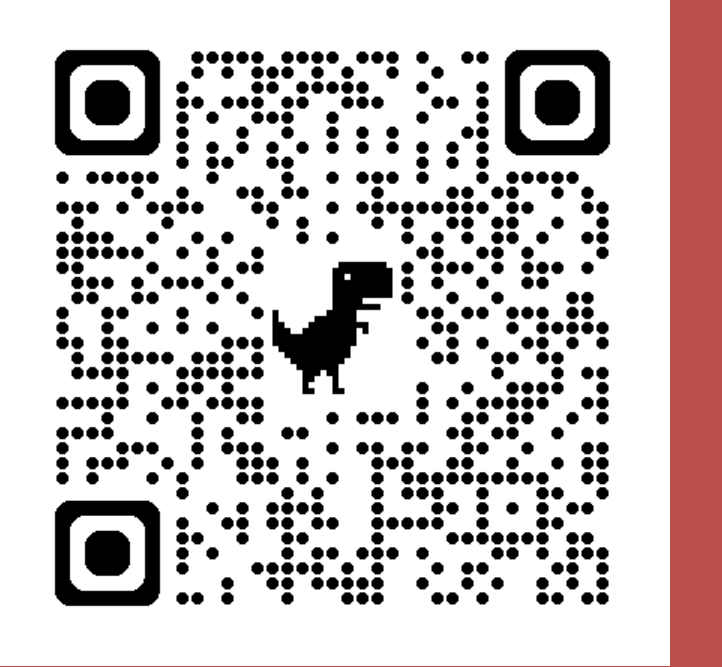


# Parameter-free Clipped Gradient Descent Meets Polyak

Yuki Takezawa<sup>1,2</sup>, Han Bao<sup>1,2</sup>, Ryoma Sato<sup>3</sup>, Kenta Niwa<sup>4</sup>, Makoto Yamada<sup>2</sup>

<sup>1</sup>Kyoto University, <sup>2</sup>OIST, <sup>3</sup>NII, <sup>4</sup>NTT Communication Science Laboratories



## Background

### ► Gradient Descent

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t). \quad (1)$$

#### Assumption ( $L$ -smoothness)

There exists a constant  $L > 0$  such that it holds that for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ ,

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|. \quad (2)$$

#### Theorem (Gradient Descent)

Assume that  $f$  is convex and  $L$ -smooth. Then, gradient descent with  $\eta_t = \frac{1}{L}$  satisfies

$$f(\bar{\mathbf{x}}) - f(\mathbf{x}^*) \leq \mathcal{O}\left(\frac{L\|\mathbf{x}_0 - \mathbf{x}^*\|^2}{T}\right). \quad (3)$$

where  $\bar{\mathbf{x}} := \frac{1}{T} \sum_{t=0}^{T-1} \mathbf{x}_t$ .

### ► Clipped Gradient Descent

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta_t \min\left\{1, \frac{c}{\|\nabla f(\mathbf{x}_t)\|}\right\} \nabla f(\mathbf{x}_t). \quad (4)$$

#### Assumption ( $(L_0, L_1)$ -smoothness)

There exists constants  $L_0 > 0$  and  $L_1 > 0$  that satisfies

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq (L_0 + L_1 \|\nabla f(\mathbf{x})\|) \|\mathbf{x} - \mathbf{y}\|, \quad (5)$$

for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$  with  $\|\mathbf{x} - \mathbf{y}\| \leq \frac{1}{L_1}$ .

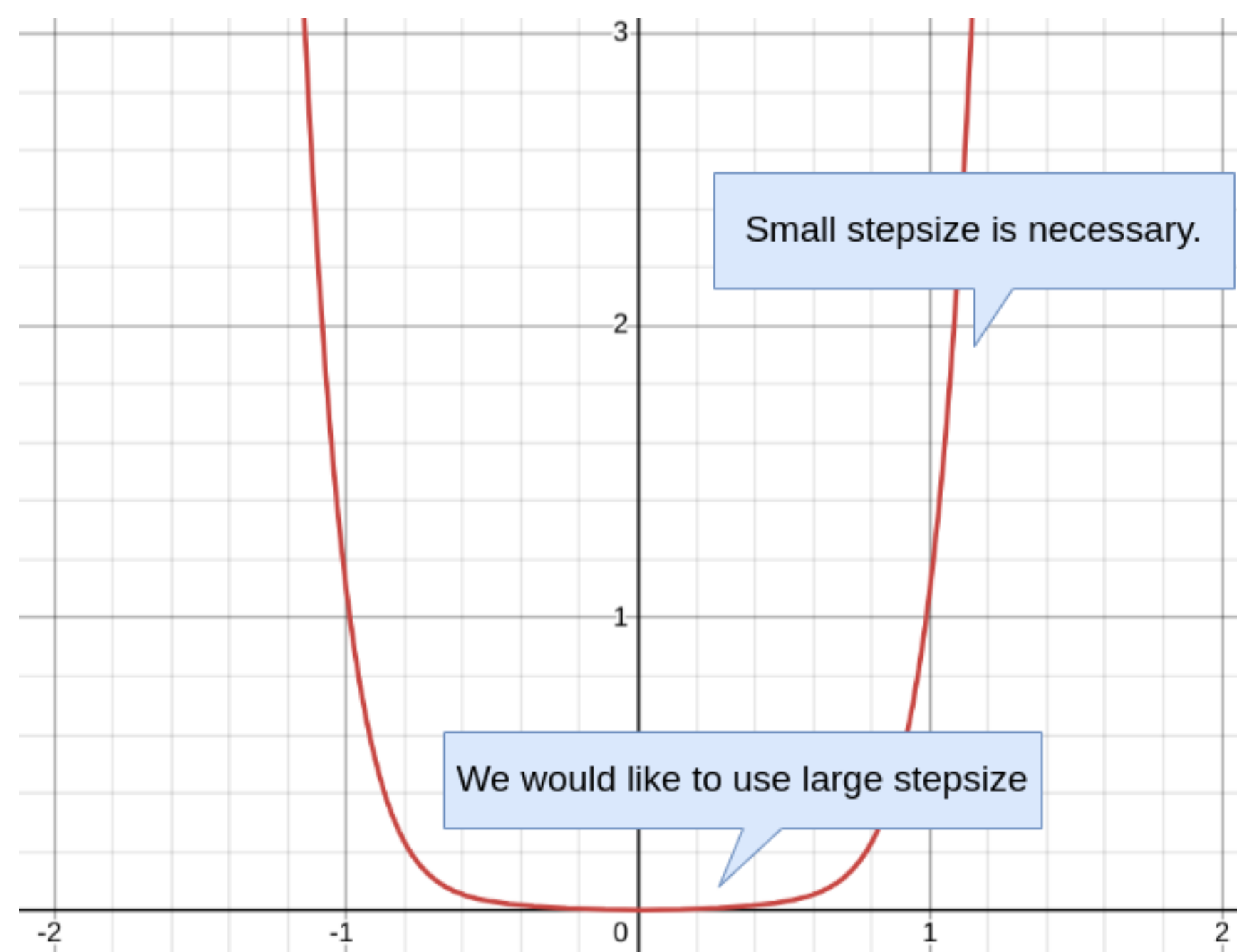
#### Theorem (Clipped Gradient Descent)

Assume that  $f$  is convex,  $L$ -smooth, and  $(L_0, L_1)$ -smooth. Then, clipped gradient descent with  $\eta_t = \frac{1}{L_0}$  and  $c = \frac{L_0}{L_1}$  satisfies

$$f(\bar{\mathbf{x}}) - f(\mathbf{x}^*) \leq \mathcal{O}\left(\frac{L_0\|\mathbf{x}_0 - \mathbf{x}^*\|^2}{T} + \frac{LL_1^2\|\mathbf{x}_0 - \mathbf{x}^*\|^4}{T^2}\right), \quad (6)$$

where  $\bar{\mathbf{x}} := \frac{1}{T} \sum_{t=0}^{T-1} \mathbf{x}_t$ .

- Since  $L_0 \ll L$  in practice, clipped gradient descent can converge faster than gradient descent.



## Contribution

*Q: Can we develop a parameter-free method whose convergence rate is asymptotically independent of  $L$  under  $(L_0, L_1)$ -smoothness?*

- We discover that Polyak stepsize can converge as fast as clipped gradient descent.
- We make Polyak stepsize parameter-free without losing the asymptotic independence of  $L$  by proposing Inexact Polyak Stepsize.

## New Convergence Result of Polyak Stepsize

### ► Polyak Stepsize

- It is well-known that Polyak stepsize allows gradient descent to converge as fast as the optimal stepsize.

$$\eta_t = \frac{f(\mathbf{x}_t) - f(\mathbf{x}^*)}{\|\nabla f(\mathbf{x}_t)\|^2}. \quad (7)$$

### ► Analysis of Polyak Stepsize under $(L_0, L_1)$ -smoothness

- Polyak stepsize can also achieve the same convergence rate as clipped gradient descent.

#### Theorem (Polyak Stepsize)

Assume that  $f$  is convex,  $L$ -smooth, and  $(L_0, L_1)$ -smooth. Then, gradient descent with Polyak stepsize satisfies

$$f(\mathbf{x}_\tau) - f(\mathbf{x}^*) \leq \mathcal{O}\left(\frac{L_0\|\mathbf{x}_0 - \mathbf{x}^*\|^2}{T} + \frac{LL_1^2\|\mathbf{x}_0 - \mathbf{x}^*\|^4}{T^2}\right),$$

where  $\tau := \arg \min_{0 \leq t \leq T} f(\mathbf{x}_t)$ .

### ► Several existing papers proposed parameter-free versions of Polyak stepsize, while they lost the fruitful property under $(L_0, L_1)$ -smoothness.

- They proposed to use the lower bound  $l^*$  instead of  $f(\mathbf{x}^*)$ .
- To prevent the stepsize from becoming too large, they make the stepsize monotonically decreasing.

$$\eta_t = \frac{f(\mathbf{x}_t) - l^*}{\|\nabla f(\mathbf{x}_t)\|^2}.$$

## Proposed Method

### Algorithm 1 Inexact Polyak Stepsize

- 1: **Input:** The number of iterations  $T$  and lower bound  $l^*$ .
- 2:  $f^{\text{best}}, \mathbf{x}^{\text{best}} \leftarrow f(\mathbf{x}_0), \mathbf{x}_0$ .
- 3: **for**  $t = 0, 1, \dots, T-1$  **do**
- 4:    $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \frac{f(\mathbf{x}_t) - l^*}{\sqrt{T}\|\nabla f(\mathbf{x}_t)\|^2} \nabla f(\mathbf{x}_t)$ .
- 5:   **if**  $f(\mathbf{x}_{t+1}) \leq f^{\text{best}}$  **then**
- 6:      $f^{\text{best}}, \mathbf{x}^{\text{best}} \leftarrow f(\mathbf{x}_{t+1}), \mathbf{x}_{t+1}$ .
- 7: **return**  $\mathbf{x}^{\text{best}}$ .

## Theorem (Inexact Polyak Stepsize)

Assume that  $f$  is convex,  $L$ -smooth, and  $(L_0, L_1)$ -smooth. Then, gradient descent with Inexact Polyak stepsize satisfies

$$\begin{aligned} & f(\mathbf{x}^{\text{best}}) - f(\mathbf{x}^*) \\ & \leq \mathcal{O}\left(\frac{L_0\|\mathbf{x}_0 - \mathbf{x}^*\|^2 + \sigma^2}{\sqrt{T}} + \frac{LL_1^2\|\mathbf{x}_0 - \mathbf{x}^*\|^4}{T} + \frac{L_1^2L\sigma^4}{L_0^2T}\right), \end{aligned}$$

where  $\sigma^2 := f(\mathbf{x}^*) - l^*$ .

- The convergence rate of Inexact Polyak Stepsize is asymptotically independent of  $L$ .
- The convergence rates of DecSPS and AdaSPS depend on  $D_T := \max_{0 \leq t \leq T} \|\mathbf{x}_t - \mathbf{x}^*\|$ , while the rate of Inexact Polyak Stepsize depends on  $\|\mathbf{x}_0 - \mathbf{x}^*\|$ .

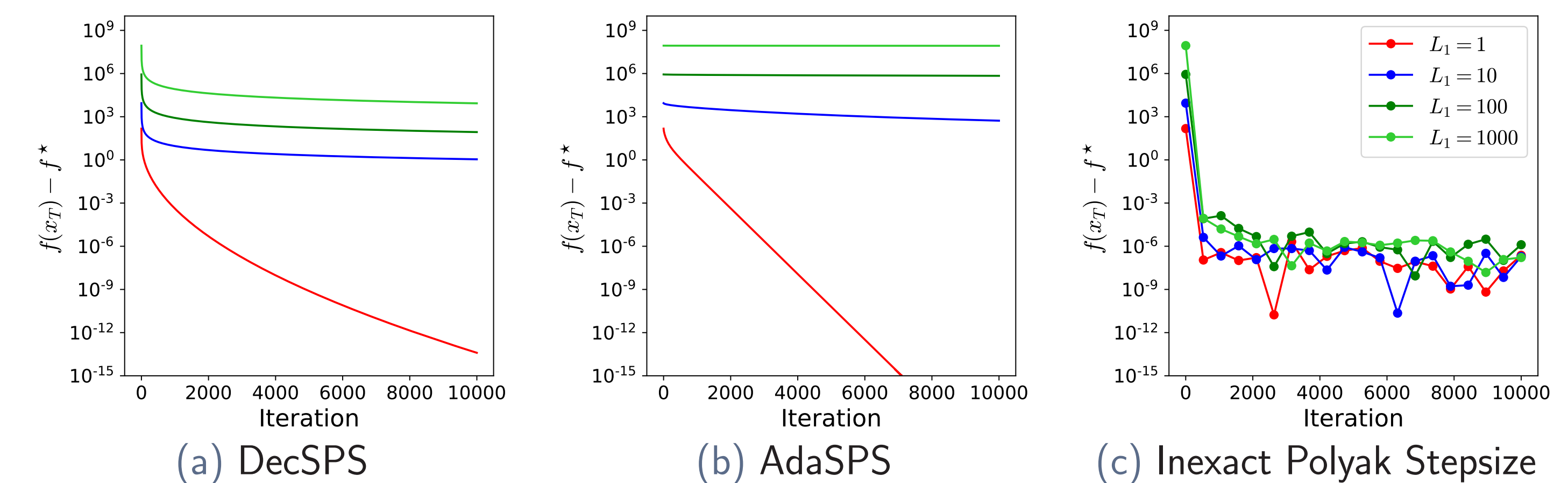
Table: Summary of convergence rates of parameter-free methods based on Polyak stepsize.

Algorithm	Convergence Rate
DecSPS (Orvieto et al., 2022)	$\mathcal{O}\left(\frac{\max\{L, \eta_0^{-1}\}D_T^2 + \sigma^2}{\sqrt{T}}\right)$
AdaSPS (Jiang et al., 2023)	$\mathcal{O}\left(\frac{LD_T^2\sigma}{\sqrt{T}} + \frac{L^2D_T^4}{T}\right)$
Inexact Polyak Stepsize (Ours)	$\mathcal{O}\left(\frac{L_0\ \mathbf{x}_0 - \mathbf{x}^*\ ^2 + \sigma^2}{\sqrt{T}} + \frac{LL_1^2\ \mathbf{x}_0 - \mathbf{x}^*\ ^4}{T} + \frac{L_1^2L\sigma^4}{L_0^2T}\right)$

## Numerical Results

### ► Synthetic Function

- The convergence behavior of Inexact Polyak stepsize is almost the same for all  $L_1$ .



### ► Neural Networks

