

Lower Complexity Bounds for Nonconvex-Strongly-Convex Bilevel Optimization with First-Order Oracles

Kaiyi Ji

Department of Computer Science and Engineering

University at Buffalo

kaiyiji@buffalo.edu

November 27, 2025

Abstract

Although upper bound guarantees for bilevel optimization have been widely studied, progress on lower bounds has been limited due to the complexity of the bilevel structure. In this work, we focus on the smooth nonconvex-strongly-convex setting and develop new hard instances that yield nontrivial lower bounds under deterministic and stochastic first-order oracle models. In the deterministic case, we prove that any first-order zero-respecting algorithm requires at least $\Omega(\kappa^{3/2}\epsilon^{-2})$ oracle calls to find an ϵ -accurate stationary point, improving the optimal lower bounds known for single-level nonconvex optimization and for nonconvex-strongly-convex min-max problems. In the stochastic case, we show that at least $\Omega(\kappa^{5/2}\epsilon^{-4})$ stochastic oracle calls are necessary, again strengthening the best known bounds in related settings. Our results expose substantial gaps between current upper and lower bounds for bilevel optimization and suggest that even simplified regimes, such as those with quadratic lower-level objectives, warrant further investigation toward understanding the optimal complexity of bilevel optimization under standard first-order oracles.

1 Introduction

In this paper, we are interested in solving the following bilevel optimization problem:

$$\begin{aligned} \min_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}) &:= f(\mathbf{x}; \mathbf{y}^*(\mathbf{x})) \\ \text{s.t. } \mathbf{y}^*(\mathbf{x}) &= \arg \min_{\mathbf{y} \in \mathcal{Y}} g(\mathbf{x}; \mathbf{y}), \end{aligned} \tag{1}$$

where $\mathcal{X} \subset \mathbb{R}^m$ and $\mathcal{Y} \subset \mathbb{R}^n$ are nonempty closed convex sets. In this paper, we study the smooth nonconvex-strongly-convex bilevel optimization setting, where the lower-level function g is smooth and strongly convex in \mathbf{y} , while the upper-level function f is smooth and potentially nonconvex. This formulation captures a variety of modern applications, including meta-learning (Rajeswaran et al., 2019), reinforcement learning (Konda & Tsitsiklis, 2000; Hong et al., 2023), robotics (Wang et al., 2024), as well as communication networks and federated learning (Ji & Ying, 2022; Tarzanagh et al., 2022; Huang et al., 2023).

Recent years have witnessed substantial progress in understanding the convergence and complexity of bilevel optimization. A broad class of works (Ji et al., 2021; Hong et al., 2023; Chen et al., 2022; Dagr  ou et al., 2022) analyzes nonconvex-strongly-convex bilevel problems under access to second-order information such as Hessian- and Jacobian-vector products. More recently, there has been growing interest in developing and analyzing *fully first-order* bilevel algorithms that avoid any second-order computations (Shen & Chen, 2023; Chen et al., 2025; Lu & Mei, 2024; Kwon et al., 2023; Liu et al., 2020).

Although upper-bound analyses for bilevel optimization have been extensively studied, progress on establishing tight *lower bounds* has been significantly slower, largely due to the complexity of the general bilevel formulation. Deriving meaningful lower bounds that reflect the dependence on condition numbers and the target accuracy ϵ requires carefully constructed hard instances; otherwise, one risks obtaining vacuous bounds no stronger than the classical single-level lower bounds. Ji & Liang (2022) establish lower bounds for strongly-convex-strongly-convex and convex-strongly-convex bilevel problems under second-order oracle access, where the hyper-objective $H(\mathbf{x})$ is assumed to be convex or strongly convex. Their results show a gap of a factor $\sqrt{\kappa}$ compared with those for min-max optimization with analogous assumptions, where κ denotes the condition number of the lower-level function. However, their analysis is restricted to the deterministic setting, and the convexity assumptions on the hyper-objective may be restrictive for general bilevel problems. More recently, Kwon et al. (2024) provide lower bounds for nonconvex-strongly-convex bilevel optimization under a so-called \mathbf{y}^* -aware stochastic first-order oracle, where the oracle returns an estimate $\hat{\mathbf{y}}$ that is ϵ -close to the exact lower level solution \mathbf{y}^* , reducing the analysis to one that resembles single-level optimization. Yet, lower bounds for standard (stochastic) first-order oracles applied directly to the upper- and lower-level functions f and g remain open. In this paper, we take a further step toward reducing this gap by developing nontrivial lower bounds for smooth nonconvex-strongly-convex bilevel optimization under standard first-order oracle models. Our main contributions are summarized below.

- **Deterministic setting.** We construct a hard instance on which no first-order zero-respecting algorithm can find an ϵ -stationary solution using fewer than $\Omega(\kappa^{3/2}\epsilon^{-2})$ first-order oracle calls for smooth nonconvex-strongly-convex bilevel problems. In comparison, the optimal lower bounds for related settings are $\Omega(\epsilon^{-2})$ for general smooth nonconvex single-level optimization (Carmon et al., 2020) and $\Omega(\sqrt{\kappa}\epsilon^{-2})$ for smooth nonconvex-strongly-convex min-max optimization (Li et al., 2021). Our result improves these bounds by factors of $\kappa^{3/2}$ and κ , respectively.

On the upper-bound side, Chen et al. (2025) propose a first-order penalty method achieving a convergence rate of order $\kappa^4\epsilon^{-2}$, which can be reduced to $\kappa^{3.5}\epsilon^{-2}$ through a naive application of Nesterov acceleration. However, even when compared with our lower bound, there remains a gap of order κ^2 , indicating substantial room for future improvements.

- **Stochastic setting.** We further construct an instance showing that no first-order zero-respecting algorithm can achieve an ϵ -stationary solution with fewer than $\Omega(\kappa^{5/2}\epsilon^{-4})$ stochastic oracle calls under bounded variance assumptions. For comparison, the lower bound for standard smooth nonconvex single-level stochastic optimization is $\Omega(\epsilon^{-4})$ (Arjevani et al., 2023), and for smooth nonconvex-strongly-convex min-max optimization it is $\Omega(\kappa^{1/3}\epsilon^{-2})$ (Li et al., 2021). Our result improves upon these by factors of $\kappa^{5/2}$ and $\kappa^{13/6}$, respectively. Compared with the $\Omega(\epsilon^{-6})$ upper bound established by Kwon et al. (2024), a notable gap still remains.
- **Implications.** Our constructions demonstrate that nontrivial lower bounds for nonconvex-strongly-convex bilevel optimization are indeed possible and are significantly stronger than the known results

for single-level and min–max problems. Nevertheless, substantial gaps persist between current upper and lower bounds, even in this restricted setting. Motivated by our findings, we suggest that closing these gaps may require first studying the simpler yet meaningful case in which the lower-level function is **quadratic**. Our lower bounds continue to apply in that regime, but obtaining tighter upper bounds in this setting remains largely unexplored and not yet well understood. We hope that the results presented in this paper offer valuable insights for future progress in this direction.

2 Related Works

Bilevel optimization algorithms. Bilevel optimization has a long history dating back to the seminal work of Bracken & McGill (1973). Early studies (Hansen et al., 1992; Shi et al., 2005) approached bilevel programs from a constrained optimization perspective, motivating the development of KKT-based reformulations and related techniques. More recently, gradient-based bilevel optimization has attracted significant attention due to its efficiency and scalability in modern machine learning applications. A major class of gradient-based approaches is the family of Approximate Implicit Differentiation (AID) methods (Domke, 2012; Liao et al., 2018; Ji et al., 2021; Dagr  ou et al., 2022; Yang et al., 2024), which compute the hypergradient via implicit differentiation and approximate the resulting linear system using iterative solvers. In contrast, Iterative Differentiation (ITD) methods (Maclaurin et al., 2015; Franceschi et al., 2017) estimate hypergradients by unrolling the lower-level optimization and applying automatic differentiation in either forward or reverse mode. Building upon these ideas, a number of stochastic bilevel algorithms have been developed using Neumann-series approximation (Chen et al., 2022; Ji et al., 2021), recursive momentum techniques (Yang et al., 2021; Guo & Yang, 2021), and variance-reduction mechanisms (Yang et al., 2021). All such methods rely on second-order information, commonly in the form of Hessian–vector or Jacobian–vector products. A comprehensive overview is provided in the survey (Liu et al., 2021a).

Recently, growing interest has shifted slightly toward designing *first-order* bilevel optimization methods that use only (stochastic) first-order oracles, thereby avoiding explicit second-order computations. Representative examples include penalty-based methods (Shen & Chen, 2023; Lu & Mei, 2024; Kwon et al., 2023; Jiang et al., 2025; Chen et al., 2025), primal–dual frameworks (Sow et al., 2022), finite-difference Hessian–vector approximation techniques (Yang et al., 2023), value-function-based approaches (Liu et al., 2020, 2021c,b), barrier-based formulations (Liu et al., 2022), and min–max optimization based methods (Lu & Mei, 2025; Wang et al., 2023). These works collectively highlight the potential of first-order bilevel algorithms to achieve competitive performance while significantly reducing computational overhead.

Upper bound analysis. A large body of work, including Ji et al. (2021); Hong et al. (2023); Chen et al. (2022), studies AID- and ITD-type algorithms for nonconvex–strongly-convex bilevel optimization. Another line of research considers cases where the lower-level objective is not strongly convex; for example, Arbel & Mairal (2022); Liu et al. (2021c) analyze settings in which the lower-level solution is characterized through a selection map (e.g., the output of a particular algorithm). For bilevel algorithms that rely solely on (stochastic) first-order oracles, Kwon et al. (2023); Chen et al. (2025) establish convergence guarantees for nonconvex–strongly-convex formulations. In addition, Shen & Chen (2023); Chen et al. (2024) study algorithms under weaker structural assumptions on the lower-level problem, extending beyond strong convexity.

Lower bound analysis. Foundational lower bounds for first-order optimization were established by Nemirovski and Nesterov and are presented in their textbooks (Nemirovsky, 1992; Nesterov et al., 2018).

A central concept in this theory is the notion of *zero-chains*, which ensure that any zero-respecting first-order method can activate coordinates only sequentially. Recent works have significantly advanced these constructions in the context of smooth nonconvex optimization (Fang et al., 2018; Carmon et al., 2020, 2021; Arjevani et al., 2023). Building upon these developments, Li et al. (2021) establish lower bounds for nonconvex-strongly-convex min-max optimization. Our work builds based on these results.

Lower bounds for bilevel optimization are relatively underexplored. Ji & Liang (2022) derive bounds for convex and strongly-convex bilevel problems using second-order oracles. More recently, Kwon et al. (2024) establish lower bounds for nonconvex-strongly-convex bilevel problems under a \mathbf{y}^* -aware stochastic oracle. In contrast, we provide lower bounds for nonconvex-strongly-convex bilevel optimization using only standard (stochastic) first-order oracles.

3 Preliminaries

Notations. We use bold lower-case letters to denote vectors and regular lower-case letters to denote scalars. For a vector $\mathbf{x} \in \mathbb{R}^d$, we use \mathbf{x}^t to denote its value at the t^{th} iteration, and x_i to denote its i th coordinate and define its support as $\text{supp}(\mathbf{x}) := \{i \mid x_i \neq 0\}$. We use $\|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^d x_i^2}$ and $\|\mathbf{x}\|_\infty = \max_{1 \leq i \leq d} |x_i|$ to denote the ℓ_2 and ℓ_∞ norms, respectively. For a matrix $M \in \mathbb{R}^{m \times n}$, we use $M_{i,j}$ to denote its (i, j) th entry. We use $\|M\|_\infty = \max_{1 \leq i \leq m} \sum_{j=1}^n |M_{i,j}|$ for the matrix infinity norm and $\|M\|_2$ for its spectral norm. For a square matrix M , we let $\text{diag}_m(M)$ denote the block diagonal matrix with m identical copies of M on the diagonal. We use standard asymptotic notation $\mathcal{O}(\cdot)$, $\Omega(\cdot)$, and $\Theta(\cdot)$.

3.1 Function Class

In this paper, we focus on the class of smooth nonconvex-strongly-convex bilevel problems that satisfy the standard assumptions used in first order bilevel optimization.

Definition 1. Given $L_f, L_g \geq \mu > 0$, $C \geq 0$ and $\Delta > 0$, define $\mathcal{F}(L_f, L_g, \mu, \Delta)$ to be the set of function pairs $\{f, g\}$ such that $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ and $g : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ for some nonempty closed convex sets $\mathcal{X} \subset \mathbb{R}^m$ and $\mathcal{Y} \subset \mathbb{R}^n$ for all $m, n \in \mathbb{N}$, which satisfy the following assumptions:

1. Functions f, g are continuously differentiable and L_f and L_g -smooth respectively, jointly in (\mathbf{x}, \mathbf{y}) over $\mathcal{X} \times \mathcal{Y}$.
2. For every $(\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{Y}$, there exists a numerical constant $C \geq 0$ such that $\|\nabla_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})\|_2 \leq C$.
3. For every $\mathbf{x} \in \mathcal{X}$, $g(\mathbf{x}, \cdot)$ is μ -strongly-convex in \mathbf{y} , that is, for any $\mathbf{y}_1, \mathbf{y}_2 \in \mathcal{Y}$,

$$g(\mathbf{x}; \mathbf{y}_1) \geq g(\mathbf{x}; \mathbf{y}_2) + \langle \nabla_{\mathbf{y}} g(\mathbf{x}; \mathbf{y}_2), \mathbf{y}_1 - \mathbf{y}_2 \rangle + \frac{\mu}{2} \|\mathbf{y}_1 - \mathbf{y}_2\|_2^2. \quad (2)$$

4. There exists a numerical constant $\rho \geq 0$ such that the second-order derivatives $\nabla_{\mathbf{x}, \mathbf{y}}^2 g$ and $\nabla_{\mathbf{y}, \mathbf{y}}^2 g$ are well-defined and ρ -Lipschitz jointly in (\mathbf{x}, \mathbf{y}) for all $(\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{Y}$.
5. $H(\mathbf{0}) - \min_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}) \leq \Delta$, where $H(\mathbf{x}) := f(\mathbf{x}; \mathbf{y}^*(\mathbf{x}))$ is the hyper-objective function.

The constants L_g , L_f , C , and ρ are all independent of the strong convexity parameter μ and the target accuracy ϵ . Note that for items 2 and 4, we only require the existence of numerical constants $C, \rho = \mathcal{O}(1)$.

3.2 Algorithm Class

We focus on algorithms that solve bilevel optimization problems using (stochastic) first order oracles. For clarity of presentation, we first define the (stochastic) first-order oracles considered in this work.

Definition 2 (Deterministic first-order oracle). *The deterministic first-order oracle of a differentiable function $f : \mathcal{X} \rightarrow \mathbb{R}$ is a mapping $O : \mathbf{x} \mapsto (f(\mathbf{x}), \nabla f(\mathbf{x}))$ for $\mathbf{x} \in \mathcal{X}$.*

Definition 3 (Stochastic first-order oracle). *The stochastic first-order oracle of a differentiable function $f : \mathcal{X} \rightarrow \mathbb{R}$ is a mapping $O : \mathbf{x} \mapsto (f(\mathbf{x}), G_f(\mathbf{x}; \xi))$ for $\mathbf{x} \in \mathcal{X}$, where ξ is a random variable satisfying $\mathbb{E}_\xi [G_f(\mathbf{x}; \xi)] = \nabla f(\mathbf{x})$ and $\mathbb{E}_\xi \|G_f(\mathbf{x}; \xi) - \nabla f(\mathbf{x})\|_2^2 \leq \sigma_f^2$.*

Note that the algorithms rely on first-order oracles for both the upper- and lower-level objectives f and g . In the stochastic setting, we assume for simplicity that the variances of the stochastic first-order oracles are identical, i.e., $\sigma_f = \sigma_g = \sigma$. We further focus on first-order bilevel algorithms that satisfy the following zero-respecting property:

Definition 4 (First-order bilevel algorithm class.). *For upper- and lower-level objective functions $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ and $g : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ and their associated first-order oracles $O_f : (\mathbf{x}, \mathbf{y}) \mapsto (f(\mathbf{x}; \mathbf{y}), \nabla f(\mathbf{x}; \mathbf{y}))$ and $O_g : (\mathbf{x}, \mathbf{y}) \mapsto (f(\mathbf{x}; \mathbf{y}), \nabla g(\mathbf{x}; \mathbf{y}))$, the $(t + 1)$ -th iterate $(\mathbf{x}^{t+1}, \mathbf{y}^{t+1})$ satisfies:*

$$\begin{aligned} \mathbf{x}^{t+1} &\in \left\{ \mathcal{P}_{\mathcal{X}}(\mathbf{u}) : \text{supp}(\mathbf{u}) \subset \bigcup_{0 \leq i \leq t} (\text{supp}(\mathbf{x}^i) \cup \text{supp}(\nabla_{\mathbf{x}} f(\mathbf{x}^i; \mathbf{y}^i)) \cup \text{supp}(\nabla_{\mathbf{x}} g(\mathbf{x}^i; \mathbf{y}^i)) \right\}; \\ \mathbf{y}^{t+1} &\in \left\{ \mathcal{P}_{\mathcal{Y}}(\mathbf{v}) : \text{supp}(\mathbf{v}) \subset \bigcup_{0 \leq i \leq t} (\text{supp}(\mathbf{y}^i) \cup \text{supp}(\nabla_{\mathbf{y}} f(\mathbf{x}^i; \mathbf{y}^i)) \cup \text{supp}(\nabla_{\mathbf{y}} g(\mathbf{x}^i; \mathbf{y}^i)) \right\}. \end{aligned} \quad (3)$$

A similar definition applies in the stochastic setting, where the gradients ∇f and ∇g are replaced by their corresponding stochastic first-order oracles.

Note that the subspaces defined in Equation (3) permit both simultaneous and alternating updates of \mathbf{x} and \mathbf{y} , thereby including single-loop and double-loop bilevel optimization algorithms. Consequently, the algorithm class introduced in Definition 4 covers all existing first-order bilevel optimization methods, including but not limited to penalty-based approaches (Shen & Chen, 2023; Lu & Mei, 2024), primal–dual methods (Sow et al., 2022), finite-difference Hessian–vector–approximation methods (Yang et al., 2023), value-function-based approaches (Liu et al., 2020, 2021c,b), and barrier-based methods (Liu et al., 2022).

4 Lower Bounds in Deterministic Setting

4.1 Useful Techniques for Lower-Bound Construction

In this paper, we focus on the bilevel optimization setting where the lower-level function $g(\mathbf{x}; \mathbf{y})$ is strongly convex in \mathbf{y} , while the upper-level function $f(\mathbf{x}; \mathbf{y})$ is smooth and possibly nonconvex. For this reason, our constructions draw on key techniques and components from the worst-case instances of Nesterov et al. (2018) for smooth strongly convex functions and Carmon et al. (2020) for smooth nonconvex functions. Their core idea is to make sure their instances satisfy the following notion of zero-chain property:

Definition 5 (Zero-chain). A function $f : \mathcal{X} \subset \mathbb{R}^d \rightarrow \mathbb{R}$ is a (first-order) zero-chain if for every $1 \leq i \leq d$,

$$\text{supp}(\mathbf{x}) := \{i : x_i \neq 0\} \subset \{1, \dots, i-1\} \implies \text{supp}(\nabla f(\mathbf{x})) \subset \{1, \dots, i\}.$$

Consider running a first-order algorithm on a zero-chain function, starting from the initialization $\mathbf{x} = 0$, and assume access to a deterministic first-order oracle. By the zero-chain property, each iteration can introduce at most one new nonzero coordinate of \mathbf{x} —that is, each iteration “activates” at most one additional coordinate. Consequently, after t iterations we must have $\text{supp}(\mathbf{x}^t) \subset \{1, \dots, t\}$. Therefore, if a good solution requires that at least T coordinates be discovered, then any deterministic first-order method must take at least T iterations, which yields a lower bound of order T on the algorithm’s complexity.

Following this strategy, [Nesterov et al. \(2018\)](#) and [Carmon et al. \(2020\)](#) provide the following key components for their constructions in strongly-convex and nonconvex settings, respectively:

- **Tri-diagonal matrix A .** Following [Nesterov et al. \(2018\)](#); [Li et al. \(2021\)](#), we use the following tri-diagonal 1-D discrete Laplacian matrix $A \in \mathbb{R}^{n \times n}$ to construct the strongly-convex lower-level instance:

$$A := \begin{bmatrix} 1 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & 1 \end{bmatrix},$$

where it is verified that A is positive semidefinite and $\|A\|_2 \leq 4$. Due to its tri-diagonal nature, it is easily verified that if $\text{supp}(x) \subset \{1, \dots, i-1\}$, then $Ax \subset \{1, \dots, i\}$. In other words, if a vector has nonzero entries only at its first $i-1$ coordinates, then multiplying it by A can activate at most one additional coordinate, namely the i -th one.

- **$\Psi(\cdot)$ and $\Phi(\cdot)$ hardness functions.** Following the construction in [Carmon et al. \(2020\)](#), we employ the component functions $\Psi(x) : \mathbb{R} \rightarrow \mathbb{R}$ and $\Phi(x) : \mathbb{R} \rightarrow \mathbb{R}$ defined below.

$$\Psi(x) := \begin{cases} 0, & x \leq \frac{1}{2}, \\ \exp\left(1 - \frac{1}{(2x-1)^2}\right), & x > \frac{1}{2}, \end{cases} \quad \text{and} \quad \Phi(x) := \sqrt{e} \int_{-\infty}^x e^{-\frac{1}{2}t^2} dt, \quad (4)$$

which possess the following key properties that will be used in our analysis.

Lemma 1 ([Carmon et al. \(2020\)](#), Lemma 1). *The functions Φ and Ψ satisfy*

1. For all $x \leq \frac{1}{2}$ and $k \in \mathbb{N}$, we have $\Psi^{(k)}(x) = 0$, where $\Psi^{(k)}$ denotes the k^{th} -order derivative.
2. For all $x \geq 1$ and $|y| < 1$, we have $\Psi(x) \Phi'(y) > 1$.
3. Both Ψ and Φ are infinitely differentiable. For all $k \in \mathbb{N}$, we have

$$\sup_x |\Psi^{(k)}(x)| \leq \exp\left(\frac{5k}{2} \log(4k)\right) \quad \text{and} \quad \sup_x |\Phi^{(k)}(x)| \leq \exp\left(\frac{3k}{2} \log \frac{3k}{2}\right).$$

4. The functions and derivatives Ψ, Ψ', Φ, Φ' are nonnegative and bounded, with

$$0 < \Psi < e, \quad 0 < \Psi' < \sqrt{\frac{54}{e}}, \quad 0 < \Phi < \sqrt{2\pi e}, \quad 0 < \Phi' < \sqrt{e}.$$

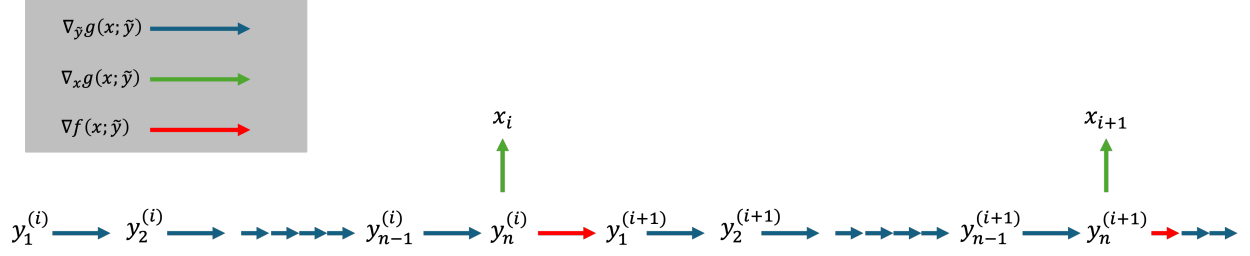


Figure 1: An illustration of the zero-chain for our constructed instance in eq. (6) for nonconvex-strongly-convex bilevel optimization.

Carmon et al. (2020) use a construction of $f(\mathbf{x}) = \sum_i [\Psi(-x_{i-1})\Phi(-x_i) - \Psi(x_{i-1})\Phi(x_i)]$, which together with $\Psi'(0) = \Psi(0) = 0$, ensures the zero-chain property that if $\mathbf{x} \subset \{1, \dots, i-1\}$, then $\nabla f(\mathbf{x}) \subset \{1, \dots, i\}$. Furthermore, as we will show later, the boundedness of Ψ , Ψ' , Φ , and Φ' is crucial for constructing a valid worst-case instance within the bilevel class $\mathcal{F}(L_f, L_g, \mu, C, \Delta)$.

4.2 Main Result: A Lower Bound on First-Order Oracle Complexity

The following theorem establishes a complexity lower bound for deterministic first-order bilevel algorithms.

Theorem 1. *For any $L_f, L_g, \mu, \Delta, \epsilon > 0$ satisfying $\kappa = L_g/\mu \geq 1$ and $\frac{\Delta}{L_f} = \mathcal{O}(1)$, there exist functions $f : \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}$ and $g : \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}$ such that $\{f, g\} \in \mathcal{F}(L_f, L_g, \mu, \Delta)$ for some $m, n \in \mathbb{N}$ with their deterministic first-order oracles. For any first-order bilevel algorithm of the form in Definition 4, in order to find an ϵ -accurate stationary point \mathbf{x} such that $\|\nabla H(\mathbf{x})\|_2 < \epsilon$, the algorithm must use at least*

$$\frac{C_0 \Delta L_f \kappa^{3/2}}{\epsilon^2} \quad (5)$$

oracle calls, where $H(\mathbf{x}) = f(\mathbf{x}; \mathbf{y}^*(\mathbf{x}))$ with $\mathbf{y}^*(\mathbf{x}) = \arg \min_{\mathbf{y}} g(\mathbf{x}; \mathbf{y})$ is the hyper-objective, and C_0 is a numerical constant.

Carmon et al. (2020) establishes a lower bound of $\Omega(1/\epsilon^2)$ for smooth nonconvex optimization, and Li et al. (2021) proves a lower bound of $\Omega(\sqrt{\kappa}/\epsilon^2)$ for smooth nonconvex-strongly-concave min-max optimization. Both results can be viewed as special cases of smooth nonconvex-strongly-convex bilevel optimization, for which we obtain in Theorem 1 a much larger lower bound of $\Omega(\kappa^{3/2}/\epsilon^2)$. This demonstrates that bilevel optimization is provably more challenging than min-max optimization. This observation is consistent with the fundamental hardness comparison for smooth strongly-convex-strongly-convex bilevel problems established in Ji & Liang (2022).

4.3 Analysis and Proof Outline for Deterministic Lower Bound

We consider the following worst-case instance. For notational simplicity, define $x_0 \equiv \frac{\lambda}{C_l M_{n,n}}$.

$$f(\mathbf{x}; \tilde{\mathbf{y}}) = \sum_{i=1}^T \frac{\lambda^2 L_f}{L} \left[\Psi\left(-\frac{C_l}{\lambda} y_n^{(i-1)}\right) \Phi\left(-\frac{C_r}{\lambda} y_1^{(i)}\right) - \Psi\left(\frac{C_l}{\lambda} y_n^{(i-1)}\right) \Phi\left(\frac{C_r}{\lambda} y_1^{(i)}\right) \right]$$

$$g(\mathbf{x}; \tilde{\mathbf{y}}) = \sum_{i=0}^T \left[\frac{L_g n^2}{2(4n^2 + 1)} (\mathbf{y}^{(i)})^\top \left(\frac{1}{n^2} I_n + A \right) \mathbf{y}^{(i)} - L_g (\mathbf{b}_x^{(i)})^\top \mathbf{y}^{(i)} \right], \quad (6)$$

where $\mathbf{x} = [x_1, \dots, x_T] \in \mathbb{R}^T$ is the upper-level variable, $\tilde{\mathbf{y}} = [\mathbf{y}^{(0)}, \mathbf{y}^{(1)}, \dots, \mathbf{y}^{(T)}]$ with each $\mathbf{y}^{(i)} \in \mathbb{R}^n$ is the lower-level variable, $y_j^{(i)}$ returns the j^{th} coordinate of $\mathbf{y}^{(i)}$, and the dimension $n = \lfloor \sqrt{\frac{L_g - \mu}{4\mu}} \rfloor$, and the design of $\mathbf{b}_x^{(i)}$ is most critical, which is given by

$$\mathbf{b}_x^{(i)} = [0, 0, \dots, x_i] = x_i \mathbf{e}_n, \quad (7)$$

where \mathbf{e}_i denotes the i^{th} standard basis vector, whose sole nonzero entry equals 1. For simple presentation, the numerical constants C_l, C_r, L and the parameter λ will be specified at a later stage.

Validation of our constructed instance. We first verify that our constructed instance belongs to the function class $\mathcal{F}(L_f, L_g, \mu, \Delta)$.

1. First, we need to verify $g(\mathbf{x}; \cdot)$ is μ -strongly convex. Since the matrix A is positive semidefinite, it can be verified that $\nabla^2 g(\mathbf{x}; \cdot) = \frac{L_g n^2}{4n^2 + 1} \text{diag}_{T+1}(A + \frac{1}{n^2})$. Let $M := \text{diag}_{T+1}\{A\}$. For any vector $\mathbf{z} \in \mathbb{R}^{n(T+1)}$, write it as a block vector $\mathbf{z} = [\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_m]$, $\mathbf{z}_i \in \mathbb{R}^n$, we have $\mathbf{z}^\top M \mathbf{z} = \sum_{i=1}^m \mathbf{z}_i^\top A \mathbf{z}_i$. Since A is positive semidefinite, each term $\mathbf{z}_i^\top A \mathbf{z}_i \geq 0$, so the sum is nonnegative. Hence $\mathbf{z}^\top M \mathbf{z} \geq 0$ for all \mathbf{z} , and therefore M is positive semidefinite. This further implies that $\|\nabla^2 g(\mathbf{x}; \cdot)\|_2 \geq \frac{L_g}{4n^2 + 1}$. Given that $n = \lfloor \sqrt{\frac{L_g - \mu}{4\mu}} \rfloor \leq \sqrt{\frac{L_g - \mu}{4\mu}}$, we have $\frac{L_g}{4n^2 + 1} \geq \mu$. This validates that $g(\mathbf{x}; \tilde{\mathbf{y}})$ is μ -strongly convex in $\tilde{\mathbf{y}}$.

2. Next, we validate the smoothness of f and g functions:

- For the lower-level function $g(\mathbf{x}; \tilde{\mathbf{y}})$, it follows from eq. (6) that $\|\nabla_{\tilde{\mathbf{y}}}^2 g(\mathbf{x}; \tilde{\mathbf{y}})\|_2 \leq \frac{L_g n^2}{4n^2 + 1} (\frac{1}{n^2} + 4) = L_g$, $\|\nabla_{\mathbf{x}, \tilde{\mathbf{y}}}^2 g(\mathbf{x}; \tilde{\mathbf{y}})\|_2 = L_g$, $\nabla_{\mathbf{x}}^2 g(\mathbf{x}; \tilde{\mathbf{y}}) = 0$ for any $\mathbf{x}, \tilde{\mathbf{y}}$, and hence $g(\mathbf{x}; \tilde{\mathbf{y}})$ is L_g -smooth.
- For the upper-level function $f(\mathbf{x}; \tilde{\mathbf{y}})$, note that $\nabla_{\tilde{\mathbf{y}}}^2 f(\mathbf{x}; \tilde{\mathbf{y}}) = \frac{L_f}{L} M$, where $M \in \mathbb{R}^{n(T+1) \times n(T+1)}$ is a tri-diagonal matrix, where the absolute value of each nonzero element is bounded by some numerical constant, due to the fact that C_r and C_l are numerical constants, and that the functions Φ and Ψ , together with their derivatives, are bounded by numerical constants, as shown in item 3 of Lemma 1. Then, we have $\|M\|_2 \leq C_M$ for some numerical constant $C_M > 0$. Thus, choosing $L = C_M$ yields $\|\nabla_{\tilde{\mathbf{y}}}^2 f(\mathbf{x}; \tilde{\mathbf{y}})\|_2 \leq L_f$. Since $f(\mathbf{x}; \tilde{\mathbf{y}})$ depends only on $\tilde{\mathbf{y}}$, it is thus L_f -smooth.

3. Next, we need to show that the gradient norm $\|\nabla_{\tilde{\mathbf{y}}} f(\mathbf{x}; \tilde{\mathbf{y}})\|_2$ is bounded by a numerical constant that is independent of both T and n . This step is particularly challenging. For example, the previous lower bound in Ji & Liang (2022) circumvents this requirement by exploiting the strong convexity of the hyper-objective to guarantee gradient boundedness during the optimization process. However, that strategy applies only to the strongly-convex–strongly-convex setting and may not extend well to nonconvex or stochastic regimes. Moreover, another lower bound in Kwon et al. (2024) sets the upper-level function as a scalar y , which ensures that the gradient norm remains bounded by a constant. For our construction in eq. (6), it can be obtained that $\|\nabla_{\tilde{\mathbf{y}}} f(\mathbf{x}; \tilde{\mathbf{y}})\|_2 = \frac{\lambda L_f}{L} \mathbf{v}$, where $\mathbf{v} \in \mathbb{R}^{n(T+1)}$ has at most $2(T+1)$ nonzero entries at coordinates $kn+1$ for $k=0, \dots, T$ and jn for $j=1, \dots, T+1$. Moreover, the absolute value of each nonzero entry is bounded by a positive numerical constant,

owing to the fact that C_r and C_l are numerical constants and that Ψ , Ψ' , Φ , and Φ' are all bounded (Lemma 1, item 3). Therefore, we have $\|v\| \leq C_0\sqrt{T}$ for some numerical constant C_0 . Thus, we have $\|\nabla_{\tilde{\mathbf{y}}} f(\mathbf{x}; \tilde{\mathbf{y}})\|_2 \leq \frac{C_0 L_f}{L} \lambda \sqrt{T}$. As will be seen later, T is chosen such that $\lambda \sqrt{T} \leq \sqrt{\frac{\Delta L}{12 L_f}}$, which, together with $\frac{\Delta}{L_f} = \mathcal{O}(1)$, implies that $\|\nabla_{\tilde{\mathbf{y}}} f(\mathbf{x}; \tilde{\mathbf{y}})\|_2 = \mathcal{O}(1)$.

Zero-chain properties and iterate subspaces. We initialize \mathbf{x} and $\tilde{\mathbf{y}}$ to be $\mathbf{0}$. Then, based on the tri-diagonal structure of A and the properties of Ψ function in Lemma 1 (item 1), it can be quickly verified from our construction in eq. (6) that

- At the first iteration, $y_n^{(0)}$ becomes activated, because $x_0 \neq 0$ and $\partial g / \partial y_n^{(0)} = -L_g x_0$. Thus, at the second iteration, $y_1^{(1)}$ becomes activated due to the zero-chain property of the $f(\mathbf{x}; \tilde{\mathbf{y}})$ function.
- Suppose the iterates have begun updating $\mathbf{y}^{(i)}$ but have not yet reached $y_n^{(i)}$ (i.e., $y_n^{(i)} = 0$) for some $i \geq 1$. This implies that $y_n^{(j)} = 0$ for all $j \geq i$. Then, based on item 1 of Lemma 1, it can be verified that for all $j \geq i$,

$$\frac{\partial f(\mathbf{x}; \tilde{\mathbf{y}})}{\partial y_1^{(j+1)}} = -\frac{C_r \lambda L_f}{L} \Psi\left(-\frac{C_l}{\lambda} y_n^{(j)}\right) \Phi'\left(-\frac{C_r}{\lambda} y_1^{(j+1)}\right) - \frac{C_r \lambda L_f}{L} \Psi\left(\frac{C_l}{\lambda} y_n^{(j)}\right) \Phi'\left(\frac{C_r}{\lambda} y_1^{(j+1)}\right) = 0,$$

which, together with the structure of the lower-level function and the condition $x_j = 0$ for all $j \geq i$, implies that $\mathbf{y}^{(j)} = \mathbf{0}$ for all $j \geq i + 1$. This property is crucial because it preserves the zero-chain structure along the sequence $\{\mathbf{y}^{(i)}\}_{i=1}^T$ and ensures that advancing from one adjacent \mathbf{y} -iterate to the next necessarily requires at least n iterations.

- Suppose the iterates have begun updating $\mathbf{y}^{(i)}$ but have not yet reached $y_n^{(i)}$ (i.e., $y_n^{(i)} = 0$) for some $i \geq 1$. Then, for all $j \geq i$, the gradient of $g(\mathbf{x}; \tilde{\mathbf{y}})$ with respect to x_j is given by $-y_n^{(j)}$. As a consequence, the coordinate x_j will not be activated until $y_n^{(j)}$ becomes activated.

Based on the above analysis, it can be derived that at any iteration $Kn + k$ with $K = 0, \dots, T - 1$ and $k = 1, \dots, n$,

$$\begin{aligned} \text{supp}(\mathbf{y}^{(i)}) &\subseteq \{1, \dots, n\}, \quad i \leq K \text{ and } i \neq 0 \\ \text{supp}(\mathbf{y}^{(K+1)}) &\subset \{1, \dots, k\} \\ \text{supp}(\mathbf{y}^{(i)}) &= \emptyset, \quad i > K + 1 \\ \text{supp}(\mathbf{x}) &\subset \{0, \dots, K\}. \end{aligned} \tag{8}$$

Accordingly, to activate all coordinates of \mathbf{x} , one must perform at least Tn iterations in total.

The overall hyper-objective function and its key properties. First, it can be verified that the lower-level solutions are given by:

$$(\mathbf{y}^{(i)})^* = \underbrace{\frac{4n^2 + 1}{n^2} \left(\frac{1}{n^2} I_n + A \right)^{-1}}_M \mathbf{b}_x^{(i)}.$$

The hyper-objective function $H(x) := f(x; \tilde{\mathbf{y}}^*)$ is then given by

$$H(\mathbf{x}) = \sum_{i=1}^T \frac{\lambda^2 L_f}{L} \left[\Psi\left(-\frac{C_l}{\lambda} M_{n,n} x_{i-1}\right) \Phi\left(-\frac{C_r}{\lambda} M_{1,n} x_i\right) - \Psi\left(\frac{C_l}{\lambda} M_{n,n} x_{i-1}\right) \Phi\left(\frac{C_r}{\lambda} M_{1,n} x_i\right) \right].$$

Note that the above definition of $H(\mathbf{x})$ involves the quantities $M_{n,n}$ and $M_{1,n}$, whose behaviors are characterized in the following lemma.

Lemma 2. *Let $A \in \mathbb{R}^{n \times n}$ be the tri-diagonal matrix*

$$A = \begin{bmatrix} 1 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & 1 \end{bmatrix}, \quad S := \left(A + \frac{1}{n^2} I_n\right)^{-1}.$$

Then for every integer $n \geq 1$,

$$cn \leq S_{1,n}, S_{n,n} \leq Cn, \quad c := 1 - \frac{\pi^2}{12}, \quad C := 1 + \frac{\pi^2}{12}. \quad (9)$$

Based on Lemma 2 and $4 \leq \frac{4n^2+1}{n^2} \leq 5$, it can be derived that $4cn \leq M_{1,n}, M_{n,n} \leq 5Cn$, where c and C are given by eq. (9). Thus, choose numerical constants C_l and C_r such that

$$\frac{C_l M_{n,n}}{n} = \frac{C_r M_{1,n}}{n} = \tilde{C}, \quad (10)$$

where $\tilde{C} = O(1)$ is a numerical constant. Then, we use the following lemma to provide a lower bound on the gradient norm when the algorithm has not yet reached the end of the chain.

Lemma 3. *If $|x_i| < \frac{\lambda}{\tilde{C}n}$ for some $i \leq T$. Then, we have $\|\nabla H(\mathbf{x})\|_2 \geq \frac{\lambda L_f \tilde{C}n}{L}$.*

The following lemma provides the bound on the optimality gap of the hyper-objective function $H(\mathbf{x})$:

Lemma 4. *The hyper-objective function $H(\mathbf{x})$ satisfies $H(\mathbf{0}) - \inf_{\mathbf{x}} H(\mathbf{x}) \leq \frac{12\lambda^2 L_f T}{L}$.*

Based on all the above auxiliary lemmas, we begin to prove our main theorem.

Proof of Theorem 1. First note that if $x_T = 0$, based on Lemma 3, we have that

$$\|\nabla H(\mathbf{x})\|_2 \geq \frac{\lambda L_f \tilde{C}n}{L}.$$

Choosing $\lambda = \frac{\epsilon L}{L_f \tilde{C}n}$ guarantees $\|\nabla H(\mathbf{x})\|_2 \geq \epsilon$. Then, we need to verify that $H(\mathbf{0}) - \inf_{\mathbf{x}} H(\mathbf{x}) \leq \Delta$. Based on Lemma 4, we havethat

$$H(\mathbf{0}) - \inf_{\mathbf{x}} H(\mathbf{x}) \leq \frac{12\lambda^2 L_f T}{L}, \quad (11)$$

which, by setting $T = \left\lfloor \frac{\Delta L}{12\lambda^2 L_f} \right\rfloor$, guarantees that $H(\mathbf{0}) - \inf_{\mathbf{x}} H(\mathbf{x}) \leq \Delta$.

Based on the subspace analysis in eq. (8), we have that $x_T = 0$ if $t < Tn$, and hence $\|\nabla H(\mathbf{x}^t)\|_2 \geq \epsilon$. Recall that $n = \left\lfloor \sqrt{\frac{L_g - \mu}{4\mu}} \right\rfloor$. Thus, to achieve an ϵ -accurate stationary solution, there are at least

$$Tn = \frac{c_0 \Delta n^3}{\epsilon^2} = \frac{\Delta L n}{12 L_f} \frac{L_f^2 \tilde{C}^2 n^2}{\epsilon^2 L^2} = \frac{C_0 \Delta L_f \kappa^{\frac{3}{2}}}{\epsilon^2} \quad (12)$$

oracle calls, where c_0 is some numerical constant. Then, the proof is complete. \square

5 Lower Bounds in Stochastic Setting

In this section, we provide a lower bound for stochastic first-order oracles. We first introduce several important definitions and lemmas from [Arjevani et al. \(2023\)](#), which serve as the foundation for our constructions in the stochastic setting.

5.1 Auxiliary Definitions and Lemmas

Following [Arjevani et al. \(2023\)](#), to establish a lower bound in the stochastic setting, we adopt the notion of a probability- p zero-chain.

Definition 6 (Probability- p zero-chain). *A function $f : \mathcal{X} \rightarrow \mathbb{R}$ with a stochastic first-order oracle $O : \mathbf{x} \mapsto (f(\mathbf{x}), G_f(\mathbf{x}; \xi))$ is a probability- p zero-chain if*

$$\text{supp}(\mathbf{x}) \subset \{1, \dots, i-1\} \implies \begin{cases} \mathbb{P}(\text{supp}(G_f(\mathbf{x}; \xi)) \not\subset \{1, \dots, i-1\}) \leq p, \\ \mathbb{P}(\text{supp}(G_f(\mathbf{x}; \xi)) \subset \{1, \dots, i\}) = 1. \end{cases}$$

The above definition implies that at each iteration, a new coordinate i becomes activated (i.e., the iterate acquires a nonzero entry at coordinate i) with probability p . The following lemma (which is an adapted version from [Li et al. \(2021\)](#)) provides a recipe for constructing a probability- p zero-chain based on a given zero-chain.

Lemma 5 (([Arjevani et al., 2023](#), Lemma 3)). *Let $f : \mathcal{X} \rightarrow \mathbb{R}$ be a zero-chain on $\mathcal{X} \subset \mathbb{R}^T$. For $\mathbf{x} \in \mathcal{X}$, let $i^*(\mathbf{x}) := \inf\{i \in [T] : x_i = 0\}$ be the next coordinate to activate. For $p \in (0, 1]$, define the stochastic gradient estimator $G_f(\mathbf{x}; \xi)$ coordinate-wisely by*

$$[G_f(\mathbf{x}, \xi)]_i := \begin{cases} \frac{\xi}{p} \nabla_i f(\mathbf{x}), & \text{if } i = i^*(\mathbf{x}), \\ \nabla_i f(\mathbf{x}), & \text{otherwise,} \end{cases}$$

where $\xi \sim \text{Bernoulli}(p)$. Suppose there exists $G < \infty$ such that $\|\nabla f(\mathbf{x})\|_\infty \leq G$ for all $\mathbf{x} \in \mathcal{X}$. Then, the oracle $O : \mathbf{x} \mapsto (f(\mathbf{x}), G_f(\mathbf{x}, \xi))$ is a stochastic first-order oracle with bounded variance $\sigma^2 \leq G^2(1-p)/p$. Moreover, f with oracle O is a probability- p zero-chain.

Lemma 5 allows us to build a probability- p zero-chain based on the zero-chain we establish in eq. (6) and Figure 1. However, as also noted by [Li et al. \(2021\)](#) for min-max optimization problems, one main challenge

lies in the unboundedness of the iterates \mathbf{x} and $\tilde{\mathbf{y}}$, such that the gradient norm of the lower-level function $\|\nabla g(\mathbf{x}; \tilde{\mathbf{y}})\|_\infty$ is unbounded. To address this challenge, [Li et al. \(2021\)](#) modify the quadratic components in their deterministic worst-case instance and introduce two bounded hypercubes as the domains for \mathbf{x} and \mathbf{y} :

$$\mathcal{C}_{R_x}^m := \{\mathbf{x} \in \mathbb{R}^m : \|\mathbf{x}\|_\infty \leq R_x\}, \quad \mathcal{C}_{R_y}^n := \{\mathbf{y} \in \mathbb{R}^n : \|\mathbf{y}\|_\infty \leq R_y\}, \quad (13)$$

where R_x and R_y are chosen so that the variance of the stochastic oracle is bounded by G . Interestingly, unlike [Li et al. \(2021\)](#), which must revise the quadratic components in their deterministic construction, we find that our deterministic instance in eq. (6) can be used directly, provided that the domain radius R_x and R_y are properly selected, as will be seen in our analysis later.

5.2 Main Result: A Lower Bound on Stochastic First-Order Oracle Complexity

The following theorem establishes a complexity lower bound for stochastic first-order bilevel algorithms.

Theorem 2. *For any $L_f, L_g, \mu, \Delta, \epsilon > 0$ satisfying $\kappa = L_g/\mu \geq 1$ and $\frac{\Delta}{L_f} = \mathcal{O}(1)$, there exist functions $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ and $g : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ such that $\{f, g\} \in \mathcal{F}(L_f, L_g, \mu, \Delta)$ for some $\mathcal{X} \subset \mathbb{R}^m$ and $\mathcal{Y} \subset \mathbb{R}^n$, and stochastic first-order oracles O for both f and g such that for any first-order bilevel algorithm of the form in Definition 4, in order to find an ϵ -accurate stationary point \mathbf{x} such that*

$$\mathbb{E} [L_h \|\mathcal{P}_{\mathcal{X}}[\mathbf{x} - (1/L_h)\nabla H(\mathbf{x})] - \mathbf{x}\|_2] < \epsilon,$$

the algorithm must use at least

$$\Omega\left(\frac{L_f^3 \Delta \kappa^{5/2} \sigma^2}{L_g^2 \epsilon^4}\right)$$

stochastic oracle calls, where L_h is the smoothness parameter of the hyper-objective $H(\mathbf{x})$.

In the stochastic setting, [Arjevani et al. \(2023\)](#) establishes a lower bound of $\Omega(1/\epsilon^4)$ for smooth nonconvex optimization, and [Li et al. \(2021\)](#) proves a lower bound of $\Omega(\kappa^{1/3}/\epsilon^4)$ for smooth nonconvex-strongly-concave min-max optimization. For smooth nonconvex-strongly-convex bilevel optimization, we obtain in Theorem 2 a significantly larger lower bound of $\Omega(\kappa^{5/2}/\epsilon^4)$. To the best of our knowledge, this is the first lower-bound result for stochastic bilevel optimization, and it shows that this setting is strictly more challenging than both smooth nonconvex optimization and smooth nonconvex-strongly-concave min-max optimization.

We note that [Kwon et al. \(2024\)](#) establish a lower bound of $\Omega(\epsilon^{-6})$ for bilevel optimization under a so-called \mathbf{y}^* -aware stochastic first order oracle with bounded variance. Their hard instance is constructed as

$$f(\mathbf{x}; y) = y, \quad g(\mathbf{x}; y) = (y - F(\mathbf{x}))^2,$$

where $\mathbf{x} \in \mathbb{R}^{\epsilon^{-2}}$, $y \in \mathbb{R}$ and $F(\mathbf{x}) = \epsilon^2 \sum_{i=1}^{\epsilon^{-2}} [\Psi(-x_{i-1})\Phi(-x_i) - \Psi(x_{i-1})\Phi(x_i)]$. It can be verified that $|F(\mathbf{x})| = \mathcal{O}(1)$, and therefore

$$\|\nabla_{\mathbf{x}, y}^2 g(\mathbf{x}; y)\|_2 = \mathcal{O}(\|F(\mathbf{x})\|_2) = \mathcal{O}(\epsilon^2 \sqrt{\epsilon^{-2}}) = \mathcal{O}(\epsilon).$$

In addition, their \mathbf{y}^* -aware oracle requires $\|y - y^*\| = \mathcal{O}(\epsilon)$, which forces $|g(\mathbf{x}; y)|$ to be of order $\mathcal{O}(\epsilon)$. These conditions can also be satisfied in our construction by choosing $L_g = \mathcal{O}(\epsilon)$, since $\|\nabla_{\mathbf{x}, \tilde{\mathbf{y}}}^2 g(\mathbf{x}; \tilde{\mathbf{y}})\|_2 = L_g$ and both \mathbf{x} and $\tilde{\mathbf{y}}$ are bounded. Under this choice, our Theorem 2 also yields a lower bound of order $\Omega(\epsilon^{-6})$.

In contrast, a more standard and practically relevant setting assumes $L_g, L_f = \Theta(1)$, independent of ϵ or the condition number κ . Under this commonly studied regime, obtaining an $\Omega(\epsilon^{-6})$ lower bound for bilevel optimization remains an open problem.

5.3 Analysis and Proof Outline for Stochastic Lower Bound

We use the following construction $\{f_{sc}(\mathbf{x}; \tilde{\mathbf{y}}), g_{sc}(\mathbf{x}; \tilde{\mathbf{y}})\}$ as the hard instance in the stochastic setting. For any $\mathbf{x} \in \mathcal{C}_{r_x \lambda/n}^T$ and $\tilde{\mathbf{y}} \in \mathcal{C}_{r_y \lambda}^{n(T+1)}$,

$$\begin{aligned} f_{sc}(\mathbf{x}; \tilde{\mathbf{y}}) &= \sum_{i=1}^T \frac{\lambda^2 L_f}{L} \left[\Psi\left(-\frac{C_l}{\lambda} y_n^{(i-1)}\right) \Phi\left(-\frac{C_r}{\lambda} y_1^{(i)}\right) - \Psi\left(\frac{C_l}{\lambda} y_n^{(i-1)}\right) \Phi\left(\frac{C_r}{\lambda} y_1^{(i)}\right) \right] \\ g_{sc}(\mathbf{x}; \tilde{\mathbf{y}}) &= \sum_{i=0}^T \left[\frac{L_g n^2}{2(4n^2 + 1)} (\mathbf{y}^{(i)})^\top \left(\frac{1}{n^2} I_n + A \right) \mathbf{y}^{(i)} - L_g (\mathbf{b}_x^{(i)})^\top \mathbf{y}^{(i)} \right], \end{aligned} \quad (14)$$

where r_x and r_y are positive numerical constants from the hypercube sizes, chosen such that $r_y \geq 10r_x$ and $r_x > \frac{1}{\tilde{C}}$, where $\tilde{C} > 0$ is the numerical constant defined in eq. (10). The constants C_l, C_r , and L are the same as in the deterministic setting. The parameter λ is selected to satisfy $\lambda\sqrt{T} = \mathcal{O}(1)$, and its exact form will be specified later. Recall that $x_0 = \frac{\lambda}{\tilde{C}n} < \frac{r_x \lambda}{n} \in \mathcal{C}_{r_x \lambda/n}^1$.

The following lemma shows that, with appropriately chosen r_x and r_y , the lower-level minimizer $\tilde{\mathbf{y}}^*$ lies within the selected bounded domain.

Lemma 6. *If $r_y \geq 10r_x$, the lower-level minimizer $\tilde{\mathbf{y}}^*$ of the instance in eq. (14) satisfies $\tilde{\mathbf{y}}^* \in \mathcal{C}_{r_y \lambda}^{n(T+1)}$.*

Building on Lemma 6, we establish the following lemma, which provides several properties of the instance in eq. (14) that will be used in the proof of the main theorem.

Lemma 7. *Suppose $r_y \geq 10r_x$, $r_x > \frac{1}{\tilde{C}}$, and $\lambda\sqrt{T} = \mathcal{O}(1)$. The functions f_{sc} and g_{sc} satisfy:*

- (a) f_{sc} and g_{sc} satisfy all items 1-4 in Definition 1.
- (b) $H_{sc}(\mathbf{0}) - \min_{\mathbf{x}} H_{sc}(\mathbf{x}) \leq \frac{12\lambda^2 L_f T}{L}$.
- (c) $H_{sc}(\mathbf{x})$ is L_h -smooth with $L_h = \frac{c_0 n^2 L_f}{L}$ for some numerical constant c_0 .
- (d) For any $(\mathbf{x}, \tilde{\mathbf{y}}) \in \mathcal{C}_{r_x \lambda/n}^T \times \mathcal{C}_{r_y \lambda}^{n(T+1)}$, we have $\|\nabla_{\tilde{\mathbf{y}}} f_{sc}(\mathbf{x}; \tilde{\mathbf{y}})\|_\infty \leq \frac{c_1 \lambda L_f}{L}$, $\|\nabla_{\tilde{\mathbf{y}}} g_{sc}(\mathbf{x}; \tilde{\mathbf{y}})\|_\infty \leq 2L_g r_y \lambda$, and $\|\nabla_{\mathbf{x}} g_{sc}(\mathbf{x}; \tilde{\mathbf{y}})\|_\infty \leq L_g r_y \lambda$, where c_1 is a numerical constant.

Similarly to Lemma 3, we then provide a lower bound of the hyper-gradient norm when the algorithm has not yet reached the end of the chain.

Lemma 8. *Suppose $r_x > \frac{1}{\tilde{C}}$. If $x_i < \frac{\lambda}{\tilde{C}n}$ for some $i \leq T$, then, we have*

$$L_h \|\mathcal{P}_{\mathcal{X}}[\mathbf{x} - (1/L_h) \nabla H_{sc}(\mathbf{x})] - \mathbf{x}\|_2 \geq \frac{c_2 L_f n \lambda}{L},$$

where $c_2 > 0$ is a numerical constant.

Based on all the above auxiliary lemmas, we begin to prove our main theorem.

Proof of Theorem 2. Based on part (d) of Lemma 7, we now construct a probability- p zero-chain following the approach of Arjevani et al. (2023), with a slight modification. In Arjevani et al. (2023), the key idea is to perturb the gradient only at the next coordinate to be activated, so that this coordinate is revealed with probability p . For our zero-chain given in eq. (8), let $i^* \in \{n+1, \dots, (T+1)n\}$ be the next coordinate to activate. Thus, we can define the stochastic gradient as follows.

- When $i^* \bmod n \neq 1$, perturb the gradients at the coordinate $i = i^*$ as $\frac{\xi}{p} \frac{\partial g_{sc}(\mathbf{x}; \tilde{\mathbf{y}})}{\partial y_i}$, where $\xi \sim \text{Bernoulli}(p)$. The gradients at all other coordinates remain unchanged and receive no perturbation.
- When $i^* \bmod n = 1$, perturb the gradients at the coordinate $i = i^*$ as $\frac{\xi}{p} \frac{\partial g_{sc}(\mathbf{x}; \tilde{\mathbf{y}})}{x_j}$, where $j = (i^* - 1)/n$ and $\xi \sim \text{Bernoulli}(p)$. The gradients at all other coordinates remain unchanged and receive no perturbation.

Note that in the above stochastic oracles, we **do not** perturb the gradients of f . It can be verified that the stochastic gradients defined above are unbiased. Using Lemma 5 together with part (d) of Lemma 7, we conclude that our construction in eq. (14), equipped with these stochastic oracles, forms a probability- p zero-chain, and the variance of the oracles is bounded by

$$c_3 L_g^2 \lambda^2 \left(\frac{1-p}{p} \right),$$

where the bound follows from (d) of Lemma 7, and c_3 is a positive numerical constant. Thus, to ensure the variance is bounded by σ^2 , it suffices to choose

$$p = \min \left\{ 1, c_3 \frac{L_g^2 \lambda^2}{\sigma^2} \right\}. \quad (15)$$

Then, based on Lemma 9 and the stochastic oracles constructed above, we have that with probability $1 - \delta$, $x_T = 0$ if

$$t \leq \frac{(n-1)T - 1 - \log(\frac{1}{\delta})}{2p}. \quad (16)$$

Based on the choice of p in eq. (15), we have

$$\frac{(n-1)T - 1 - \log(\frac{1}{\delta})}{2p} \geq \frac{((n-1)T - 1 - \log(\frac{1}{\delta}))\sigma^2}{2c_3 L_g^2 \lambda^2},$$

which, together with eq. (16), yields that with probability $1 - \delta$, $x_T = 0$ for all

$$t \leq \frac{((n-1)T - 1 - \log(\frac{1}{\delta}))\sigma^2}{2c_3 L_g^2 \lambda^2}.$$

This, with Lemma 8, implies that with probability $1 - \delta$, $x_T = 0$ for all $t \leq \frac{((n-1)T - 1 - \log(\frac{1}{\delta}))\sigma^2}{2c_3 L_g^2 \lambda^2}$, and hence

$$L_h \|\mathcal{P}_{\mathcal{X}}[\mathbf{x}^t - (1/L_h)\nabla H_{sc}(\mathbf{x}^t)] - \mathbf{x}^t\|_2 \geq \frac{c_2 L_f n \lambda}{L},$$

which, by setting $\lambda = \frac{2L\epsilon}{c_2 L_f n}$, yields that $L_h \|\mathcal{P}_{\mathcal{X}}[\mathbf{x}^t - (1/L_h)\nabla H_{sc}(\mathbf{x}^t)] - \mathbf{x}^t\|_2 \geq 2\epsilon$. Set $\delta = \frac{1}{2}$. Then, for all $t \leq \frac{((n-1)T - 1 - \log(\frac{1}{\delta}))\sigma^2}{2c_3 L_g^2 \lambda^2}$,

$$\mathbb{E} [L_h \|\mathcal{P}_{\mathcal{X}}[\mathbf{x}^t - (1/L_h)\nabla H_{sc}(\mathbf{x}^t)] - \mathbf{x}^t\|_2] \geq \frac{1}{2}(2\epsilon) = \epsilon.$$

Based on (b) of Lemma 7, we have $\frac{12\lambda^2 L_f T}{L} = \Delta$, which implies that $T = \frac{\Delta L}{12\lambda^2 L_f}$. Thus, to achieve an ϵ -accurate stationary point, the algorithm must use at least

$$\Omega\left(\frac{nT\sigma^2}{L_g^2\lambda^2}\right) = \Omega\left(\frac{n\Delta\sigma^2}{L_f L_g^2\lambda^4}\right) = \Omega\left(\frac{n^5 L_f^3 \Delta\sigma^2}{L_g^2\epsilon^4}\right),$$

which, together with $n = \sqrt{\kappa}$, finishes the proof. \square

6 Conclusion and Future Works

In this work, we developed new hard instances that establish improved lower bounds for smooth nonconvex and strongly convex bilevel optimization under both deterministic and stochastic first order oracle models. Our results demonstrate that bilevel optimization is fundamentally more challenging than classical single-level and min-max formulations, and they reveal significant separations between the best known upper and lower bounds. These findings highlight that the current theoretical understanding of bilevel optimization is still far from complete.

There are several promising directions for future research. First, even for the simplified and practically meaningful setting in which the lower level function is quadratic, the optimal complexity remains open. We suggest that closing these gaps may require first studying this simpler yet meaningful quadratic setting. Moreover, our constructions suggest that sharper lower and upper bounds may be obtained by designing algorithms that exploit higher-order structure of the lower-level function. Second, closing the large gaps between the existing upper bounds and our lower bounds, especially the gap of order κ^2 in the deterministic case and the dependence on ϵ in the stochastic case, represents an important challenge. Third, another compelling direction is to investigate whether an $\Omega(\epsilon^{-6})$ lower bound can be achieved under the standard regime where the smoothness constants L_f and L_g are $\Theta(1)$, a question that remains unresolved. Finally, extending the lower bound framework to broader variants of bilevel optimization, including settings with constraints, approximate inner solvers, or distributed architectures, may deepen the understanding of the fundamental limits of bilevel learning.

Overall, we hope that the insights developed in this paper serve as a starting point for further studies toward a complete theory of the computational complexity of bilevel optimization.

A Proofs for Deterministic Lower Bound

A.1 Proof of Lemma 2

It is straightforward to verify that $c \leq S_{1,1} = 1 \leq C$, so the claim holds for $n = 1$. For the remainder of the proof, we assume $n \geq 2$. Set $s := 1/n^2$. The eigenpairs of A are

$$\mu_k = 2\left(1 - \cos\left(\frac{(k-1)\pi}{n}\right)\right), \quad k = 1, \dots, n,$$

with orthonormal eigenvectors

$$q_1(j) = \frac{1}{\sqrt{n}}, \quad q_k(j) = \sqrt{\frac{2}{n}} \cos\left(\frac{(k-1)(j-\frac{1}{2})\pi}{n}\right), \quad k \geq 2.$$

Thus

$$A = Q\Lambda Q^\top, \quad \Lambda = \text{diag}(\mu_1, \dots, \mu_n), \quad Q = [q_1 \dots q_n].$$

Hence

$$S = (A + sI_n)^{-1} = Q(\Lambda + sI_n)^{-1}Q^\top = \sum_{k=1}^n \frac{1}{\mu_k + s} q_k q_k^\top,$$

so we can express $S_{i,j}$ as

$$S_{i,j} = \sum_{k=1}^n \frac{q_k(i)q_k(j)}{\mu_k + s}.$$

Note that $\frac{q_1(i)q_1(j)}{s} = \frac{1/n}{1/n^2} = n$. Thus, we have

$$S_{i,j} = n + R_{i,j}, \quad R_{i,j} := \sum_{k=2}^n \frac{q_k(i)q_k(j)}{\mu_k + s}.$$

Higher eigenmodes. Because $|q_k(\cdot)| \leq \sqrt{2/n}$,

$$|q_k(i)q_k(j)| \leq \frac{2}{n}.$$

Also for $k \geq 2$,

$$\mu_k = 2(1 - \cos(\frac{(k-1)\pi}{n})) \geq \frac{4(k-1)^2}{n^2},$$

so

$$\frac{1}{\mu_k + s} \leq \frac{n^2}{4(k-1)^2}.$$

Thus

$$|R_{i,j}| \leq \sum_{k=2}^n \frac{2}{n} \cdot \frac{n^2}{4(k-1)^2} = \frac{n}{2} \sum_{m=1}^{n-1} \frac{1}{m^2} \leq \frac{\pi^2}{12} n.$$

Final bounds.

$$S_{n,n} = n + R_{n,n}, \quad R_{n,n} \geq 0, \quad S_{1,n} = n + R_{1,n}, \quad |R_{1,n}| \leq \frac{\pi^2}{12} n.$$

Hence for all $n \geq 2$,

$$\left(1 - \frac{\pi^2}{12}\right)n \leq S_{1,n}, \quad S_{n,n} \leq \left(1 + \frac{\pi^2}{12}\right)n.$$

Then, the proof is complete.

A.2 Proof of Lemma 3

Note that $x_0 = \frac{\lambda}{C_l M_{n,n}} = \frac{\lambda}{\tilde{C}_n}$. Since $|x_0| \geq \frac{\lambda}{\tilde{C}_n}$ and $|x_i| < \frac{\lambda}{\tilde{C}_n}$, we can find some $0 < j \leq i$ such that $|x_{j-1}| \geq \frac{\lambda}{\tilde{C}_n}$ and $|x_j| < \frac{\lambda}{\tilde{C}_n}$. Thus, look at

$$\begin{aligned} \frac{\partial H(\mathbf{x})}{\partial x_j} = & -\frac{\lambda L_f \tilde{C}_n}{L} \left[\Psi\left(-\frac{\tilde{C}_n}{\lambda} x_{j-1}\right) \Phi'\left(-\frac{\tilde{C}_n}{\lambda} x_j\right) + \Psi\left(\frac{\tilde{C}_n}{\lambda} x_{j-1}\right) \Phi'\left(\frac{\tilde{C}_n}{\lambda} x_j\right) \right] \\ & -\frac{\lambda L_f \tilde{C}_n}{L} \left[\Psi'\left(-\frac{\tilde{C}_n}{\lambda} x_j\right) \Phi\left(-\frac{\tilde{C}_n}{\lambda} x_{j+1}\right) + \Psi'\left(\frac{\tilde{C}_n}{\lambda} x_j\right) \Phi\left(\frac{\tilde{C}_n}{\lambda} x_{j+1}\right) \right], \end{aligned}$$

which, in conjunction with Lemma 1 (items 2 and 4), implies that

$$\|\nabla H(\mathbf{x})\|_2 \geq \left| \frac{\partial H(\mathbf{x})}{\partial x_j} \right| \geq \frac{\lambda L_f \tilde{C}_n}{L}.$$

Then, the proof is complete.

A.3 Proof of Lemma 4

First note that

$$H(\mathbf{0}) = \frac{\lambda^2 L_f}{L} \left[\left(\Psi\left(-\frac{C_l}{\lambda} M_{n,n} x_0\right) - \Psi\left(\frac{C_l}{\lambda} M_{n,n} x_0\right) \right) \Phi(0) \right] \leq 0, \quad (17)$$

where the inequality follows because $\frac{C_l}{\lambda} M_{n,n} x_0 \geq 0$ and from the definitions of Ψ and Φ functions in eq. (4). Furthermore, based on Lemma 1 (item 4), we have that

$$H(\mathbf{x}) \geq -\frac{\lambda^2 L_f}{L} \sum_{i=1}^T \Psi\left(\frac{C_l}{\lambda} M_{n,n} x_{i-1}\right) \Phi\left(\frac{C_r}{\lambda} M_{1,n} x_i\right) \geq -\frac{12\lambda^2 L_f T}{L}, \quad (18)$$

which, combined with $H(\mathbf{0}) \leq 0$, implies that

$$H(\mathbf{0}) - \inf_{\mathbf{x}} H(\mathbf{x}) \leq \frac{12\lambda^2 L_f T}{L},$$

which finishes the proof.

B Proofs for Stochastic Lower Bound

B.1 Auxiliary Lemmas

For a probability- p zero-chain, at each iteration, a new coordinate is discovered with probability at most p . Therefore, it takes at least $1/p$ steps in expectation to activate a new coordinate. The following lemma, adapted from Arjevani et al. (2023); Li et al. (2021), shows that at least $\Omega(T/p)$ iterations are required to reach the end of a probability- p zero-chain.

Lemma 9 ((Arjevani et al., 2023, Lemma 1)). *Let $f : \mathcal{X} \rightarrow \mathbb{R}$, where $\mathcal{X} \subset \mathbb{R}^T$ satisfies $\text{supp}(P_{\mathcal{X}}(\mathbf{x})) = \text{supp}(\mathbf{x})$, $\forall \mathbf{x} \in \mathbb{R}^T$, and suppose f is a probability- p zero-chain with a stochastic first-order oracle. Then, for any first-order algorithm, with probability at least $1 - \delta$, the T -th coordinate of \mathbf{x} at the t^{th} iteration, satisfies*

$$x_T^t = 0, \quad \forall t \leq \frac{T - \log(1/\delta)}{2p}.$$

Lemma 10. *Recall $S := (A + \frac{1}{n^2}I_n)^{-1}$. For every $i = 1, \dots, n$,*

$$S_{1,n} \leq S_{i,n} \leq S_{n,n}.$$

Proof. Define $B := A + \frac{1}{n^2}I_n$, and $\mathbf{v} \in \mathbb{R}^n$ the last column of S , i.e., $\mathbf{v} := S_{\cdot,n}$, $v_i := S_{i,n}$, $i = 1, \dots, n$. Since \mathbf{v} is the last column of $S = B^{-1}$, it solves the linear system

$$B\mathbf{v} = \mathbf{e}_n.$$

Writing this componentwise, we obtain

$$\begin{aligned} (1 + n^{-2})v_1 - v_2 &= 0, \\ -v_{i-1} + (2 + n^{-2})v_i - v_{i+1} &= 0, \quad i = 2, \dots, n-1, \\ -v_{n-1} + (1 + n^{-2})v_n &= 1. \end{aligned}$$

Define the forward differences

$$d_i := v_{i+1} - v_i, \quad i = 1, \dots, n-1.$$

We next derive a system of equations for $\mathbf{d} = (d_1, \dots, d_{n-1})^\top$. For $i = 2, \dots, n-2$, subtracting the equation at index i from that at index $i+1$ gives

$$(-v_i + (2 + n^{-2})v_{i+1} - v_{i+2}) - (-v_{i-1} + (2 + n^{-2})v_i - v_{i+1}) = 0,$$

which can be rewritten as

$$-d_{i-1} + (2 + n^{-2})d_i - d_{i+1} = 0, \quad i = 2, \dots, n-2.$$

From the first equation, we obtain

$$(1 + n^{-2})v_1 - v_2 = 0 \implies (2 + n^{-2})d_1 - d_2 = 0.$$

From the last equation, we obtain

$$-v_{n-1} + (1 + n^{-2})v_n = 1 \implies -d_{n-2} + (2 + n^{-2})d_{n-1} = 1.$$

Thus, the vector $\mathbf{d} = (d_1, \dots, d_{n-1})^\top$ satisfies a tri-diagonal linear system

$$B'\mathbf{d} = \mathbf{e}_{n-1},$$

where $B' \in \mathbb{R}^{(n-1) \times (n-1)}$ is the symmetric tri-diagonal matrix

$$B' = \begin{bmatrix} 2 + n^{-2} & -1 & & & \\ -1 & 2 + n^{-2} & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 + n^{-2} \end{bmatrix}.$$

The matrix B' is strictly diagonally dominant with positive diagonal entries and nonpositive off-diagonal entries, hence an irreducible M -matrix (Berman & Plemmons, 1994). It is therefore positive definite and its inverse is entrywise nonnegative:

$$(B')^{-1} \geq 0 \quad (\text{entrywise}).$$

Since $\mathbf{d} = (B')^{-1} \mathbf{e}_{n-1}$, we obtain

$$d_i \geq 0, \quad i = 1, \dots, n-1.$$

Equivalently,

$$v_{i+1} - v_i = d_i \geq 0 \implies v_1 \leq v_2 \leq \dots \leq v_n.$$

The inequalities $S_{1,n} \leq S_{i,n} \leq S_{n,n}$ follow immediately from this monotonicity. \square

B.2 Proof of Lemma 6

Note that the minimizers $(\mathbf{y}^{(i)})^*, i = 0, \dots, T$ take the forms of

$$(\mathbf{y}^{(i)})^* = \underbrace{\frac{4n^2 + 1}{n^2} \left(\frac{1}{n^2} I_n + A \right)^{-1}}_M \mathbf{b}_x^{(i)}.$$

Combining Lemma 2 and Lemma 10, we have for all $i = 1, \dots, n$,

$$4cn \leq M_{i,n} \leq 5Cn,$$

where $c = 1 - \frac{\pi^2}{12}$ and $C = 1 + \frac{\pi^2}{12}$. Thus, we have

$$\|(\mathbf{y}^{(i)})^*\|_\infty \leq 5Cn|x_i| \leq 5Cr_x\lambda < 10r_x\lambda < r_y\lambda,$$

which finishes the proof.

B.3 Proof of Lemma 7

The proof of (a) is identical to the deterministic case. The proof of (b) follows the same reasoning as in Lemma 4. To establish (c), recall that

$$H_{sc}(\mathbf{x}) = \sum_{i=1}^T \frac{\lambda^2 L_f}{L} \left[\Psi\left(-\frac{C_l}{\lambda} M_{n,n} x_{i-1}\right) \Phi\left(-\frac{C_r}{\lambda} M_{1,n} x_i\right) - \Psi\left(\frac{C_l}{\lambda} M_{n,n} x_{i-1}\right) \Phi\left(\frac{C_r}{\lambda} M_{1,n} x_i\right) \right].$$

Then one can verify that $\nabla^2 H_{sc}(\mathbf{x})$ is a tri-diagonal matrix whose entries are all of order $\mathcal{O}\left(\frac{L_f n^2}{L}\right)$. Consequently, $\|\nabla^2 H_{sc}(\mathbf{x})\|_2 = \mathcal{O}\left(\frac{L_f n^2}{L}\right)$. To prove (d), note that each coordinate of $\nabla_{\tilde{\mathbf{y}}} f_{sc}(\mathbf{x}; \tilde{\mathbf{y}})$ takes an order of $\mathcal{O}\left(\frac{\lambda L_f}{L}\right)$, and hence $\|\nabla_{\tilde{\mathbf{y}}} f_{sc}(\mathbf{x}; \tilde{\mathbf{y}})\|_\infty = \mathcal{O}\left(\frac{\lambda L_f}{L}\right)$.

For $\nabla_{\tilde{\mathbf{y}}} g_{sc}(\mathbf{x}; \tilde{\mathbf{y}})$, note that

$$\begin{aligned} \left\| \frac{\partial g_{sc}(\mathbf{x}; \tilde{\mathbf{y}})}{\partial \mathbf{y}^{(i)}} \right\|_\infty &= \left\| \frac{L_g n^2}{4n^2 + 1} \left(\frac{1}{n^2} I_n + A \right) \mathbf{y}^{(i)} - L_g \mathbf{b}_x^{(i)} \right\|_\infty \\ &\leq \frac{L_g}{4} \left\| \frac{1}{n^2} I_n + A \right\|_\infty \left\| \mathbf{y}^{(i)} \right\|_\infty + L_g |x_i| \\ &\leq \frac{L_g}{4} \left(\frac{1}{n^2} + 4 \right) r_y \lambda + \frac{L_g r_x \lambda}{n} \\ &\leq \frac{5}{4} L_g r_y \lambda + \frac{1}{10} L_g r_y \lambda \leq 2 L_g r_y \lambda, \end{aligned}$$

which holds for all $i = 0, \dots, T$. This implies that $\|\nabla_{\tilde{\mathbf{y}}} g_{sc}(\mathbf{x}; \tilde{\mathbf{y}})\|_\infty \leq 2 L_g r_y \lambda$.

For $\|\nabla_{\mathbf{x}} g_{sc}(\mathbf{x}; \tilde{\mathbf{y}})\|_\infty$, note that

$$\left| \frac{\partial g_{sc}(\mathbf{x}; \tilde{\mathbf{y}})}{\partial x_i} \right| = L_g |y_n^{(i)}| \leq L_g r_y \lambda,$$

which yields that $\|\nabla_{\mathbf{x}} g_{sc}(\mathbf{x}; \tilde{\mathbf{y}})\|_\infty \leq L_g r_y \lambda$. Then, the proof is complete.

B.4 Proof of Lemma 8

Note that $x_0 = \frac{\lambda}{C_l M_{n,n}} = \frac{\lambda}{\tilde{C}n}$. Since $|x_0| \geq \frac{\lambda}{\tilde{C}n}$ and $|x_i| < \frac{\lambda}{\tilde{C}n}$, we can find some $0 < j \leq i$ such that $|x_{j-1}| \geq \frac{\lambda}{\tilde{C}n}$ and $|x_j| < \frac{\lambda}{\tilde{C}n}$. Thus, look at

$$\begin{aligned} \frac{\partial H_{sc}(\mathbf{x})}{\partial x_j} &= -\frac{\lambda L_f \tilde{C}n}{L} \left[\Psi\left(-\frac{\tilde{C}n}{\lambda} x_{j-1}\right) \Phi'\left(-\frac{\tilde{C}n}{\lambda} x_j\right) + \Psi\left(\frac{\tilde{C}n}{\lambda} x_{j-1}\right) \Phi'\left(\frac{\tilde{C}n}{\lambda} x_j\right) \right] \\ &\quad - \frac{\lambda L_f \tilde{C}n}{L} \left[\Psi'\left(-\frac{\tilde{C}n}{\lambda} x_j\right) \Phi\left(-\frac{\tilde{C}n}{\lambda} x_{j+1}\right) + \Psi'\left(\frac{\tilde{C}n}{\lambda} x_j\right) \Phi\left(\frac{\tilde{C}n}{\lambda} x_{j+1}\right) \right]. \end{aligned} \quad (19)$$

Then, if $|x_j - (1/L_h) \frac{\partial H_{sc}(\mathbf{x})}{\partial x_j}| \leq r_x \lambda / n$, then we have

$$L_h \|\mathcal{P}_{\mathcal{X}}[\mathbf{x} - (1/L_h) \nabla H_{sc}(\mathbf{x})] - \mathbf{x}\|_2 \geq \left| \frac{\partial H_{sc}(\mathbf{x})}{\partial x_j} \right| \geq \frac{\lambda L_f \tilde{C}n}{L}.$$

Otherwise, i.e., $|x_j - (1/L_h) \frac{\partial H_{sc}(\mathbf{x})}{\partial x_j}| > r_x \lambda / n$, we have

$$\begin{aligned} L_h \|\mathcal{P}_{\mathcal{X}}[\mathbf{x} - (1/L_h) \nabla H_{sc}(\mathbf{x})] - \mathbf{x}\|_2 &\geq L_h \left| \mathcal{P}_{\mathcal{C}_{r_x \lambda / n}^1} \left[x_j - (1/L_h) \frac{\partial H_{sc}(\mathbf{x})}{\partial x_j} \right] - x_j \right| \\ &\geq L_h \left(\frac{r_x \lambda}{n} - |x_j| \right) \geq L_h \left(\frac{r_x \lambda}{n} - \frac{\lambda}{\tilde{C}n} \right) \\ &\stackrel{(i)}{\geq} \frac{c_0 n^2 L_f}{L} \left(r_x - \frac{1}{\tilde{C}} \right) \frac{\lambda}{n} = \frac{c_0 L_f}{L} \left(r_x - \frac{1}{\tilde{C}} \right) n \lambda, \end{aligned}$$

where (i) follows from (c) of Lemma 7. Combining the above two cases completes the proof.

References

- Michael Arbel and Julien Mairal. Non-convex bilevel games with critical point selection maps. *Advances in Neural Information Processing Systems*, 35:8013–8026, 2022.
- Yossi Arjevani, Yair Carmon, John C Duchi, Dylan J Foster, Nathan Srebro, and Blake Woodworth. Lower bounds for non-convex stochastic optimization. *Mathematical Programming*, 199(1):165–214, 2023.
- Abraham Berman and Robert J Plemmons. *Nonnegative matrices in the mathematical sciences*. SIAM, 1994.
- Jerome Bracken and James T McGill. Mathematical programs with optimization problems in the constraints. *Operations Research*, 21(1):37–44, 1973.
- Yair Carmon, John C Duchi, Oliver Hinder, and Aaron Sidford. Lower bounds for finding stationary points i. *Mathematical Programming*, 184(1):71–120, 2020.
- Yair Carmon, John C Duchi, Oliver Hinder, and Aaron Sidford. Lower bounds for finding stationary points ii: first-order methods. *Mathematical Programming*, 185(1):315–355, 2021.
- Lesi Chen, Jing Xu, and Jingzhao Zhang. On finding small hyper-gradients in bilevel optimization: Hardness results and improved analysis. In *Conference on Learning Theory*, pp. 947–980. PMLR, 2024.
- Lesi Chen, Yaohua Ma, and Jingzhao Zhang. Near-optimal nonconvex-strongly-convex bilevel optimization with fully first-order oracles. *Journal of Machine Learning Research*, 26(109):1–56, 2025.
- Tianyi Chen, Yuejiao Sun, Quan Xiao, and Wotao Yin. A single-timescale method for stochastic bilevel optimization. In *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, 2022.
- Mathieu Dagr  ou, Pierre Ablin, Samuel Vaiter, and Thomas Moreau. A framework for bilevel optimization that enables stochastic and global variance reduction algorithms. *Advances in Neural Information Processing Systems*, 35:26698–26710, 2022.
- Justin Domke. Generic methods for optimization-based modeling. In *Artificial Intelligence and Statistics*, pp. 318–326. PMLR, 2012.
- Cong Fang, Chris Junchi Li, Zhouchen Lin, and Tong Zhang. Spider: Near-optimal non-convex optimization via stochastic path-integrated differential estimator. *Advances in neural information processing systems*, 31, 2018.
- Luca Franceschi, Michele Donini, Paolo Frasconi, and Massimiliano Pontil. Forward and reverse gradient-based hyperparameter optimization. In *International Conference on Machine Learning (ICML)*, pp. 1165–1173, 2017.
- Zhishuai Guo and Tianbao Yang. Randomized stochastic variance-reduced methods for stochastic bilevel optimization. *arXiv preprint arXiv:2105.02266*, 2021.
- Pierre Hansen, Brigitte Jaumard, and Gilles Savard. New branch-and-bound rules for linear bilevel programming. *SIAM Journal on Scientific and Statistical Computing*, 13(5):1194–1217, 1992.

- Mingyi Hong, Hoi-To Wai, Zhaoran Wang, and Zhuoran Yang. A two-timescale framework for bilevel optimization: Complexity analysis and application to actor-critic. *SIAM Journal on Optimization*, 33(1): 147–180, 2023.
- Minhui Huang, Dewei Zhang, and Kaiyi Ji. Achieving linear speedup in non-iid federated bilevel learning. In *International conference on machine learning*, pp. 14039–14059. PMLR, 2023.
- Kaiyi Ji and Yingbin Liang. Lower bounds and acceleration algorithms for bilevel optimization. *Journal of machine learning research*, 2022.
- Kaiyi Ji and Lei Ying. Network utility maximization with general and unknown utility functions: A distributed, data-driven bilevel optimization approach. *Submitted*, 2022.
- Kaiyi Ji, Junjie Yang, and Yingbin Liang. Bilevel optimization: Convergence analysis and enhanced design. In *International Conference on Machine Learning*, pp. 4882–4892. PMLR, 2021.
- Liuyuan Jiang, Quan Xiao, Lisha Chen, and Tianyi Chen. Beyond value functions: Single-loop bilevel optimization under flatness conditions. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025.
- Vijay R Konda and John N Tsitsiklis. Actor-critic algorithms. In *Advances in neural information processing systems*, pp. 1008–1014, 2000.
- Jeongyeol Kwon, Dohyun Kwon, Stephen Wright, and Robert Nowak. A fully first-order method for stochastic bilevel optimization. pp. 18083–18113, 2023.
- Jeongyeol Kwon, Dohyun Kwon, and Hanbaek Lyu. On the complexity of first-order methods in stochastic bilevel optimization. *arXiv preprint arXiv:2402.07101*, 2024.
- Haochuan Li, Yi Tian, Jingzhao Zhang, and Ali Jadbabaie. Complexity lower bounds for nonconvex-strongly-concave min-max optimization. *Advances in Neural Information Processing Systems (NeurIPS)*, 34: 1792–1804, 2021.
- Renjie Liao, Yuwen Xiong, Ethan Fetaya, Lisa Zhang, KiJung Yoon, Xaq Pitkow, Raquel Urtasun, and Richard Zemel. Reviving and improving recurrent back-propagation. In *International Conference on Machine Learning*, pp. 3082–3091. PMLR, 2018.
- Bo Liu, Mao Ye, Stephen Wright, Peter Stone, and Qiang Liu. Bome! bilevel optimization made easy: A simple first-order approach. *Advances in Neural Information Processing Systems*, 35:17248–17262, 2022.
- Risheng Liu, Pan Mu, Xiaoming Yuan, Shangzhi Zeng, and Jin Zhang. A generic first-order algorithmic framework for bi-level programming beyond lower-level singleton. In *International Conference on Machine Learning*, 2020.
- Risheng Liu, Jiaxin Gao, Jin Zhang, Deyu Meng, and Zhouchen Lin. Investigating bi-level optimization for learning and vision from a unified perspective: A survey and beyond. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(12):10045–10067, 2021a.

- Risheng Liu, Xuan Liu, Xiaoming Yuan, Shangzhi Zeng, and Jin Zhang. A value-function-based interior-point method for non-convex bi-level optimization. In *International Conference on Machine Learning*, 2021b.
- Risheng Liu, Yaohua Liu, Shangzhi Zeng, and Jin Zhang. Towards gradient-based bilevel optimization with non-convex followers and beyond. *Advances in Neural Information Processing Systems*, 34:8662–8675, 2021c.
- Zhaosong Lu and Sanyou Mei. First-order penalty methods for bilevel optimization. *SIAM Journal on Optimization*, 34(2):1937–1969, 2024.
- Zhaosong Lu and Sanyou Mei. Solving bilevel optimization via sequential minimax optimization. *arXiv preprint arXiv:2511.07398*, 2025.
- Dougal Maclaurin, David Duvenaud, and Ryan Adams. Gradient-based hyperparameter optimization through reversible learning. In *International conference on machine learning*, pp. 2113–2122. PMLR, 2015.
- Arkadi S Nemirovsky. Information-based complexity of linear operator equations. *Journal of Complexity*, 8(2):153–175, 1992.
- Yurii Nesterov et al. *Lectures on convex optimization*, volume 137. Springer, 2018.
- Aravind Rajeswaran, Chelsea Finn, Sham M Kakade, and Sergey Levine. Meta-learning with implicit gradients. In *Advances in Neural Information Processing Systems*, pp. 113–124, 2019.
- Han Shen and Tianyi Chen. On penalty-based bilevel gradient descent method. In *Proceedings of the 40th International Conference on Machine Learning*, pp. 30992–31015, 2023.
- Chenggen Shi, Jie Lu, and Guangquan Zhang. An extended kuhn–tucker approach for linear bilevel programming. *Applied Mathematics and Computation*, 162(1):51–63, 2005.
- Daouda Sow, Kaiyi Ji, Ziwei Guan, and Yingbin Liang. A primal-dual approach to bilevel optimization with multiple inner minima. *arXiv preprint arXiv:2203.01123*, 2022.
- Davoud Ataee Tarzanagh, Mingchen Li, Christos Thrampoulidis, and Samet Oymak. Fednest: Federated bilevel, minimax, and compositional optimization. In *International Conference on Machine Learning*, pp. 21146–21179. PMLR, 2022.
- Chen Wang, Kaiyi Ji, Junyi Geng, Zhongqiang Ren, Taimeng Fu, Fan Yang, Yifan Guo, Haonan He, Xiangyu Chen, Zitong Zhan, et al. Imperative learning: A self-supervised neural-symbolic learning framework for robot autonomy. *arXiv preprint arXiv:2406.16087*, 2024.
- Xiaoyu Wang, Rui Pan, Renjie Pi, and Tong Zhang. Effective bilevel optimization via minimax reformulation. *arXiv preprint arXiv:2305.13153*, 2023.
- Junjie Yang, Kaiyi Ji, and Yingbin Liang. Provably faster algorithms for bilevel optimization. *Advances in Neural Information Processing Systems*, 34:13670–13682, 2021.
- Yifan Yang, Peiyao Xiao, and Kaiyi Ji. Achieving $\mathcal{O}(\epsilon^{-1.5})$ complexity in Hessian/Jacobian-free stochastic bilevel optimization. *arXiv preprint arXiv:2312.03807*, 2023.

Yifan Yang, Hao Ban, Minhui Huang, Shiqian Ma, and Kaiyi Ji. Tuning-free bilevel optimization: New algorithms and convergence analysis. In *The Thirteenth International Conference on Learning Representations*, 2024.