

# NON-CONVERGENCE OF PROBABILISTIC DIRECT SEARCH

ICNONLA, TAIYUAN, CHINA

---

Cunxin Huang, joint work with Zaikun Zhang

July 19, 2023

Department of Applied Mathematics  
The Hong Kong Polytechnic University

# DERIVATIVE-FREE OPTIMIZATION

## Derivative-free optimization

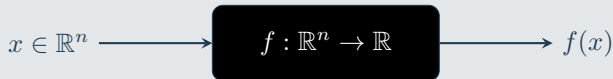
Derivative-free optimization (DFO) studies optimization algorithms that **do not use derivatives** (only use function evaluations).

# DERIVATIVE-FREE OPTIMIZATION

## Derivative-free optimization

Derivative-free optimization (DFO) studies optimization algorithms that **do not use derivatives** (only use function evaluations).

**Example** (Typical case:  $f$  is a blackbox)

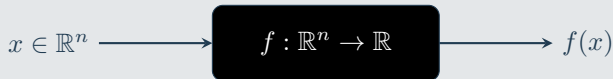


# DERIVATIVE-FREE OPTIMIZATION

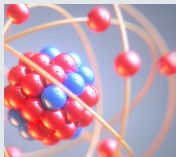
## Derivative-free optimization

Derivative-free optimization (DFO) studies optimization algorithms that **do not use derivatives** (only use function evaluations).

Example (Typical case:  $f$  is a blackbox)



## Applications



Nuclear Physics



Machine Learning



Cosmology

In this talk, we solve the **unconstrained** problem

$$\min_{x \in \mathbb{R}^n} f(x),$$

where we follow standard setup that

- $\nabla f$  is **Lipschitz continuous** with constant  $\nu$ , cannot be evaluated,
- $f$  is bounded below,
- the evaluation of  $f$  is expensive.

## Direct-search methods

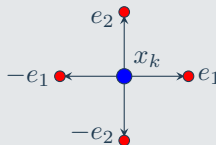
Sample points based on a [finite direction set](#) and make decisions by comparing values.

## Direct-search methods

Sample points based on a [finite direction set](#) and make decisions by comparing values.

### Example (typical direction set in $\mathbb{R}^2$ )

$$D = \{e_1, -e_1, e_2, -e_2\}$$

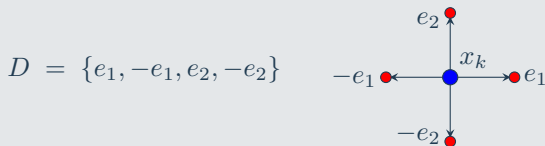


# DIRECT-SEARCH METHODS

## Direct-search methods

Sample points based on a [finite direction set](#) and make decisions by comparing values.

### Example (typical direction set in $\mathbb{R}^2$ )



In this talk, we assume direction set always be a set of unit vectors in  $\mathbb{R}^n$ .



---

**Algorithm 1:** Direct Search with sufficient decrease

---

---

**Algorithm 1:** Direct Search with sufficient decrease

---

**Input:**  $x_0 \in \mathbb{R}^n$ ,  $\alpha_0 \in (0, \infty)$ ,  $0 < \theta < 1 \leq \gamma$ .

---

**Algorithm 1:** Direct Search with sufficient decrease

---

**Input:**  $x_0 \in \mathbb{R}^n$ ,  $\alpha_0 \in (0, \infty)$ ,  $0 < \theta < 1 \leq \gamma$ .

**for**  $k = 0, 1, \dots$  **do**

---

**Algorithm 1:** Direct Search with sufficient decrease

---

**Input:**  $x_0 \in \mathbb{R}^n$ ,  $\alpha_0 \in (0, \infty)$ ,  $0 < \theta < 1 \leq \gamma$ .

**for**  $k = 0, 1, \dots$  **do**

    Select  $D_k \subset \mathbb{R}^n$ .



---

**Algorithm 1:** Direct Search with sufficient decrease

---

**Input:**  $x_0 \in \mathbb{R}^n$ ,  $\alpha_0 \in (0, \infty)$ ,  $0 < \theta < 1 \leq \gamma$ .

**for**  $k = 0, 1, \dots$  **do**

    Select  $D_k \subset \mathbb{R}^n$ .

**if**  $f(x_k + \alpha_k d) < f(x_k) - \alpha_k^2$  for some  $d \in \mathcal{D}_k$  **then**

---

**Algorithm 1:** Direct Search with sufficient decrease

---

**Input:**  $x_0 \in \mathbb{R}^n$ ,  $\alpha_0 \in (0, \infty)$ ,  $0 < \theta < 1 \leq \gamma$ .

**for**  $k = 0, 1, \dots$  **do**

    Select  $D_k \subset \mathbb{R}^n$ .

**if**  $f(x_k + \alpha_k d) < f(x_k) - \alpha_k^2$  for some  $d \in \mathcal{D}_k$  **then**

        Expand step size, and move to that point

---

## Algorithm 1: Direct Search with sufficient decrease

---

**Input:**  $x_0 \in \mathbb{R}^n$ ,  $\alpha_0 \in (0, \infty)$ ,  $0 < \theta < 1 \leq \gamma$ .

**for**  $k = 0, 1, \dots$  **do**

    Select  $D_k \subset \mathbb{R}^n$ .

**if**  $f(x_k + \alpha_k d) < f(x_k) - \alpha_k^2$  for some  $d \in \mathcal{D}_k$  **then**

        Expand step size, and move to that point

        Set  $\alpha_{k+1} = \gamma \alpha_k$  and  $x_{k+1} = x_k + \alpha_k d$ .

**else**

---

## Algorithm 1: Direct Search with sufficient decrease

---

**Input:**  $x_0 \in \mathbb{R}^n$ ,  $\alpha_0 \in (0, \infty)$ ,  $0 < \theta < 1 \leq \gamma$ .

**for**  $k = 0, 1, \dots$  **do**

    Select  $D_k \subset \mathbb{R}^n$ .

**if**  $f(x_k + \alpha_k d) < f(x_k) - \alpha_k^2$  for some  $d \in \mathcal{D}_k$  **then**

        Expand step size, and move to that point

        Set  $\alpha_{k+1} = \gamma \alpha_k$  and  $x_{k+1} = x_k + \alpha_k d$ .

**else**

        Shrink step size, and stand still



---

## Algorithm 1: Direct Search with sufficient decrease

---

**Input:**  $x_0 \in \mathbb{R}^n$ ,  $\alpha_0 \in (0, \infty)$ ,  $0 < \theta < 1 \leq \gamma$ .

**for**  $k = 0, 1, \dots$  **do**

    Select  $D_k \subset \mathbb{R}^n$ .

**if**  $f(x_k + \alpha_k d) < f(x_k) - \alpha_k^2$  for some  $d \in \mathcal{D}_k$  **then**

        Expand step size, and move to that point

        Set  $\alpha_{k+1} = \gamma \alpha_k$  and  $x_{k+1} = x_k + \alpha_k d$ .

**else**

        Shrink step size, and stand still

        Set  $\alpha_{k+1} = \theta \alpha_k$  and  $x_{k+1} = x_k$ .

---

---

## Algorithm 1: Direct Search with sufficient decrease

---

**Input:**  $x_0 \in \mathbb{R}^n$ ,  $\alpha_0 \in (0, \infty)$ ,  $0 < \theta < 1 \leq \gamma$ .

**for**  $k = 0, 1, \dots$  **do**

    Select  $D_k \subset \mathbb{R}^n$ .

**if**  $f(x_k + \alpha_k d) < f(x_k) - \alpha_k^2$  for some  $d \in \mathcal{D}_k$  **then**

        Expand step size, and move to that point

        Set  $\alpha_{k+1} = \gamma \alpha_k$  and  $x_{k+1} = x_k + \alpha_k d$ .

**else**

        Shrink step size, and stand still

        Set  $\alpha_{k+1} = \theta \alpha_k$  and  $x_{k+1} = x_k$ .

---

Everything almost clear, except “select  $D_k$ ”.

## “GOOD” DIRECTION SET

Question 1: what is a “good”  $D_k$ ?

## “GOOD” DIRECTION SET

Question 1: what is a “good”  $D_k$ ?

Include at least one **descent direction**, i.e. exists a  $d \in D_k$  s.t.,  $-g_k^\top d > 0$ .

## “GOOD” DIRECTION SET

Question 1: what is a “good”  $D_k$ ?

Include at least one **descent direction**, i.e. exists a  $d \in D_k$  s.t.,  $-g_k^\top d > 0$ .

Question 2: how to choose a “good”  $D_k$  when we don't know  $g_k$ ?

## “GOOD” DIRECTION SET

Question 1: what is a “good”  $D_k$ ?

Include at least one **descent direction**, i.e. exists a  $d \in D_k$  s.t.,  $-g_k^\top d > 0$ .

Question 2: how to choose a “good”  $D_k$  when we don't know  $g_k$ ?

Natural idea: choose a  $D_k$  s.t., **for any  $v \in \mathbb{R}^n$ , there exists at least a  $d \in D_k$  satisfying  $d^\top v > 0$ .**

## “GOOD” DIRECTION SET

Question 1: what is a “good”  $D_k$ ?

Include at least one **descent direction**, i.e. exists a  $d \in D_k$  s.t.,  $-g_k^\top d > 0$ .

Question 2: how to choose a “good”  $D_k$  when we don't know  $g_k$ ?

Natural idea: choose a  $D_k$  s.t., **for any  $v \in \mathbb{R}^n$ , there exists at least a  $d \in D_k$  satisfying  $d^\top v > 0$ .**

**Positive Spanning Set (PSS)**

## “GOOD” DIRECTION SET

Question 1: what is a “good”  $D_k$ ?

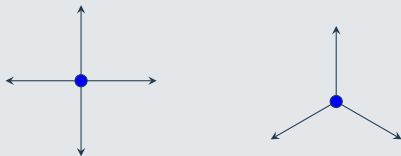
Include at least one **descent direction**, i.e. exists a  $d \in D_k$  s.t.,  $-g_k^\top d > 0$ .

Question 2: how to choose a “good”  $D_k$  when we don’t know  $g_k$ ?

Natural idea: choose a  $D_k$  s.t., **for any  $v \in \mathbb{R}^n$ , there exists at least a  $d \in D_k$  satisfying  $d^\top v > 0$ .**

### Positive Spanning Set (PSS)

Example (typical PSS in  $\mathbb{R}^2$ )





## “GOOD” DIRECTION SET

Question 3: how “good” is a PSS?

# “GOOD” DIRECTION SET

Question 3: how “good” is a PSS?

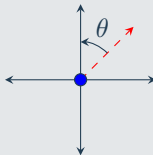
## Cosine measure

Define the cosine measure,  $\text{cm}(D)$ , for a finite set of nonzero vectors,  $D \in \mathbb{R}^n$ , as

$$\text{cm}(D) = \min_{v \in \mathbb{R}^n \setminus \{0\}} \text{cm}(D, v) = \min_{v \in \mathbb{R}^n \setminus \{0\}} \max_{d \in D} \frac{d^\top v}{\|d\| \|v\|}.$$

## Example

$$\text{cm}(D) = \cos \theta = \frac{\sqrt{2}}{2}$$



## Theorem

*Under standard assumptions, if in Algorithm 1, there exists a positive constant  $\kappa$  such that*

$$\text{cm}(D_k) \geq \kappa \quad \text{for all } k \geq 0,$$

*then we have*

$$\liminf_{k \rightarrow \infty} \|g_k\| = 0.$$

## Theorem

*Under standard assumptions, if in Algorithm 1, there exists a positive constant  $\kappa$  such that*

$$\text{cm}(D_k) \geq \kappa \quad \text{for all } k \geq 0,$$

*then we have*

$$\liminf_{k \rightarrow \infty} \|g_k\| = 0.$$

Each direction set  $D_k$  should be a “good enough” PSS.



Bad news: if  $D$  is a PSS in  $\mathbb{R}^n$ , then  $|D| \geq n + 1$ .

Complexity:  $\mathcal{O}(n^2 \varepsilon^{-2})$

## RANDOMIZED METHOD

Bad news: if  $D$  is a PSS in  $\mathbb{R}^n$ , then  $|D| \geq n + 1$ .

Complexity:  $\mathcal{O}(n^2 \varepsilon^{-2})$

Goal: use randomized methods to reduce  $|D|$  from  $\mathcal{O}(n)$  to  $\mathcal{O}(1)$ .

## RANDOMIZED METHOD

Bad news: if  $D$  is a PSS in  $\mathbb{R}^n$ , then  $|D| \geq n + 1$ .

Complexity:  $\mathcal{O}(n^2 \varepsilon^{-2})$

Goal: use randomized methods to reduce  $|D|$  from  $\mathcal{O}(n)$  to  $\mathcal{O}(1)$ .

---

### Algorithm 2: Probabilistic Direct Search

---

**Input:**  $x_0 \in \mathbb{R}^n$ ,  $\alpha_0 \in (0, \infty)$ ,  $0 < \theta < 1 \leq \gamma$ .

**for**  $k = 0, 1, \dots$  **do**

    Randomly generate  $\mathfrak{D}_k \subset \mathbb{R}^n$ .

**if**  $f(x_k + \alpha_k d) < f(x_k) - \rho(\alpha_k)$  for some  $d \in \mathfrak{D}_k$  **then**

        Set  $\alpha_{k+1} = \gamma \alpha_k$  and  $x_{k+1} = x_k + \alpha_k d$ .

**else**

        Set  $\alpha_{k+1} = \theta \alpha_k$  and  $x_{k+1} = x_k$ .

---



## RANDOMIZED METHOD

Bad news: if  $D$  is a PSS in  $\mathbb{R}^n$ , then  $|D| \geq n + 1$ .

Complexity:  $\mathcal{O}(n^2 \varepsilon^{-2})$

Goal: use randomized methods to reduce  $|D|$  from  $\mathcal{O}(n)$  to  $\mathcal{O}(1)$ .

---

### Algorithm 2: Probabilistic Direct Search

---

**Input:**  $x_0 \in \mathbb{R}^n$ ,  $\alpha_0 \in (0, \infty)$ ,  $0 < \theta < 1 \leq \gamma$ .

**for**  $k = 0, 1, \dots$  **do**

    Randomly generate  $\mathfrak{D}_k \subset \mathbb{R}^n$ .

**if**  $f(x_k + \alpha_k d) < f(x_k) - \rho(\alpha_k)$  for some  $d \in \mathfrak{D}_k$  **then**

        Set  $\alpha_{k+1} = \gamma \alpha_k$  and  $x_{k+1} = x_k + \alpha_k d$ .

**else**

        Set  $\alpha_{k+1} = \theta \alpha_k$  and  $x_{k+1} = x_k$ .

---

All the randomness comes from  $\mathfrak{D}_k$

## RANDOMIZED METHOD

Bad news: if  $D$  is a PSS in  $\mathbb{R}^n$ , then  $|D| \geq n + 1$ .

Complexity:  $\mathcal{O}(n^2 \varepsilon^{-2})$

Goal: use randomized methods to reduce  $|D|$  from  $\mathcal{O}(n)$  to  $\mathcal{O}(1)$ .

---

### Algorithm 2: Probabilistic Direct Search

---

**Input:**  $x_0 \in \mathbb{R}^n$ ,  $\alpha_0 \in (0, \infty)$ ,  $0 < \theta < 1 \leq \gamma$ .

**for**  $k = 0, 1, \dots$  **do**

    Randomly generate  $\mathfrak{D}_k \subset \mathbb{R}^n$ .

**if**  $f(x_k + \alpha_k d) < f(x_k) - \rho(\alpha_k)$  for some  $d \in \mathfrak{D}_k$  **then**

        Set  $\alpha_{k+1} = \gamma \alpha_k$  and  $x_{k+1} = x_k + \alpha_k d$ .

**else**

        Set  $\alpha_{k+1} = \theta \alpha_k$  and  $x_{k+1} = x_k$ .

---

All the randomness comes from  $\mathfrak{D}_k$

Notations for random variables or random vectors

$D_k \Rightarrow \mathfrak{D}_k$ ,  $d \Rightarrow \mathfrak{d}$ ,  $x_k \Rightarrow X_k$ ,  $\alpha_k \Rightarrow A_k$ ,  $g_k \Rightarrow G_k$



Actually, what we need is not  $\text{cm}(D_k) \geq \kappa$ , but  $\text{cm}(D_k, -G_k) \geq \kappa$ .

Actually, what we need is not  $\text{cm}(D_k) \geq \kappa$ , but  $\text{cm}(D_k, -G_k) \geq \kappa$ .

### Definition ( $p$ -probabilistically $\kappa$ -descent)

$(\mathfrak{D}_k)_{k \geq 0}$  is said to be  $p$ -probabilistically  $\kappa$ -descent, if

$$\mathbb{P}(\text{cm}(\mathfrak{D}_k, -G_k) \geq \kappa \mid \mathfrak{D}_0, \dots, \mathfrak{D}_{k-1}) \geq p \quad \text{for each } k \geq 0.$$

Intuition: high probability that  $\mathfrak{D}_k$  is “always good enough”,  
no matter what happened before.

## Theorem (Gratton et al. 2015)

*Under standard assumptions, if in Algorithm 2, there exists a **positive constant  $\kappa$**  such that  $(\mathfrak{D}_k)_{k \geq 0}$  is  $p$ -probabilistically  $\kappa$ -descent with  $p = \log \theta / \log(\gamma^{-1} \theta)$ , then we have*

$$\mathbb{P} \left( \liminf_{k \rightarrow \infty} \|G_k\| = 0 \right) = 1.$$

Complexity:  $\mathcal{O}(mn\varepsilon^{-2})$  with extremely high probability, where  $m \ll n$ .

1. Why and how of “non-convergence”:  
ideas and main results
2. Tightness of non-convergence: counterexample
3. Conclusion

WHY AND HOW OF  
“NON-CONVERGENCE”:  
IDEAS AND MAIN RESULTS

---



Recall “Converge or Diverge? A Story of Adam” by Prof. Luo

Recall “Converge or Diverge? A Story of Adam” by Prof. Luo

What will happen if  $(\mathfrak{D}_k)_{k \geq 0}$  fails to meet  $p$ -probabilistically  $\kappa$ -descent?  
Or worse  $\mathfrak{D}_k$  is always “bad”?

Can we obtain some kind of “non-convergence”?

### $p$ -probabilistically ascent

$(\mathfrak{D}_k)_{k \in \mathbb{N}}$  is said to be  $p$ -probabilistically ascent if it satisfies

$$\mathbb{P}(\text{cm}(\mathfrak{D}_k, -G_k) \leq 0 \mid \mathfrak{D}_0, \dots, \mathfrak{D}_{k-1}) \geq p \quad \text{for each } k \geq 0.$$

### Observation

For any convex function, if  $\text{cm}(D_k, -g_k) \leq 0$ , then no descent direction in  $D_k$ , leading to shrinking of step size.

### $p$ -probabilistically ascent

$(\mathfrak{D}_k)_{k \in \mathbb{N}}$  is said to be  $p$ -probabilistically ascent if it satisfies

$$\mathbb{P}(\text{cm}(\mathfrak{D}_k, -G_k) \leq 0 \mid \mathfrak{D}_0, \dots, \mathfrak{D}_{k-1}) \geq p \quad \text{for each } k \geq 0.$$

### Observation

For any convex function, if  $\text{cm}(D_k, -g_k) \leq 0$ , then no descent direction in  $D_k$ , leading to shrinking of step size.

If we show non-convergence for convex functions, then we cannot expect convergence for others.

## NAIVE IDEA OF NON-CONVERGENCE

$\mathfrak{D}_k$  always bad

## NAIVE IDEA OF NON-CONVERGENCE

$\mathfrak{D}_k$  always bad



$A_k$  always shrinks

## NAIVE IDEA OF NON-CONVERGENCE

$\mathfrak{D}_k$  always bad



$A_k$  always shrinks



$$\sum_{k=0}^{\infty} A_k < \infty \text{ a.s. ?}$$

## NAIVE IDEA OF NON-CONVERGENCE

$\mathfrak{D}_k$  always bad



$A_k$  always shrinks



$$\sum_{k=0}^{\infty} A_k < \infty \text{ a.s. ?}$$



$$x_o \notin \bar{\mathcal{B}}(x^*, \sum_{k=0}^{\infty} A_k)$$



## NAIVE IDEA OF NON-CONVERGENCE

$\mathfrak{D}_k$  always bad



$A_k$  always shrinks



$$\sum_{k=0}^{\infty} A_k < \infty \text{ a.s. ?}$$



$$x_o \notin \bar{\mathcal{B}}(x^*, \sum_{k=0}^{\infty} A_k)$$



Non-convergence:

$$\mathbb{P}(\liminf_k \text{dist}(X_k, x^*) = 0) < 1$$

- Define indicator function  $Y_k = \mathbb{1}_{\{\text{cm}(\mathfrak{D}_k, -G_k) > 0\}}$   
Indicator for “good” event

- Define indicator function  $Y_k = \mathbb{1}_{\{\text{cm}(\mathfrak{D}_k, -G_k) > 0\}}$   
Indicator for “good” event
- $A_{k+1} \leq \gamma^{Y_k} \theta^{1-Y_k} A_k$ , when  $f$  is convex

- Define indicator function  $Y_k = \mathbb{1}_{\{\text{cm}(\mathfrak{D}_k, -G_k) > 0\}}$   
Indicator for “good” event
- $A_{k+1} \leq \gamma^{Y_k} \theta^{1-Y_k} A_k$ , when  $f$  is convex
- $A_k \leq \alpha_0 \prod_{i=0}^{k-1} \gamma^{Y_i} \theta^{1-Y_i} =: \alpha_0 U_k$

- Define indicator function  $Y_k = \mathbb{1}_{\{\text{cm}(\mathfrak{D}_k, -G_k) > 0\}}$   
Indicator for “good” event
- $A_{k+1} \leq \gamma^{Y_k} \theta^{1-Y_k} A_k$ , when  $f$  is convex
- $A_k \leq \alpha_0 \prod_{i=0}^{k-1} \gamma^{Y_i} \theta^{1-Y_i} =: \alpha_0 U_k$
- $\sum_{k=1}^{\infty} A_k \leq \alpha_0 \sum_{k=1}^{\infty} U_k$

## Assumption

$$\mathbb{P}(Y_k = 0 \mid Y_0, \dots, Y_{k-1}) \geq p > p_0 \quad \text{for each } k \geq 0,$$

where  $p_0 = \log \gamma / \log(\theta^{-1} \gamma)$ .

## Assumption

$$\mathbb{P}(Y_k = 0 \mid Y_0, \dots, Y_{k-1}) \geq p > p_0 \quad \text{for each } k \geq 0,$$

where  $p_0 = \log \gamma / \log(\theta^{-1} \gamma)$ .

## Result

1.

$$\mathbb{P}\left(\sum_{k=1}^{\infty} U_k < \infty\right) = 1$$

2.

$$\mathbb{P}\left(\sum_{k=1}^{\infty} U_k < \zeta\right) > 0 \iff \zeta > \frac{\theta}{1 - \theta}$$

## Assumption

$$\mathbb{P}(Y_k = 0 \mid Y_0, \dots, Y_{k-1}) \geq p > p_0 \quad \text{for each } k \geq 0,$$

where  $p_0 = \log \gamma / \log(\theta^{-1} \gamma)$ .

## Result

1.

$$\mathbb{P}\left(\sum_{k=1}^{\infty} U_k < \infty\right) = 1$$

2.

$$\mathbb{P}\left(\sum_{k=1}^{\infty} U_k < \zeta\right) > 0 \iff \zeta > \frac{\theta}{1 - \theta}$$

## Observation

- Easy to achieve non-convergence:  $\theta/(1 - \theta)$  is fixed number.
- Cannot further improve:  $\sum_{k=1}^{\infty} U_k \geq \theta/(1 - \theta)$ .





Space of  $(\mathfrak{D}_k)_{k \geq 0}$

$\mathfrak{D}_k$  is measurable map from  $\Omega$  to  $2^{\mathbb{R}^n}$

Red: Convergence    Blue: Non-convergence



Space of  $(\mathfrak{D}_k)_{k \geq 0}$

$\mathfrak{D}_k$  is measurable map from  $\Omega$  to  $2^{\mathbb{R}^n}$

Red: Convergence    Blue: Non-convergence

Too much “shadow area”, not satisfied!

Question: what if we only need  $\sum_{k=1}^{\infty} U_k < \infty$ ?

Question: what if we only need  $\sum_{k=1}^{\infty} U_k < \infty$ ?  
Only need  $\mathfrak{D}_k$  to be “bad” eventually.

### Assumption

$$\liminf_{k \rightarrow \infty} \mathbb{P}(Y_k = 0 \mid Y_0, \dots, Y_{k-1}) > p_0.$$

Question: what if we only need  $\sum_{k=1}^{\infty} U_k < \infty$ ?  
Only need  $\mathfrak{D}_k$  to be “bad” eventually.

## Assumption

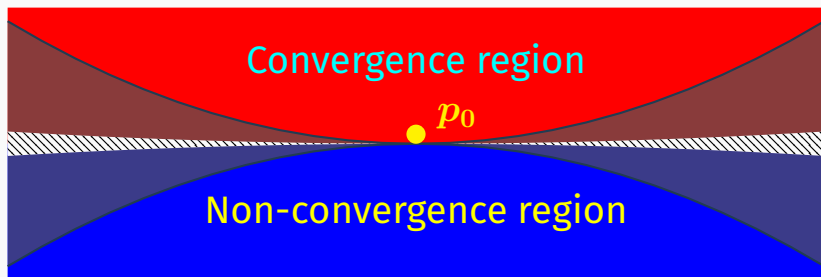
$$\liminf_{k \rightarrow \infty} \mathbb{P}(Y_k = 0 \mid Y_0, \dots, Y_{k-1}) > p_0.$$

## Result

$$\mathbb{P}\left(\sum_{k=1}^{\infty} U_k < \infty\right) = 1.$$

Also use “liminf”-type assumption in convergence theory

Also use “liminf”-type assumption in convergence theory



## TIGHTNESS OF NON-CONVERGENCE: COUNTEREXAMPLE

---



## CAN $p_0$ BE INCLUDED?

Natural question:

can  $p_0$  be included in the non-convergence region?

## CAN $p_0$ BE INCLUDED?

Natural question:

can  $p_0$  be included in the non-convergence region?

Answer: NO

## CAN $p_0$ BE INCLUDED?

Natural question:

can  $p_0$  be included in the non-convergence region?

Answer: NO

### Example

We assume

- $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be gradient Lipschitz and strongly convex,
- $\theta = 1/2$  and  $\gamma = 2$ ,
- $\mathbb{P}(\mathfrak{d}_k = G_k/\|G_k\|) = \mathbb{P}(\mathfrak{d}_k = -G_k/\|G_k\|) = 1/2$ ,

then we have

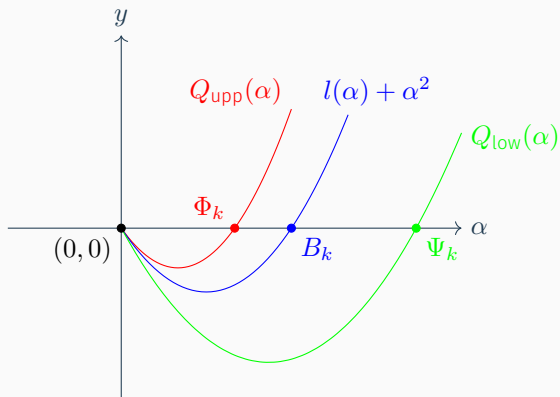
$$\mathbb{P} \left( \lim_{k \rightarrow \infty} \|G_k\| = 0 \right) = 1.$$

$$\bullet \quad l(\alpha) = f(X_k - \alpha G_k / \|G_k\|) - f(X_k)$$

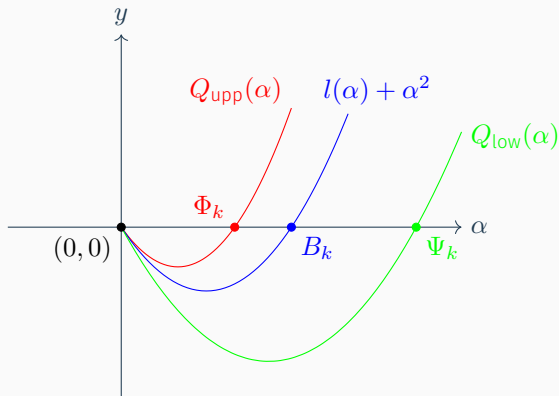
- $l(\alpha) = f(X_k - \alpha G_k / \|G_k\|) - f(X_k)$
- Sufficient decrease condition becomes  $l(\alpha) + \alpha^2 < 0$

- $l(\alpha) = f(X_k - \alpha G_k / \|G_k\|) - f(X_k)$
- Sufficient decrease condition becomes  $l(\alpha) + \alpha^2 < 0$
- $l(\alpha) + \alpha^2$  is both upper and lower bounded by quadratic functions

## SKETCH OF PROOF



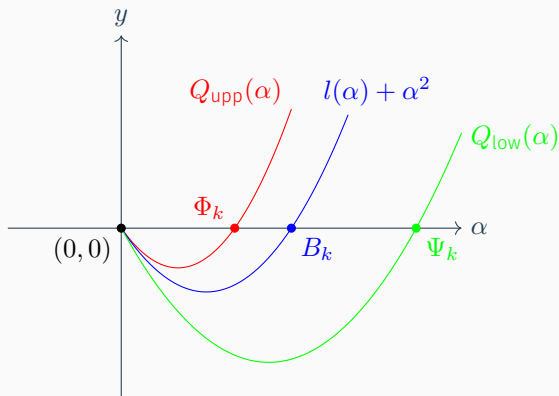
## SKETCH OF PROOF



$$\cdot C_1 \|G_k\| = \Phi_k \leq B_k \leq \Psi_k = C_2 \|G_k\|$$

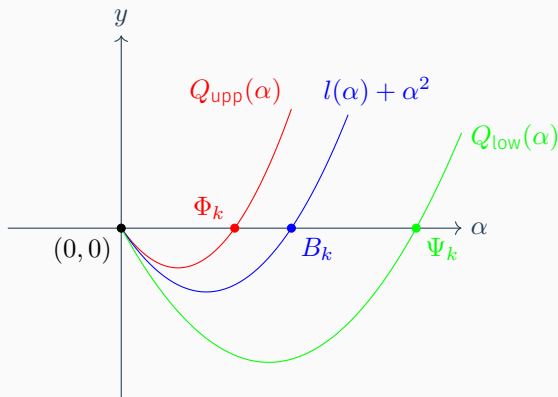


## SKETCH OF PROOF



- $C_1 \|G_k\| = \Phi_k \leq B_k \leq \Psi_k = C_2 \|G_k\|$
- Just to prove  $\mathbb{P}(\liminf_k B_k = 0) = 1$

## SKETCH OF PROOF



- $C_1 \|G_k\| = \Phi_k \leq B_k \leq \Psi_k = C_2 \|G_k\|$
- Just to prove  $\mathbb{P}(\liminf_k B_k = 0) = 1$
- Known that  $A_k(\omega) \rightarrow 0 \forall \omega \in \Omega$ , just to prove  $\mathbb{P}(A_k \geq B_k \text{ i.o.}) = 1$

Let  $S_k = \log_2(A_k/\alpha_0)$  and  $R_k = \log_2(B_k/\alpha_0)$ . Then we abstract the following problem.

# ANALYSIS OF A RANDOM WALK

Let  $S_k = \log_2(A_k/\alpha_0)$  and  $R_k = \log_2(B_k/\alpha_0)$ . Then we abstract the following problem.

Assume that  $(S_k)_{k \geq 0}$  and  $(R_k)_{k \geq 0}$  are two stochastic processes such that

- $S_0 = 0$  and  $(R_k)_{k \geq 0}$  uniformly bounded,

- 

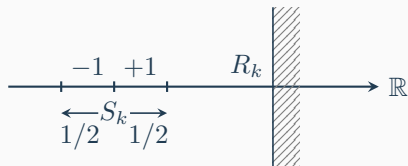
$$\mathbb{P}(S_{k+1} = S_k - 1 \mid S_1, \dots, S_k) = \frac{1}{2} \mathbb{1}_{\{S_k < R_k\}} + \mathbb{1}_{\{S_k \geq R_k\}},$$

$$\mathbb{P}(S_{k+1} = S_k + 1 \mid S_1, \dots, S_k) = \frac{1}{2} \mathbb{1}_{\{S_k < R_k\}}.$$

Then,  $\mathbb{P}(S_k \geq R_k \text{ i.o.}) = 1$ ?

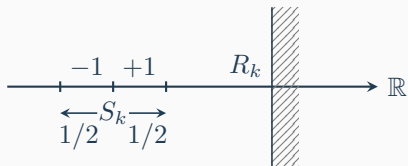
## GRAPHICAL EXPLANATION

When  $S_k < R_k$ ,

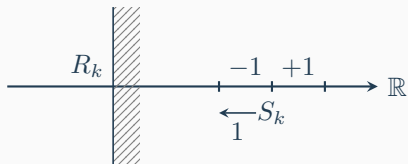


## GRAPHICAL EXPLANATION

When  $S_k < R_k$ ,

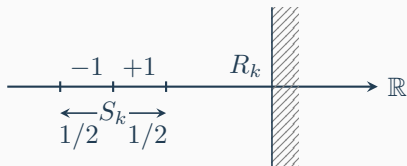


When  $S_k \geq R_k$ ,

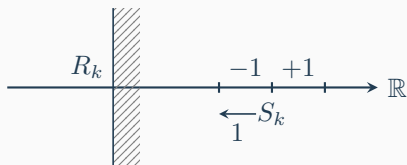


## GRAPHICAL EXPLANATION

When  $S_k < R_k$ ,



When  $S_k \geq R_k$ ,



Does  $S_k$  go beyond the “wall”  $R_k$  i.o. with probability 1?

- Construct  $\tilde{S}_k = S_k + \sum_{i=0}^{k-1} T_i \mathbb{1}_{\{S_k \geq R_k\}}$ ,  
where  $T_0, T_1, \dots$  are i.i.d.,  $\mathbb{P}(T_i = 0) = \mathbb{P}(T_i = 2) = 1/2$



- Construct  $\tilde{S}_k = S_k + \sum_{i=0}^{k-1} T_i \mathbb{1}_{\{S_k \geq R_k\}}$ ,  
 where  $T_0, T_1, \dots$  are i.i.d.,  $\mathbb{P}(T_i = 0) = \mathbb{P}(T_i = 2) = 1/2$   
 Intuition: Add the probability of  $S_k$  going right in the  
 case  $S_k \geq R_k$

- Construct  $\tilde{S}_k = S_k + \sum_{i=0}^{k-1} T_i \mathbb{1}_{\{S_k \geq R_k\}}$ ,  
where  $T_0, T_1, \dots$  are i.i.d.,  $\mathbb{P}(T_i = 0) = \mathbb{P}(T_i = 2) = 1/2$   
Intuition: Add the probability of  $S_k$  going right in the  
case  $S_k \geq R_k$
- $(\tilde{S}_k)_{k \geq 0}$  is simple random walk  $\Rightarrow \limsup_k \tilde{S}_k = \infty$  w.p. 1

- Construct  $\tilde{S}_k = S_k + \sum_{i=0}^{k-1} T_i \mathbb{1}_{\{S_k \geq R_k\}}$ ,  
where  $T_0, T_1, \dots$  are i.i.d.,  $\mathbb{P}(T_i = 0) = \mathbb{P}(T_i = 2) = 1/2$   
Intuition: Add the probability of  $S_k$  going right in the  
case  $S_k \geq R_k$
- $(\tilde{S}_k)_{k \geq 0}$  is simple random walk  $\Rightarrow \limsup_k \tilde{S}_k = \infty$  w.p. 1
- $(R_k)_{k \geq 0}$  uniformly bounded  $\Rightarrow$  so does  $(S_k)_{k \geq 0}$

- Construct  $\tilde{S}_k = S_k + \sum_{i=0}^{k-1} T_i \mathbb{1}_{\{S_k \geq R_k\}}$ ,  
where  $T_0, T_1, \dots$  are i.i.d.,  $\mathbb{P}(T_i = 0) = \mathbb{P}(T_i = 2) = 1/2$   
Intuition: Add the probability of  $S_k$  going right in the  
case  $S_k \geq R_k$
- $(\tilde{S}_k)_{k \geq 0}$  is simple random walk  $\Rightarrow \limsup_k \tilde{S}_k = \infty$  w.p. 1
- $(R_k)_{k \geq 0}$  uniformly bounded  $\Rightarrow$  so does  $(S_k)_{k \geq 0}$
- $\sum_{i=0}^{\infty} T_i \mathbb{1}_{\{S_k \geq R_k\}} = \infty$  w.p. 1

## CONCLUSION

---

# CONCLUSION

In this talk:

# CONCLUSION

In this talk:

- establish non-convergence theory for probabilistic direct search,

# CONCLUSION

In this talk:

- establish non-convergence theory for probabilistic direct search,
- distinguish convergence region and non-convergence region,



In this talk:

- establish non-convergence theory for probabilistic direct search,
- distinguish convergence region and non-convergence region,
- construct one counterexample to show the tightness of non-convergence region boundary.

Future work:

- find estimation or lower bound for the probability of non-convergence,
- establish non-convergence theory for other models.

Thank you!

## REFERENCES I

- ▶ Biviano, A. et al. (2013). “CLASH-VLT: the mass, velocity-anisotropy, and pseudo-phase-space density profiles of the  $z = 0.44$  galaxy cluster MACS J1206.2-0847”. *A&A* 558, A1:1–A1:22.
- ▶ Conn, A. R., Scheinberg, K., and Vicente, L. N. (2009). *Introduction to Derivative-Free Optimization*. Vol. 8. MOS-SIAM Ser. Optim. Philadelphia: SIAM.
- ▶ Fermi, E. and Metropolis, N. (1952). *Numerical solution of a minimum problem*. Tech. rep. Alamos National Laboratory, Los Alamos, USA.
- ▶ Ghanbari, H. and Scheinberg, K. (2017). “Black-box optimization in machine learning with trust region based derivative free algorithm”. *arXiv:1703.06925*.
- ▶ Gratton, S. et al. (2015). “Direct search based on probabilistic descent”. *SIAM J. Optim.* 25, pp. 1515–1541.