ZEROTH-ORDER NON-CONVEX SMOOTH OPTIMIZATION: Local minimax rates

Yining Wang, CMU

joint work with Sivaraman Balakrishnan and Aarti Singh

BACKGROUND

- ► Optimization: $\min_{x \in \mathcal{X}} f(x)$
- ► Classical setting (first-order):
 - * f is known (e.g., a likelihood function or an NN objective)
 - * $\nabla f(x)$ can be evaluated, or unbiasedly approximated.
- Zeroth-order setting:
 - * f is unknown, or very complicated.
 - * $\nabla f(x)$ is unknown, or very difficult to evaluate.
 - * f(x) can be evaluated, or unbiasedly approximated.

BACKGROUND

- ► Hyper-parameter tuning
 - * f maps hyper-parameter θ to system performance r
 - * f is essentially unknown
- ► Experimental design
 - * f maps experimental setting (pressure, temperature, etc.) to synthesized material quality.
- Communication-efficient optimization
 - * Data defining the objective scattered throughout machines
 - * Communicating abla f(x) is expensive, but f(x) is ok.

- ► Compact domain $\mathcal{X} = [0, 1]^d$
- ► Objective function $f : \mathcal{X} \to \mathbb{R}$
 - * f belongs to the Holder class of order α
 - * f may be non-convex
- ► Query model: adaptive $x_1, x_2, \cdots, x_n \in \mathcal{X}$ * $y_t = f(x_t) + \xi_t \quad \xi_t \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$
- ► Goal: minimize $f(\hat{x}_n) \inf_{x \in \mathcal{X}} f(x)$

- ► Compact domain $\mathcal{X} = [0, 1]^d$
- ► Objective function $f : \mathcal{X} \to \mathbb{R}$
 - * f belongs to the Holder class of order α
 - * f may be non-convex
- ► Query model: adaptive $x_1, x_2, \cdots, x_n \in \mathcal{X}$ * $y_t = f(x_t) + \xi_t \quad \xi_t \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$
- ► Goal: minimize $f(\hat{x}_n) \inf_{x \in \mathcal{X}} f(x)$

- ► Compact domain $\mathcal{X} = [0, 1]^d$
- ► Objective function $f : \mathcal{X} \to \mathbb{R}$
 - * f belongs to the Holder class of order α
- $\|f^{(\alpha)}\|_{\infty} \le M$

- * f may be non-convex
- ► Query model: adaptive $x_1, x_2, \cdots, x_n \in \mathcal{X}$ * $y_t = f(x_t) + \xi_t \quad \xi_t \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$
- ► Goal: minimize $f(\hat{x}_n) \inf_{x \in \mathcal{X}} f(x)$

- ► Compact domain $\mathcal{X} = [0, 1]^d$
- ► Objective function $f : \mathcal{X} \to \mathbb{R}$
 - * f belongs to the Holder class of order $\alpha \longrightarrow \|f^{(\alpha)}\|_{\infty} \leq M$
 - * f may be non-convex
- ► Query model: adaptive $x_1, x_2, \cdots, x_n \in \mathcal{X}$ * $y_t = f(x_t) + \xi_t \quad \xi_t \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$
- ► Goal: minimize $f(\hat{x}_n) \inf_{x \in \mathcal{X}} f(x)$

- ► Compact domain $\mathcal{X} = [0, 1]^d$
- ► Objective function $f : \mathcal{X} \to \mathbb{R}$
 - * f belongs to the Holder class of order $\alpha \longrightarrow \|f^{(\alpha)}\|_{\infty} \leq M$
 - * f may be non-convex
- ► Query model: adaptive $x_1, x_2, \cdots, x_n \in \mathcal{X}$ * $y_t = f(x_t) + \xi_t \quad \xi_t \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$
- ► Goal: minimize $f(\hat{x}_n) \inf_{x \in \mathcal{X}} f(x)$

- ► Compact domain $\mathcal{X} = [0, 1]^d$
- ► Objective function $f : \mathcal{X} \to \mathbb{R}$

* f belongs to the Holder class of order $\alpha \longrightarrow \|f^{(\alpha)}\|_{\infty} \leq M$

- * f may be non-convex
- ► Query model: adaptive $x_1, x_2, \cdots, x_n \in \mathcal{X}$ * $y_t = f(x_t) + \xi_t \quad \xi_t \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$

► Goal: minimize $f(\hat{x}_n) - \inf_{x \in \mathcal{X}} f(x)$

 $=: \mathfrak{L}(\hat{x}_n; f)$

- ► Compact domain $\mathcal{X} = [0, 1]^d$
- ► Objective function $f : \mathcal{X} \to \mathbb{R}$
 - * f belongs to the Holder class of order $\alpha \longrightarrow \|f^{(\alpha)}\|_{\infty} \leq M$
 - * f may be non-convex
- ► Query model: adaptive $x_1, x_2, \cdots, x_n \in \mathcal{X}$ * $y_t = f(x_t) + \xi_t \quad \xi_t \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$
- ► Goal: minimize $f(\hat{x}_n) \inf_{x \in \mathcal{X}} f(x)$ =: $\mathfrak{L}(\hat{x}_n; f)$

Uniform sampling + nonparametric reconstruction



Uniform sampling + nonparametric reconstruction



Uniform sampling + nonparametric reconstruction



- Uniform sampling + nonparametric reconstruction
 - * Classical Non-parametric analysis

$$\|\widehat{f}_n - f\|_{\infty} = \widetilde{O}_P\left(n^{-\alpha/(2\alpha+d)}\right)$$

- * Implies optimization error: $f(\hat{x}_n) f^* \leq 2 \|\hat{f}_n f\|_{\infty}$
- ► Can we do better?
- ► NO!

$$\inf_{\hat{x}_n} \sup_{f \in \Sigma^{\alpha}(M)} \mathbb{E}_f \left[\mathfrak{L}(\hat{x}_n; f) \right] \gtrsim n^{-\alpha/(2\alpha+d)}$$

- Uniform sampling + nonparametric reconstruction
 - * Classical Non-parametric analysis

$$\|\widehat{f}_n - f\|_{\infty} = \widetilde{O}_P\left(n^{-\alpha/(2\alpha+d)}\right)$$

* Implies optimization error: $f(\hat{x}_n) - f^* \leq 2 \|\hat{f}_n - f\|_{\infty}$

► Can we do better? No! Intuitions:

$$h_n \sim n^{-1/(2\alpha+d)}$$

$$h_n^{\alpha} \sim n^{-\alpha/(2\alpha+d)}$$

LOCAL RESULTS

- Characterize error for functions "near" a reference function f₀
- ► What is the error rate for f close to f_0 that is ...
 - * a constant function?
 - * strongly convex?
 - * has regular level sets?

* ...

Can an algorithm achieve instance-optimal error, without knowing f₀?

NOTATIONS

► Some definitions

- * Level set: $L_f(\epsilon) := \{x \in \mathcal{X} : f(x) \le f^* + \epsilon\}$
- * Distribution function: $\mu_f(\epsilon) := \operatorname{vol}(L_f(\epsilon))$



- Some definitions
 - * Level set: $L_f(\epsilon) := \{x \in \mathcal{X} : f(x) \le f^* + \epsilon\}$
 - * Distribution function: $\mu_f(\epsilon) := \operatorname{vol}(L_f(\epsilon))$
- ► Regularity condition (A1):
 - * # of δ -radius balls needed to cover $L_f(\epsilon) \simeq 1 + \mu_f(\epsilon)/\delta^d$



Regular level-set

- Some definitions
 - * Level set: $L_f(\epsilon) := \{x \in \mathcal{X} : f(x) \le f^* + \epsilon\}$
 - * Distribution function: $\mu_f(\epsilon) := \operatorname{vol}(L_f(\epsilon))$
- ► Regularity condition (A1):
 - * # of δ -radius balls needed to cover $L_f(\epsilon) \simeq 1 + \mu_f(\epsilon)/\delta^d$



Irregular level-set

- Some definitions
 - * Level set: $L_f(\epsilon) := \{x \in \mathcal{X} : f(x) \le f^* + \epsilon\}$
 - * Distribution function: $\mu_f(\epsilon) := \operatorname{vol}(L_f(\epsilon))$
- ► Regularity condition (A1):
 - * # of δ -radius balls needed to cover $L_f(\epsilon) \simeq 1 + \mu_f(\epsilon)/\delta^d$



Irregular level-set

Some definitions

- * Level set: $L_f(\epsilon) := \{x \in \mathcal{X} : f(x) \le f^* + \epsilon\}$
- * Distribution function: $\mu_f(\epsilon) := \operatorname{vol}(L_f(\epsilon))$
- ► Regularity condition (A2): * $\mu_f(\epsilon \log n) \le \mu_f(\epsilon) \times O(\log^{\gamma} n)$



- Some definitions
 - * Level set: $L_f(\epsilon) := \{x \in \mathcal{X} : f(x) \le f^* + \epsilon\}$
 - * Distribution function: $\mu_f(\epsilon) := \operatorname{vol}(L_f(\epsilon))$
- ► Regularity condition (A2): * $\mu_f(\epsilon \log n) \le \mu_f(\epsilon) \times O(\log^{\gamma} n)$



► Main result on local upper bound:

THEOREM 1. Suppose regularity conditions hold. There exists an algorithm such that for sufficiently large n,

$$\sup_{f \in \Sigma^{\alpha}(M)} \Pr_{f} \left[\mathfrak{L}(\hat{x}_{n}; f) \geq C \varepsilon_{n}(f) \log^{c} n \right] \leq 1/4$$

where $\varepsilon_{n}(f) := \sup \left\{ \epsilon > 0 : \epsilon^{-(2+d/\alpha)} \mu_{f}(\epsilon) \geq n \right\}$

► Main result on local upper bound:

THEOREM 1. Suppose regularity conditions hold. There exists an algorithm such that for sufficiently large n,

$$\sup_{f \in \Sigma^{\alpha}(M)} \Pr_{f} \left[\mathfrak{L}(\hat{x}_{n}; f) \geq C \varepsilon_{n}(f) \log^{c} n \right] \leq 1/4$$

where $\varepsilon_{n}(f) := \sup \left\{ \epsilon > 0 : \epsilon^{-(2+d/\alpha)} \mu_{f}(\epsilon) \geq n \right\}$

► Main result on local upper bound:

THEOREM 1. Suppose regularity conditions hold. There exists an algorithm such that for sufficiently large *n*,

$$\sup_{f \in \Sigma^{\alpha}(M)} \Pr_{f} \left[\mathfrak{L}(\hat{x}_{n}; f) \geq C \varepsilon_{n}(f) \log^{c} n \right] \leq 1/4$$

where $\varepsilon_{n}(f) := \sup \left\{ \epsilon > 0 : \epsilon^{-(2+d/\alpha)} \mu_{f}(\epsilon) \geq n \right\}$

► Main result on local upper bound:

THEOREM 1. Suppose regularity conditions hold. There exists an algorithm such that for sufficiently large n,

$$\sup_{f \in \Sigma^{\alpha}(M)} \Pr \left[\mathfrak{L}(\hat{x}_{n}; f) \geq C \varepsilon_{n}(f) \log^{c} n \right] \leq 1/4$$

where $\varepsilon_{n}(f) := \sup \left\{ \epsilon > 0 : \epsilon^{-(2+d/\alpha)} \mu_{f}(\epsilon) \geq n \right\}$

Adaptivity:

The algo does not know *f*.

► Main result on local upper bound:

THEOREM 1. Suppose regularity conditions hold. There exists an algorithm such that for sufficiently large n,

$$\sup_{f \in \Sigma^{\alpha}(M)} \Pr \left[\mathfrak{L}(\hat{x}_{n}; f) \geq C \varepsilon_{n}(f) \log^{c} n \right] \leq 1/4$$

where $\varepsilon_{n}(f) := \sup \left\{ \epsilon > 0 : \epsilon^{-(2+d/\alpha)} \mu_{f}(\epsilon) \geq n \right\}$

Adaptivity:

The algo does not know *f*.

► Main result on local upper bound:

THEOREM 1. Suppose regularity conditions hold. There exists an algorithm such that for sufficiently large n, $\sup_{f \in \Sigma^{\alpha}(M)} \Pr_{f} \left[\mathfrak{L}(\hat{x}_{n}; f) \geq C \varepsilon_{n}(f) \log^{c} n \right] \leq 1/4$

where
$$\varepsilon_n(f) := \sup\left\{\epsilon > 0 : e^{-(2+d/\alpha)}\mu_f(\epsilon) \ge n\right\}$$

Adaptivity:

The algo does not know *f*.

► Main result on local upper bound:

THEOREM 1. Suppose regularity conditions hold. There
exists an algorithm such that for sufficiently large n, $\sup_{f \in \Sigma^{\alpha}(M)} \Pr_{f} [\mathfrak{L}(\hat{x}_{n}; f) \geq C \varepsilon_{n}(f) \log^{c} n] \leq 1/4$ where $\varepsilon_{n}(f) := \sup \left\{ \epsilon > 0 : e^{-(2+d/\alpha)} \mu_{f}(\epsilon) \geq n \right\}$ Adaptivity:Instance dependent:

The algo does not know *f*.

Error rate depends on f

► Main result on local upper bound:

THEOREM 1. Suppose regularity conditions hold. There exists an algorithm such that for sufficiently large n,

$$\sup_{f \in \Sigma^{\alpha}(M)} \Pr_{f} \left[\mathfrak{L}(\hat{x}_{n}; f) \geq C \varepsilon_{n}(f) \log^{c} n \right] \leq 1/4$$

where $\varepsilon_{n}(f) := \sup \left\{ \epsilon > 0 : \epsilon^{-(2+d/\alpha)} \mu_{f}(\epsilon) \geq n \right\}$

► Example 1: polynomial growth $\mu_f(\epsilon) \asymp \epsilon^\beta, \beta \ge 0$

$$\varepsilon_n(f) \asymp n^{-\alpha/(2\alpha + d - \alpha\beta)}$$

Much faster than the "baseline" rate $n^{-\alpha/(2\alpha+d)}$

► Main result on local upper bound:

THEOREM 1. Suppose regularity conditions hold. There exists an algorithm such that for sufficiently large n,

$$\sup_{f \in \Sigma^{\alpha}(M)} \Pr_{f} \left[\mathfrak{L}(\hat{x}_{n}; f) \geq C \varepsilon_{n}(f) \log^{c} n \right] \leq 1/4$$

where $\varepsilon_{n}(f) := \sup \left\{ \epsilon > 0 : \epsilon^{-(2+d/\alpha)} \mu_{f}(\epsilon) \geq n \right\}$

► Example 2: constant function $f \equiv 0$

$$\varepsilon_n(f) \asymp n^{-\alpha/(2\alpha+d)}$$

This is the worst case function, matching global rate $n^{-\alpha/(2\alpha+d)}$

► Main result on local upper bound:

THEOREM 1. Suppose regularity conditions hold. There exists an algorithm such that for sufficiently large n,

$$\sup_{f \in \Sigma^{\alpha}(M)} \Pr_{f} \left[\mathfrak{L}(\hat{x}_{n}; f) \geq C \varepsilon_{n}(f) \log^{c} n \right] \leq 1/4$$

where $\varepsilon_{n}(f) := \sup \left\{ \epsilon > 0 : \epsilon^{-(2+d/\alpha)} \mu_{f}(\epsilon) \geq n \right\}$

► Example 3: strongly convex $f: \mu_f(\epsilon) \asymp \epsilon^{d/2}$

$$\varepsilon_n(f) \asymp n^{-1/2}$$

Match the classical zeroth-order convex rate $n^{-1/2}$ up to log terms

► Main result on local lower bound:

```
THEOREM 2. Suppose f_0 satisfies regularity conditions.

Then we have

\inf_{\hat{x}_n} \sup_{\substack{f \in \Sigma^{\alpha}(M), \\ \|f - f_0\|_{\infty} \lesssim \varepsilon_n(f_0)}} \mathbb{E}_f \left[ \mathfrak{L}(\hat{x}_n; f) \right] \gtrsim \varepsilon_n(f_0)

where \varepsilon_n(f) := \sup \left\{ \epsilon > 0 : \epsilon^{-(2+d/\alpha)} \mu_f(\epsilon) \ge n \right\}
```

full knowledge of reference:

The algo. knows f_0

► Main result on local lower bound:

THEOREM 2. Suppose f_0 satisfies regularity conditions. Then we have $\inf_{\hat{x}_n} \sup_{\substack{f \in \Sigma^{\alpha}(M), \\ \|f - f_0\|_{\infty} \lesssim \varepsilon_n(f_0)}} \mathbb{E}_f \left[\mathfrak{L}(\hat{x}_n; f) \right] \gtrsim \varepsilon_n(f_0)$

where $\varepsilon_n(f) := \sup\left\{\epsilon > 0 : \epsilon^{-(2+d/\alpha)} \mu_f(\epsilon) \ge n\right\}$

full knowledge of reference:

The algo. knows f_0

► Main result on local lower bound:

THEOREM 2. Suppose f_0 satisfies regularity conditions. Then we have

$$\inf_{\substack{\hat{x}_n \\ \|f-f_0\|_{\infty} \lesssim \varepsilon_n(f_0)}} \mathbb{E}_f \left[\mathfrak{L}(\hat{x}_n; f) \right] \gtrsim \varepsilon_n(f_0)$$
where $\varepsilon_n(f) := \sup \left\{ \epsilon > 0 : \epsilon^{-(2+d/\alpha)} \mu_f(\epsilon) \ge n \right\}$

full knowledge of reference:

The algo. knows f_0

► Main result on local lower bound:

THEOREM 2. Suppose f_0 satisfies regularity conditions. Then we have

$$\inf_{\substack{\hat{x}_n \\ \|f-f_0\|_{\infty} \lesssim \varepsilon_n(f_0)}} \mathbb{E}_f \left[\mathfrak{L}(\hat{x}_n; f) \right] \gtrsim \varepsilon_n(f_0)$$
where $\varepsilon_n(f) := \sup \left\{ \epsilon > 0 : \epsilon^{-(2+d/\alpha)} \mu_f(\epsilon) \ge n \right\}$

full knowledge of reference:

The algo. knows f_0

► Main result on local lower bound:

THEOREM 2. Suppose f_0 satisfies regularity conditions. Then we have

$$\inf_{\substack{\hat{x}_n \\ \|f \in \Sigma^{\alpha}(M), \\ \|f - f_0\|_{\infty} \leq \varepsilon_n(f_0)}} \mathbb{E}_f \left[\mathfrak{L}(\hat{x}_n; f) \right] \gtrsim \varepsilon_n(f_0)$$
where $\varepsilon_n(f) := \sup \left\{ \epsilon > 0 : \epsilon^{-(2+d/\alpha)} \mu_f(\epsilon) \ge n \right\}$

full knowledge of reference:

The algo. knows f_0

► Main result on local lower bound:

THEOREM 2. Suppose f_0 satisfies regularity conditions. Then we have

 $\inf_{\hat{x}_n} \sup_{\substack{f \in \Sigma^{\alpha}(M), \\ \|f - f_0\|_{\infty} \lesssim \varepsilon_n(f_0)}} \mathbb{E}_f \left[\mathfrak{L}(\hat{x}_n; f) \right] \gtrsim \varepsilon_n(f_0)$ where $\varepsilon_n(f) \coloneqq \sup \left\{ \epsilon > 0 : \epsilon^{-(2+d/\alpha)} \mu_f(\epsilon) \ge n \right\}$

full knowledge of reference:

The algo. knows f_0

► Main result on local lower bound:

THEOREM 2. Suppose f_0 satisfies regularity conditions. Then we have

 $\inf_{\substack{\hat{x}_n \\ \|f \in \Sigma^{\alpha}(M), \\ \|f - f_0\|_{\infty} \leq \varepsilon_n(f_0)}} \mathbb{E}_f \left[\mathfrak{L}(\hat{x}_n; f) \right] \gtrsim \varepsilon_n(f_0)$ where $\varepsilon_n(f) \coloneqq \sup \left\{ \epsilon > 0 : \epsilon^{-(2+d/\alpha)} \mu_f(\epsilon) \ge n \right\}$

Local minimality:

estimation error of f close to the reference f_0

- ► Our algorithm: "successive rejection"
 - * Step I: uniform sampling and build confidence intervals (CI)
 - * Step 2: remove sub-optimal points
 - * Step 3: uniform sample in the remaining points.

- ► Our algorithm: "successive rejection"
 - * Step I: uniform sampling and build confidence intervals (CI)
 - * Step 2: remove sub-optimal points
 - * Step 3: uniform sample in the remaining points.



- ► Our algorithm: "successive rejection"
 - * Step I: uniform sampling and build confidence intervals (CI)
 - * Step 2: remove sub-optimal points
 - * Step 3: uniform sample in the remaining points.



- ► Our algorithm: "successive rejection"
 - * Step 1: uniform sampling and build confidence intervals (CI)
 - * Step 2: remove sub-optimal points
 - * Step 3: uniform sample in the remaining points.



- ► Our algorithm: "successive rejection"
 - * Step 1: uniform sampling and build confidence intervals (CI)
 - * Step 2: remove sub-optimal points
 - * Step 3: uniform sample in the remaining points.
- ► Key observation between iterations:

$$S_{\tau} \subseteq L_f(\varepsilon) \Longrightarrow S_{\tau+1} \subseteq L_f(\varepsilon/2)$$

Until $\varepsilon \sim \varepsilon_n(f) \times \log^c n$

► An $O(\log n)$ number of iterations suffice.

► Step 1: constructing "packings" on $L_{f_0}(\epsilon_n)$



 $h_n \asymp \epsilon_n^{1/\alpha}$ Discrepancy in ball: $2\epsilon_n$

Must identify the ball w.h.p.

► Step 1: constructing "packings" on $L_{f_0}(\epsilon_n)$



Resembles

Bandit Pure Exploration

$$\mu_1, \mu_2, \cdots, \mu_H \in \mathbb{R}$$
$$\mu_i = 2\epsilon_n; \mu_{-i} = 0$$

Identify the non-zero arm









TAKE-HOME MESSAGES

- (Noisy) zeroth-order optimization of smooth functions is in general difficult
 - * As difficult as estimating the function in sup-norm.
- The optimal convergence rates exhibit significant gaps locally for different objective functions
 - * Local minimax rate mostly dictated by level set growth;
 - * The constant function is the hardest example;
 - * Strongly convex functions do **not** exhibit curse of dim.
- ► A successive-rejection type algorithm is near-optimal.

FUTURE DIRECTIONS

- Are the regularity conditions absolutely necessary?
 - * Can the level sets of f be irregular?
 - * Can the volumes of level sets of *f* grow heterogeneously?
- Are there more computationally efficient algorithms?
 - * Key challenge: avoiding creating sup-norm CIs explicitly.
- ► Log factors: are they removable? (conjecture: yes!)
 - * Active queries methods for nonparametric estimation / bandit pure exploration do not have log terms.

QUESTIONS