# Avoid Overclaims:
# Summary of Complexity Bounds for Algorithms in Minimization and Minimax Optimization

Siqi Zhang

SME, NJU

Yifan Hu

EPFL & ETH Zürich

# Introduction

- Review the complexity bounds of first-order methods in optimization in different convexity/smoothness scenarios

- Two problem settings
  - Minimization
  - Minimax Optimization

$$\min_{x \in \mathcal{X}} f(x) \qquad \min_{x \in \mathcal{X}} \left[ f(x) \triangleq \max_{y \in \mathcal{Y}} g(x,y) \right]$$

- Three stochastic settings
  - Deterministic (general)
  - Finite-sum
  - Stochastic optimization

$$\min_{x \in \mathcal{X}} f(x) \qquad \min_{x \in \mathcal{X}} f(x) \triangleq \frac{1}{n} \sum_{i=1}^{n} f_i(x)$$

$$\min_{x \in \mathcal{X}} f(x) \triangleq \mathbb{E}_{\xi \sim \mathcal{D}}[f(x; \xi)]$$

# Prior Arts

- Sebastian Ruder's Blog on common optimization algorithms in ML

- Several monographs reviewed optimization algorithms in various settings

- Popular repository tracking nonconvex optimization research by Ju Sun

- We focus on the SOTA upper and lower bounds in various settings

# Oracle Complexity Framework



Function Class

Algorithm and Oracle Class: Information

Convergence Criteria

Complexity Analysis Framework

A Hard Problem Instance

Algorithm (Oracle-based Method)

Lower Complexity Bound

Upper Complexity Bound

# Oracle Complexity Framework (cont.)

# Main Results

- Case 1-1: Deterministic Minimization

- Case 1-2: Finite-sum and Stochastic Minimization

- Case 2-1: (S)C-(S)C Deterministic Minimax Optimization

- Case 2-2: (S)C-(S)C Finite-sum and Stochastic Minimax Optimization

- Case 2-3: NC-(S)C Deterministic Minimax Optimization

- Case 2-4: NC-(S)C Finite-sum and Stochastic Minimax Optimization

- SOTA upper and lower bounds comparison

# Main Results

**Case 1-1: Deterministic Minimization**

| Problem Type | Measure | Lower Bound | Upper Bound | Reference (LB) | Reference (UB) [5] |
|---|---|---|---|---|---|
| $L$-Smooth Convex | Optimality gap | $\Omega\left(\sqrt{L\epsilon^{-1}}\right)$ | ✓ | [25], Theorem 2.1.7 | [25], Theorem 2.2.2 |
| $L$-Smooth $\mu$-SC | Optimality gap | $\Omega\left(\sqrt{\kappa}\log\frac{1}{\epsilon}\right)$ | ✓ | [25], Theorem 2.1.13 | [25], Theorem 2.2.2 |
| NS $G$-Lip Cont. Convex | Optimality gap | $\Omega(G^2\epsilon^{-2})$ | ✓ | [25], Theorem 3.2.1 | [25], Theorem 3.2.2 |
| NS $G$-Lip Cont. $\mu$-SC | Optimality gap | $\Omega(G^2(\mu\epsilon)^{-1})$ | ✓ | [25], Theorem 3.2.5 | [7], Theorem 3.9 [6] |
| $L$-Smooth Convex (function case) | Stationarity | $\Omega\left(\sqrt{\Delta L\epsilon^{-1}}\right)$ | ✓ (within logarithmic) | [26], Theorem 1 | [26], Appendix A.1 |
| $L$-Smooth Convex (domain case) | Stationarity | $\Omega\left(\sqrt{DL}\epsilon^{-\frac{1}{2}}\right)$ | ✓ | [26], Theorem 1 | [27] Section 6.5 |
| $L$-Smooth NC | Stationarity | $\Omega(\Delta L\epsilon^{-2})$ | ✓ | [20], Theorem 1 | [28], Theorem 10.15 |
| NS $G$-Lip Cont. $\rho$-WC | Near-stationarity | Unknown | $\mathcal{O}(\epsilon^{-4})$ | / | [29], Corollary 2.2 |
| $L$-Smooth $\mu$-PL | Optimality gap | $\Omega\left(\kappa\log\frac{1}{\epsilon}\right)$ | ✓ | [30], Theorem 3 | [24], Theorem 1 |

**Remark:**

1. References: [25] [7] [26] [20] [27] [28] [29] [30] [24]

2. $\kappa \triangleq L/\mu \geq 1$ is called the condition number, which can be very large in many applications, e.g., the optimal regularization parameter choice in statistical learning can lead to $\kappa = \Omega(\sqrt{n})$ where $n$ is the sample size [31].

3. The PL condition is a popular assumption in nonconvex optimization, generalizing the strong convexity condition. Based on the summary above, we can find that both smooth strongly convex and smooth PL condition optimization problems have established the optimal complexities (i.e., UB matches LB). However, the LB in the PL case is strictly larger than that of the SC case. Thus, regarding the worst-case complexity, we can say that the PL case is "strictly harder" than the strongly convex case.

# Future Directions

- Richer Problem Structure
  - Bilevel Optimization
  - Compositional Stochastic Optimization
  - Conditional Stochastic Optimization
  - Performative Prediction
  - Contextual Stochastic Optimization
  - Distributionally Robust Optimization
  - ......

- Various optimization problems arising from ML & OR, which come with more involved problem structure and complicated landscape characteristics
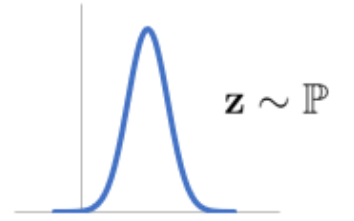
**Deterministic Optimization**

$$\min_{\beta} \; h_\beta(\mathbf{z})$$

$\bullet$ $\mathbf{z}$

**Stochastic Optimization**

$$\inf_{\beta} \mathbb{E}^{\mathbb{P}}[h_\beta(\mathbf{z})]$$

$\mathbf{z} \sim \mathbb{P}$

**Robust Optimization**

$$\min_{\beta} \max_{\mathbf{z} \in \mathcal{Z}} \; h_\beta(\mathbf{z})$$

$\mathbf{z} \in \mathcal{Z}$

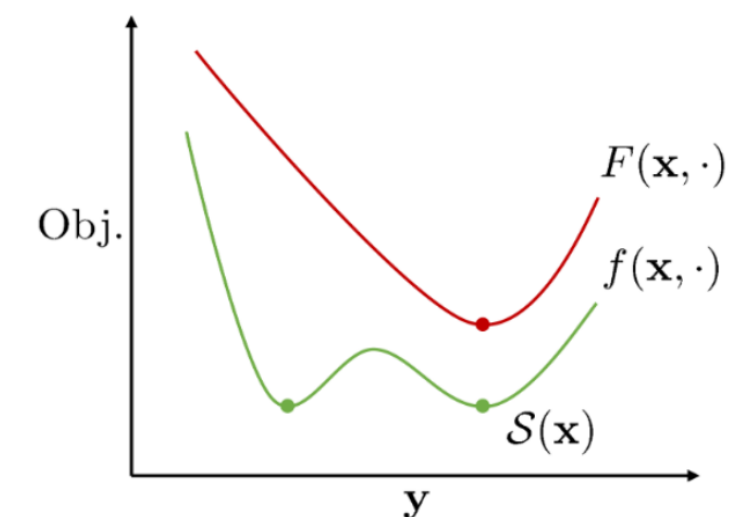**Distributionally Robust Optimization**

$$\inf_{\beta} \sup_{\mathbb{P} \in \Omega} \mathbb{E}^{\mathbb{P}}[h_\beta(\mathbf{z})]$$

$\mathbf{z} \sim \mathbb{P}$

$\mathbb{P} \in \Omega$

$$\min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \mathbf{y}), \; s.t. \; \mathbf{y} \in \mathcal{S}(\mathbf{x})$$
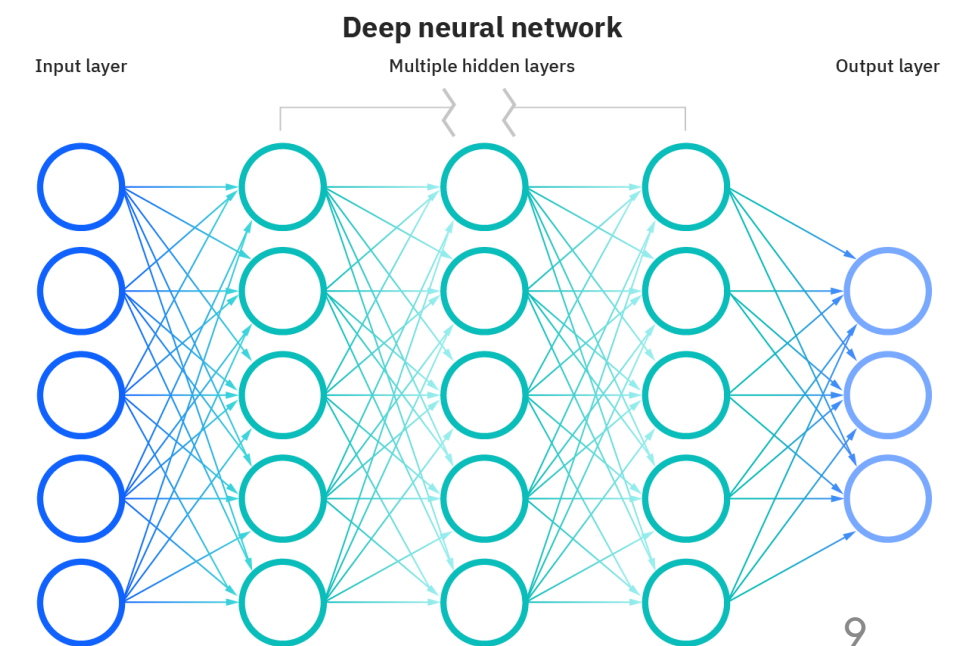
$$\mathcal{S}(\mathbf{x}) := \arg\min_{\mathbf{y} \in \mathcal{Y}} f(\mathbf{x}, \mathbf{y})$$

Obj.

$F(\mathbf{x}, \cdot)$

$f(\mathbf{x}, \cdot)$

$\mathcal{S}(\mathbf{x})$

$\mathbf{y}$

# Future Directions (cont.)

- Landscape Analysis
  - Hidden Convexity
  - PL/KL Conditions
  - Relaxed Smoothness
  - Low Rank Structure

- Beyond Classical Oracle Model
  - Average-case complexity
  - Arithmetic complexity
  - Communication complexity in distributed optimization
  - Long stepsize in first-order methods and achieve a faster convergence rate

- Unified Lower Bounds
  - Lower bound valid for any given dimension
  - Information theoretic-based lower bounds

**Deep neural network**

Input layer     Multiple hidden layers     Output layer

# Feedback Appreciated!

- Possibility of overlooking certain relevant works, subtle technical conditions, or potential inaccuracies in interpreting the literature

- Don't hesitate to send emails to bring them to our attention!

- Constructive feedback, corrections, and suggestions are highly appreciated.