcal{P}} [\mathbf{g}(\zeta) \mathbf{g}(\zeta)^\top]$ (see (13)). Let ℓ be any symmetric, quasiconvex, and lower semicontinuous function, and let \mathbf{w}_k denote any sequence of estimators, and $\mathbf{w}_{\mathbf{v}}^*$ denote the primal and (active) dual solution of the perturbed problem (C.4) for small enough \mathbf{v} , i.e., $\mathbf{w}_{\mathbf{v}}^* = s(\mathbf{v})$. Note that the active set is consistent with the unperturbed active set \mathcal{A}^* . For any finite subset $\mathcal{W}_0 \subseteq \mathbb{R}^{d+m+|\mathcal{A}_{\mathbf{t}}^*|+|\mathcal{A}_{\mathbf{u}}^*|}$, we define $c = \max_{\mathbf{v} \in \mathcal{W}_0} \|\mathbf{v}\|$. Since $\mathbf{w}_{\mathbf{v}}^* = s(\mathbf{v})$ and $\Gamma_k(\mathbf{v}) = s(\mathbf{v}/\sqrt{k})$, we have

$$\liminf_{k \rightarrow \infty} \sup_{\|\mathbf{v}\| \leq c/\sqrt{k}} \mathbb{E}_{\mathcal{P}_{\mathbf{v}}^{\otimes k}} [\ell(\sqrt{k}(\mathbf{w}_k - \mathbf{w}_{\mathbf{v}}^*))] \geq \liminf_{k \rightarrow \infty} \max_{\mathbf{v} \in \mathcal{W}_0} \mathbb{E}_{\mathcal{P}_{k,\mathbf{v}}} [\ell(\sqrt{k}(\mathbf{w}_k - \Gamma_k(\mathbf{v})))] .$$

Taking the supremum over all finite $\mathcal{W}_0 \subseteq \mathbb{R}^{d+m+|\mathcal{A}_{\mathbf{t}}^*|+|\mathcal{A}_{\mathbf{u}}^*|}$ (equivalently, letting $c \rightarrow \infty$), and applying the Hájek and Le Cam's local minimax theorem in Theorem B.3, we obtain

$$\lim_{c \rightarrow \infty} \liminf_{k \rightarrow \infty} \sup_{\|\mathbf{v}\| \leq c/\sqrt{k}} \mathbb{E}_{\mathcal{P}_{\mathbf{v}}^{\otimes k}} [\ell(\sqrt{k}(\mathbf{w}_k - \mathbf{w}_{\mathbf{v}}^*))] \geq \mathbb{E} [\ell(Z_\gamma)],$$

where $Z_\gamma \sim \mathcal{N}(\mathbf{0}, (\mathbf{H}^*)^{-1} \Sigma^* (\Sigma^* + \gamma \mathbf{I})^{-1} \Sigma^* (\mathbf{H}^*)^{-1})$ for any $\gamma > 0$. Since $\lim_{\gamma \rightarrow 0} \Sigma^* (\Sigma^* + \gamma \mathbf{I})^{-1} \Sigma^* (\mathbf{H}^*)^{-1} = \Sigma^* (\Sigma^* + \gamma \mathbf{I})^{-1} \Sigma^* (\mathbf{H}^*)^{-1}$ for any $\gamma > 0$. Since $\lim_{\gamma \rightarrow 0} \Sigma^* (\Sigma^* + \gamma \mathbf{I})^{-1} \Sigma^* (\mathbf{H}^*)^{-1} = \Sigma^* (\Sigma^* + \gamma \mathbf{I})^{-1} \Sigma^* (\mathbf{H}^*)^{-1}$ for any $\gamma > 0$.

$\gamma \mathbf{I})^{-1} \Sigma^* = \Sigma^*$, the right-hand side converges to $\mathbb{E}[\ell(Z)]$ with $Z \sim \mathcal{N}(\mathbf{0}, \Omega^*)$. For the left-hand side, it has been shown in [Duchi and Ruan, 2021](#) that \mathcal{P}_v with $\|v\| \leq c/\sqrt{k}$ implies $\mathcal{P}_v \in \mathcal{B}(\mathcal{P}; c'/k)$ for some $c' > 0$. This establishes the claimed lower bound.

C.3 Proof of Lemma 3.7

We first state a preliminary lemma. The proof is omitted, as it follows directly from the observation that $\mathbf{x} + \alpha \Delta \mathbf{x} = (1 - \alpha)\mathbf{x} + \alpha(\mathbf{x} + \Delta \mathbf{x})$ and that the box constraint set is convex.

Lemma C.3. *Suppose $\ell \leq \mathbf{x} \leq \mathbf{u}$ and $\ell \leq \mathbf{x} + \Delta \mathbf{x} \leq \mathbf{u}$. Then, for any $\alpha \in [0, 1]$, $\ell \leq \mathbf{x} + \alpha \Delta \mathbf{x} \leq \mathbf{u}$. In particular, if $\ell \leq \mathbf{x}_0 \leq \mathbf{u}$, then the entire SSQP sequence $\{\mathbf{x}_k\}$ satisfies $\ell \leq \mathbf{x}_k \leq \mathbf{u}$.*

The KKT conditions for the SQP subproblem at \mathbf{x}_k with true gradient ∇f_k show that there exist some dual multipliers $(\lambda_k^{\text{sub}}, \mu_k^{\text{sub}})$ satisfying

$$\begin{aligned} \nabla f_k + \bar{\mathbf{B}}_k \Delta \mathbf{x}_k + \nabla \mathbf{c}_k^\top \lambda_k^{\text{sub}} - \mu_{1,k}^{\text{sub}} + \mu_{2,k}^{\text{sub}} &= \mathbf{0}, \\ \theta_k \mathbf{c}_k + \nabla \mathbf{c}_k \Delta \mathbf{x}_k &= \mathbf{0}, \quad \ell \leq \mathbf{x}_k + \Delta \mathbf{x}_k \leq \mathbf{u}, \\ \mu_{1,k}^{\text{sub}\top} (\ell - \mathbf{x}_k - \Delta \mathbf{x}_k) &= 0, \quad \mu_{2,k}^{\text{sub}\top} (\mathbf{x}_k + \Delta \mathbf{x}_k - \mathbf{u}) = 0, \\ \mu_{1,k}^{\text{sub}}, \mu_{2,k}^{\text{sub}} &\geq \mathbf{0}. \end{aligned} \tag{C.6}$$

Multiplying both sides of the first equality by $\Delta \mathbf{x}_k$, we obtain

$$\begin{aligned} \nabla f_k^\top \Delta \mathbf{x}_k + \Delta \mathbf{x}_k^\top \bar{\mathbf{B}}_k \Delta \mathbf{x}_k &= -\Delta \mathbf{x}_k^\top \nabla \mathbf{c}_k^\top \lambda_k^{\text{sub}} + \Delta \mathbf{x}_k^\top \mu_{1,k}^{\text{sub}} - \Delta \mathbf{x}_k^\top \mu_{2,k}^{\text{sub}} \\ &\stackrel{\text{(C.6)}}{=} \theta_k \lambda_k^{\text{sub}\top} \mathbf{c}_k - \mu_{1,k}^{\text{sub}\top} (\mathbf{x}_k - \ell) + \mu_{2,k}^{\text{sub}\top} (\mathbf{x}_k - \mathbf{u}) \\ &\leq \theta_k \lambda_k^{\text{sub}\top} \mathbf{c}_k \leq \|\lambda_k^{\text{sub}}\| \|\mathbf{c}_k\| \leq M_{\text{dual}} \|\mathbf{c}_k\|, \end{aligned} \tag{C.7}$$

where the third inequality comes from $\mu_{1,k}^{\text{sub}}, \mu_{2,k}^{\text{sub}} \geq \mathbf{0}$ and $\ell \leq \mathbf{x}_k \leq \mathbf{u}$ (by Lemma C.3). By the definitions (21) and (22), we have

$$\Delta \varphi_\rho^{\text{loc}}(\mathbf{x}_k, \Delta \mathbf{x}_k, \bar{\mathbf{B}}_k) = -\nabla f_k^\top \Delta \mathbf{x}_k - \frac{1}{2} \Delta \mathbf{x}_k^\top \bar{\mathbf{B}}_k \Delta \mathbf{x}_k + \rho \theta_k \|\mathbf{c}_k\|. \tag{C.8}$$

Thus, to satisfy $\Delta \varphi_\rho^{\text{loc}}(\mathbf{x}_k, \Delta \mathbf{x}_k, \bar{\mathbf{B}}_k) \geq \frac{1}{2} \Delta \mathbf{x}_k^\top \bar{\mathbf{B}}_k \Delta \mathbf{x}_k + \nu \rho \theta_k \|\mathbf{c}_k\|$, it suffices to satisfy

$$\nabla f_k^\top \Delta \mathbf{x}_k + \Delta \mathbf{x}_k^\top \bar{\mathbf{B}}_k \Delta \mathbf{x}_k \leq (1 - \nu) \rho \theta_k \|\mathbf{c}_k\|.$$

Combining with (C.7), we know the above display holds if

$$M_{\text{dual}} \|\mathbf{c}_k\| \leq (1 - \nu) \rho \theta_k \|\mathbf{c}_k\|,$$

which is further implied by the stated condition $\rho \geq \bar{\rho} := \frac{M_{\text{dual}}}{(1-\nu)\tau\bar{\theta}}$ since $\theta_k \geq \tau\bar{\theta}$ by Assumption

3.4. This completes the proof.

C.4 Proof of Theorem 3.8

To simplify notation, we slightly abuse the symbol $\mathbb{E}_k[\cdot]$ throughout the proof to denote the conditional expectation given the randomness of \mathbf{x}_0 through \mathbf{x}_k . Since $\bar{\alpha}_k \leq \alpha_k + \psi\alpha_k^p$ and $\theta_k \leq 1$, we may, without loss of generality, assume that $\bar{\alpha}_k\theta_k \leq 1$ (otherwise, we can simply restrict our analysis for $k \geq k_0$ for some deterministic k_0). For any $\rho > 0$, it follows from Assumptions 3.4, 3.5, and 3.6 that

$$\begin{aligned} & \varphi_\rho(\mathbf{x}_{k+1}) - \varphi_\rho(\mathbf{x}_k) \\ &= f(\mathbf{x}_k + \bar{\alpha}_k \bar{\Delta} \mathbf{x}_k) - f(\mathbf{x}_k) + \rho(\|\mathbf{c}(\mathbf{x}_k + \bar{\alpha}_k \bar{\Delta} \mathbf{x}_k)\| - \|\mathbf{c}_k\|) \\ &\leq \bar{\alpha}_k \nabla f_k^\top \bar{\Delta} \mathbf{x}_k + \frac{\kappa_{\nabla f}}{2} \bar{\alpha}_k^2 \|\bar{\Delta} \mathbf{x}_k\|^2 + \rho \left(\|\mathbf{c}_k + \bar{\alpha}_k \nabla \mathbf{c}_k \bar{\Delta} \mathbf{x}_k\| - \|\mathbf{c}_k\| + \frac{\kappa_{\nabla c}}{2} \bar{\alpha}_k^2 \|\bar{\Delta} \mathbf{x}_k\|^2 \right) \\ &= \bar{\alpha}_k \nabla f_k^\top \bar{\Delta} \mathbf{x}_k - \bar{\alpha}_k \rho \theta_k \|\mathbf{c}_k\| + \frac{\kappa_{\nabla f} + \rho \kappa_{\nabla c}}{2} \bar{\alpha}_k^2 \|\bar{\Delta} \mathbf{x}_k\|^2 \\ &= \bar{\alpha}_k \nabla f_k^\top \bar{\Delta} \mathbf{x}_k + \bar{\alpha}_k \nabla f_k^\top (\bar{\Delta} \mathbf{x}_k - \Delta \mathbf{x}_k) - \bar{\alpha}_k \rho \theta_k \|\mathbf{c}_k\| + \frac{\kappa_{\nabla f} + \rho \kappa_{\nabla c}}{2} \bar{\alpha}_k^2 \|\bar{\Delta} \mathbf{x}_k\|^2 \\ &\stackrel{(\text{C.8})}{=} -\bar{\alpha}_k \Delta \varphi_\rho^{\text{loc}}(\mathbf{x}_k, \Delta \mathbf{x}_k, \bar{\mathbf{B}}_k) - \frac{\bar{\alpha}_k}{2} \Delta \mathbf{x}_k^\top \bar{\mathbf{B}}_k \Delta \mathbf{x}_k + \bar{\alpha}_k \nabla f_k^\top (\bar{\Delta} \mathbf{x}_k - \Delta \mathbf{x}_k) + \frac{\kappa_{\nabla f} + \rho \kappa_{\nabla c}}{2} \bar{\alpha}_k^2 \|\bar{\Delta} \mathbf{x}_k\|^2 \\ &\leq -\bar{\alpha}_k \Delta \varphi_\rho^{\text{loc}}(\mathbf{x}_k, \Delta \mathbf{x}_k, \bar{\mathbf{B}}_k) + \bar{\alpha}_k \nabla f_k^\top (\bar{\Delta} \mathbf{x}_k - \Delta \mathbf{x}_k) + \frac{\kappa_{\nabla f} + \rho \kappa_{\nabla c}}{2} \bar{\alpha}_k^2 \|\bar{\Delta} \mathbf{x}_k\|^2 \\ &\leq -\bar{\alpha}_k \Delta \varphi_\rho^{\text{loc}}(\mathbf{x}_k, \Delta \mathbf{x}_k, \bar{\mathbf{B}}_k) + \bar{\alpha}_k M_{\nabla f} \|\bar{\Delta} \mathbf{x}_k - \Delta \mathbf{x}_k\| + \frac{\kappa_{\nabla f} + \bar{\rho} \kappa_{\nabla c}}{2} \bar{\alpha}_k^2 M_{\ell, u}^2, \end{aligned}$$

where the last inequality holds for some constants $M_{\nabla f}, M_{\ell, u}^2 > 0$ since box constraints form a compact set. Setting $\rho = \bar{\rho}$, we apply Lemma 3.7, Assumptions 3.4 and 3.5, and obtain

$$\Delta \varphi_{\bar{\rho}}^{\text{loc}}(\mathbf{x}_k, \Delta \mathbf{x}_k, \bar{\mathbf{B}}_k) \geq \frac{1}{2} \Delta \mathbf{x}_k^\top \bar{\mathbf{B}}_k \Delta \mathbf{x}_k + \nu \bar{\rho} \tau \bar{\theta} \|\mathbf{c}_k\| \geq \frac{\kappa_1}{2} \|\Delta \mathbf{x}_k\|^2 + \nu \bar{\rho} \tau \bar{\theta} \|\mathbf{c}_k\|.$$

Combining the above two displays and taking conditional expectation over \mathbf{x}_0 to \mathbf{x}_k , we obtain

$$\mathbb{E}_k \left[\varphi_{\bar{\rho}}(\mathbf{x}_{k+1}) - \min_{\ell \leq \mathbf{x} \leq \mathbf{u}} \varphi_{\bar{\rho}}(\mathbf{x}) \right] \leq \left\{ \varphi_{\bar{\rho}}(\mathbf{x}_k) - \min_{\ell \leq \mathbf{x} \leq \mathbf{u}} \varphi_{\bar{\rho}}(\mathbf{x}) \right\} - \alpha_k \mathbb{E}_k \left[\frac{\kappa_1}{2} \|\Delta \mathbf{x}_k\|^2 + \nu \bar{\rho} \tau \bar{\theta} \|\mathbf{c}_k\| \right]$$

$$+ M_{\nabla f} (\alpha_k + \psi \alpha_k^p) \mathbb{E}_k[\|\bar{\Delta} \mathbf{x}_k - \Delta \mathbf{x}_k\|] + \frac{\kappa_{\nabla f} + \bar{\rho} \kappa_{\nabla c}}{2} M_{\ell, \mathbf{u}}^2 (\alpha_k + \psi \alpha_k^p)^2, \quad \forall k \geq 0. \quad (\text{C.9})$$

Now, we introduce the following lemma, which is proved in Appendix D.3.

Lemma C.4. *Under the conditions of Theorem 3.8, we have $\sum_{k=0}^{\infty} \alpha_k \mathbb{E}[\|\bar{\Delta} \mathbf{x}_k - \Delta \mathbf{x}_k\|] < \infty$.*

With the above lemma, we directly have $\mathbb{E}[\sum_{k=0}^{\infty} \alpha_k \mathbb{E}_k[\|\bar{\Delta} \mathbf{x}_k - \Delta \mathbf{x}_k\|]] = \sum_{k=0}^{\infty} \alpha_k \mathbb{E}[\|\bar{\Delta} \mathbf{x}_k - \Delta \mathbf{x}_k\|] < \infty$, which implies that $\sum_{k=0}^{\infty} (\alpha_k + \psi \alpha_k^p) \mathbb{E}_k[\|\bar{\Delta} \mathbf{x}_k - \Delta \mathbf{x}_k\|] < \infty$ almost surely. From (C.9), together with the condition $\sum_{k=0}^{\infty} \alpha_k^2 < \infty$, we apply the Robbins-Siegmund theorem (Robbins and Siegmund, 1971) and obtain $\sum_{k=0}^{\infty} \alpha_k (\|\Delta \mathbf{x}_k\|^2 + \|\mathbf{c}_k\|) < \infty$. This completes the proof by noting that $\sum_{k=0}^{\infty} \alpha_k = \infty$.

C.5 Proof of Theorem 3.9

Let us denote the objective function of Problem (23) by

$$G(\boldsymbol{\lambda}, \boldsymbol{\mu}; \mathbf{x}_k) = \|\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}_k, \boldsymbol{\lambda}, \boldsymbol{\mu})\|^2 + \|\boldsymbol{\mu}_1 \odot (\boldsymbol{\ell} - \mathbf{x}_k)\|^2 + \|\boldsymbol{\mu}_2 \odot (\mathbf{x}_k - \mathbf{u})\|^2.$$

By the definition of the KKT residual $\mathbf{R}(\mathbf{x}_k, \boldsymbol{\lambda}_k^*, \boldsymbol{\mu}_k^*)$ in (20), it suffices to show $G(\boldsymbol{\lambda}_k^*, \boldsymbol{\mu}_k^*; \mathbf{x}_k) \rightarrow 0$ and $\|\mathbf{c}_k\| \rightarrow 0$ as $k \rightarrow \infty$ almost surely.

• **Convergence of $G(\boldsymbol{\lambda}_k^*, \boldsymbol{\mu}_k^*; \mathbf{x}_k)$.** Recall from Assumption 3.5 that $(\boldsymbol{\lambda}_k^{\text{sub}}, \boldsymbol{\mu}_k^{\text{sub}})$ denotes the dual multipliers of the SQP subproblem (17) with the true gradient ∇f_k . It follows from the KKT conditions in (C.6) that

$$\begin{aligned} G(\boldsymbol{\lambda}_k^*, \boldsymbol{\mu}_k^*; \mathbf{x}_k) &\leq G(\boldsymbol{\lambda}_k^{\text{sub}}, \boldsymbol{\mu}_k^{\text{sub}}; \mathbf{x}_k) \leq \|\mathbf{B}_k \Delta \mathbf{x}_k\|^2 + \|\boldsymbol{\mu}_{1,k}^{\text{sub}} \odot \Delta \mathbf{x}_k\|^2 + \|\boldsymbol{\mu}_{2,k}^{\text{sub}} \odot \Delta \mathbf{x}_k\|^2 \\ &\leq (\kappa_2^2 + 2M_{\text{dual}}^2) \|\Delta \mathbf{x}_k\|^2, \end{aligned} \quad (\text{C.10})$$

which implies from Theorem 3.8 that $\liminf_{k \rightarrow \infty} G(\boldsymbol{\lambda}_k^*, \boldsymbol{\mu}_k^*; \mathbf{x}_k) = 0$. Suppose $\limsup_{k \rightarrow \infty} G(\boldsymbol{\lambda}_k^*, \boldsymbol{\mu}_k^*; \mathbf{x}_k) > 0$, we can find a sufficiently small number $\varepsilon > 0$ and two infinite sequences $\{m_i, n_i\}$ with $m_i < n_i$, such that

$$G(\boldsymbol{\lambda}_{m_i}^*, \boldsymbol{\mu}_{m_i}^*; \mathbf{x}_{m_i}) > 2\varepsilon, \quad \|\Delta \mathbf{x}_k\| \geq \sqrt{\frac{\varepsilon}{\kappa_2^2 + 2M_{\text{dual}}^2}} \text{ for } k \in [m_i, n_i), \quad \|\Delta \mathbf{x}_{n_i}\| \leq \sqrt{\frac{\varepsilon}{\kappa_2^2 + 2M_{\text{dual}}^2}}.$$

We can always find such sequences since $G(\boldsymbol{\lambda}_{m_i}^*, \boldsymbol{\mu}_{m_i}^*; \mathbf{x}_{m_i}) > 2\varepsilon$ implies $\|\Delta \mathbf{x}_{m_i}\| \geq \sqrt{\frac{2\varepsilon}{\kappa_2^2 + 2M_{\text{dual}}^2}}$;

and due to $\liminf_{k \rightarrow \infty} \|\Delta \mathbf{x}_k\| = 0$, there must exist $n_i > m_i$ such that $\|\Delta \mathbf{x}_{n_i}\| \leq \sqrt{\frac{\varepsilon}{\kappa_2^2 + 2M_{\text{dual}}^2}}$. Let

$$\tilde{G}(\boldsymbol{\lambda}, \boldsymbol{\mu}; \mathbf{x}) = G(\boldsymbol{\lambda}, \boldsymbol{\mu}; \mathbf{x}) + \frac{\varepsilon}{6M_{\text{dual}}^2} \|(\boldsymbol{\lambda}, \boldsymbol{\mu})\|^2$$

and denote $\mathbf{w}_k = (\tilde{\boldsymbol{\lambda}}_k, \tilde{\boldsymbol{\mu}}_k) \in \arg \min_{\boldsymbol{\lambda}, \boldsymbol{\mu} \geq \mathbf{0}} \tilde{G}(\boldsymbol{\lambda}, \boldsymbol{\mu}; \mathbf{x}_k)$. By the above definition of \tilde{G} and the construction of the sequences $\{m_i, n_i\}$, we have

$$\begin{aligned} \tilde{G}(\tilde{\boldsymbol{\lambda}}_{m_i}, \tilde{\boldsymbol{\mu}}_{m_i}; \mathbf{x}_{m_i}) &\geq G(\tilde{\boldsymbol{\lambda}}_{m_i}, \tilde{\boldsymbol{\mu}}_{m_i}; \mathbf{x}_{m_i}) \geq G(\boldsymbol{\lambda}_{m_i}^*, \boldsymbol{\mu}_{m_i}^*; \mathbf{x}_{m_i}) \geq 2\varepsilon, \\ \tilde{G}(\tilde{\boldsymbol{\lambda}}_{n_i}, \tilde{\boldsymbol{\mu}}_{n_i}; \mathbf{x}_{n_i}) &\leq \tilde{G}(\boldsymbol{\lambda}_{n_i}^{\text{sub}}, \boldsymbol{\mu}_{n_i}^{\text{sub}}; \mathbf{x}_{n_i}) \stackrel{\text{(C.10)}}{\leq} (\kappa_2^2 + 2M_{\text{dual}}^2) \|\Delta \mathbf{x}_{n_i}\|^2 + \frac{\varepsilon}{3} \leq \frac{4}{3}\varepsilon. \end{aligned} \quad (\text{C.11})$$

Next, we analyze the consecutive difference of \tilde{G} . In particular, let us write $\tilde{G}(\boldsymbol{\lambda}, \boldsymbol{\mu}; \mathbf{x}_k)$ into a quadratic form of $\mathbf{w} = (\boldsymbol{\lambda}, \boldsymbol{\mu}) = (\boldsymbol{\lambda}, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2)$:

$$\tilde{G}(\boldsymbol{\lambda}, \boldsymbol{\mu}; \mathbf{x}_k) = \|\nabla f_k\|^2 + \mathbf{q}_k^\top \mathbf{w} + \frac{1}{2} \mathbf{w}^\top \nabla^2 \tilde{G}_k \mathbf{w},$$

where

$$\mathbf{q}_k = \begin{pmatrix} 2\nabla \mathbf{c}_k \nabla f_k \\ -2\nabla f_k \\ 2\nabla f_k \end{pmatrix}$$

and

$$\nabla^2 \tilde{G}_k = \begin{pmatrix} 2\nabla \mathbf{c}_k \nabla \mathbf{c}_k^\top & -2\nabla \mathbf{c}_k & 2\nabla \mathbf{c}_k \\ -2\nabla \mathbf{c}_k^\top & 2\mathbf{I} + 2\text{diag}^2(\boldsymbol{\ell} - \mathbf{x}_k) & -2\mathbf{I} \\ 2\nabla \mathbf{c}_k^\top & -2\mathbf{I} & 2\mathbf{I} + 2\text{diag}^2(\mathbf{x}_k - \mathbf{u}) \end{pmatrix} + \frac{\varepsilon}{3M_{\text{dual}}^2} \mathbf{I}.$$

Then, we have for any $k \geq 0$,

$$\begin{aligned} &|\tilde{G}(\tilde{\boldsymbol{\lambda}}_{k+1}, \tilde{\boldsymbol{\mu}}_{k+1}; \mathbf{x}_{k+1}) - \tilde{G}(\tilde{\boldsymbol{\lambda}}_k, \tilde{\boldsymbol{\mu}}_k; \mathbf{x}_k)| \\ &\leq \left| \mathbf{q}_{k+1}^\top \mathbf{w}_{k+1} + \frac{1}{2} \mathbf{w}_{k+1}^\top \nabla^2 \tilde{G}_{k+1} \mathbf{w}_{k+1} - \mathbf{q}_k^\top \mathbf{w}_k - \frac{1}{2} \mathbf{w}_k^\top \nabla^2 \tilde{G}_k \mathbf{w}_k \right| + \left| \|\nabla f_{k+1}\|^2 - \|\nabla f_k\|^2 \right| \\ &\leq \left| \mathbf{q}_{k+1}^\top \mathbf{w}_{k+1} + \frac{1}{2} \mathbf{w}_{k+1}^\top \nabla^2 \tilde{G}_{k+1} \mathbf{w}_{k+1} - \mathbf{q}_k^\top \mathbf{w}_{k+1} - \frac{1}{2} \mathbf{w}_{k+1}^\top \nabla^2 \tilde{G}_k \mathbf{w}_{k+1} \right| \\ &\quad + \left| \mathbf{q}_k^\top \mathbf{w}_{k+1} + \frac{1}{2} \mathbf{w}_{k+1}^\top \nabla^2 \tilde{G}_k \mathbf{w}_{k+1} - \mathbf{q}_k^\top \mathbf{w}_k - \frac{1}{2} \mathbf{w}_k^\top \nabla^2 \tilde{G}_k \mathbf{w}_k \right| + \left| \|\nabla f_{k+1}\|^2 - \|\nabla f_k\|^2 \right| \\ &\leq \|\mathbf{w}_{k+1}\| \|\mathbf{q}_{k+1} - \mathbf{q}_k\| + \frac{1}{2} \|\mathbf{w}_{k+1}\|^2 \|\nabla^2 \tilde{G}_{k+1} - \nabla^2 \tilde{G}_k\| + \|\mathbf{q}_k\| \|\mathbf{w}_{k+1} - \mathbf{w}_k\| \\ &\quad + \frac{1}{2} (\|\mathbf{w}_{k+1}\| + \|\mathbf{w}_k\|) \|\nabla^2 \tilde{G}_k\| \|\mathbf{w}_{k+1} - \mathbf{w}_k\| + \|\nabla f_{k+1} - \nabla f_k\| (\|\nabla f_{k+1}\| + \|\nabla f_k\|). \end{aligned} \quad (\text{C.12})$$

We provide the bound for each of term in the above display. We note that for any $k \geq 0$,

$$\frac{\varepsilon}{6M_{\text{dual}}^2} \|\mathbf{w}_k\|^2 \leq \tilde{G}(\tilde{\boldsymbol{\lambda}}_k, \tilde{\boldsymbol{\mu}}_k; \mathbf{x}_k)$$

$$\leq G(\boldsymbol{\lambda}_k^{\text{sub}}, \boldsymbol{\mu}_k^{\text{sub}}; \mathbf{x}_k) + \frac{\varepsilon}{6M_{\text{dual}}^2} \|(\boldsymbol{\lambda}_k^{\text{sub}}, \boldsymbol{\mu}_k^{\text{sub}})\|^2 \stackrel{\text{(C.10)}}{\leq} (\kappa_2^2 + 2M_{\text{dual}}^2) \|\Delta \mathbf{x}_k\|^2 + \frac{\varepsilon}{3}.$$

Thus, we have

$$\|\mathbf{w}_k\| \leq \sqrt{\frac{6M_{\text{dual}}^2 (\kappa_2^2 + 2M_{\text{dual}}^2) M_{\ell, \mathbf{u}}^2}{\varepsilon}} + \sqrt{2}M_{\text{dual}}. \quad (\text{C.13})$$

Furthermore, by the smoothness of the objective $f(\mathbf{x})$ and the constraints $\mathbf{c}(\mathbf{x})$, as well as the boundedness of $\nabla \mathbf{c}, \nabla f, \bar{\Delta} \mathbf{x}$, we know there exists a constant $M_{\mathbf{q}, \nabla^2 \tilde{G}} > 0$ such that for $k \geq 0$,

$$\max\{\|\mathbf{q}_{k+1} - \mathbf{q}_k\|, \|\nabla^2 \tilde{G}_{k+1} - \nabla^2 \tilde{G}_k\|\} \leq M_{\mathbf{q}, \nabla^2 \tilde{G}} \cdot \alpha_k. \quad (\text{C.14})$$

Finally, we note that for any $k \geq 0$ and any vector $\mathbf{w} = (\boldsymbol{\lambda}, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2) \in \mathbb{R}^m \times \mathbb{R}^d \times \mathbb{R}^d$,

$$\frac{1}{2} \mathbf{w}^\top \left(\nabla^2 \tilde{G}_k - \frac{\varepsilon}{3M_{\text{dual}}^2} \mathbf{I} \right) \mathbf{w} = \|\nabla \mathbf{c}_k^\top \boldsymbol{\lambda} - \boldsymbol{\mu}_1 + \boldsymbol{\mu}_2\|^2 + \|(\mathbf{x}_k - \boldsymbol{\ell}) \odot \boldsymbol{\mu}_1\|^2 + \|(\mathbf{x}_k - \mathbf{u}) \odot \boldsymbol{\mu}_2\|^2 \geq 0.$$

We combine the above display with (C.13) and (C.14), apply Lemma B.5, and have $\|\mathbf{w}_{k+1} - \mathbf{w}_k\| = \mathcal{O}\left(\frac{\alpha_k}{\varepsilon^{3/2}}\right)$, where we have omitted universal deterministic constants that are independent of ε . Combining this result with (C.13) and (C.14), and plugging into (C.12), we obtain

$$|\tilde{G}(\tilde{\boldsymbol{\lambda}}_{k+1}, \tilde{\boldsymbol{\mu}}_{k+1}; \mathbf{x}_{k+1}) - \tilde{G}(\tilde{\boldsymbol{\lambda}}_k, \tilde{\boldsymbol{\mu}}_k; \mathbf{x}_k)| \leq M_G \frac{\alpha_k}{\varepsilon^2}, \quad \forall k \geq 0,$$

for some universal constant $M_G > 0$ that is independent of α_k and ε . Combining the above display with (C.11), we obtain

$$\begin{aligned} \frac{2}{3} \varepsilon &\leq \tilde{G}(\tilde{\boldsymbol{\lambda}}_{m_i}, \tilde{\boldsymbol{\mu}}_{m_i}; \mathbf{x}_{m_i}) - \tilde{G}(\tilde{\boldsymbol{\lambda}}_{n_i}, \tilde{\boldsymbol{\mu}}_{n_i}; \mathbf{x}_{n_i}) \\ &\leq \sum_{k=m_i}^{n_i-1} \left| \tilde{G}(\tilde{\boldsymbol{\lambda}}_k, \tilde{\boldsymbol{\mu}}_k; \mathbf{x}_k) - \tilde{G}(\tilde{\boldsymbol{\lambda}}_{k+1}, \tilde{\boldsymbol{\mu}}_{k+1}; \mathbf{x}_{k+1}) \right| \leq \frac{M_G}{\varepsilon^2} \sum_{k=m_i}^{n_i-1} \alpha_k. \end{aligned}$$

Summing up both sides from $i = 1$ to ∞ , we have

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

On the other hand, since $\|\Delta \mathbf{x}_k\| \geq \sqrt{\frac{\varepsilon}{\kappa_2^2 + 2M_{\text{dual}}^2}}$ for $m_i \leq k \leq n_i - 1$, we know that

$$\sum_{i=1}^{\infty} \sum_{k=m_i}^{n_i-1} \alpha_k \leq \frac{\kappa_2^2 + 2M_{\text{dual}}^2}{\varepsilon} \sum_{i=1}^{\infty} \sum_{k=m_i}^{n_i-1} \alpha_k \|\Delta \mathbf{x}_k\|^2 \leq \frac{\kappa_2^2 + 2M_{\text{dual}}^2}{\varepsilon} \sum_{k=0}^{\infty} \alpha_k \|\Delta \mathbf{x}_k\|^2 < \infty,$$

where the last inequality is due to the proof of Theorem 3.8. This derives the contradiction and

we complete the proof of $\lim_{k \rightarrow \infty} G(\boldsymbol{\lambda}_k^*, \boldsymbol{\mu}_k^*; \mathbf{x}_k) = 0$.

• **Convergence of $\|\mathbf{c}_k\|$.** We follow the same proof strategy. By Theorem 3.8 we know that $\liminf_{k \rightarrow \infty} \|\mathbf{c}_k\| = 0$. Suppose $\limsup_{k \rightarrow \infty} \|\mathbf{c}_k\| > 0$, then we can find a sufficiently small number $\varepsilon > 0$ and two infinite sequences $\{m_i, n_i\}$ with $m_i < n_i$, such that

$$\|\mathbf{c}_{m_i}\| > 2\varepsilon, \quad \|\mathbf{c}_k\| \geq \varepsilon \text{ for } k \in [m_i, n_i), \quad \|\mathbf{c}_{n_i}\| \leq \varepsilon.$$

It follows from the construction of the sequences that

$$\varepsilon \leq \|\mathbf{c}_{m_i}\| - \|\mathbf{c}_{n_i}\| = \sum_{k=m_i}^{n_i-1} (\|\mathbf{c}_k\| - \|\mathbf{c}_{k+1}\|) \leq \sum_{k=m_i}^{n_i-1} \|\mathbf{c}_k - \mathbf{c}_{k+1}\| \leq M_c \sum_{k=m_i}^{n_i-1} \alpha_k$$

for some universal constant $M_c > 0$ that is independent of α_k and ε . Multiplying both sides by ε and noting the fact that $\|\mathbf{c}_k\| \geq \varepsilon$ for $m_i \leq k < n_i$, we have

$$\varepsilon^2 \leq M_c \sum_{k=m_i}^{n_i-1} \alpha_k \|\mathbf{c}_k\|,$$

which implies that $\infty \leq \sum_{i=1}^{\infty} \sum_{k=m_i}^{n_i-1} \alpha_k \|\mathbf{c}_k\| \leq \sum_{k=0}^{\infty} \alpha_k \|\mathbf{c}_k\| < \infty$. This leads to the contradiction and we complete the proof.

C.6 Proof of Lemma 4.2

• **Proof of (a).** By Assumption 4.1, we know LICQ holds at \mathbf{x}^* . Thus, when \mathbf{x}_k is close to \mathbf{x}^* , the rows of $\tilde{\mathbf{J}}_k = (\nabla \mathbf{c}_k; -\mathbf{I}_{\mathcal{A}_\ell^*}; \mathbf{I}_{\mathcal{A}_u^*})$ are linearly independent. To show $\theta_k = 1$ for k large enough, it suffices to show that there exists $\mathbf{z}_k \in \mathbb{R}^d$ such that $\mathbf{c}_k + \nabla \mathbf{c}_k \mathbf{z}_k = \mathbf{0}$ and $\ell \leq \mathbf{x}_k + \mathbf{z}_k \leq \mathbf{u}$.

To this end, let us define the linear system at \mathbf{x}_k as

$$\begin{pmatrix} \nabla \mathbf{c}_k \\ -\mathbf{I}_{\mathcal{A}_\ell^*} \\ \mathbf{I}_{\mathcal{A}_u^*} \end{pmatrix} \mathbf{z}_k = - \begin{pmatrix} \mathbf{c}_k \\ [\ell - \mathbf{x}_k]_{\mathcal{A}_\ell^*} \\ [\mathbf{x}_k - \mathbf{u}]_{\mathcal{A}_u^*} \end{pmatrix}.$$

By the full row-rankness of $\tilde{\mathbf{J}}_k$, the above linear system has a solution given by

$$\mathbf{z}_k = -\tilde{\mathbf{J}}_k^\top (\tilde{\mathbf{J}}_k \tilde{\mathbf{J}}_k^\top)^{-1} \begin{pmatrix} \mathbf{c}_k \\ [\ell - \mathbf{x}_k]_{\mathcal{A}_\ell^*} \\ [\mathbf{x}_k - \mathbf{u}]_{\mathcal{A}_u^*} \end{pmatrix}.$$

Clearly, the above \mathbf{z}_k satisfies $\mathbf{c}_k + \nabla \mathbf{c}_k \mathbf{z}_k = \mathbf{0}$ and $[\mathbf{l}]_i \leq [\mathbf{x}_k + \mathbf{z}_k]_i \leq [\mathbf{u}]_i$ for all $i \in \mathcal{A}^* = \mathcal{A}_\ell^* \cup \mathcal{A}_u^*$. Furthermore, since $\mathbf{z}_k \rightarrow \mathbf{0}$ as $\mathbf{x}_k \rightarrow \mathbf{x}^*$, we also have $[\mathbf{l}]_i \leq [\mathbf{x}_k + \mathbf{z}_k]_i \leq [\mathbf{u}]_i$ for all $i \in (\mathcal{A}^*)^-$. This implies that $\mathbf{z}_k \in \Omega(\mathbf{x}_k; \theta_k)$ with $\theta_k = 1$ and we complete the proof.

• **Proof of (b).** Let us consider \mathbf{x}_k such that $\|\mathbf{x}_k - \mathbf{x}^*\|_{\max} = \max_{i \in [d]} |[\mathbf{x}_k - \mathbf{x}^*]_i| \leq 0.5\epsilon$. For any

$i \in \mathcal{A}_\ell^*$, we have $[\mathbf{x}_k - \boldsymbol{\ell}]_i \leq \|\mathbf{x}_k - \mathbf{x}^*\|_{\max} \leq 0.5\epsilon$, implying that $\mathcal{A}_\ell^* \subseteq \mathcal{A}_{\ell,k}(\epsilon)$ by the definition of (25). On the other hand, for any $i \in \mathcal{A}_{\ell,k}(\epsilon)$, we have $[\mathbf{x}^* - \boldsymbol{\ell}]_i \leq |[\mathbf{x}_k - \mathbf{x}^*]_i| + [\mathbf{x}_k - \boldsymbol{\ell}]_i \leq \|\mathbf{x}_k - \mathbf{x}^*\|_{\max} + [\mathbf{x}_k - \boldsymbol{\ell}]_i \leq 1.5\epsilon < [\mathbf{x}^* - \boldsymbol{\ell}]_j$, for any $j \in (\mathcal{A}^*)^-$, where the last inequality is by the definition (24). Thus, we know $i \in \mathcal{A}^*$. By the same derivations, we also have $[\mathbf{x}^* - \boldsymbol{\ell}]_i < [\mathbf{u} - \boldsymbol{\ell}]_i$, which implies $i \in \mathcal{A}_\ell^*$. Therefore, we have $\mathcal{A}_{\ell,k}(\epsilon) = \mathcal{A}_\ell^*$ as long as $\|\mathbf{x}_k - \mathbf{x}^*\|_{\max} \leq 0.5\epsilon$. Similar arguments hold for showing that $\mathcal{A}_{\mathbf{u},k}(\epsilon) = \mathcal{A}_{\mathbf{u}}^*$, and we complete the proof.

• **Proof of (c) and (d).** We prove (c) and (d) jointly. We need the following supporting lemma, which is proved in Appendix D.6.

Lemma C.5. *Under Assumptions 3.6 and 4.1, let $\beta_k = \iota_2(k+1)^{-b_2}$ satisfy $\iota_2 > 0, b_2 \in (0.5, 1]$. We have $\bar{\mathbf{g}}_k \rightarrow \nabla f^*$ as $k \rightarrow \infty$ almost surely.*

We will first show that the subproblem (17) admits a solution $\bar{\Delta}\mathbf{x}_k$ near the origin for k large enough, with $\bar{\Delta}\mathbf{x}_k \rightarrow 0$ as $k \rightarrow \infty$ and satisfying $\mathcal{A}_\ell(\mathbf{x}_k + \bar{\Delta}\mathbf{x}_k) = \mathcal{A}_\ell^*$ and $\mathcal{A}_{\mathbf{u}}(\mathbf{x}_k + \bar{\Delta}\mathbf{x}_k) = \mathcal{A}_{\mathbf{u}}^*$. We then show that the subproblem solution is also unique in a neighborhood of the origin.

Let us consider the following linear system:

$$\begin{pmatrix} \bar{\mathbf{B}}_k & \nabla \mathbf{c}_k^\top & -\mathbf{I}_{\mathcal{A}_{\ell,k}(\epsilon)}^\top & \mathbf{I}_{\mathcal{A}_{\mathbf{u},k}(\epsilon)}^\top \\ \nabla \mathbf{c}_k & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ -\mathbf{I}_{\mathcal{A}_{\ell,k}(\epsilon)} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{I}_{\mathcal{A}_{\mathbf{u},k}(\epsilon)} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \bar{\Delta}\mathbf{x}_k \\ \bar{\boldsymbol{\lambda}}_k^{\text{sub}} \\ [\bar{\boldsymbol{\mu}}_{1,k}^{\text{sub}}]_{\mathcal{A}_{\ell,k}(\epsilon)} \\ [\bar{\boldsymbol{\mu}}_{2,k}^{\text{sub}}]_{\mathcal{A}_{\mathbf{u},k}(\epsilon)} \end{pmatrix} = - \begin{pmatrix} \bar{\mathbf{g}}_k \\ \mathbf{c}_k \\ [\boldsymbol{\ell} - \mathbf{x}_k]_{\mathcal{A}_{\ell,k}(\epsilon)} \\ [\mathbf{x}_k - \mathbf{u}]_{\mathcal{A}_{\mathbf{u},k}(\epsilon)} \end{pmatrix}. \quad (\text{C.15})$$

By the construction of $\bar{\mathbf{B}}_k$, the local LICQ in Assumption 4.1, and the properties $\mathcal{A}_{\ell,k}(\epsilon) = \mathcal{A}_\ell^*$ and $\mathcal{A}_{\mathbf{u},k}(\epsilon) = \mathcal{A}_{\mathbf{u}}^*$ proved in (b), we know the above system has a (unique) solution. Let us rewrite the system as

$$\begin{pmatrix} \bar{\mathbf{B}}_k & \nabla \mathbf{c}_k^\top & -\mathbf{I}_{\mathcal{A}_{\ell,k}(\epsilon)}^\top & \mathbf{I}_{\mathcal{A}_{\mathbf{u},k}(\epsilon)}^\top \\ \nabla \mathbf{c}_k & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ -\mathbf{I}_{\mathcal{A}_{\ell,k}(\epsilon)} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{I}_{\mathcal{A}_{\mathbf{u},k}(\epsilon)} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \bar{\Delta}\mathbf{x}_k \\ \bar{\boldsymbol{\lambda}}_k^{\text{sub}} - \boldsymbol{\lambda}^* \\ [\bar{\boldsymbol{\mu}}_{1,k}^{\text{sub}} - \boldsymbol{\mu}_1^*]_{\mathcal{A}_{\ell,k}(\epsilon)} \\ [\bar{\boldsymbol{\mu}}_{2,k}^{\text{sub}} - \boldsymbol{\mu}_2^*]_{\mathcal{A}_{\mathbf{u},k}(\epsilon)} \end{pmatrix} = - \begin{pmatrix} \bar{\mathbf{g}}_k + \nabla \mathbf{c}_k^\top \boldsymbol{\lambda}^* - \mathbf{I}_{\mathcal{A}_{\ell,k}(\epsilon)}^\top [\boldsymbol{\mu}_1^*]_{\mathcal{A}_{\ell,k}(\epsilon)} + \mathbf{I}_{\mathcal{A}_{\mathbf{u},k}(\epsilon)}^\top [\boldsymbol{\mu}_2^*]_{\mathcal{A}_{\mathbf{u},k}(\epsilon)} \\ \mathbf{c}_k \\ [\boldsymbol{\ell} - \mathbf{x}_k]_{\mathcal{A}_{\ell,k}(\epsilon)} \\ [\mathbf{x}_k - \mathbf{u}]_{\mathcal{A}_{\mathbf{u},k}(\epsilon)} \end{pmatrix}. \quad (\text{C.16})$$

By the convergence of \mathbf{x}_k , $\mathcal{A}_{\ell,k}(\epsilon)$, $\mathcal{A}_{\mathbf{u},k}(\epsilon)$, and Lemma C.5, we know the right hand side of (C.16)

converges to 0 because it is the KKT conditions at $(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$, implying that $(\bar{\Delta}\mathbf{x}_k; \bar{\boldsymbol{\lambda}}_k^{\text{sub}} - \boldsymbol{\lambda}^*; [\bar{\boldsymbol{\mu}}_{1,k}^{\text{sub}} - \boldsymbol{\mu}_1^*]_{\mathcal{A}_{\ell,k}(\epsilon)}; [\bar{\boldsymbol{\mu}}_{2,k}^{\text{sub}} - \boldsymbol{\mu}_2^*]_{\mathcal{A}_{u,k}(\epsilon)}) \rightarrow 0$. Now, let us show that $\mathcal{A}_{\ell}(\mathbf{x}_k + \bar{\Delta}\mathbf{x}_k) = \mathcal{A}_{\ell,k}(\epsilon) = \mathcal{A}_{\ell}^*$ and $\mathcal{A}_u(\mathbf{x}_k + \bar{\Delta}\mathbf{x}_k) = \mathcal{A}_{u,k}(\epsilon) = \mathcal{A}_u^*$ when $\|\bar{\Delta}\mathbf{x}_k\|_{\max} \leq \epsilon$. By the system (C.16), we have $\mathcal{A}_{\ell,k}(\epsilon) \subseteq \mathcal{A}_{\ell}(\mathbf{x}_k + \bar{\Delta}\mathbf{x}_k)$. On the other hand, for any $i \in \mathcal{A}_{\ell}(\mathbf{x}_k + \bar{\Delta}\mathbf{x}_k)$, we have $|[\mathbf{x}_k]_i - [\boldsymbol{\ell}]_i| \leq |[\bar{\Delta}\mathbf{x}_k]_i| \leq \epsilon$, suggesting that $i \in \mathcal{A}_{\ell,k}(\epsilon)$ and hence $\mathcal{A}_{\ell}(\mathbf{x}_k + \bar{\Delta}\mathbf{x}_k) \subseteq \mathcal{A}_{\ell,k}(\epsilon)$. This leads to $\mathcal{A}_{\ell,k}(\epsilon) = \mathcal{A}_{\ell}(\mathbf{x}_k + \bar{\Delta}\mathbf{x}_k)$. Similarly, we can show $\mathcal{A}_{u,k}(\epsilon) = \mathcal{A}_u(\mathbf{x}_k + \bar{\Delta}\mathbf{x}_k)$. Next, we complement $[\bar{\boldsymbol{\mu}}_k^{\text{sub}}]_{\mathcal{A}_k(\epsilon)} = ([\bar{\boldsymbol{\mu}}_{1,k}^{\text{sub}}]_{\mathcal{A}_{\ell,k}(\epsilon)}, [\bar{\boldsymbol{\mu}}_{2,k}^{\text{sub}}]_{\mathcal{A}_{u,k}(\epsilon)})$ by defining $[\bar{\boldsymbol{\mu}}_k^{\text{sub}}]_{\mathcal{A}_k^-(\epsilon)} = \mathbf{0}$, and aim to show that $(\bar{\Delta}\mathbf{x}_k, \bar{\boldsymbol{\lambda}}_k^{\text{sub}}, \bar{\boldsymbol{\mu}}_k^{\text{sub}})$ is the local solution of the SSQP subproblem (17). By (C.15) and the complement $[\bar{\boldsymbol{\mu}}_k^{\text{sub}}]_{\mathcal{A}_k^-(\epsilon)} = \mathbf{0}$, we have

$$\begin{aligned} \bar{\mathbf{g}}_k + \bar{\mathbf{B}}_k \bar{\Delta}\mathbf{x}_k + \nabla \mathbf{c}_k^{\top} \bar{\boldsymbol{\lambda}}_k^{\text{sub}} - \bar{\boldsymbol{\mu}}_{1,k}^{\text{sub}} + \bar{\boldsymbol{\mu}}_{2,k}^{\text{sub}} &= \mathbf{0}, \\ \bar{\boldsymbol{\mu}}_{1,k}^{\text{sub}\top} (\boldsymbol{\ell} - \mathbf{x}_k - \bar{\Delta}\mathbf{x}_k) &= \mathbf{0}, \\ \bar{\boldsymbol{\mu}}_{2,k}^{\text{sub}\top} (\mathbf{x}_k + \bar{\Delta}\mathbf{x}_k - \mathbf{u}) &= \mathbf{0}. \end{aligned} \tag{C.17}$$

Furthermore, since $[\boldsymbol{\mu}^*]_{\mathcal{A}_k(\epsilon)} > 0$ by strict complementarity condition and $[\bar{\boldsymbol{\mu}}_k^{\text{sub}} - \boldsymbol{\mu}^*]_{\mathcal{A}_k(\epsilon)} \rightarrow 0$ as $k \rightarrow \infty$, we have $[\bar{\boldsymbol{\mu}}_k^{\text{sub}}]_{\mathcal{A}_k(\epsilon)} > 0$ for k large enough. This verifies the KKT conditions of the subproblem (17). The conditions of LICQ, SOSC, and strict complementarity of the subproblem are trivial due to the facts that $\mathcal{A}_{\ell}(\mathbf{x}_k + \bar{\Delta}\mathbf{x}_k) = \mathcal{A}_{\ell,k}(\epsilon)$ and $\mathcal{A}_u(\mathbf{x}_k + \bar{\Delta}\mathbf{x}_k) = \mathcal{A}_{u,k}(\epsilon)$.

Finally, we show the uniqueness of the solution $\bar{\Delta}\mathbf{x}_k$ in the neighborhood of $\mathbf{0}$. Suppose there exists another solution $\bar{\Delta}\mathbf{x}'_k$ with $\|\bar{\Delta}\mathbf{x}'_k\|_{\max} \leq \epsilon$. We know for any $i \in \mathcal{A}_{\ell}(\mathbf{x}_k + \bar{\Delta}\mathbf{x}'_k)$, $[\mathbf{x}_k - \boldsymbol{\ell}]_i \leq \|\bar{\Delta}\mathbf{x}'_k\|_{\max} \leq \epsilon$, suggesting that $i \in \mathcal{A}_{\ell,k}(\epsilon)$ and hence $\mathcal{A}_{\ell}(\mathbf{x}_k + \bar{\Delta}\mathbf{x}'_k) \subseteq \mathcal{A}_{\ell,k}(\epsilon)$. Similarly, we can show $\mathcal{A}_u(\mathbf{x}_k + \bar{\Delta}\mathbf{x}'_k) \subseteq \mathcal{A}_{u,k}(\epsilon)$. This suggests that the LICQ of the subproblem holds at $\bar{\Delta}\mathbf{x}'_k$ when \mathbf{x}_k is close to \mathbf{x}^* , as $\mathbf{J}_k = (\nabla \mathbf{c}_k; -\mathbf{I}_{\mathcal{A}_{\ell,k}(\epsilon)}; \mathbf{I}_{\mathcal{A}_{u,k}(\epsilon)})$ has full row-rank when \mathbf{x}_k is close to \mathbf{x}^* . Then, the KKT condition of the subproblem at $\bar{\Delta}\mathbf{x}'_k$ implies the existence of unique dual multipliers $(\bar{\boldsymbol{\lambda}}_k^{\text{sub}'}, \bar{\boldsymbol{\mu}}_k^{\text{sub}'})$ such that

$$\bar{\mathbf{g}}_k + \bar{\mathbf{B}}_k \bar{\Delta}\mathbf{x}'_k + \nabla \mathbf{c}_k^{\top} \bar{\boldsymbol{\lambda}}_k^{\text{sub}'} - \mathbf{I}_{\mathcal{A}_{\ell}(\mathbf{x}_k + \bar{\Delta}\mathbf{x}'_k)}^{\top} [\bar{\boldsymbol{\mu}}_{1,k}^{\text{sub}'}]_{\mathcal{A}_{\ell}(\mathbf{x}_k + \bar{\Delta}\mathbf{x}'_k)} + \mathbf{I}_{\mathcal{A}_u(\mathbf{x}_k + \bar{\Delta}\mathbf{x}'_k)}^{\top} [\bar{\boldsymbol{\mu}}_{2,k}^{\text{sub}'}]_{\mathcal{A}_u(\mathbf{x}_k + \bar{\Delta}\mathbf{x}'_k)} = \mathbf{0}.$$

Comparing the above display with (C.17), we obtain

$$\mathbf{J}_k^\top [\bar{\boldsymbol{\lambda}}_k^{\text{sub}} - \bar{\boldsymbol{\lambda}}_k^{\text{sub}'}; [\bar{\boldsymbol{\mu}}_{1,k}^{\text{sub}} - \bar{\boldsymbol{\mu}}_{1,k}^{\text{sub}'}]_{\mathcal{A}_{\ell,k}(\epsilon)}; [\bar{\boldsymbol{\mu}}_{2,k}^{\text{sub}} - \bar{\boldsymbol{\mu}}_{2,k}^{\text{sub}'}]_{\mathcal{A}_{u,k}(\epsilon)}] = \bar{\mathbf{B}}_k(\bar{\Delta}\mathbf{x}'_k - \bar{\Delta}\mathbf{x}_k).$$

By the full row-rankness of \mathbf{J}_k as implied by LICQ at \mathbf{x}^* , we know $\|[\bar{\boldsymbol{\mu}}_k^{\text{sub}} - \bar{\boldsymbol{\mu}}_k^{\text{sub}'}]_{\mathcal{A}_k(\epsilon)}\| \leq \text{const} \cdot (\|\bar{\Delta}\mathbf{x}'_k\| + \|\bar{\Delta}\mathbf{x}_k\|)$ for some $\text{const} > 0$. Thus, using the fact that $[\bar{\boldsymbol{\mu}}_k^{\text{sub}}]_{\mathcal{A}_k(\epsilon)} > 0$, we can choose the neighborhood of the origin small enough such that if $\bar{\Delta}\mathbf{x}'_k, \bar{\Delta}\mathbf{x}_k$ are in the neighborhood, we have $[\bar{\boldsymbol{\mu}}_k^{\text{sub}'}]_{\mathcal{A}_k(\epsilon)} > 0$ as well, further implying that $\mathcal{A}_{\ell,k}(\epsilon) \subseteq \mathcal{A}_{\ell}(\mathbf{x}_k + \bar{\Delta}\mathbf{x}'_k)$ and $\mathcal{A}_{u,k}(\epsilon) \subseteq \mathcal{A}_u(\mathbf{x}_k + \bar{\Delta}\mathbf{x}'_k)$ due to the complementarity conditions of the subproblem. This leads to $\mathcal{A}_{\ell}(\mathbf{x}_k + \bar{\Delta}\mathbf{x}'_k) = \mathcal{A}_{\ell,k}(\epsilon)$ and $\mathcal{A}_{\ell}(\mathbf{x}_k + \bar{\Delta}\mathbf{x}'_k) = \mathcal{A}_{u,k}(\epsilon)$. However, by the KKT condition of the SSQP subproblem, the solution $\bar{\Delta}\mathbf{x}'_k$ must also follow the linear system (C.15), which admits a unique solution. Thus, the solution of SSQP subproblem is unique near $\mathbf{0}$. The argument about $\Delta\mathbf{x}_k$ follows the same analysis by noting that $\nabla f_k \rightarrow \nabla f^*$ almost surely. This completes the proof of (c) and (d).

C.7 Proof of Lemma 4.4

We first prove the convergence of $(\mathbf{x}_k, \boldsymbol{\lambda}_k, \boldsymbol{\mu}_k)$. By the proof of Lemma 4.2(c) in Appendix C.6, we know $\bar{\boldsymbol{\lambda}}_k^{\text{sub}} - \boldsymbol{\lambda}^* \rightarrow \mathbf{0}$ and $[\bar{\boldsymbol{\mu}}_k^{\text{sub}} - \boldsymbol{\mu}^*]_{\mathcal{A}_k(\epsilon)} \rightarrow \mathbf{0}$ as $k \rightarrow \infty$ almost surely. Let ϵ satisfy (24) and we have $\mathcal{A}_{\ell,k}(\epsilon) = \mathcal{A}^*$ for k large enough by Lemma 4.2(b). Since $[\bar{\boldsymbol{\mu}}_k^{\text{sub}}]_{\mathcal{A}_k^-(\epsilon)} = \mathbf{0} = \boldsymbol{\mu}_{(\mathcal{A}^*)-}^*$, we have $\bar{\boldsymbol{\mu}}_k^{\text{sub}} - \boldsymbol{\mu}^* \rightarrow \mathbf{0}$ as $k \rightarrow \infty$ almost surely. With this result, we now prove $(\boldsymbol{\lambda}_k, \boldsymbol{\mu}_k) \rightarrow (\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$.

We take $\boldsymbol{\lambda}_k$ as an example. By the update (18), we have

$$\begin{aligned} \boldsymbol{\lambda}_{k+1} - \boldsymbol{\lambda}^* &= \boldsymbol{\lambda}_k - \boldsymbol{\lambda}^* + \bar{\alpha}_k (\bar{\boldsymbol{\lambda}}_k^{\text{sub}} - \boldsymbol{\lambda}_k) = (1 - \bar{\alpha}_k)(\boldsymbol{\lambda}_k - \boldsymbol{\lambda}^*) + \bar{\alpha}_k(\bar{\boldsymbol{\lambda}}_k^{\text{sub}} - \boldsymbol{\lambda}^*) \\ &= \prod_{i=0}^k (1 - \bar{\alpha}_i)(\boldsymbol{\lambda}_0 - \boldsymbol{\lambda}^*) + \sum_{i=0}^k \prod_{j=i+1}^k (1 - \bar{\alpha}_j) \bar{\alpha}_i (\bar{\boldsymbol{\lambda}}_i^{\text{sub}} - \boldsymbol{\lambda}^*). \end{aligned}$$

To proceed, we present the following lemma, which is proved in Appendix D.7.

Lemma C.6. *Suppose $\beta_k \geq 0$ satisfies $\beta_k \rightarrow 0$ and $\sum_{k=0}^{\infty} \beta_k = \infty$, then for any scalar a and sequence $e_k \rightarrow 0$, we have $\sum_{i=0}^k \prod_{j=i+1}^k (1 - \beta_j) \beta_i e_i + a \prod_{j=0}^k (1 - \beta_j) \rightarrow 0$ as $k \rightarrow \infty$.*

By Lemma C.6 and the fact that $\sum_{k=0}^{\infty} \bar{\alpha}_k = \infty$, we know $\boldsymbol{\lambda}_k \rightarrow \boldsymbol{\lambda}^*$ almost surely. Following

the same derivations, we obtain $\boldsymbol{\mu}_k \rightarrow \boldsymbol{\mu}^*$. Now, we consider the convergence of $\bar{\mathbf{B}}_k$. By the construction of $\boldsymbol{\Delta}_k$, the facts that $\mathcal{A}_{\ell,k}(\epsilon) = \mathcal{A}_\ell^*$ and $\mathcal{A}_{u,k}(\epsilon) = \mathcal{A}_u^*$, and the conditions of LICQ and SOSC at $(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$, it suffices to show $\bar{\mathbf{Q}}_k \rightarrow \nabla^2 f^*$ almost surely. This is an analogy of Lemma C.5 and we summarize it to the following lemma. See Appendix D.8 for the proof.

Lemma C.7. *Under Assumptions 4.1 and 4.3, let $\gamma_k = \iota_3(k+1)^{-b_3}$ satisfy $\iota_3 > 0, b_3 \in (0.5, 1]$. We have $\bar{\mathbf{Q}}_k \rightarrow \nabla^2 f^*$ as $k \rightarrow \infty$ almost surely.*

The above lemma directly implies $\bar{\mathbf{B}}_k \rightarrow \nabla_{\mathbf{x}}^2 \mathcal{L}^*$ and we complete the proof.

C.8 Proof of Theorem 4.6

We introduce some notation first. We let

$$\mathbf{H}_k = \begin{pmatrix} \bar{\mathbf{B}}_k & \bar{\mathbf{J}}_k^\top \\ \bar{\mathbf{J}}_k & \mathbf{0} \end{pmatrix} = \begin{pmatrix} \bar{\mathbf{B}}_k & \nabla \mathbf{c}_k^\top & -\mathbf{I}_{\mathcal{A}_\ell^*}^\top & \mathbf{I}_{\mathcal{A}_u^*}^\top \\ \nabla \mathbf{c}_k & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ -\mathbf{I}_{\mathcal{A}_\ell^*} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{I}_{\mathcal{A}_u^*} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix}, \quad \bar{\nabla}_{\mathbf{w}} \mathcal{L}_k = \begin{pmatrix} \bar{\mathbf{g}}_k + \nabla \mathbf{c}_k^\top \boldsymbol{\lambda}_k - \mathbf{I}_{\mathcal{A}_\ell^*}^\top [\boldsymbol{\mu}_{1,k}]_{\mathcal{A}_\ell^*} + \mathbf{I}_{\mathcal{A}_u^*}^\top [\boldsymbol{\mu}_{2,k}]_{\mathcal{A}_u^*} \\ \mathbf{c}_k \\ [\boldsymbol{\ell} - \mathbf{x}_k]_{\mathcal{A}_\ell^*} \\ [\mathbf{x}_k - \mathbf{u}]_{\mathcal{A}_u^*} \end{pmatrix},$$

and $\nabla_{\mathbf{w}} \mathcal{L}_k$ is defined in the same way as $\bar{\nabla}_{\mathbf{w}} \mathcal{L}_k$ except that $\bar{\mathbf{g}}_k$ is replaced by the true gradient ∇f_k . We further define the step direction

$$\bar{\Delta} \mathbf{w}_k = \begin{pmatrix} \bar{\Delta} \mathbf{x}_k \\ \bar{\Delta} \boldsymbol{\lambda}_k \\ [\bar{\Delta} \boldsymbol{\mu}_{1,k}]_{\mathcal{A}_\ell^*} \\ [\bar{\Delta} \boldsymbol{\mu}_{2,k}]_{\mathcal{A}_u^*} \end{pmatrix} = -\mathbf{H}_k^{-1} \bar{\nabla}_{\mathbf{w}} \mathcal{L}_k \quad \text{and} \quad [\bar{\Delta} \boldsymbol{\mu}_k]_{(\mathcal{A}^*)^-} = -[\boldsymbol{\mu}_k]_{(\mathcal{A}^*)^-}. \quad (\text{C.18})$$

By the KKT conditions of the SSQP subproblem (C.15), the facts from Lemma 4.2 that $\mathcal{A}_\ell(\mathbf{x}_k + \bar{\Delta} \mathbf{x}_k) = \mathcal{A}_{\ell,k}(\epsilon) = \mathcal{A}_\ell^*$ and $\mathcal{A}_u(\mathbf{x}_k + \bar{\Delta} \mathbf{x}_k) = \mathcal{A}_{u,k}(\epsilon) = \mathcal{A}_u^*$, and using the relation of $(\bar{\boldsymbol{\lambda}}_k^{\text{sub}}, \bar{\boldsymbol{\mu}}_k^{\text{sub}}) = (\boldsymbol{\lambda}_k + \bar{\Delta} \boldsymbol{\lambda}_k, \boldsymbol{\mu}_k + \bar{\Delta} \boldsymbol{\mu}_k)$, we know the update of the iterates when k is large enough can be written by utilizing (C.18) as

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \bar{\alpha}_k \bar{\Delta} \mathbf{w}_k \quad \text{and} \quad [\boldsymbol{\mu}_{k+1}]_{(\mathcal{A}^*)^-} = (1 - \bar{\alpha}_k) \cdot [\boldsymbol{\mu}_k]_{(\mathcal{A}^*)^-}. \quad (\text{C.19})$$

Let $\Upsilon > 0$ be a constant and we propose the following conditions:

- (a) $\max\{\beta_k, \bar{\alpha}_k\} \leq 1; \|\mathbf{H}_k^{-1} - (\mathbf{H}^*)^{-1}\| \leq \Upsilon \|\mathbf{H}_k - \mathbf{H}^*\|; \|\mathbf{H}_k^{-1}\| \leq \Upsilon.$
- (b) $\max\{\|\nabla f_k - \nabla f^*\|, \|\nabla^2 f_k - \nabla^2 f^*\|\} \leq \Upsilon \|\mathbf{x}_k - \mathbf{x}^*\|; \|\nabla f_k - \nabla f_{k-1}\| \leq \Upsilon \|\mathbf{x}_k - \mathbf{x}_{k-1}\|.$
- (c) $\max\{\|\mathbf{w}_k - \mathbf{w}^*\|, \|\nabla f_k - \nabla f^*\|\} \leq \frac{1}{\Upsilon}; \max\{\|\nabla_{\mathbf{w}} \mathcal{L}_k\|, \|\nabla_{\mathbf{w}}^2 \mathcal{L}_k - \nabla_{\mathbf{w}}^2 \mathcal{L}^*\|\} \leq \Upsilon \|\mathbf{w}_k - \mathbf{w}^*\|.$

- (d) $\|\nabla_{\mathbf{w}} \mathcal{L}_k - \mathbf{H}^*(\mathbf{w}_k - \mathbf{w}^*)\| \leq \Upsilon \|\mathbf{w}_k - \mathbf{w}^*\|^2$; $\|\mathbf{w}_k - \mathbf{w}^* - \mathbf{H}_k^{-1} \nabla_{\mathbf{w}} \mathcal{L}_k\| \leq \|\mathbf{w}_k - \mathbf{w}^*\|/\Upsilon$.
- (e) $\mathcal{A}_{\ell}(\mathbf{x}_k + \bar{\Delta} \mathbf{x}_k) = \mathcal{A}_{\ell,k}(\epsilon) = \mathcal{A}_{\ell}^*$; $\mathcal{A}_{\mathbf{u}}(\mathbf{x}_k + \bar{\Delta} \mathbf{x}_k) = \mathcal{A}_{\mathbf{u},k}(\epsilon) = \mathcal{A}_{\mathbf{u}}^*$.
- (f) The subproblem solution $(\bar{\Delta} \mathbf{x}_k, \bar{\boldsymbol{\lambda}}_k^{\text{sub}}, \bar{\boldsymbol{\mu}}_k^{\text{sub}}) = (\bar{\Delta} \mathbf{x}_k, \boldsymbol{\lambda}_k + \bar{\Delta} \boldsymbol{\lambda}_k, \boldsymbol{\mu}_k + \bar{\Delta} \boldsymbol{\mu}_k)$ is from (C.18).

For any $k_0 \geq 0$, we also define the stopping time

$$\tau_{k_0} := \inf_j \{j \geq k_0 : \text{any of the above conditions does not hold at } j\text{-th iteration}\}. \quad (\text{C.20})$$

By Lemmas 4.2, 4.4, Assumption 4.1, and choosing Υ to be large enough, we know for any run of the method, there exists a (potentially random) $\tilde{k}_0 < \infty$ such that $\tau_{k_0} = \infty$, $\forall k_0 \geq \tilde{k}_0$.

With the above preparation definitions, we first focus on the inactive dual components. For any run of the method, we know from (C.19) that there exists a (potentially random) $\tilde{k}_0 < \infty$ such that for all $k \geq \tilde{k}_0$,

$$\begin{aligned} \|[\boldsymbol{\mu}_k]_{(\mathcal{A}^*)-}\| &= \prod_{j=\tilde{k}_0}^{k-1} (1 - \bar{\alpha}_j) \cdot \|[\boldsymbol{\mu}_{\tilde{k}_0}]_{(\mathcal{A}^*)-}\| \leq \prod_{j=0}^{\tilde{k}_0-1} (1 - \alpha_j)^{-1} \|[\boldsymbol{\mu}_{\tilde{k}_0}]_{(\mathcal{A}^*)-}\| \cdot \prod_{j=0}^{k-1} (1 - \alpha_j) \\ &\leq \prod_{j=0}^{\tilde{k}_0-1} (1 - \alpha_j)^{-1} \|[\boldsymbol{\mu}_{\tilde{k}_0}]_{(\mathcal{A}^*)-}\| \exp\left(-\sum_{j=0}^{k-1} \alpha_j\right). \end{aligned}$$

Note that

$$\sum_{i=0}^{k-1} \alpha_k \geq \iota_1 \int_1^{k+1} x^{-b_1} dx = \begin{cases} \frac{\iota_1}{1-b_1} ((k+1)^{1-b_1} - 1) & \text{if } b_1 < 1, \\ \iota_1 \log(k+1) & \text{if } b_1 = 1. \end{cases} \quad (\text{C.21})$$

Combining the above two displays, we know

$$\|[\boldsymbol{\mu}_k]_{(\mathcal{A}^*)-}\| = \begin{cases} o(k^{-b}) \text{ for any } b > 0 & \text{if } b_1 < 1, \\ \mathcal{O}(k^{-\iota_1}) & \text{if } b_1 = 1. \end{cases}$$

We now focus on the primal and active dual components. For notational simplicity, we denote $\bar{\Delta} \mathbf{w}_k = (\bar{\Delta} \mathbf{x}_k, \bar{\boldsymbol{\lambda}}_k^{\text{sub}} - \boldsymbol{\lambda}_k, [\bar{\boldsymbol{\mu}}_k^{\text{sub}} - \boldsymbol{\mu}_k]_{\mathcal{A}^*})$ for any $k \geq 0$, while noting from (C.18) that we only have $\bar{\Delta} \mathbf{w}_k = -\mathbf{H}_k^{-1} \bar{\nabla}_{\mathbf{w}} \mathcal{L}_k$ for $k_0 \leq k < \tau_{k_0}$ with any given $k_0 \geq 0$. For any $k \geq 0$, we have

$$\begin{aligned} &\mathbf{w}_{k+1} - \mathbf{w}^* \\ &= \mathbf{w}_k - \mathbf{w}^* + \bar{\alpha}_k \bar{\Delta} \mathbf{w}_k = \mathbf{w}_k - \mathbf{w}^* + \alpha_k \bar{\Delta} \mathbf{w}_k + (\bar{\alpha}_k - \alpha_k) \bar{\Delta} \mathbf{w}_k \\ &= \mathbf{w}_k - \mathbf{w}^* - \alpha_k (\mathbf{H}^*)^{-1} \nabla_{\mathbf{w}} \mathcal{L}_k + \alpha_k (\mathbf{H}^*)^{-1} (\nabla_{\mathbf{w}} \mathcal{L}_k - \bar{\nabla}_{\mathbf{w}} \mathcal{L}_k) + \alpha_k ((\mathbf{H}^*)^{-1} - \mathbf{H}_k^{\dagger}) \bar{\nabla}_{\mathbf{w}} \mathcal{L}_k \end{aligned}$$

$$\begin{aligned}
& + \alpha_k(\mathbf{H}_k^\dagger \bar{\nabla} \mathbf{w} \mathcal{L}_k + \bar{\Delta} \mathbf{w}_k) + (\bar{\alpha}_k - \alpha_k) \bar{\Delta} \mathbf{w}_k \\
& = \mathbf{w}_k - \mathbf{w}^* - \alpha_k \mathbf{H}_k^\dagger \nabla \mathbf{w} \mathcal{L}_k + \alpha_k (\mathbf{H}^*)^{-1} (\nabla \mathbf{w} \mathcal{L}_k - \bar{\nabla} \mathbf{w} \mathcal{L}_k) + \alpha_k ((\mathbf{H}^*)^{-1} - \mathbf{H}_k^\dagger) (\bar{\nabla} \mathbf{w} \mathcal{L}_k - \nabla \mathbf{w} \mathcal{L}_k) \\
& + \alpha_k(\mathbf{H}_k^\dagger \bar{\nabla} \mathbf{w} \mathcal{L}_k + \bar{\Delta} \mathbf{w}_k) + (\bar{\alpha}_k - \alpha_k) \bar{\Delta} \mathbf{w}_k \\
& = (1 - \alpha_k)(\mathbf{w}_k - \mathbf{w}^*) - \alpha_k \mathbf{H}_k^\dagger (\nabla \mathbf{w} \mathcal{L}_k - \mathbf{H}^*(\mathbf{w}_k - \mathbf{w}^*)) - \alpha_k (\mathbf{H}_k^\dagger \mathbf{H}^* - \mathbf{I})(\mathbf{w}_k - \mathbf{w}^*) \\
& - \alpha_k (\mathbf{H}_k^\dagger - (\mathbf{H}^*)^{-1}) (\bar{\nabla} \mathbf{w} \mathcal{L}_k - \nabla \mathbf{w} \mathcal{L}_k) + \alpha_k (\mathbf{H}_k^\dagger \bar{\nabla} \mathbf{w} \mathcal{L}_k + \bar{\Delta} \mathbf{w}_k) + (\bar{\alpha}_k - \alpha_k) \bar{\Delta} \mathbf{w}_k \\
& - \alpha_k (\mathbf{H}^*)^{-1} (\bar{\nabla} \mathbf{w} \mathcal{L}_k - \nabla \mathbf{w} \mathcal{L}_k).
\end{aligned}$$

Applying the above recursion from 0 to k , we have

$$\begin{aligned}
& \mathbf{w}_{k+1} - \mathbf{w}^* - \prod_{i=0}^k (1 - \alpha_i) (\mathbf{w}_0 - \mathbf{w}^*) \\
& = - \sum_{i=0}^k \prod_{j=i+1}^k (1 - \alpha_j) \alpha_i \left\{ \mathbf{H}_i^\dagger (\nabla \mathbf{w} \mathcal{L}_i - \mathbf{H}^*(\mathbf{w}_i - \mathbf{w}^*)) + (\mathbf{H}_i^\dagger \mathbf{H}^* - \mathbf{I})(\mathbf{w}_i - \mathbf{w}^*) \right\} \\
& - \sum_{i=0}^k \prod_{j=i+1}^k (1 - \alpha_j) \alpha_i \left\{ (\mathbf{H}_i^\dagger - (\mathbf{H}^*)^{-1}) (\bar{\nabla} \mathbf{w} \mathcal{L}_i - \nabla \mathbf{w} \mathcal{L}_i) - (\mathbf{H}_i^\dagger \bar{\nabla} \mathbf{w} \mathcal{L}_i + \bar{\Delta} \mathbf{w}_i) - \frac{\bar{\alpha}_i - \alpha_i}{\alpha_i} \bar{\Delta} \mathbf{w}_i \right\} \\
& - \sum_{i=0}^k \prod_{j=i+1}^k (1 - \alpha_j) \alpha_i (\mathbf{H}^*)^{-1} (\bar{\nabla} \mathbf{w} \mathcal{L}_i - \nabla \mathbf{w} \mathcal{L}_i) =: -\mathcal{J}_{1,k} - \mathcal{J}_{2,k} - \mathcal{J}_{3,k}. \tag{C.22}
\end{aligned}$$

First, for the term $\prod_{i=0}^k (1 - \alpha_i) (\mathbf{w}_0 - \mathbf{w}^*)$ on the left, we note from (C.21) that $\prod_{i=0}^k (1 - \alpha_i) (\mathbf{w}_0 - \mathbf{w}^*) = o(\sqrt{\alpha_k})$ as long as $\iota_1 > 0.5$ when $b_1 = 1$. Next, we aim to show $\mathcal{J}_{1,k} = o_p(\sqrt{\alpha_k})$, $\mathcal{J}_{2,k} = o_p(\sqrt{\alpha_k})$, and asymptotic normality of $\mathcal{J}_{3,k}$. Throughout the proof, we use $o_p(\cdot)$ and $\mathcal{O}_p(\cdot)$ to denote the standard small- and big- O notation in the probability sense. We state the following preparation lemma.

Lemma C.8 (Lemma C.3 in Na, 2025). *Consider a sequence of random variables $\{X_k\}_{k=0}^\infty$ and a sequence of events $\{\mathcal{A}_k\}_{k=0}^\infty$. Let $\tau_{k_0} = \inf\{k \geq k_0 : \mathcal{A}_k \text{ happens}\}$ be the first index k after k_0 such that \mathcal{A}_k happens. Suppose that for each realization of the sequence, there exists a (potentially random) $\tilde{k}_0 < \infty$ such that $\tau_{\tilde{k}_0} = \infty$. Also, for the sequence $\alpha_k = \iota_1(k+1)^{-b_1}$ with $b_1 \in (0, 1]$ and $\iota_1 > 0.5\mathbf{1}_{b_1=1}$, suppose there exists a deterministic $\bar{k}_0 > 0$ such that for any fixed $k_0 \geq \bar{k}_0$,*

$X_k \mathbf{1}_{\tau_{k_0} > k} = o_p(\sqrt{\alpha_k})$. Then, we have

$$\sum_{i=0}^k \prod_{j=i+1}^k (1 - \alpha_j) \alpha_i X_i = o_p(\sqrt{\alpha_k}).$$

By the above lemma and (C.22), we know $\mathcal{J}_{1,k} = o_p(\sqrt{\alpha_k})$ and $\mathcal{J}_{2,k} = o_p(\sqrt{\alpha_k})$ as long as

$$\|\mathbf{H}_k^\dagger (\nabla_{\mathbf{w}} \mathcal{L}_k - \mathbf{H}^*(\mathbf{w}_k - \mathbf{w}^*))\| \mathbf{1}_{\tau_{k_0} > k} = o_p(\sqrt{\alpha_k}), \quad (\text{C.23a})$$

$$\|(\mathbf{H}_k^\dagger \mathbf{H}^* - \mathbf{I})(\mathbf{w}_k - \mathbf{w}^*)\| \mathbf{1}_{\tau_{k_0} > k} = o_p(\sqrt{\alpha_k}), \quad (\text{C.23b})$$

$$\|(\mathbf{H}_k^\dagger - (\mathbf{H}^*)^{-1})(\bar{\nabla}_{\mathbf{w}} \mathcal{L}_k - \nabla_{\mathbf{w}} \mathcal{L}_k)\| \mathbf{1}_{\tau_{k_0} > k} = o_p(\sqrt{\alpha_k}), \quad (\text{C.23c})$$

$$\|(\mathbf{H}_k^\dagger \bar{\nabla}_{\mathbf{w}} \mathcal{L}_k + \bar{\Delta} \mathbf{w}_k)\| \mathbf{1}_{\tau_{k_0} > k} = o_p(\sqrt{\alpha_k}), \quad (\text{C.23d})$$

$$\frac{\bar{\alpha}_k - \alpha_k}{\alpha_k} \|\bar{\Delta} \mathbf{w}_k\| \mathbf{1}_{\tau_{k_0} > k} = o_p(\sqrt{\alpha_k}). \quad (\text{C.23e})$$

For (C.23a), we apply the definition of τ_{k_0} in (C.20) and have

$$\|\mathbf{H}_k^\dagger (\nabla_{\mathbf{w}} \mathcal{L}_k - \mathbf{H}^*(\mathbf{w}_k - \mathbf{w}^*))\| \mathbf{1}_{\tau_{k_0} > k} \leq \Upsilon^2 \|\mathbf{w}_k - \mathbf{w}^*\|^2 \mathbf{1}_{\tau_{k_0} > k}.$$

For (C.23b), we have

$$\begin{aligned} \|(\mathbf{H}_k^\dagger \mathbf{H}^* - \mathbf{I})(\mathbf{w}_k - \mathbf{w}^*)\| \mathbf{1}_{\tau_{k_0} > k} &\stackrel{(\text{C.20})}{\leq} \Upsilon \|\mathbf{H}_k - \mathbf{H}^*\| \|\mathbf{w}_k - \mathbf{w}^*\| \mathbf{1}_{\tau_{k_0} > k} \\ &\leq 0.5 \Upsilon \|\mathbf{H}_k - \mathbf{H}^*\|^2 \mathbf{1}_{\tau_{k_0} > k} + 0.5 \Upsilon \|\mathbf{w}_k - \mathbf{w}^*\|^2 \mathbf{1}_{\tau_{k_0} > k}. \end{aligned}$$

For (C.23c), we have

$$\begin{aligned} \|(\mathbf{H}_k^\dagger - (\mathbf{H}^*)^{-1})(\bar{\nabla}_{\mathbf{w}} \mathcal{L}_k - \nabla_{\mathbf{w}} \mathcal{L}_k)\| \mathbf{1}_{\tau_{k_0} > k} &\stackrel{(\text{C.20})}{\leq} \Upsilon \|\mathbf{H}_k - \mathbf{H}^*\| \|\bar{\mathbf{g}}_k - \nabla f_k\| \mathbf{1}_{\tau_{k_0} > k} \\ &\leq 0.5 \Upsilon \|\mathbf{H}_k - \mathbf{H}^*\|^2 \mathbf{1}_{\tau_{k_0} > k} + 0.5 \Upsilon \|\bar{\mathbf{g}}_k - \nabla f_k\|^2 \mathbf{1}_{\tau_{k_0} > k}. \end{aligned}$$

For (C.23d), it is trivial due to (C.18). For (C.23e), we have

$$\begin{aligned} \frac{\bar{\alpha}_k - \alpha_k}{\alpha_k} \|\bar{\Delta} \mathbf{w}_k\| \mathbf{1}_{\tau_{k_0} > k} &\leq \psi \alpha_k^{p-1} \|\bar{\Delta} \mathbf{w}_k\| \mathbf{1}_{\tau_{k_0} > k} \stackrel{(\text{C.18})}{\leq} \Upsilon \psi \alpha_k^{p-1} \|\bar{\nabla}_{\mathbf{w}} \mathcal{L}_k\| \mathbf{1}_{\tau_{k_0} > k} \\ &\stackrel{(\text{C.20})}{\leq} \Upsilon \psi \alpha_k^{p-1} \|\bar{\mathbf{g}}_k - \nabla f_k\| \mathbf{1}_{\tau_{k_0} > k} + \Upsilon^2 \psi \alpha_k^{p-1} \|\mathbf{w}_k - \mathbf{w}^*\| \mathbf{1}_{\tau_{k_0} > k}. \end{aligned}$$

Combining the above four displays, we know (C.23) is implied by showing that

$$\begin{aligned} (\|\mathbf{w}_k - \mathbf{w}^*\|^2 + \|\bar{\mathbf{g}}_k - \nabla f_k\|^2 + \|\mathbf{H}_k - \mathbf{H}^*\|^2) \mathbf{1}_{\tau_{k_0} > k} &= o_p(\sqrt{\alpha_k}), \\ \alpha_k^{p-1} (\|\mathbf{w}_k - \mathbf{w}^*\| + \|\bar{\mathbf{g}}_k - \nabla f_k\|) \mathbf{1}_{\tau_{k_0} > k} &= o_p(\sqrt{\alpha_k}). \end{aligned} \quad (\text{C.24})$$

We present the following two lemmas to corroborate the above result (C.24). See Appendices

D.9 and D.11 for the proofs, respectively.

Lemma C.9. *Under Assumption 3.6 and suppose $\alpha_k = \iota_1(k+1)^{-b_1}$, $\beta_k = \iota_2(k+1)^{-b_2}$ satisfy $\iota_1, \iota_2 > 0$, $0 < b_2 < b_1 \leq 1$, and $\iota_1 > 0.5$ if $b_1 = 1$. There exists a deterministic integer $\bar{k}_0 > 0$ such that for any $k_0 \geq \bar{k}_0$, the following result holds for some deterministic constant $\Upsilon(k_0) > 0$:*

$$\mathbb{E}[(\|\mathbf{w}_k - \mathbf{w}^*\|^2 + \|\bar{\mathbf{g}}_k - \nabla f_k\|^2)\mathbf{1}_{\tau_{k_0} > k}] \leq \Upsilon(k_0)\beta_k, \quad \forall k \geq k_0.$$

Lemma C.10. *Under Assumptions 3.6, 4.3 and suppose $\alpha_k = \iota_1(k+1)^{-b_1}$, $\beta_k = \iota_2(k+1)^{-b_2}$, $\gamma_k = \iota_3(k+1)^{-b_3}$ satisfy $\iota_1, \iota_2, \iota_3 > 0$, $b_1 \in (0, 1]$, $b_2 \in (0.5b_1, b_1)$, $b_3 \in (0.5b_1, 1]$, and $\iota_1 > 0.5$ if $b_1 = 1$ and $\iota_3 > 0.25b_1$ if $b_3 = 1$. There exists a deterministic integer $\bar{k}_0 > 0$ such that for any $k_0 \geq \bar{k}_0$, the following result holds for some deterministic constant $\Upsilon(k_0) > 0$:*

$$\mathbb{E}[\|\mathbf{H}_k - \mathbf{H}^*\|^2 \mathbf{1}_{\tau_{k_0} > k}] / (\Upsilon(k_0)\sqrt{\alpha_k}) \rightarrow 0 \quad \text{as } k \rightarrow \infty.$$

With Lemmas C.9 and C.10, and applying the conditions that $b_2 > 0.5 \geq 0.5b_1$ and $p > 1.5 - 0.5b_2/b_1$, we know (C.24) holds, leading to $\mathcal{J}_{1,k} = o_p(\sqrt{\alpha_k})$ and $\mathcal{J}_{2,k} = o_p(\sqrt{\alpha_k})$. Finally, we present the following lemma demonstrating the asymptotic normality property of $\mathcal{J}_{3,k}$. See Appendix D.12 for the proof.

Lemma C.11. *Under Assumptions 4.1, 4.5 and suppose $\alpha_k = \iota_1(k+1)^{-b_1}$, $\beta_k = \iota_2(k+1)^{-b_2}$ satisfy $\iota_1, \iota_2 > 0$, $b_1 \in ((2 - 2\delta)/(2 + \delta), 1]$, $b_2 \in (0, b_1)$, and $\iota_1 > 2/3$ if $b_1 = 1$. Then, we have $1/\sqrt{\alpha_k} \cdot \mathcal{J}_{3,k} \xrightarrow{d} \mathcal{N}(\mathbf{0}, \eta \cdot \mathbf{\Omega}^*)$, where η and $\mathbf{\Omega}^*$ are defined in Theorem 4.6.*

Noting that $\alpha_k/\bar{\alpha}_k \rightarrow 1$ almost surely since $p > 1.5 - 0.5b_2/b_1$ and $b_1 > b_2 > 0$ imply $p > 1$, we combine Lemma C.11 with (C.22) and complete the proof.

C.9 Proof of Theorem 4.8

By Assumption 4.1 and Lemma 4.4, it suffices to show the convergence of $\text{Cov}(\{\nabla F(\mathbf{x}_i; \zeta_i)\}_{i=0}^k)$.

In particular, we have

$$\|\text{Cov}(\{\nabla F(\mathbf{x}_i; \zeta_i)\}_{i=0}^k) - \text{Cov}(\nabla F(\mathbf{x}^*; \zeta))\|$$

$$\begin{aligned}
&\leq \left\| \frac{1}{k+1} \sum_{i=0}^k \nabla F(\mathbf{x}_i; \zeta_i) \nabla F(\mathbf{x}_i; \zeta_i)^\top - \mathbb{E} [\nabla F(\mathbf{x}^*; \zeta) \nabla F(\mathbf{x}^*; \zeta)^\top] \right\| \\
&\quad + \left\| \left(\frac{1}{k+1} \sum_{i=0}^k \nabla F(\mathbf{x}_i; \zeta_i) \right) \left(\frac{1}{k+1} \sum_{i=0}^k \nabla F(\mathbf{x}_i; \zeta_i) \right)^\top - \nabla f^* \nabla f^{*\top} \right\|. \quad (\text{C.25})
\end{aligned}$$

We analyze the two terms on the right hand side separately. For the first term in (C.25), we have

$$\begin{aligned}
&\left\| \frac{1}{k+1} \sum_{i=0}^k \nabla F(\mathbf{x}_i; \zeta_i) \nabla F(\mathbf{x}_i; \zeta_i)^\top - \mathbb{E} [\nabla F(\mathbf{x}^*; \zeta) \nabla F(\mathbf{x}^*; \zeta)^\top] \right\| \\
&\leq \left\| \frac{1}{k+1} \sum_{i=0}^k \nabla F(\mathbf{x}_i; \zeta_i) \nabla F(\mathbf{x}_i; \zeta_i)^\top - \mathbb{E}_i [\nabla F(\mathbf{x}_i; \zeta_i) \nabla F(\mathbf{x}_i; \zeta_i)^\top] \right\| \\
&\quad + \left\| \frac{1}{k+1} \sum_{i=0}^k \mathbb{E}_i [\nabla F(\mathbf{x}_i; \zeta_i) \nabla F(\mathbf{x}_i; \zeta_i)^\top] - \mathbb{E} [\nabla F(\mathbf{x}^*; \zeta) \nabla F(\mathbf{x}^*; \zeta)^\top] \right\|.
\end{aligned}$$

The first term converges to zero almost surely by bounded 4-th moment condition of $\nabla F(\mathbf{x}_k; \zeta_k)$ and the strong law of large numbers for square integrable martingale (Duflo, 1997, Theorem 1.3.15); the second term converges to zero almost surely by Assumptions 4.1, 4.5, and the fact that $x_k \rightarrow 0 \Rightarrow \frac{1}{k+1} \sum_{i=0}^k x_i \rightarrow 0$. Similarly, for the second term in (C.25), we just note that

$$\left\| \frac{1}{k+1} \sum_{i=0}^k \nabla F(\mathbf{x}_i; \zeta_i) - \nabla f^* \right\| \leq \left\| \frac{1}{k+1} \sum_{i=0}^k (\nabla F(\mathbf{x}_i; \zeta_i) - \nabla f_i) \right\| + \left\| \frac{1}{k+1} \sum_{i=0}^k (\nabla f_i - \nabla f^*) \right\|,$$

and it converges to zero almost surely due to the same reasons as above. This completes the proof.

D Proofs of Supporting Lemmas

D.1 Proof of Lemma C.1

We adapt the proof from (Davis et al., 2024, Lemma B.1). Consider the normalizing constant $C(\mathbf{v}) = 1 + \int h(\mathbf{v}^\top \mathbf{g}(\zeta)) d\mathcal{P}(\zeta)$. Since h is C^3 -smooth with bounded derivatives and $\mathbf{g} \in \mathcal{G}$ satisfies $\mathbb{E} \|\mathbf{g}(\zeta)\|^2 < \infty$, the dominated convergence theorem ensures $C(\mathbf{v})$ is twice continuously differentiable with $\nabla C(\mathbf{v}) = \int h'(\mathbf{v}^\top \mathbf{g}(\zeta)) \mathbf{g}(\zeta) d\mathcal{P}(\zeta)$, $\nabla^2 C(\mathbf{v}) = \int h''(\mathbf{v}^\top \mathbf{g}(\zeta)) \mathbf{g}(\zeta) \mathbf{g}(\zeta)^\top d\mathcal{P}(\zeta)$. Similarly, $\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w})$ is C^1 -smooth with $\nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{w}) = \mathbb{E}_{\zeta \sim \mathcal{P}} [\nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{w}; \zeta)]$. Define $\tilde{\nabla}_{\mathbf{w}} \mathcal{L}_{\mathbf{v}}(\mathbf{w}) = \int h(\mathbf{v}^\top \mathbf{g}(\zeta)) \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}; \zeta) d\mathcal{P}(\zeta)$. By the dominated convergence theorem, $\tilde{\nabla}_{\mathbf{w}} \mathcal{L}_{\mathbf{v}}(\mathbf{w})$ is contin-

uously differentiable in both \mathbf{v} and \mathbf{w} , and its partial derivatives are given by

$$\begin{aligned}\nabla_{\mathbf{v}}(\tilde{\nabla}_{\mathbf{w}}\mathcal{L}_{\mathbf{v}}(\mathbf{w})) &= \int h'(\mathbf{v}^\top \mathbf{g}(\zeta)) \nabla_{\mathbf{w}}\mathcal{L}(\mathbf{w}; \zeta) \mathbf{g}(\zeta)^\top d\mathcal{P}(\zeta), \\ \nabla_{\mathbf{w}}(\tilde{\nabla}_{\mathbf{w}}\mathcal{L}_{\mathbf{v}}(\mathbf{w})) &= \int h(\mathbf{v}^\top \mathbf{g}(\zeta)) \nabla_{\mathbf{w}}^2\mathcal{L}(\mathbf{w}; \zeta) d\mathcal{P}(\zeta),\end{aligned}\tag{D.1}$$

both of which are continuous in (\mathbf{w}, \mathbf{v}) . Thus, in a small neighborhood of $(\mathbf{w}, \mathbf{v}) = (\mathbf{w}^*, \mathbf{0})$, we apply $h(t) = t$ and $\mathbb{E}_{\zeta \sim \mathcal{P}}[\mathbf{g}(\zeta)] = \mathbf{0}$, and know $\nabla_{\mathbf{w}}\mathcal{L}_{\mathbf{v}}(\mathbf{w}) = \tilde{\nabla}_{\mathbf{w}}\mathcal{L}_{\mathbf{v}}(\mathbf{w})/C(\mathbf{v})$ is C^1 -smooth. Finally, the first identity in (C.5) directly follows from the fact that $\nabla_{\mathbf{w}}\mathcal{L}_{\mathbf{v}}(\mathbf{w})|_{\mathbf{v}=\mathbf{0}} = \nabla_{\mathbf{w}}\mathcal{L}(\mathbf{w})$ for any \mathbf{w} . For the second identity, we note that $\nabla C(\mathbf{0}) = \mathbf{0}$ and $\nabla^2 C(\mathbf{0}) = \mathbf{0}$, implying $C(\mathbf{v}) = 1 + o(\|\mathbf{v}\|^2)$. Thus, $\nabla_{\mathbf{v}\mathbf{w}}\mathcal{L}_{\mathbf{v}}(\mathbf{w})|_{\mathbf{v}=\mathbf{0}} = \nabla_{\mathbf{v}}(\tilde{\nabla}_{\mathbf{w}}\mathcal{L}_{\mathbf{v}}(\mathbf{w}))|_{\mathbf{v}=\mathbf{0}} \stackrel{(\text{D.1})}{=} \mathbb{E}_{\zeta \sim \mathcal{P}}[\nabla_{\mathbf{w}}\mathcal{L}(\mathbf{w}; \zeta) \mathbf{g}(\zeta)^\top]$. This completes the proof.

D.2 Proof of Lemma C.2

We adapt the proof from (Davis et al., 2024, Lemma B.2). Define the linearized mapping $\Psi(\mathbf{w}) = \nabla_{\mathbf{w}}\mathcal{L}(\mathbf{w}^*) + \nabla_{\mathbf{w}}^2\mathcal{L}(\mathbf{w}^*)(\mathbf{w} - \mathbf{w}^*) = \mathbf{H}^*(\mathbf{w} - \mathbf{w}^*)$, which is invertible and is C^1 -smooth. By Lemma C.1, we can equivalently express $\Psi(\mathbf{w}) = \nabla_{\mathbf{w}}\mathcal{L}_{\mathbf{v}}(\mathbf{w}^*)|_{\mathbf{v}=\mathbf{0}} + \nabla_{\mathbf{w}}^2\mathcal{L}_{\mathbf{v}}(\mathbf{w}^*)|_{\mathbf{v}=\mathbf{0}}(\mathbf{w} - \mathbf{w}^*)$. By the implicit function theorem (Dontchev and Rockafellar, 2009, Theorem 2D.6), the mapping $S(\mathbf{v})$ admits a single-valued localization $s(\mathbf{v})$ around \mathbf{w}^* and $\mathbf{v} = \mathbf{0}$ with $\nabla s(\mathbf{0}) = -\nabla\sigma(\mathbf{0}) \cdot \nabla_{\mathbf{v}\mathbf{w}}\mathcal{L}_{\mathbf{v}}(\mathbf{w}^*)|_{\mathbf{v}=\mathbf{0}} = -(\mathbf{H}^*)^{-1} \cdot \mathbb{E}_{\zeta \sim \mathcal{P}}[\nabla_{\mathbf{w}}\mathcal{L}(\mathbf{w}^*; \zeta) \mathbf{g}(\zeta)^\top]$. This completes the proof.

D.3 Proof of Lemma C.4

We present two preparation lemmas that are proved in Appendices D.4 and D.5, respectively.

Lemma D.1. *Let $\alpha_k = \iota_1(k+1)^{-b_1}$, $\beta_k = \iota_2(k+1)^{-b_2}$ for $\iota_1, \iota_2, b_1, b_2 > 0$ with $b_2 < \min\{b_1, 1\}$.*

Under Assumption 3.6 and the Lipschitz continuity of ∇f (cf. Assumption 3.5), we have

$$\mathbb{E}[\|\bar{\mathbf{g}}_k - \nabla f_k\|^2] = \mathcal{O}(\beta_k + \alpha_k^2/\beta_k^2).$$

Lemma D.2. *Suppose $\mathbf{B}_k \succeq \kappa_1 \mathbf{I}$, then*

$$\|\bar{\Delta}\mathbf{x}_k - \Delta\mathbf{x}_k\| \leq \kappa_1^{-1} \|\bar{\mathbf{g}}_k - \nabla f_k\|,$$

where $\bar{\Delta}\mathbf{x}_k$ is the solution of (17) and $\Delta\mathbf{x}_k$ is the solution with $\bar{\mathbf{g}}_k$ replaced by ∇f_k .

Now, we apply the results in Lemmas [D.1](#) and [D.2](#) and obtain

$$\mathbb{E}[\|\bar{\Delta}\mathbf{x}_k - \Delta\mathbf{x}_k\|] \leq \kappa_1^{-1} \mathbb{E}[\|\bar{\mathbf{g}}_k - \nabla f_k\|] \leq \kappa_1^{-1} \{\mathbb{E}[\|\bar{\mathbf{g}}_k - \nabla f_k\|^2]\}^{1/2} = \mathcal{O}\left(\sqrt{\beta_k} + \alpha_k/\beta_k\right).$$

By the setup of α_k, β_k , it follows that $\sum_{k=0}^{\infty} \alpha_k \mathbb{E}[\|\bar{\Delta}\mathbf{x}_k - \Delta\mathbf{x}_k\|] < \infty$ if $b_1 + 0.5b_2 > 1$ and $2b_1 - b_2 > 1$. Equivalently, this requires $b_1 \in (0.75, 1]$ and $b_2 \in (2 - 2b_1, 2b_1 - 1)$. We complete the proof.

D.4 Proof of Lemma [D.1](#)

From the update rule of $\bar{\mathbf{g}}_k$ in [\(15\)](#), we have

$$\begin{aligned} \bar{\mathbf{g}}_k - \nabla f_k &= \beta_k(\nabla F(\mathbf{x}_k; \zeta_k) - \nabla f_k) + (1 - \beta_k)(\bar{\mathbf{g}}_{k-1} - \nabla f_{k-1}) + (1 - \beta_k)(\nabla f_{k-1} - \nabla f_k) \\ &= \beta_k(\nabla F(\mathbf{x}_k; \zeta_k) - \nabla f_k) + (1 - \beta_k)\{\beta_{k-1}(\nabla F(\mathbf{x}_{k-1}; \zeta_{k-1}) - \nabla f_{k-1}) \\ &\quad + (1 - \beta_{k-1})(\bar{\mathbf{g}}_{k-2} - \nabla f_{k-2}) + (1 - \beta_{k-1})(\nabla f_{k-2} - \nabla f_{k-1})\} + (1 - \beta_k)(\nabla f_{k-1} - \nabla f_k) \\ &= \sum_{i=0}^k \left(\prod_{j=i+1}^k (1 - \beta_j) \right) \beta_i(\nabla F(\mathbf{x}_i; \zeta_i) - \nabla f_i) + \sum_{i=0}^k \left(\prod_{j=i}^k (1 - \beta_j) \right) (\nabla f_{i-1} - \nabla f_i) \\ &:= \mathcal{W}_{1,k} + \mathcal{W}_{2,k}, \end{aligned} \tag{D.2}$$

where we denote $\nabla f_{-1} = \bar{\mathbf{g}}_{-1}$ in the second last equality. By Lemma [B.4](#) and noting that

$\alpha_k \leq \bar{\alpha}_k \leq \alpha_k + \psi\alpha_k^p$, we have

$$\|\mathcal{W}_{2,k}\| \leq \prod_{j=0}^k (1 - \beta_j)(\|\bar{\mathbf{g}}_{-1}\| + \|\nabla f_0\|) + \sum_{i=1}^k \left(\prod_{j=i}^k (1 - \beta_j) \right) \bar{\alpha}_{i-1} M_{\ell, \mathbf{u}} = \mathcal{O}(\alpha_k/\beta_k).$$

Moreover, by Assumption [3.6](#) and Lemma [B.4](#), we have

$$\begin{aligned} \mathbb{E}[\|\mathcal{W}_{1,k}\|^2] &= \sum_{i=0}^k \left(\prod_{j=i+1}^k (1 - \beta_j) \right)^2 \beta_i^2 \mathbb{E}[\|\nabla F(\mathbf{x}_i; \zeta_i) - \nabla f_i\|^2] \\ &\leq \sigma_g^2 \sum_{i=0}^k \left(\prod_{j=i+1}^k (1 - \beta_j) \right)^2 \beta_i^2 = \mathcal{O}(\beta_k). \end{aligned}$$

Combining the above three displays together, we complete the proof.

D.5 Proof of Lemma D.2

The subproblem (17) at \mathbf{x}_k with the averaged gradient $\bar{\mathbf{g}}_k$ can be written as

$$\min_{\bar{\Delta}\mathbf{x} \in \Omega(\mathbf{x}_k; \theta_k)} \frac{1}{2} \|\bar{\Delta}\mathbf{x} + \mathbf{B}_k^{-1} \bar{\mathbf{g}}_k\|_{\mathbf{B}_k}^2,$$

where $\|\mathbf{x}\|_{\mathbf{B}_k}^2 = \mathbf{x}^\top \mathbf{B}_k \mathbf{x}$. By the optimality condition using variational inequalities, we have

$$\langle \Delta\mathbf{x}_k - \bar{\Delta}\mathbf{x}_k, -\mathbf{B}_k^{-1} \bar{\mathbf{g}}_k - \bar{\Delta}\mathbf{x}_k \rangle_{\mathbf{B}_k} \leq 0.$$

Similarly, since $\Delta\mathbf{x}_k$ is the solution with true gradient ∇f_k , it satisfies

$$\langle \bar{\Delta}\mathbf{x}_k - \Delta\mathbf{x}_k, -\mathbf{B}_k^{-1} \nabla f_k - \Delta\mathbf{x}_k \rangle_{\mathbf{B}_k} \leq 0.$$

Adding the two inequalities yields

$$\begin{aligned} 0 &\geq \langle \Delta\mathbf{x}_k - \bar{\Delta}\mathbf{x}_k, -\mathbf{B}_k^{-1} \bar{\mathbf{g}}_k - \bar{\Delta}\mathbf{x}_k + \mathbf{B}_k^{-1} \nabla f_k + \Delta\mathbf{x}_k \rangle_{\mathbf{B}_k} \\ &= \|\Delta\mathbf{x}_k - \bar{\Delta}\mathbf{x}_k\|_{\mathbf{B}_k}^2 + \langle \bar{\Delta}\mathbf{x}_k - \Delta\mathbf{x}_k, \bar{\mathbf{g}}_k - \nabla f_k \rangle \\ &\geq \|\bar{\Delta}\mathbf{x}_k - \Delta\mathbf{x}_k\|_{\mathbf{B}_k}^2 - \|\bar{\Delta}\mathbf{x}_k - \Delta\mathbf{x}_k\| \cdot \|\bar{\mathbf{g}}_k - \nabla f_k\|. \end{aligned}$$

Since $\mathbf{B}_k \succeq \kappa_1 \mathbf{I}$, we know $\|\bar{\Delta}\mathbf{x}_k - \Delta\mathbf{x}_k\|_{\mathbf{B}_k}^2 \geq \kappa_1 \|\bar{\Delta}\mathbf{x}_k - \Delta\mathbf{x}_k\|^2$. This completes the proof.

D.6 Proof of Lemma C.5

We follow the proof of Lemma D.1 in Appendix D.4 and slightly abuse the notation of $\mathcal{W}_{1,k}$ and $\mathcal{W}_{2,k}$ there. In particular, we revisit the definition of $\bar{\mathbf{g}}_k$ and have

$$\begin{aligned} \bar{\mathbf{g}}_k - \nabla f^\star &= (1 - \beta_k)(\bar{\mathbf{g}}_{k-1} - \nabla f^\star) + \beta_k(\nabla F(\mathbf{x}_k; \zeta_k) - \nabla f_k) + \beta_k(\nabla f_k - \nabla f^\star) \\ &= \sum_{i=0}^k \left(\prod_{j=i+1}^k (1 - \beta_j) \right) \beta_i(\nabla F(\mathbf{x}_i; \zeta_i) - \nabla f_i) \\ &\quad + \sum_{i=0}^k \left(\prod_{j=i+1}^k (1 - \beta_j) \right) \beta_i(\nabla f_i - \nabla f^\star) + \prod_{j=0}^k (1 - \beta_j)(\bar{\mathbf{g}}_{-1} - \nabla f^\star) \\ &:= \mathcal{W}_{1,k} + \mathcal{W}_{2,k}. \end{aligned}$$

Applying Lemma C.6 and the fact that $\nabla f_i \rightarrow \nabla f^*$, we immediately have $\mathcal{W}_{2,k} \rightarrow \mathbf{0}$ as $k \rightarrow \infty$.

For the term $\mathcal{W}_{1,k}$, we rewrite it as

$$\mathcal{W}_{1,k} = (1 - \beta_k)\mathcal{W}_{1,k-1} + \beta_k(\nabla F(\mathbf{x}_k; \zeta_k) - \nabla f_k).$$

Thus, for large enough k , we have from Assumption 3.6 that

$$\begin{aligned} \mathbb{E}_k[\|\mathcal{W}_{1,k}\|^2] &= (1 - \beta_k)^2\|\mathcal{W}_{1,k-1}\|^2 + \beta_k^2\mathbb{E}_k[\|\nabla F(\mathbf{x}_k; \zeta_k) - \nabla f_k\|^2] \\ &\leq (1 - \beta_k)\|\mathcal{W}_{1,k-1}\|^2 + \beta_k^2\sigma_g^2. \end{aligned}$$

Since $\sum_{k=0}^{\infty} \beta_k = \infty$ and $\sum_{k=0}^{\infty} \beta_k^2 < \infty$, by Robbins-Siegmund theorem (Robbins and Siegmund, 1971), we know $\|\mathcal{W}_{1,k}\|^2$ is convergent and $\sum_{k=0}^{\infty} \beta_k\|\mathcal{W}_{1,k-1}\|^2 < \infty$ almost surely. This further implies $\|\mathcal{W}_{1,k}\| \rightarrow 0$ almost surely and we complete the proof.

D.7 Proof of Lemma C.6

For any $\epsilon > 0$, there exists $i' > 0$ such that $|e_i| \leq \epsilon$ and $\beta_i \in [0, 1]$ for any $i \geq i'$. Thus, for $k \geq i'$,

$$\left| a \prod_{j=0}^k (1 - \beta_j) \right| \leq |a| \prod_{j=0}^{i'-1} |1 - \beta_j| \cdot \prod_{j=i'}^k (1 - \beta_j) \leq |a| \prod_{j=0}^{i'-1} |1 - \beta_j| \cdot \exp\left(-\sum_{j=i'}^k \beta_j\right),$$

and

$$\begin{aligned} \left| \sum_{i=0}^k \prod_{j=i+1}^k (1 - \beta_j) \beta_i e_i \right| &\leq \sum_{i=0}^{i'-1} \prod_{j=i+1}^k |1 - \beta_j| |\beta_i e_i| + \sum_{i=i'}^k \prod_{j=i+1}^k (1 - \beta_j) |\beta_i e_i| \\ &\leq \prod_{j=i'}^k (1 - \beta_j) \cdot \sum_{i=0}^{i'-1} \prod_{j=i+1}^{i'-1} |1 - \beta_j| |\beta_i e_i| + \epsilon \sum_{i=i'}^k \prod_{j=i+1}^k (1 - \beta_j) \beta_i \\ &= \prod_{j=i'}^k (1 - \beta_j) \cdot \sum_{i=0}^{i'-1} \prod_{j=i+1}^{i'-1} |1 - \beta_j| |\beta_i e_i| + \epsilon \left\{ 1 - \prod_{j=i'}^k (1 - \beta_j) \right\} \\ &\leq \exp\left(-\sum_{j=i'}^k \beta_j\right) \cdot \sum_{i=0}^{i'-1} \prod_{j=i+1}^{i'-1} |1 - \beta_j| |\beta_i e_i| + \epsilon. \end{aligned}$$

Since $\sum_{i=0}^{\infty} \beta_i = \infty$, we can find $k' \geq i'$ large enough such that $\exp(-\sum_{j=i'}^k \beta_j) \cdot \{|a| \prod_{j=0}^{i'-1} |1 - \beta_j| + \sum_{i=0}^{i'-1} \prod_{j=i+1}^{i'-1} |1 - \beta_j| |\beta_i e_i|\} \leq \epsilon$ for any $k \geq k'$. Then, we obtain for $k \geq k'$ that

$$\left| a \prod_{j=0}^k (1 - \beta_j) \right| + \left| \sum_{i=0}^k \prod_{j=i+1}^k (1 - \beta_j) \beta_i e_i \right| \leq 2\epsilon.$$

This shows $a \prod_{j=0}^k (1 - \beta_j) + \sum_{i=0}^k \prod_{j=i+1}^k (1 - \beta_j) \beta_i e_i \rightarrow 0$ as $k \rightarrow \infty$. We complete the proof.

D.8 Proof of Lemma C.7

We follow the proof of Lemma C.5 in Appendix D.6. By the update rule (15), we have

$$\begin{aligned}
\bar{\mathbf{Q}}_k - \nabla^2 f^\star &= (1 - \gamma_k)(\bar{\mathbf{Q}}_{k-1} - \nabla^2 f^\star) + \gamma_k(\nabla^2 F(\mathbf{x}_k; \zeta_k) - \nabla^2 f_k) + \gamma_k(\nabla^2 f_k - \nabla^2 f^\star) \\
&= \sum_{i=0}^k \left(\prod_{j=i+1}^k (1 - \gamma_j) \right) \gamma_i (\nabla^2 F(\mathbf{x}_i; \zeta_i) - \nabla^2 f_i) \\
&\quad + \sum_{i=0}^k \left(\prod_{j=i+1}^k (1 - \gamma_j) \right) \gamma_i (\nabla^2 f_i - \nabla^2 f^\star) + \prod_{j=0}^k (1 - \gamma_j) (\bar{\mathbf{Q}}_{-1} - \nabla^2 f^\star) \\
&:= \mathcal{V}_{1,k} + \mathcal{V}_{2,k}.
\end{aligned}$$

Applying Lemma C.6 and the fact that $\nabla^2 f_i \rightarrow \nabla^2 f^\star$, we immediately have $\mathcal{V}_{2,k} \rightarrow \mathbf{0}$ as $k \rightarrow \infty$.

For the term $\mathcal{V}_{1,k}$, we rewrite it as

$$\mathcal{V}_{1,k} = (1 - \gamma_k)\mathcal{V}_{1,k-1} + \gamma_k(\nabla^2 F(\mathbf{x}_k; \zeta_k) - \nabla^2 f_k).$$

Thus, for large enough k , we have from Assumption 4.3 that

$$\begin{aligned}
\mathbb{E}_k[\|\mathcal{V}_{1,k}\|_F^2] &= (1 - \gamma_k)^2 \|\mathcal{V}_{1,k-1}\|_F^2 + \gamma_k^2 \mathbb{E}_k[\|\nabla^2 F(\mathbf{x}_k; \zeta_k) - \nabla^2 f_k\|_F^2] \\
&\leq (1 - \gamma_k) \|\mathcal{V}_{1,k-1}\|_F^2 + \gamma_k^2 d\sigma_H^2.
\end{aligned}$$

Since $\sum_{k=0}^\infty \gamma_k = \infty$ and $\sum_{k=0}^\infty \gamma_k^2 < \infty$, by Robbins-Siegmund theorem (Robbins and Siegmund, 1971), we know $\|\mathcal{V}_{1,k}\|_F^2$ is convergent and $\sum_{k=0}^\infty \gamma_k \|\mathcal{V}_{1,k-1}\|_F^2 < \infty$ almost surely. This further implies $\|\mathcal{V}_{1,k}\|_F \rightarrow 0$ almost surely and we complete the proof.

D.9 Proof of Lemma C.9

We introduce a lemma that shows the recursive relation between $\|\mathbf{w}_k - \mathbf{w}^\star\|^2$ and $\|\bar{\mathbf{g}}_k - \nabla f_k\|^2$.

See Appendix D.10 for the proof.

Lemma D.3. *Under Assumption 3.6 and suppose $\alpha_k = \iota_1(k+1)^{-b_1}$, $\beta_k = \iota_2(k+1)^{-b_2}$ satisfy $\iota_1, \iota_2, b_1 > 0$, $b_2 \in (0, 1)$. There exists a deterministic integer $\bar{k}_0 > 0$ such that for any $k_0 \geq \bar{k}_0$, the following result holds for some deterministic constant $\Upsilon_1(k_0) > 0$: for any $k \geq k_0$,*

$$\mathbb{E}[\|\mathbf{w}_{k+1} - \mathbf{w}^\star\|^2 \mathbf{1}_{\tau_{k_0} > k+1}]$$

$$\begin{aligned}
&\leq \{1 - 2(1 - 2/\Upsilon)\alpha_k\} \mathbb{E}[\|\mathbf{w}_k - \mathbf{w}^*\|^2 \mathbf{1}_{\tau_{k_0} > k}] + (2 + \psi)\Upsilon^3 \alpha_k \mathbb{E}[\|\bar{\mathbf{g}}_k - \nabla f_k\|^2 \mathbf{1}_{\tau_{k_0} > k}], \\
&\mathbb{E}[\|\bar{\mathbf{g}}_{k+1} - \nabla f_{k+1}\|^2 \mathbf{1}_{\tau_{k_0} > k+1}] \\
&\leq \Upsilon_1(k_0) \left\{ \beta_k + \left(\sum_{i=k_0}^k \prod_{j=i+1}^k (1 - \beta_j) \alpha_i \left\{ \mathbb{E}[(\|\bar{\mathbf{g}}_i - \nabla f_i\|^2 + \|\mathbf{w}_i - \mathbf{w}^*\|^2) \mathbf{1}_{\tau_{k_0} > i}] \right\}^{1/2} \right)^2 \right\},
\end{aligned}$$

where $\Upsilon \geq 2$ is defined in the stopping time τ_{k_0} in (C.20).

With the above recursion, we claim that for any $q \geq 0$, there exists a deterministic integer $\bar{k}_0 > 0$ such that for any $k_0 \geq \bar{k}_0$, we have for some constant $\Upsilon_2(k_0) > 0$ that

$$\mathbb{E}[(\|\mathbf{w}_k - \mathbf{w}^*\|^2 + \|\bar{\mathbf{g}}_k - \nabla f_k\|^2) \mathbf{1}_{\tau_{k_0} > k}] \leq \Upsilon_2(k_0) (\beta_k + (\alpha_k/\beta_k)^{2q}), \quad \forall k \geq k_0. \quad (\text{D.3})$$

We prove the above statement by induction. When $q = 0$, we apply the definition of τ_{k_0} in (C.20) and have $\mathbb{E}[\|\mathbf{w}_k - \mathbf{w}^*\|^2 \mathbf{1}_{\tau_{k_0} > k}] \leq 1/\Upsilon$ for any $k \geq k_0$. Furthermore, by the recursive form of $\bar{\mathbf{g}}_k$ in (15) and Assumption 3.6, for any $k \geq k_0$,

$$\begin{aligned}
\mathbb{E}[\|\bar{\mathbf{g}}_k\|^2 \mathbf{1}_{\tau_{k_0} > k}] &\leq (1 - \beta_k) \mathbb{E}[\|\bar{\mathbf{g}}_{k-1}\|^2 \mathbf{1}_{\tau_{k_0} > k-1}] + \beta_k \mathbb{E}[\|\nabla F(\mathbf{x}_k; \zeta_k)\|^2 \mathbf{1}_{\tau_{k_0} > k}] \\
&\leq (1 - \beta_k) \mathbb{E}[\|\bar{\mathbf{g}}_{k-1}\|^2 \mathbf{1}_{\tau_{k_0} > k-1}] + 3\beta_k (\sigma_g^2 + \|\nabla f_k - \nabla f^*\|^2 \mathbf{1}_{\tau_{k_0} > k} + \|\nabla f^*\|^2) \\
&\stackrel{(\text{C.20})}{\leq} (1 - \beta_k) \mathbb{E}[\|\bar{\mathbf{g}}_{k-1}\|^2 \mathbf{1}_{\tau_{k_0} > k-1}] + 3\beta_k (\sigma_g^2 + 1/\Upsilon + \|\nabla f^*\|^2).
\end{aligned}$$

The above display, together with the fact that $\mathbb{E}[\|\bar{\mathbf{g}}_{k_0}\|^2] \leq \Upsilon_3(k_0)$ for some constant $\Upsilon_3(k_0)$ due to Assumption 3.6, implies that $\mathbb{E}[\|\bar{\mathbf{g}}_k - \nabla f_k\|^2 \mathbf{1}_{\tau_{k_0} > k}] \leq \Upsilon_4(k_0)$, $\forall k \geq k_0$ for some constant $\Upsilon_4(k_0)$. This verifies (D.3) for $q = 0$. Suppose that (D.3) holds for $q \geq 0$, then we establish the results for $q + 1$. In particular, we apply Lemmas D.3 and B.4 and have (\lesssim hides deterministic constants that may depend on k_0)

$$\begin{aligned}
&\mathbb{E}[\|\bar{\mathbf{g}}_{k+1} - \nabla f_{k+1}\|^2 \mathbf{1}_{\tau_{k_0} > k+1}] \\
&\lesssim \beta_k + \left(\sum_{i=k_0}^k \prod_{j=i+1}^k (1 - \beta_j) \alpha_i \left\{ \mathbb{E}[(\|\bar{\mathbf{g}}_i - \nabla f_i\|^2 + \|\mathbf{w}_i - \mathbf{w}^*\|^2) \mathbf{1}_{\tau_{k_0} > i}] \right\}^{1/2} \right)^2 \\
&\lesssim \beta_k + \left(\sum_{i=k_0}^k \prod_{j=i+1}^k (1 - \beta_j) \alpha_i \left\{ \beta_i + \left(\frac{\alpha_i}{\beta_i} \right)^{2q} \right\}^{1/2} \right)^2 \lesssim \beta_k + \left(\frac{\alpha_k}{\beta_k} \right)^{2(q+1)}.
\end{aligned}$$

Moreover, we still apply Lemma D.3 and have for any $b > 0$ and $k \geq k_0$,

$$\begin{aligned}
& \mathbb{E}[\|\mathbf{w}_{k+1} - \mathbf{w}^*\|^2 \mathbf{1}_{\tau_{k_0} > k+1}] \\
& \leq \{1 - 2(1 - 2/\Upsilon)\alpha_k\} \mathbb{E}[\|\mathbf{w}_k - \mathbf{w}^*\|^2 \mathbf{1}_{\tau_{k_0} > k}] + (2 + \psi)\Upsilon^3 \alpha_k \mathbb{E}[\|\bar{\mathbf{g}}_k - \nabla f_k\|^2 \mathbf{1}_{\tau_{k_0} > k}] \\
& \leq (2 + \psi)\Upsilon^3 \sum_{i=k_0}^k \prod_{j=i+1}^k \{1 - 2(1 - 2/\Upsilon)\alpha_j\} \alpha_i \mathbb{E}[\|\bar{\mathbf{g}}_i - \nabla f_i\|^2 \mathbf{1}_{\tau_{k_0} > i}] \\
& \quad + \prod_{j=k_0}^k \{1 - 2(1 - 2/\Upsilon)\alpha_j\} \mathbb{E}[\|\mathbf{w}_{k_0} - \mathbf{w}^*\|^2 \mathbf{1}_{\tau_{k_0} > k_0}] \\
& \stackrel{(\text{C.21})}{\lesssim} \sum_{i=k_0}^k \prod_{j=i+1}^k \{1 - 2(1 - 2/\Upsilon)\alpha_j\} \alpha_i \mathbb{E}[\|\bar{\mathbf{g}}_i - \nabla f_i\|^2 \mathbf{1}_{\tau_{k_0} > i}] + (k^{-b} \mathbf{1}_{b_1 < 1} + k^{-2(1-2/\Upsilon)\iota_1} \mathbf{1}_{b_1=1}) \\
& \lesssim \sum_{i=0}^k \prod_{j=i+1}^k \{1 - 2(1 - 2/\Upsilon)\alpha_j\} \alpha_i \left\{ \beta_i + \left(\frac{\alpha_i}{\beta_i} \right)^{2(q+1)} \right\} + (k^{-b} \mathbf{1}_{b_1 < 1} + k^{-2(1-2/\Upsilon)\iota_1} \mathbf{1}_{b_1=1}) \\
& \lesssim \beta_k + (\alpha_k/\beta_k)^{2(q+1)},
\end{aligned}$$

where the last inequality is because when $b_1 = 1$, we have $2(1 - 2/\Upsilon)\iota_1 > 1 > b_2$ by choosing Υ is large enough. This completes the induction step of (D.3) and further completes the proof.

D.10 Proof of Lemma D.3

By the local update (C.19), we know for any $k \geq k_0$,

$$\begin{aligned}
& \|\mathbf{w}_{k+1} - \mathbf{w}^*\|^2 = \|\mathbf{w}_k - \mathbf{w}^* + \bar{\alpha}_k \bar{\Delta} \mathbf{w}_k\|^2 = \|\mathbf{w}_k - \mathbf{w}^* - \bar{\alpha}_k \mathbf{H}_k^{-1} \bar{\nabla} \mathbf{w} \mathcal{L}_k\|^2 \\
& = \|\mathbf{w}_k - \mathbf{w}^* - \bar{\alpha}_k \mathbf{H}_k^{-1} \nabla \mathbf{w} \mathcal{L}_k - \bar{\alpha}_k \mathbf{H}_k^{-1} (\bar{\nabla} \mathbf{w} \mathcal{L}_k - \nabla \mathbf{w} \mathcal{L}_k)\|^2 \\
& = \|\mathbf{w}_k - \mathbf{w}^* - \bar{\alpha}_k \mathbf{H}_k^{-1} \nabla \mathbf{w} \mathcal{L}_k\|^2 + \bar{\alpha}_k^2 \|\mathbf{H}_k^{-1} (\bar{\nabla} \mathbf{w} \mathcal{L}_k - \nabla \mathbf{w} \mathcal{L}_k)\|^2 \\
& \quad - 2\bar{\alpha}_k \langle \mathbf{w}_k - \mathbf{w}^* - \bar{\alpha}_k \mathbf{H}_k^{-1} \nabla \mathbf{w} \mathcal{L}_k, \mathbf{H}_k^{-1} (\bar{\nabla} \mathbf{w} \mathcal{L}_k - \nabla \mathbf{w} \mathcal{L}_k) \rangle \\
& \leq \left(1 + \frac{1.5\bar{\alpha}_k}{\Upsilon}\right) \|\mathbf{w}_k - \mathbf{w}^* - \bar{\alpha}_k \mathbf{H}_k^{-1} \nabla \mathbf{w} \mathcal{L}_k\|^2 + \left(\bar{\alpha}_k^2 + \frac{\Upsilon \bar{\alpha}_k}{1.5}\right) \|\mathbf{H}_k^{-1} (\bar{\nabla} \mathbf{w} \mathcal{L}_k - \nabla \mathbf{w} \mathcal{L}_k)\|^2.
\end{aligned}$$

For the first term on the right hand side, we have

$$\begin{aligned}
& \|\mathbf{w}_k - \mathbf{w}^* - \bar{\alpha}_k \mathbf{H}_k^{-1} \nabla \mathbf{w} \mathcal{L}_k\|^2 \mathbf{1}_{\tau_{k_0} > k} \\
& = (\|\mathbf{w}_k - \mathbf{w}^*\|^2 - 2\bar{\alpha}_k \langle \mathbf{w}_k - \mathbf{w}^*, \mathbf{H}_k^{-1} \nabla \mathbf{w} \mathcal{L}_k \rangle + \bar{\alpha}_k^2 \|\mathbf{H}_k^{-1} \nabla \mathbf{w} \mathcal{L}_k\|^2) \mathbf{1}_{\tau_{k_0} > k} \\
& = ((1 - 2\bar{\alpha}_k) \|\mathbf{w}_k - \mathbf{w}^*\|^2 + 2\bar{\alpha}_k \langle \mathbf{w}_k - \mathbf{w}^*, \mathbf{w}_k - \mathbf{w}^* - \mathbf{H}_k^{-1} \nabla \mathbf{w} \mathcal{L}_k \rangle + \bar{\alpha}_k^2 \|\mathbf{H}_k^{-1} \nabla \mathbf{w} \mathcal{L}_k\|^2) \mathbf{1}_{\tau_{k_0} > k}
\end{aligned}$$

$$\stackrel{\text{(C.20)}}{\leq} \left(1 - 2\bar{\alpha}_k + \frac{2\bar{\alpha}_k}{\Upsilon} + \Upsilon^4 \bar{\alpha}_k^2\right) \|\mathbf{w}_k - \mathbf{w}^*\|^2 \mathbf{1}_{\tau_{k_0} > k}.$$

For the second term on the right hand side, we have

$$\|\mathbf{H}_k^{-1}(\bar{\nabla}_{\mathbf{w}} \mathcal{L}_k - \nabla_{\mathbf{w}} \mathcal{L}_k)\|^2 \mathbf{1}_{\tau_{k_0} > k} \stackrel{\text{(C.20)}}{\leq} \Upsilon^2 \|\bar{\mathbf{g}}_k - \nabla f_k\|^2 \mathbf{1}_{\tau_{k_0} > k}.$$

Combining the above three displays and noting that $\mathbf{1}_{\tau_{k_0} > k+1} \leq \mathbf{1}_{\tau_{k_0} > k}$, we obtain

$$\begin{aligned} \|\mathbf{w}_{k+1} - \mathbf{w}^*\|^2 \mathbf{1}_{\tau_{k_0} > k+1} &\leq \left(1 + \frac{1.5\bar{\alpha}_k}{\Upsilon}\right) \left(1 - 2\bar{\alpha}_k + \frac{2\bar{\alpha}_k}{\Upsilon} + \Upsilon^4 \bar{\alpha}_k^2\right) \|\mathbf{w}_k - \mathbf{w}^*\|^2 \mathbf{1}_{\tau_{k_0} > k} \\ &\quad + \Upsilon^2 \left(\bar{\alpha}_k^2 + \frac{\Upsilon \bar{\alpha}_k}{1.5}\right) \|\bar{\mathbf{g}}_k - \nabla f_k\|^2 \mathbf{1}_{\tau_{k_0} > k} \\ &\leq \left(1 - 2\bar{\alpha}_k + \frac{3.5\bar{\alpha}_k}{\Upsilon} + \Upsilon^4 \bar{\alpha}_k^2 + \frac{3\bar{\alpha}_k^2}{\Upsilon^2} + 1.5\Upsilon^3 \bar{\alpha}_k^3\right) \|\mathbf{w}_k - \mathbf{w}^*\|^2 \mathbf{1}_{\tau_{k_0} > k} \\ &\quad + \Upsilon^2 \left(\bar{\alpha}_k^2 + \frac{\Upsilon \bar{\alpha}_k}{1.5}\right) \|\bar{\mathbf{g}}_k - \nabla f_k\|^2 \mathbf{1}_{\tau_{k_0} > k}. \end{aligned}$$

To simplify the above display, we let \bar{k}_0 be large enough (with a deterministic threshold) such that

$$\Upsilon^4 \bar{\alpha}_k + \frac{3\bar{\alpha}_k}{\Upsilon^2} + 1.5\Upsilon^3 \bar{\alpha}_k^2 \leq \frac{0.5}{\Upsilon} \quad \text{and} \quad \left(\bar{\alpha}_k + \frac{\Upsilon}{1.5}\right) \bar{\alpha}_k \leq (2 + \psi) \Upsilon \alpha_k, \quad \forall k \geq \bar{k}_0. \quad (\text{D.4})$$

This is achievable since $b_1 > 0$ and $\limsup_{k \rightarrow \infty} \bar{\alpha}_k / \alpha_k \leq 1 + \psi$. Then, we obtain

$$\begin{aligned} &\|\mathbf{w}_{k+1} - \mathbf{w}^*\|^2 \mathbf{1}_{\tau_{k_0} > k+1} \\ &\leq \{1 - 2(1 - 2/\Upsilon)\bar{\alpha}_k\} \|\mathbf{w}_k - \mathbf{w}^*\|^2 \mathbf{1}_{\tau_{k_0} > k} + (2 + \psi) \Upsilon^3 \alpha_k \|\bar{\mathbf{g}}_k - \nabla f_k\|^2 \mathbf{1}_{\tau_{k_0} > k} \\ &\leq \{1 - 2(1 - 2/\Upsilon)\alpha_k\} \|\mathbf{w}_k - \mathbf{w}^*\|^2 \mathbf{1}_{\tau_{k_0} > k} + (2 + \psi) \Upsilon^3 \alpha_k \|\bar{\mathbf{g}}_k - \nabla f_k\|^2 \mathbf{1}_{\tau_{k_0} > k}, \end{aligned}$$

where the last inequality uses $\Upsilon \geq 2$. This shows the first argument of the lemma. For the second

argument of the lemma, we recall (D.2) from the proof in Appendix D.4 and have

$$\begin{aligned} \bar{\mathbf{g}}_{k+1} - \nabla f_{k+1} &= \sum_{i=0}^{k+1} \prod_{j=i+1}^{k+1} (1 - \beta_j) \beta_i (\nabla F(\mathbf{x}_i; \zeta_i) - \nabla f_i) + \sum_{i=0}^{k+1} \prod_{j=i}^{k+1} (1 - \beta_j) (\nabla f_{i-1} - \nabla f_i) \\ &:= \mathcal{W}_{1,k+1} + \mathcal{W}_{2,k+1}. \end{aligned} \quad (\text{D.5})$$

By Assumption 3.6 and Lemma B.4, we have

$$\begin{aligned} \mathbb{E}[\|\mathcal{W}_{1,k+1}\|^2 \mathbf{1}_{\tau_{k_0} > k+1}] &\leq \mathbb{E}[\|\mathcal{W}_{1,k+1}\|^2] \\ &= \sum_{i=0}^{k+1} \prod_{j=i+1}^{k+1} (1 - \beta_j)^2 \beta_i^2 \mathbb{E}[\|\nabla F(\mathbf{x}_i; \zeta_i) - \nabla f_i\|^2] = \mathcal{O}(\beta_k). \end{aligned} \quad (\text{D.6})$$

For the term $\mathcal{W}_{2,k+1}$, we use $\beta_k \leq 1$ for $k \geq k_0$ and separate the sum in $\mathcal{W}_{2,k+1}$ into two terms:

$$\begin{aligned}
& \|\mathcal{W}_{2,k+1}\| \mathbf{1}_{\tau_{k_0} > k+1} \\
& \leq \sum_{i=0}^{k_0} \prod_{j=i}^{k+1} |1 - \beta_j| \|\nabla f_{i-1} - \nabla f_i\| + \sum_{i=k_0+1}^{k+1} \prod_{j=i}^{k+1} (1 - \beta_j) \|\nabla f_{i-1} - \nabla f_i\| \mathbf{1}_{\tau_{k_0} > k+1} \\
& \stackrel{\text{(C.20)}}{\leq} \prod_{j=0}^{k+1} |1 - \beta_j| \sum_{i=0}^{k_0} \prod_{j=0}^{i-1} |1 - \beta_j|^{-1} \|\nabla f_{i-1} - \nabla f_i\| + \Upsilon \sum_{i=k_0+1}^{k+1} \prod_{j=i}^{k+1} (1 - \beta_j) \bar{\alpha}_{i-1} \|\bar{\Delta} \mathbf{x}_{i-1}\| \mathbf{1}_{\tau_{k_0} > k+1} \\
& \leq \Upsilon_1(k_0) \prod_{j=0}^{k+1} |1 - \beta_j| + \Upsilon \sum_{i=k_0}^k \prod_{j=i+1}^{k+1} (1 - \beta_j) \bar{\alpha}_i \|\bar{\Delta} \mathbf{x}_i\| \mathbf{1}_{\tau_{k_0} > i} \\
& \stackrel{\text{(C.20)}}{\leq} \Upsilon_1(k_0) \prod_{j=0}^{k+1} |1 - \beta_j| + (3 + 1.5\psi) \Upsilon \sum_{i=k_0}^k \prod_{j=i+1}^k (1 - \beta_j) \alpha_i \|\bar{\Delta} \mathbf{x}_i\| \mathbf{1}_{\tau_{k_0} > i}, \tag{D.7}
\end{aligned}$$

where the second last inequality holds for some deterministic constant $\Upsilon_1(k_0) > 0$ due to the boundedness of ∇f_i . For the first term on the right hand side, we apply (C.21) and know that $\prod_{j=0}^{k+1} |1 - \beta_j| = o(\beta_k)$. For the second term on the right hand side, we have for $k_0 \leq i < \tau_{k_0}$,

$$\begin{aligned}
\|\bar{\Delta} \mathbf{x}_i\|^2 &= \|\mathbf{H}_i^{-1} \bar{\nabla}_{\mathbf{w}} \mathcal{L}_i\|^2 \stackrel{\text{(C.20)}}{\leq} 2\Upsilon^2 (\|\bar{\nabla}_{\mathbf{w}} \mathcal{L}_i - \nabla_{\mathbf{w}} \mathcal{L}_i\|^2 + \|\nabla_{\mathbf{w}} \mathcal{L}_i\|^2) \\
&\stackrel{\text{(C.20)}}{\leq} 2\Upsilon^4 (\|\bar{\mathbf{g}}_i - \nabla f_i\|^2 + \|\mathbf{w}_i - \mathbf{w}^*\|^2). \tag{D.8}
\end{aligned}$$

Therefore, we have

$$\begin{aligned}
& \mathbb{E} \left[\left(\sum_{i=k_0}^k \prod_{j=i+1}^k (1 - \beta_j) \alpha_i \|\bar{\Delta} \mathbf{x}_i\| \mathbf{1}_{\tau_{k_0} > i} \right)^2 \right] \leq \left(\sum_{i=k_0}^k \prod_{j=i+1}^k (1 - \beta_j) \alpha_i \left\{ \mathbb{E}[\|\bar{\Delta} \mathbf{x}_i\|^2 \mathbf{1}_{\tau_{k_0} > i}] \right\}^{1/2} \right)^2 \\
& \stackrel{\text{(D.8)}}{\leq} 2\Upsilon^4 \left(\sum_{i=k_0}^k \prod_{j=i+1}^k (1 - \beta_j) \alpha_i \left\{ \mathbb{E}[(\|\bar{\mathbf{g}}_i - \nabla f_i\|^2 + \|\mathbf{w}_i - \mathbf{w}^*\|^2) \mathbf{1}_{\tau_{k_0} > i}] \right\}^{1/2} \right)^2. \tag{D.9}
\end{aligned}$$

Finally, combining (D.5), (D.6), (D.7), (D.9), we complete the proof.

D.11 Proof of Lemma C.10

We follow the proof in Appendix D.8 and have

$$\begin{aligned}
\bar{\mathbf{Q}}_k - \nabla^2 f^* &= \sum_{i=0}^k \prod_{j=i+1}^k (1 - \gamma_j) \gamma_i (\nabla^2 F(\mathbf{x}_i; \zeta_i) - \nabla^2 f_i) \\
&+ \sum_{i=0}^k \prod_{j=i+1}^k (1 - \gamma_j) \gamma_i (\nabla^2 f_i - \nabla^2 f^*) + \prod_{j=0}^k (1 - \gamma_j) (\bar{\mathbf{Q}}_{-1} - \nabla^2 f^*) := \mathcal{V}_{1,k} + \mathcal{V}_{2,k}.
\end{aligned}$$

We analyze each term separately. For the term $\mathcal{V}_{1,k}$, we apply Assumption 4.3 and have

$$\mathbb{E}[\|\mathcal{V}_{1,k}\|_F^2] = \sum_{i=0}^k \prod_{j=i+1}^k (1-\gamma_j)^2 \gamma_i^2 \mathbb{E}[\|\nabla^2 F(\mathbf{x}_i; \zeta_i) - \nabla^2 f_i\|_F^2] = \begin{cases} \mathcal{O}(\gamma_k) & \text{if } b_3 < 1 \\ o(\sqrt{\alpha_k}) & \text{if } b_3 = 1 \end{cases} = o(\sqrt{\alpha_k}),$$

where the second last equality uses Lemma B.4 and the condition that $\iota_3 > 0.25b_1$ if $b_3 = 1$; and

the last equality uses the condition $b_3 > 0.5b_1$. For the term $\mathcal{V}_{2,k}$, we apply (C.21) and have

$$\prod_{j=0}^k |1 - \gamma_j|^2 = \begin{cases} o(\gamma_k) & \text{if } b_3 < 1 \\ \mathcal{O}\left(\frac{1}{k^{2\iota_3}}\right) & \text{if } b_3 = 1 \end{cases} = o(\sqrt{\alpha_k}). \quad (\text{D.10})$$

Next, we analyze the first term in $\mathcal{V}_{2,k}$. We have for some constants $\Upsilon_1(k_0), \Upsilon_2(k_0) > 0$ that

$$\begin{aligned} & \mathbb{E} \left[\left\| \sum_{i=0}^k \prod_{j=i+1}^k (1 - \gamma_j) \gamma_i (\nabla^2 f_i - \nabla^2 f^*) \right\|^2 \mathbf{1}_{\tau_{k_0} > k} \right] \\ & \leq 2 \mathbb{E} \left[\left\| \sum_{i=0}^{k_0} \prod_{j=i+1}^k (1 - \gamma_j) \gamma_i (\nabla^2 f_i - \nabla^2 f^*) \right\|^2 + \left\| \sum_{i=k_0}^k \prod_{j=i+1}^k (1 - \gamma_j) \gamma_i (\nabla^2 f_i - \nabla^2 f^*) \right\|^2 \mathbf{1}_{\tau_{k_0} > k} \right] \\ & \stackrel{(\text{D.10})}{\leq} o(\Upsilon_1(k_0)\sqrt{\alpha_k}) + 2 \left(\sum_{i=k_0}^k \prod_{j=i+1}^k (1 - \gamma_j) \gamma_i \left\{ \mathbb{E}[\|\nabla f_i - \nabla f^*\|^2 \mathbf{1}_{\tau_{k_0} > i}] \right\}^{1/2} \right)^2 \\ & \stackrel{(\text{C.20})}{\leq} o(\Upsilon_1(k_0)\sqrt{\alpha_k}) + 2\Upsilon^2 \left(\sum_{i=k_0}^k \prod_{j=i+1}^k (1 - \gamma_j) \gamma_i \left\{ \mathbb{E}[\|\mathbf{x}_i - \mathbf{x}^*\|^2 \mathbf{1}_{\tau_{k_0} > i}] \right\}^{1/2} \right)^2 = o(\Upsilon_2(k_0)\sqrt{\alpha_k}), \end{aligned}$$

where the last equality is due to Lemmas B.4, C.9, and the facts that

(a) $b_3 \in (0, 1)$: We apply $b_2 > 0.5b_1$ and have

$$\sum_{i=k_0}^k \prod_{j=i+1}^k (1 - \gamma_j) \gamma_i \sqrt{\beta_i} = O(\sqrt{\beta_k}) = o(\alpha_k^{0.25}).$$

(b) $b_3 = 1$: We apply $b_2 > 0.5b_1$ and $\iota_3 > 0.25b_1 \Leftrightarrow 1 - \frac{0.5b_1}{b_2} \cdot \frac{0.5b_2}{\iota_3} > 0$, and have

$$\frac{1}{\alpha_k^{0.25}} \sum_{i=k_0}^k \prod_{j=i+1}^k (1 - \gamma_j) \gamma_i \sqrt{\beta_i} = \frac{(\sqrt{\beta_k})^{0.5b_1/b_2}}{\alpha_k^{0.25}} \frac{1}{(\sqrt{\beta_k})^{0.5b_1/b_2}} \sum_{i=k_0}^k \prod_{j=i+1}^k (1 - \gamma_j) \gamma_i \sqrt{\beta_i} = o(1).$$

Thus, we have shown that $\mathbb{E}[\|\bar{\mathbf{Q}}_k - \nabla^2 f^*\|^2 \mathbf{1}_{\tau_{k_0} > k}] = o(\Upsilon_3(k_0)\sqrt{\alpha_k})$ for some constant $\Upsilon_3(k_0) >$

0. Finally, by noting from (C.20) that

$$\begin{aligned} \|\mathbf{H}_k - \mathbf{H}^*\| \mathbf{1}_{\tau_{k_0} > k} & \leq (\|\bar{\mathbf{Q}}_k - \nabla^2 f^*\| + \|\nabla_{\mathbf{w}}^2 \mathcal{L}_k - \nabla_{\mathbf{w}}^2 \mathcal{L}^*\| + \|\nabla^2 f_k - \nabla^2 f^*\|) \mathbf{1}_{\tau_{k_0} > k} \\ & \leq (\|\bar{\mathbf{Q}}_k - \nabla^2 f^*\| + 2\Upsilon \|\mathbf{w}_k - \mathbf{w}^*\|) \mathbf{1}_{\tau_{k_0} > k}, \end{aligned}$$

and applying Lemma C.9, we complete the proof.

D.12 Proof of Lemma C.11

Recalling the definition of $\mathcal{J}_{3,k}$ from (C.22), we have

$$\begin{aligned}
\mathcal{J}_{3,k} &= \sum_{i=0}^k \prod_{j=i+1}^k (1 - \alpha_j) \alpha_i (\mathbf{H}^*)^{-1} (\bar{\nabla}_{\mathbf{w}} \mathcal{L}_i - \nabla_{\mathbf{w}} \mathcal{L}_i) \\
&\stackrel{(13)}{=} \sum_{i=0}^k \prod_{j=i+1}^k (1 - \alpha_j) \alpha_i \begin{pmatrix} \nabla_{\mathbf{x}}^2 \mathcal{L}^* & (\mathbf{J}^*)^\top \\ \mathbf{J}^* & \mathbf{0} \end{pmatrix}^{-1} \begin{pmatrix} \bar{\mathbf{g}}_i - \nabla f_i \\ \mathbf{0} \end{pmatrix} \\
&\stackrel{(D.2)}{=} \sum_{i=0}^k \prod_{j=i+1}^k (1 - \alpha_j) \alpha_i \begin{pmatrix} \nabla_{\mathbf{x}}^2 \mathcal{L}^* & (\mathbf{J}^*)^\top \\ \mathbf{J}^* & \mathbf{0} \end{pmatrix}^{-1} \left\{ \begin{pmatrix} \mathcal{W}_{1,i} \\ \mathbf{0} \end{pmatrix} + \begin{pmatrix} \mathcal{W}_{2,i} \\ \mathbf{0} \end{pmatrix} \right\} = \mathcal{J}_{3,k}^{(1)} + \mathcal{J}_{3,k}^{(2)}.
\end{aligned}$$

We first analyze $\mathcal{J}_{3,k}^{(2)}$ and aim to show that $\mathcal{J}_{3,k}^{(2)} = o_p(\sqrt{\alpha_k})$. By Lemma C.8, it suffices to show $\mathcal{W}_{2,k} \mathbf{1}_{\tau_{k_0} > k} = o_p(\sqrt{\alpha_k})$. From the definition of $\mathcal{W}_{2,k}$ in (D.2) and Lemma C.8, this result is further implied by

$$\|\nabla f_k - \nabla f_{k+1}\| \mathbf{1}_{\tau_{k_0} > k+1} \stackrel{(C.20)}{\leq} \Upsilon \bar{\alpha}_k \|\bar{\Delta} \mathbf{x}_k\| \mathbf{1}_{\tau_{k_0} > k} \leq o_p(\beta_k \sqrt{\alpha_k}).$$

By (D.8), Lemma C.9, and the fact that $\alpha_k \sqrt{\beta_k} = o(\beta_k \sqrt{\alpha_k})$, we obtain the above desired result.

Next, we analyze $\mathcal{J}_{3,k}^{(1)}$. By the definition of $\mathcal{W}_{1,k}$ in (D.2), we have

$$\begin{aligned}
\mathcal{J}_{3,k}^{(1)} &= \sum_{i=0}^k \prod_{j=i+1}^k (1 - \alpha_j) \alpha_i \begin{pmatrix} \nabla_{\mathbf{x}}^2 \mathcal{L}^* & (\mathbf{J}^*)^\top \\ \mathbf{J}^* & \mathbf{0} \end{pmatrix}^{-1} \begin{pmatrix} \mathcal{W}_{1,i} \\ \mathbf{0} \end{pmatrix} \\
&= \sum_{i=0}^k \prod_{j=i+1}^k (1 - \alpha_j) \alpha_i \sum_{h=0}^i \prod_{l=h+1}^i (1 - \beta_l) \beta_h \underbrace{\begin{pmatrix} \nabla_{\mathbf{x}}^2 \mathcal{L}^* & (\mathbf{J}^*)^\top \\ \mathbf{J}^* & \mathbf{0} \end{pmatrix}^{-1} \begin{pmatrix} \nabla F(\mathbf{x}_h; \zeta_h) - \nabla f_h \\ \mathbf{0} \end{pmatrix}}_{\phi_h^*} \\
&= \sum_{h=0}^k \underbrace{\sum_{i=h}^k \prod_{j=i+1}^k (1 - \alpha_j) \alpha_i \prod_{l=h+1}^i (1 - \beta_l) \beta_h}_{a_{h,k}} \phi_h^* = \sum_{h=0}^k a_{h,k} \phi_h^*.
\end{aligned}$$

We first claim that $\mathbb{E}_h[\phi_h^* \phi_h^{*\top}] \rightarrow \mathbf{\Omega}^*$ almost surely as $h \rightarrow \infty$. In fact, by Assumptions 4.1, 4.5, and the smoothness of $\nabla f(\mathbf{x})$, we have

$$\begin{aligned}
&\|\mathbb{E}_h[\nabla F(\mathbf{x}_h; \zeta_h) \nabla F(\mathbf{x}_h; \zeta_h)^\top - \nabla f_h \nabla f_h^\top] - \mathbb{E}[\nabla F(\mathbf{x}^*; \zeta) \nabla F(\mathbf{x}^*; \zeta)^\top - \nabla f^* \nabla f^{*\top}]\| \\
&\leq \|\mathbb{E}_h[\nabla F(\mathbf{x}_h; \zeta_h) \nabla F(\mathbf{x}_h; \zeta_h)^\top] - \mathbb{E}[\nabla F(\mathbf{x}^*; \zeta) \nabla F(\mathbf{x}^*; \zeta)^\top]\| \\
&\quad + (\|\nabla f_h\| + \|\nabla f^*\|) \|\nabla f_h - \nabla f^*\| \rightarrow 0 \quad \text{as } h \rightarrow \infty.
\end{aligned}$$

With this result, we then analyze the conditional variance process. We have

$$\begin{aligned}
& \frac{1}{\alpha_k} \sum_{h=0}^k a_{h,k}^2 \mathbb{E}_h[\phi_h^* \phi_h^{*\top}] \\
&= \frac{1}{\alpha_k} \sum_{h=0}^k \sum_{i=h}^k \sum_{i'=h}^k \prod_{j=i+1}^k (1-\alpha_j) \alpha_i \prod_{l=h+1}^i (1-\beta_l) \beta_h \prod_{j'=i'+1}^k (1-\alpha_{j'}) \alpha_{i'} \prod_{l'=h+1}^{i'} (1-\beta_{l'}) \beta_h \mathbb{E}_h[\phi_h^* \phi_h^{*\top}] \\
&= \frac{1}{\alpha_k} \sum_{i=0}^k \sum_{i'=0}^k \prod_{j=i+1}^k (1-\alpha_j) \alpha_i \prod_{j'=i'+1}^k (1-\alpha_{j'}) \alpha_{i'} \sum_{h=0}^{\min\{i,i'\}} \prod_{l=h+1}^i (1-\beta_l) \prod_{l'=h+1}^{i'} (1-\beta_{l'}) \beta_h^2 \mathbb{E}_h[\phi_h^* \phi_h^{*\top}] \\
&= \frac{2}{\alpha_k} \sum_{i=0}^k \sum_{i'=0}^i \prod_{j=i+1}^k (1-\alpha_j) \alpha_i \prod_{j'=i'+1}^k (1-\alpha_{j'}) \alpha_{i'} \sum_{h=0}^{i'} \prod_{l=h+1}^i (1-\beta_l) \prod_{l'=h+1}^{i'} (1-\beta_{l'}) \beta_h^2 \mathbb{E}_h[\phi_h^* \phi_h^{*\top}] \\
&\quad - \frac{1}{\alpha_k} \sum_{i=0}^k \prod_{j=i+1}^k (1-\alpha_j)^2 \alpha_i^2 \sum_{h=0}^i \prod_{l=h+1}^i (1-\beta_l)^2 \beta_h^2 \mathbb{E}_h[\phi_h^* \phi_h^{*\top}] \\
&= \frac{2}{\alpha_k} \sum_{i=0}^k \prod_{j=i+1}^k (1-\alpha_j)^2 \alpha_i^2 \sum_{i'=0}^i \prod_{j'=i'+1}^i (1-\alpha_{j'}) (1-\beta_{j'}) \alpha_{i'} \sum_{h=0}^{i'} \prod_{l'=h+1}^{i'} (1-\beta_{l'})^2 \beta_h^2 \mathbb{E}_h[\phi_h^* \phi_h^{*\top}] \\
&\quad - \frac{1}{\alpha_k} \sum_{i=0}^k \prod_{j=i+1}^k (1-\alpha_j)^2 \alpha_i^2 \sum_{h=0}^i \prod_{l=h+1}^i (1-\beta_l)^2 \beta_h^2 \mathbb{E}_h[\phi_h^* \phi_h^{*\top}].
\end{aligned}$$

We apply Lemma B.4 and note that

$$\begin{aligned}
& \lim_{i \rightarrow \infty} \frac{1}{\beta_i} \sum_{h=0}^i \prod_{l=h+1}^i (1-\beta_l)^2 \beta_h^2 \mathbb{E}_h[\phi_h^* \phi_h^{*\top}] = 0.5 \Omega^*, \\
& \lim_{i \rightarrow \infty} \frac{1}{\alpha_i} \sum_{i'=0}^i \prod_{j'=i'+1}^i (1-\alpha_{j'}) (1-\beta_{j'}) \alpha_{i'} \beta_{i'} = 1, \\
& \lim_{k \rightarrow \infty} \frac{1}{\alpha_k} \sum_{i=0}^k \prod_{j=i+1}^k (1-\alpha_j)^2 \alpha_i^2 = \eta := \begin{cases} 0.5, & \text{if } b_1 \in (0, 1), \\ \frac{\iota_1}{2\iota_1-1}, & \text{if } b_1 = 1, \end{cases} \\
& \lim_{k \rightarrow \infty} \frac{1}{\alpha_k} \sum_{i=0}^k \prod_{j=i+1}^k (1-\alpha_j)^2 \alpha_i^2 \beta_i = 0.
\end{aligned}$$

Combining the above two displays, we obtain almost surely,

$$\lim_{k \rightarrow \infty} \frac{1}{\alpha_k} \sum_{h=0}^k a_{h,k}^2 \mathbb{E}_h[\phi_h^* \phi_h^{*\top}] = \eta \cdot \Omega^*. \quad (\text{D.11})$$

We next verify the Lindeberg condition. It is equivalent to showing that for any $\epsilon > 0$,

$$\lim_{k \rightarrow \infty} \frac{1}{\alpha_k} \sum_{h=0}^k a_{h,k}^2 \mathbb{E}_h [\|\phi_h^*\|^2 \cdot \mathbf{1}_{\|a_{h,k} \phi_h^*\| \geq \epsilon \sqrt{\alpha_k}}] \leq \lim_{k \rightarrow \infty} \frac{1}{\epsilon^\delta \alpha_k^{1+0.5\delta}} \sum_{h=0}^k a_{h,k}^{2+\delta} \mathbb{E}_h [\|\phi_h^*\|^{2+\delta}] = 0.$$

By Assumption 4.5, it suffices to show

$$\frac{1}{\alpha_k^{1+0.5\delta}} \sum_{h=0}^k a_{h,k}^{2+\delta} \leq \frac{k+1}{\alpha_k^{1+0.5\delta}} \left(\frac{1}{k+1} \sum_{h=0}^k a_{h,k}^3 \right)^{\frac{2+\delta}{3}} = \left(\frac{\sum_{h=0}^k a_{h,k}^3}{\alpha_k^{\frac{3+1.5\delta}{2+\delta} + \frac{1-\delta}{b_1(2+\delta)}}} \right)^{\frac{2+\delta}{3}} \rightarrow 0. \quad (\text{D.12})$$

In particular, we have

$$\begin{aligned} \sum_{h=0}^k a_{h,k}^3 &= \sum_{h=0}^k \sum_{i=h}^k \sum_{i'=h}^k \sum_{i''=h}^k \prod_{j=i+1}^k (1-\alpha_j) \alpha_i \prod_{l=h+1}^i (1-\beta_l) \beta_h \prod_{j'=i'+1}^k (1-\alpha_{j'}) \alpha_{i'} \prod_{l'=h+1}^{i'} (1-\beta_{l'}) \beta_{i'} \\ &\quad \prod_{j''=i''+1}^k (1-\alpha_{j''}) \alpha_{i''} \prod_{l''=h+1}^{i''} (1-\beta_{l''}) \beta_{i''} \\ &= \sum_{i=0}^k \sum_{i'=0}^k \sum_{i''=0}^k \prod_{j=i+1}^k (1-\alpha_j) \alpha_i \prod_{j'=i'+1}^k (1-\alpha_{j'}) \alpha_{i'} \prod_{j''=i''+1}^k (1-\alpha_{j''}) \alpha_{i''} \\ &\quad \sum_{h=0}^{\min\{i,i',i''\}} \prod_{l=h+1}^i (1-\beta_l) \prod_{l'=h+1}^{i'} (1-\beta_{l'}) \prod_{l''=h+1}^{i''} (1-\beta_{l''}) \beta_h^3 \\ &\leq 6 \sum_{i=0}^k \sum_{i'=0}^i \sum_{i''=0}^{i'} \prod_{j=i+1}^k (1-\alpha_j) \alpha_i \prod_{j'=i'+1}^k (1-\alpha_{j'}) \alpha_{i'} \prod_{j''=i''+1}^k (1-\alpha_{j''}) \alpha_{i''} \\ &\quad \sum_{h=0}^{i''} \prod_{l=h+1}^i (1-\beta_l) \prod_{l'=h+1}^{i'} (1-\beta_{l'}) \prod_{l''=h+1}^{i''} (1-\beta_{l''}) \beta_h^3 \quad (i \geq i' \geq i'') \\ &= 6 \sum_{i=0}^k \prod_{j=i+1}^k (1-\alpha_j)^3 \alpha_i \sum_{i'=0}^i \prod_{j'=i'+1}^i (1-\alpha_{j'})^2 (1-\beta_{j'}) \alpha_{i'} \\ &\quad \sum_{i''=0}^{i'} \prod_{j''=i''+1}^{i'} (1-\alpha_{j''}) (1-\beta_{j''})^2 \alpha_{i''} \sum_{h=0}^{i''} \prod_{l''=h+1}^{i''} (1-\beta_{l''})^3 \beta_h^3. \end{aligned}$$

We apply Lemma B.4 and note that

$$\begin{aligned} \lim_{i'' \rightarrow \infty} \frac{1}{\beta_{i''}^2} \sum_{h=0}^{i''} \prod_{l''=h+1}^{i''} (1-\beta_{l''})^3 \beta_h^3 &= \frac{1}{3}, \\ \lim_{i' \rightarrow \infty} \frac{1}{\alpha_{i'} \beta_{i'}} \sum_{i''=0}^{i'} \prod_{j''=i''+1}^{i'} (1-\alpha_{j''}) (1-\beta_{j''})^2 \alpha_{i''} \beta_{i''}^2 &= \frac{1}{2}, \\ \lim_{i \rightarrow \infty} \frac{1}{\alpha_i^2} \sum_{i'=0}^i \prod_{j'=i'+1}^i (1-\alpha_{j'})^2 (1-\beta_{j'}) \alpha_{i'}^2 \beta_{i'} &= 1, \\ \lim_{k \rightarrow \infty} \frac{1}{\alpha_k^{\frac{3+1.5\delta}{2+\delta} + \frac{1-\delta}{b_1(2+\delta)}}} \sum_{i=0}^k \prod_{j=i+1}^k (1-\alpha_j)^3 \alpha_i^3 &= 0, \end{aligned}$$

where the last equality holds due to the following facts:

(a) $b_1 \in (0, 1)$: We apply $b_1 > (1 - \delta)/(1 + 0.5\delta)$ and have

$$\frac{1}{\alpha_k^{\frac{3+1.5\delta}{2+\delta} + \frac{1-\delta}{b_1(2+\delta)}}} \sum_{i=0}^k \prod_{j=i+1}^k (1 - \alpha_j)^3 \alpha_i^3 = \mathcal{O} \left(\frac{\alpha_k^2}{\alpha_k^{\frac{3+1.5\delta}{2+\delta} + \frac{1-\delta}{b_1(2+\delta)}}} \right) \rightarrow 0.$$

(b) $b_1 = 1$: We apply $\iota_1 > 2/3$ and have

$$\frac{1}{\alpha_k^{\frac{4+0.5\delta}{2+\delta}}} \sum_{i=0}^k \prod_{j=i+1}^k (1 - \alpha_j)^3 \alpha_i^3 = \mathcal{O} \left(\frac{\alpha_k^2}{\alpha_k^{\frac{4+0.5\delta}{2+\delta}}} \right) \rightarrow 0.$$

Thus, we have verified the Lindeberg condition. By the central limit theorem of martingale arrays (Hall and Heyde, 2014, Corollary 3.1), the results (D.11) and (D.12) lead to $1/\sqrt{\alpha_k} \cdot \mathcal{J}_{3,k}^{(1)} \xrightarrow{d} \mathcal{N}(\mathbf{0}, \eta \cdot \boldsymbol{\Omega}^*)$. This completes the proof.

E Boundedness of Dual Multipliers

Theorem E.1 shows that if the iteration sequence \mathbf{x}_k generated by our method converges to a feasible point $\bar{\mathbf{x}}$ satisfying EGMFCQ (Definition 2.1), then the corresponding Lagrangian dual multipliers $(\boldsymbol{\lambda}_k^{\text{sub}}, \boldsymbol{\mu}_k^{\text{sub}})$ associated with the subproblem (17) with the exact true gradient ∇f_k are bounded. This result validates and provides theoretical justification for the boundedness assumption on the dual variables in Assumption 3.5.

Theorem E.1. *Suppose $\bar{\mathbf{x}}$ is feasible for Problem (1) and satisfies EGMFCQ. If $\lim_{k \rightarrow \infty} \mathbf{x}_k = \bar{\mathbf{x}}$, then the corresponding dual multipliers $(\boldsymbol{\lambda}_k^{\text{sub}}, \boldsymbol{\mu}_k^{\text{sub}})$ of the SQP subproblem (17) are bounded.*

Proof. We prove it by contradiction. Suppose that there exists a sequence $\{(\mathbf{x}_k, \bar{\mathbf{B}}_k, \boldsymbol{\lambda}_k^{\text{sub}}, \boldsymbol{\mu}_k^{\text{sub}})\}$ such that $\mathbf{x}_k \rightarrow \bar{\mathbf{x}}$, $\|(\boldsymbol{\lambda}_k^{\text{sub}}, \boldsymbol{\mu}_k^{\text{sub}})\| \rightarrow \infty$ and $\kappa_1 \mathbf{I} \preceq \bar{\mathbf{B}}_k \preceq \kappa_2 \mathbf{I}$. Recall that $\Delta \mathbf{x}_k$ and $(\boldsymbol{\lambda}_k^{\text{sub}}, \boldsymbol{\mu}_k^{\text{sub}})$ denote the primal and dual solutions of the SQP subproblem (17) at \mathbf{x}_k with the true gradient ∇f_k , satisfying the following KKT conditions:

$$\begin{aligned} \nabla f_k + \bar{\mathbf{B}}_k \Delta \mathbf{x}_k + \nabla \mathbf{c}_k^\top \boldsymbol{\lambda}_k^{\text{sub}} - \boldsymbol{\mu}_{1,k}^{\text{sub}} + \boldsymbol{\mu}_{2,k}^{\text{sub}} &= \mathbf{0}, \\ \theta_k \mathbf{c}_k + \nabla \mathbf{c}_k \Delta \mathbf{x}_k &= \mathbf{0}, \quad \boldsymbol{\ell} \leq \mathbf{x}_k + \Delta \mathbf{x}_k \leq \mathbf{u}, \quad \boldsymbol{\mu}_{1,k}^{\text{sub}}, \boldsymbol{\mu}_{2,k}^{\text{sub}} \geq \mathbf{0}, \\ \boldsymbol{\mu}_{1,k}^{\text{sub}\top} (\boldsymbol{\ell} - \mathbf{x}_k - \Delta \mathbf{x}_k) &= 0, \quad \boldsymbol{\mu}_{2,k}^{\text{sub}\top} (\mathbf{x}_k + \Delta \mathbf{x}_k - \mathbf{u}) = 0. \end{aligned} \tag{E.1}$$

Note that the sequences $(\boldsymbol{\lambda}_k^{\text{sub}}, \boldsymbol{\mu}_k^{\text{sub}})/\|(\boldsymbol{\lambda}_k^{\text{sub}}, \boldsymbol{\mu}_k^{\text{sub}})\|$ and $\Delta \mathbf{x}_k$ are bounded. Without loss of gen-

erality, we assume $(\boldsymbol{\lambda}_k^{\text{sub}}, \boldsymbol{\mu}_k^{\text{sub}})/\|(\boldsymbol{\lambda}_k^{\text{sub}}, \boldsymbol{\mu}_k^{\text{sub}})\| \rightarrow (\bar{\boldsymbol{\lambda}}, \bar{\boldsymbol{\mu}})$ and $\Delta \mathbf{x}_k \rightarrow \Delta \bar{\mathbf{x}}$. Dividing the two sides of (E.1) by $\|(\boldsymbol{\lambda}_k^{\text{sub}}, \boldsymbol{\mu}_k^{\text{sub}})\|$, taking the limit $k \rightarrow \infty$, and applying the feasibility of $\bar{\mathbf{x}}$, we have

$$\begin{aligned} \nabla \mathbf{c}(\bar{\mathbf{x}})^\top \bar{\boldsymbol{\lambda}} - \bar{\boldsymbol{\mu}}_1 + \bar{\boldsymbol{\mu}}_2 &= \mathbf{0}, & \nabla \mathbf{c}(\bar{\mathbf{x}}) \Delta \bar{\mathbf{x}} &= \mathbf{0}, \\ \bar{\boldsymbol{\mu}}_1^\top (\boldsymbol{\ell} - \bar{\mathbf{x}}) &= \bar{\boldsymbol{\mu}}_1^\top \Delta \bar{\mathbf{x}}, & \bar{\boldsymbol{\mu}}_2^\top (\bar{\mathbf{x}} - \mathbf{u}) &= -\bar{\boldsymbol{\mu}}_2^\top \Delta \bar{\mathbf{x}}. \end{aligned} \quad (\text{E.2})$$

The above four equalities imply

$$\bar{\boldsymbol{\mu}}_1^\top (\boldsymbol{\ell} - \bar{\mathbf{x}}) + \bar{\boldsymbol{\mu}}_2^\top (\bar{\mathbf{x}} - \mathbf{u}) = 0. \quad (\text{E.3})$$

Since $\bar{\boldsymbol{\mu}}_1, \bar{\boldsymbol{\mu}}_2 \geq \mathbf{0}$, we deduce from (E.3) that $[\bar{\boldsymbol{\mu}}_1]_i > 0$ only if $[\bar{\mathbf{x}}]_i = [\boldsymbol{\ell}]_i$ and $[\bar{\boldsymbol{\mu}}_2]_i > 0$ only if $[\bar{\mathbf{x}}]_i = [\mathbf{u}]_i$. The EGMFCQ condition at $\bar{\mathbf{x}}$ (cf. Definition 2.1) implies that there exists $\mathbf{z} \in \mathbb{R}^d$ such that $\mathbf{c}(\bar{\mathbf{x}}) + \nabla \mathbf{c}(\bar{\mathbf{x}}) \mathbf{z} = \mathbf{0}$, $[\mathbf{z}]_i > 0$ if $[\bar{\mathbf{x}}]_i = [\boldsymbol{\ell}]_i$, and $[\mathbf{z}]_i < 0$ if $[\bar{\mathbf{x}}]_i = [\mathbf{u}]_i$. Then, we have $-\mathbf{z}^\top \bar{\boldsymbol{\mu}}_1 + \mathbf{z}^\top \bar{\boldsymbol{\mu}}_2 < 0$ if $\bar{\mathbf{x}}$ is on the boundary of the box constraints. Multiplying $-\mathbf{z}$ on both sides of the first equality in (E.2), we have

$$0 = -\mathbf{z}^\top (\nabla \mathbf{c}(\bar{\mathbf{x}})^\top \bar{\boldsymbol{\lambda}} - \bar{\boldsymbol{\mu}}_1 + \bar{\boldsymbol{\mu}}_2) = \mathbf{c}(\bar{\mathbf{x}})^\top \bar{\boldsymbol{\lambda}} + \mathbf{z}^\top \bar{\boldsymbol{\mu}}_1 - \mathbf{z}^\top \bar{\boldsymbol{\mu}}_2 = \mathbf{z}^\top \bar{\boldsymbol{\mu}}_1 - \mathbf{z}^\top \bar{\boldsymbol{\mu}}_2,$$

which contradicts $-\mathbf{z}^\top \bar{\boldsymbol{\mu}}_1 + \mathbf{z}^\top \bar{\boldsymbol{\mu}}_2 < 0$. On the other hand, if $\bar{\mathbf{x}}$ is in the interior of the box constraints, then $\bar{\boldsymbol{\mu}}_1 = \bar{\boldsymbol{\mu}}_2 = \mathbf{0}$. Together with the first equality of (E.2), the linear independence of the rows of $\nabla \mathbf{c}(\bar{\mathbf{x}})$ shows $\bar{\boldsymbol{\lambda}} = \mathbf{0}$, which contradicts to the fact that $\|(\bar{\boldsymbol{\lambda}}, \bar{\boldsymbol{\mu}})\| = 1$. This completes the proof. ■