

New Perspectives on the Polyak Stepsize: Surrogate Functions and Negative Results

Francesco Orabona and Ryan D'Orazio

Workshop on Regret, Optimization, and Games, 2025

Motivation

- This presentation is about better understanding previous results on Polyak stepsize

Motivation

- This presentation is about better understanding previous results on Polyak stepsize
- Polyak stepsize adapts to lipschitzness, strong convexity, smoothness, sharpness [e.g., Hazan&Kakade, 2019]

Motivation

- This presentation is about better understanding previous results on Polyak stepsize
- Polyak stepsize adapts to lipschitzness, strong convexity, smoothness, sharpness [e.g., Hazan&Kakade, 2019]
- But, **why?**

Motivation

- This presentation is about better understanding previous results on Polyak stepsize
- Polyak stepsize adapts to lipschitzness, strong convexity, smoothness, sharpness [e.g., Hazan&Kakade, 2019]
- But, **why**?
- In optimization for ML, a convergence rate is rarely predictive of reality, so the **why** is actually more important than a non-informative theorem

Motivation

- This presentation is about better understanding previous results on Polyak stepsize
- Polyak stepsize adapts to lipschitzness, strong convexity, smoothness, sharpness [e.g., Hazan&Kakade, 2019]
- But, **why**?
- In optimization for ML, a convergence rate is rarely predictive of reality, so the **why** is actually more important than a non-informative theorem
- In this talk I'll present a way to *truly understand and explain* the behaviour of the Polyak stepsize

Motivation

- This presentation is about better understanding previous results on Polyak stepsize
- Polyak stepsize adapts to lipschitzness, strong convexity, smoothness, sharpness [e.g., Hazan&Kakade, 2019]
- But, **why**?
- In optimization for ML, a convergence rate is rarely predictive of reality, so the **why** is actually more important than a non-informative theorem
- In this talk I'll present a way to *truly understand and explain* the behaviour of the Polyak stepsize
- No really new rates, but many negative results

The Challenge of Stepsize Selection

Gradient Descent (GD)

The foundational first-order optimization algorithm:

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

The Challenge of Stepsize Selection

Gradient Descent (GD)

The foundational first-order optimization algorithm:

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

- The stepsize (or learning rate) η_t is critical

The Challenge of Stepsize Selection

Gradient Descent (GD)

The foundational first-order optimization algorithm:

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

- The stepsize (or learning rate) η_t is critical
- Tuning η_t often requires knowledge of problem parameters (e.g., smoothness constant L , distance to optimum) that are unknown

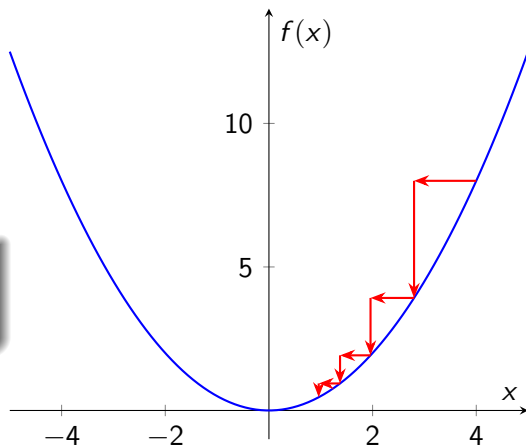
Gradient Descent on Smooth f : Step Size Too Small

Function: $f(x) = \frac{1}{2}x^2$

- Curvature (Smoothness)
 $L = 1$
- Step size $\eta = 0.3$

Observation

The algorithm converges, but takes many small steps



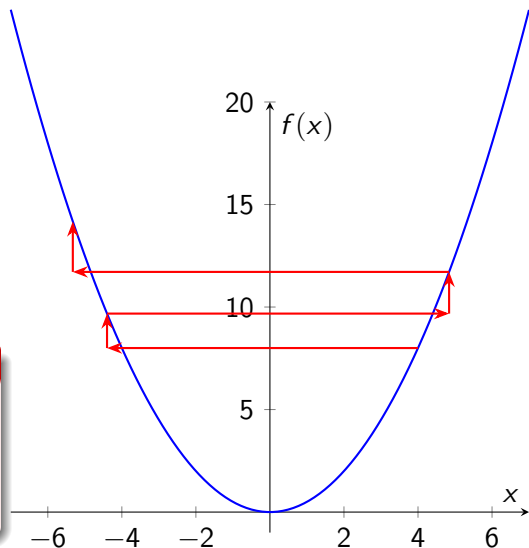
Gradient Descent on Smooth f : Step Size Too Large

Function: $f(x) = \frac{1}{2}x^2$

- Curvature (Smoothness)
 $L = 1$
- Step size $\eta = 2.1$
- This is larger than the divergence threshold
 $\eta = 2/L = 2$

Observation

Each step overshoots the minimum by a larger amount, and the iterates move further away from the solution



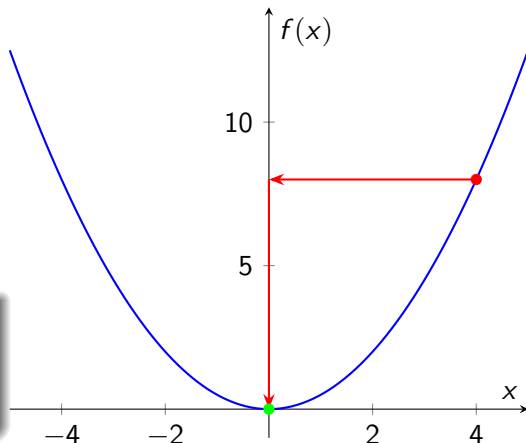
Gradient Descent on Smooth f : Optimal Constant Step Size

Function: $f(x) = \frac{1}{2}x^2$

- Curvature (Smoothness)
 $L = 1$
- Optimal step size
 $\eta = 1/L = 1$

Observation

This stepsize maximizes the worst-case decrease of the function



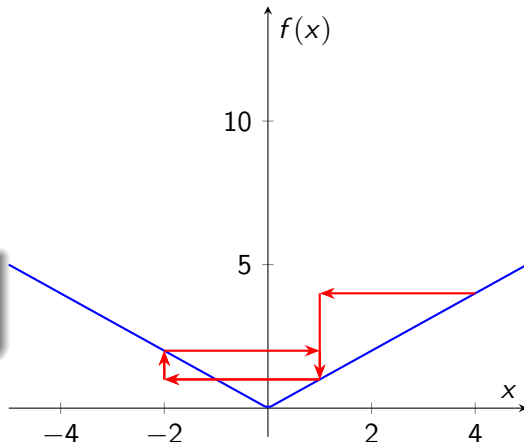
Gradient Descent on Non-Smooth f

Function: $f(x) = |x|$

- Step size $\eta = 3$

Observation

The learning rate is too large,
it will oscillate



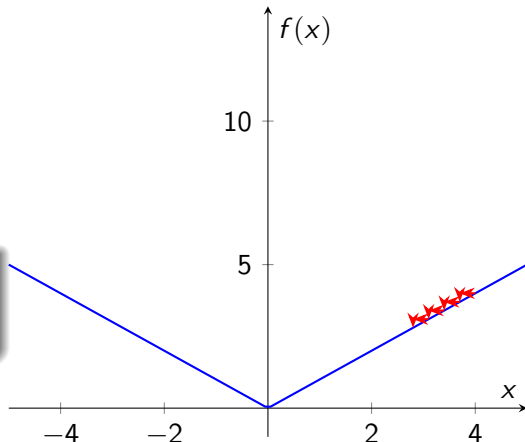
Gradient Descent on Non-Smooth f

Function: $f(x) = |x|$

- Step size $\eta = .3$

Observation

The learning rate is too small, it will converge very slowly



Gradient Descent on Non-Smooth f

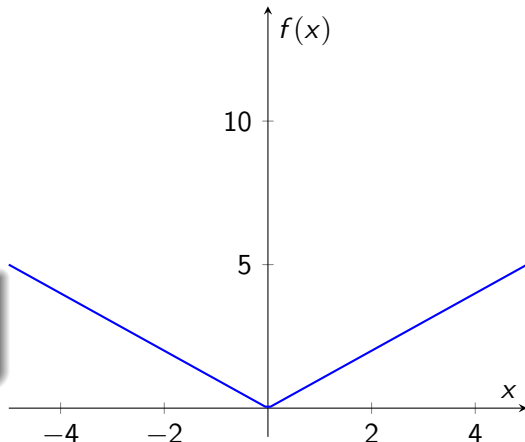
Function: $f(x) = |x|$

- Optimal step size

$$\eta = \frac{\|x_1 - x^*\|}{\sqrt{T}}$$

Observation

The optimal learning rate depends on where you start



The “Magic” of the Polyak Stepsize

Proposed by Boris Polyak in 1969, the stepsize is defined as

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \frac{f(\mathbf{x}_t) - f^*}{\|\mathbf{g}_t\|_2^2} \mathbf{g}_t,$$

where $\mathbf{g}_t \in \partial f(\mathbf{x}_t)$ is a subgradient and $f^* = \min_{\mathbf{x}} f(\mathbf{x})$

The “Magic” of the Polyak Stepsize

Proposed by Boris Polyak in 1969, the stepsize is defined as

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \frac{f(\mathbf{x}_t) - f^*}{\|\mathbf{g}_t\|_2^2} \mathbf{g}_t,$$

where $\mathbf{g}_t \in \partial f(\mathbf{x}_t)$ is a subgradient and $f^* = \min_{\mathbf{x}} f(\mathbf{x})$

Remarkable Adaptivity

A single update rule achieves near-optimal rates for:

- Non-smooth convex functions: $\mathcal{O}(1/\sqrt{T})$
- Smooth convex functions: $\mathcal{O}(1/T)$
- Smooth & strongly convex functions: Linear convergence

...all without knowing smoothness or strong convexity constants!

The Central Research Question

Despite a resurgence of interest and many new variants [e.g., Rolinek&Martius, NeurIPS'18; Berrada et al., ICML'20; Loizou et al., AISTATS'21; Prazeres&Oberman, 2021], a fundamental question remains:

What makes the Polyak stepsize so adaptive, and when can it fail?

The Central Research Question

Despite a resurgence of interest and many new variants [e.g., Rolinek&Martius, NeurIPS'18; Berrada et al., ICML'20; Loizou et al., AISTATS'21; Prazeres&Oberman, 2021], a fundamental question remains:

What makes the Polyak stepsize so adaptive, and when can it fail?

Our Contributions

- 1 A new, unified perspective: Polyak's method is just **Gradient Descent on a surrogate function**

The Central Research Question

Despite a resurgence of interest and many new variants [e.g., Rolinek&Martius, NeurIPS'18; Berrada et al., ICML'20; Loizou et al., AISTATS'21; Prazeres&Oberman, 2021], a fundamental question remains:

What makes the Polyak stepsize so adaptive, and when can it fail?

Our Contributions

- 1 A new, unified perspective: Polyak's method is just **Gradient Descent on a surrogate function**
- 2 This surrogate is **always locally smooth** and we know the **smoothness constant**

The Central Research Question

Despite a resurgence of interest and many new variants [e.g., Rolinek&Martius, NeurIPS'18; Berrada et al., ICML'20; Loizou et al., AISTATS'21; Prazeres&Oberman, 2021], a fundamental question remains:

What makes the Polyak stepsize so adaptive, and when can it fail?

Our Contributions

- 1 A new, unified perspective: Polyak's method is just **Gradient Descent on a surrogate function**
- 2 This surrogate is **always locally smooth** and we know the **smoothness constant**
- 3 We use this framework to analyze a general family of Polyak-like algorithms

The Central Research Question

Despite a resurgence of interest and many new variants [e.g., Rolinek&Martius, NeurIPS'18; Berrada et al., ICML'20; Loizou et al., AISTATS'21; Prazeres&Oberman, 2021], a fundamental question remains:

What makes the Polyak stepsize so adaptive, and when can it fail?

Our Contributions

- ➊ A new, unified perspective: Polyak's method is just **Gradient Descent on a surrogate function**
- ➋ This surrogate is **always locally smooth** and we know the **smoothness constant**
- ➌ We use this framework to analyze a general family of Polyak-like algorithms
- ➍ We prove several **negative results**, showing that the non-convergence seen in some analyses is real, not an artifact

Polyak Stepsize as GD on a Surrogate

Let f be convex with minimizer \mathbf{x}^*

Polyak Stepsize as GD on a Surrogate

Let f be convex with minimizer \mathbf{x}^*

Instead of minimizing $f(\mathbf{x})$, consider minimizing a new surrogate function:

$$\phi(\mathbf{x}) = \frac{1}{2} (f(\mathbf{x}) - f(\mathbf{x}^*))^2$$

Polyak Stepsize as GD on a Surrogate

Let f be convex with minimizer \mathbf{x}^*

Instead of minimizing $f(\mathbf{x})$, consider minimizing a new surrogate function:

$$\phi(\mathbf{x}) = \frac{1}{2} (f(\mathbf{x}) - f(\mathbf{x}^*))^2$$

Key Insight

The subgradient of $\phi(\mathbf{x})$ is $\nabla\phi(\mathbf{x}) = (f(\mathbf{x}) - f^*)\mathbf{g}_x$, where $\mathbf{g}_x \in \partial f(\mathbf{x})$

A subgradient step on $\phi(\mathbf{x})$ with stepsize $\eta = \frac{1}{\|\mathbf{g}_x\|_2^2}$ is:

$$\mathbf{x} - \eta \nabla\phi(\mathbf{x}) = \mathbf{x} - \frac{1}{\|\mathbf{g}_x\|_2^2} (f(\mathbf{x}) - f^*)\mathbf{g}_x$$

This is exactly the Polyak update!

A New Notion of Local Curvature

So, the Polyak stepsize is GD on ϕ with stepsize $\eta'_t = 1/\|\mathbf{g}_t\|_2^2$. But why is this a good stepsize?

A New Notion of Local Curvature

So, the Polyak stepsize is GD on ϕ with stepsize $\eta'_t = 1/\|\mathbf{g}_t\|_2^2$. But why is this a good stepsize?

Definition (Local Star Upper Curvature - LSUC)

A function ϕ has $\lambda_{\mathbf{y}}$ -LSUC around \mathbf{y} if

$$\phi(\mathbf{x}^*) - \langle \nabla \phi(\mathbf{y}), \mathbf{x}^* - \mathbf{y} \rangle - \frac{1}{2\lambda_{\mathbf{y}}} \|\nabla \phi(\mathbf{y})\|_2^2 \geq \phi(\mathbf{y})$$

A New Notion of Local Curvature

So, the Polyak stepsize is GD on ϕ with stepsize $\eta'_t = 1/\|\mathbf{g}_t\|_2^2$. But why is this a good stepsize?

Definition (Local Star Upper Curvature - LSUC)

A function ϕ has $\lambda_{\mathbf{y}}$ -LSUC around \mathbf{y} if

$$\phi(\mathbf{x}^*) - \langle \nabla \phi(\mathbf{y}), \mathbf{x}^* - \mathbf{y} \rangle - \frac{1}{2\lambda_{\mathbf{y}}} \|\nabla \phi(\mathbf{y})\|_2^2 \geq \phi(\mathbf{y})$$

- L -smooth functions are L -LSUC everywhere, in fact convex smooth functions satisfy

$$f(\mathbf{x}) - \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle - \frac{1}{2L} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|^2 \geq f(\mathbf{y}), \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d$$

- This is a local smoothness-like condition

The Local Curvature is the Source of Adaptivity

Theorem 1: Curvature of the Polyak Surrogate

For any f convex, the surrogate $\phi(\mathbf{x}) = \frac{1}{2}(f(\mathbf{x}) - f^\star)^2$ is $\|\mathbf{g}_\mathbf{x}\|_2^2$ -LSUC around any \mathbf{x}

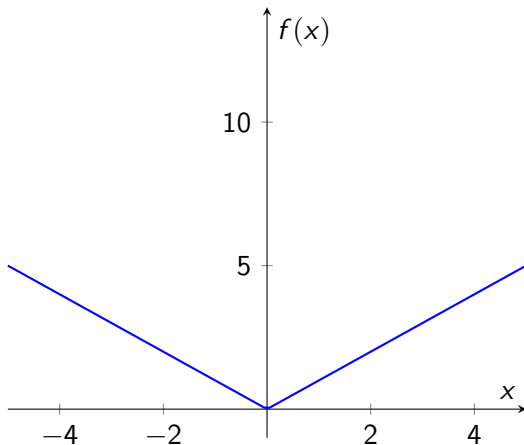
The Local Curvature is the Source of Adaptivity

Theorem 1: Curvature of the Polyak Surrogate

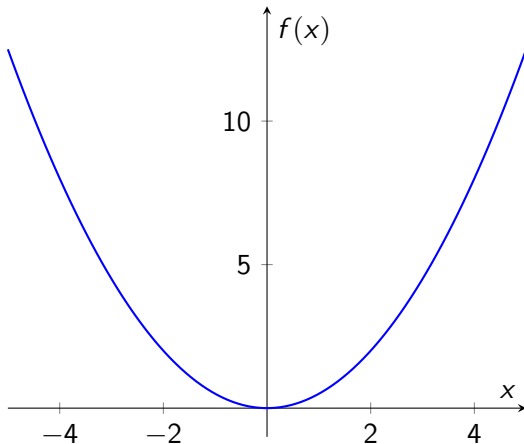
For any f convex, the surrogate $\phi(\mathbf{x}) = \frac{1}{2}(f(\mathbf{x}) - f^\star)^2$ is $\|\mathbf{g}_\mathbf{x}\|_2^2$ -LSUC around any \mathbf{x}

- **This is the magic!** The surrogate ϕ is *always* “locally smooth”
- The adaptive stepsize $1/\|\mathbf{g}_t\|_2^2$ is simply the inverse of this local curvature constant!

Non-smooth: Choice of Stepsize is Difficult



Smooth: Choice of Stepsize is Easy, *Knowing Smoothness*



Recovering Convergence Rates

Using this perspective, we can easily derive convergence guarantees

Lemma (One-step Progress)

Using stepsize $\eta'_t = 1/\lambda_{\mathbf{x}_t} = 1/\|\mathbf{g}_t\|_2^2$ on ϕ gives:

$$\eta'_t (\phi(\mathbf{x}_t) - \phi(\mathbf{x}^\star)) \leq \frac{1}{2} \|\mathbf{x}_t - \mathbf{x}^\star\|_2^2 - \frac{1}{2} \|\mathbf{x}_{t+1} - \mathbf{x}^\star\|_2^2$$

Recovering Convergence Rates

Using this perspective, we can easily derive convergence guarantees

Lemma (One-step Progress)

Using stepsize $\eta'_t = 1/\lambda_{\mathbf{x}_t} = 1/\|\mathbf{g}_t\|_2^2$ on ϕ gives:

$$\eta'_t (\phi(\mathbf{x}_t) - \phi(\mathbf{x}^*)) \leq \frac{1}{2} \|\mathbf{x}_t - \mathbf{x}^*\|_2^2 - \frac{1}{2} \|\mathbf{x}_{t+1} - \mathbf{x}^*\|_2^2$$

Summing over T steps:

$$\sum_{t=1}^T \eta'_t \phi(\mathbf{x}_t) \leq \frac{1}{2} \|\mathbf{x}_1 - \mathbf{x}^*\|_2^2$$

Recovering Convergence Rates

Using this perspective, we can easily derive convergence guarantees

Lemma (One-step Progress)

Using stepsize $\eta'_t = 1/\lambda_{\mathbf{x}_t} = 1/\|\mathbf{g}_t\|_2^2$ on ϕ gives:

$$\eta'_t (\phi(\mathbf{x}_t) - \phi(\mathbf{x}^*)) \leq \frac{1}{2} \|\mathbf{x}_t - \mathbf{x}^*\|_2^2 - \frac{1}{2} \|\mathbf{x}_{t+1} - \mathbf{x}^*\|_2^2$$

Summing over T steps:

$$\sum_{t=1}^T \eta'_t \phi(\mathbf{x}_t) \leq \frac{1}{2} \|\mathbf{x}_1 - \mathbf{x}^*\|_2^2$$

- If f is G -Lipschitz: $\sum \eta'_t \geq T/G^2 \Rightarrow \phi(\bar{\mathbf{x}}_T) = \mathcal{O}(1/T)$
- If f is L -self-bounded: We recover $\phi(\bar{\mathbf{x}}_T) = \mathcal{O}(1/T^2)$
- If f is sharp: We recover linear convergence

From the Surrogate to the Original Function

- We can easily convert a rate on the surrogate to a rate on the original function by inverting the surrogate
- For example for smooth losses we have

$$f(\mathbf{x}) - f(\mathbf{x}^*) = \sqrt{2\phi(\mathbf{x})} = \sqrt{\mathcal{O}(1/T^2)} = \mathcal{O}(1/T)$$

Not a Completely New Idea

- Gower et al. [ArXiv'21] showed that the stochastic Polyak stepsize can be casted as online convex optimization problem on surrogate losses
- The adversarial nature of online convex optimization means that it is not possible to say that we are minimizing a specific function
- For the same reason, they need slightly stronger assumptions

A Family of Surrogates

We can generalize this idea beyond knowing f^*

General Surrogate

Consider $\psi(\mathbf{x}) = \frac{1}{2}h^2(\mathbf{x})$, where $h : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$ is convex

A Family of Surrogates

We can generalize this idea beyond knowing f^\star

General Surrogate

Consider $\psi(\mathbf{x}) = \frac{1}{2}h^2(\mathbf{x})$, where $h : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$ is convex

Examples of $h(\mathbf{x})$:

- Original Polyak: $h(\mathbf{x}) = f(\mathbf{x}) - f^\star$
- Unknown f^\star : $h(\mathbf{x}) = (f(\mathbf{x}) - c)_+$ for some estimate c
- Stochastic Variants, for example, SPS_+ [Garrigos et al., 2023]:
 $h(\mathbf{x}, \xi) = |f(\mathbf{x}, \xi) - f(\mathbf{x}^\star, \xi)|_+$

A Family of Surrogates

We can generalize this idea beyond knowing f^*

General Surrogate

Consider $\psi(\mathbf{x}) = \frac{1}{2}h^2(\mathbf{x})$, where $h : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$ is convex

Examples of $h(\mathbf{x})$:

- Original Polyak: $h(\mathbf{x}) = f(\mathbf{x}) - f^*$
- Unknown f^* : $h(\mathbf{x}) = (f(\mathbf{x}) - c)_+$ for some estimate c
- Stochastic Variants, for example, SPS₊ [Garrigos et al., 2023]:
 $h(\mathbf{x}, \xi) = |f(\mathbf{x}, \xi) - f(\mathbf{x}^*, \xi)|_+$

Problem

If the minimum of h is not zero (i.e., $h(\mathbf{x}^*) > 0$), the surrogate only has *approximate* local curvature. This leads to convergence to a **neighborhood**, not the true optimum.

The Stochastic Setting

Consider minimizing $F(\mathbf{x}) = \mathbb{E}_{\xi \sim D}[f(\mathbf{x}, \xi)]$.

Stochastic Polyak variants use a surrogate based on a single sample ξ_t : $\frac{1}{2}h(\mathbf{x}, \xi_t)$, and update with

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta_t h(\mathbf{x}_t, \xi_t) \mathbf{g}_t, \text{ where } \mathbf{g}_t \in \partial h(\mathbf{x}, \xi_t)$$

The Stochastic Setting

Consider minimizing $F(\mathbf{x}) = \mathbb{E}_{\xi \sim D}[f(\mathbf{x}, \xi)]$.

Stochastic Polyak variants use a surrogate based on a single sample ξ_t : $\frac{1}{2}h(\mathbf{x}, \xi_t)$, and update with

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta_t h(\mathbf{x}_t, \xi_t) \mathbf{g}_t, \text{ where } \mathbf{g}_t \in \partial h(\mathbf{x}, \xi_t)$$

A Fundamental Mismatch

The algorithm is effectively minimizing the expectation of the *surrogate*:

$$\mathbb{E}_{\xi \sim D} \left[\frac{1}{2} h^2(\mathbf{x}, \xi) \right]$$

This is generally **not** the same as minimizing the original objective $F(\mathbf{x})$!

$$\operatorname{argmin}_{\mathbf{x}} \mathbb{E}[h^2(\mathbf{x}, \xi)] \neq \operatorname{argmin}_{\mathbf{x}} \mathbb{E}[f(\mathbf{x}, \xi)]$$

The Stochastic Setting

Consider minimizing $F(\mathbf{x}) = \mathbb{E}_{\xi \sim D}[f(\mathbf{x}, \xi)]$.

Stochastic Polyak variants use a surrogate based on a single sample ξ_t : $\frac{1}{2}h(\mathbf{x}, \xi_t)$, and update with

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta_t h(\mathbf{x}_t, \xi_t) \mathbf{g}_t, \text{ where } \mathbf{g}_t \in \partial h(\mathbf{x}, \xi_t)$$

A Fundamental Mismatch

The algorithm is effectively minimizing the expectation of the *surrogate*:

$$\mathbb{E}_{\xi \sim D} \left[\frac{1}{2} h^2(\mathbf{x}, \xi) \right]$$

This is generally **not** the same as minimizing the original objective $F(\mathbf{x})$!

$$\operatorname{argmin}_{\mathbf{x}} \mathbb{E}[h^2(\mathbf{x}, \xi)] \neq \operatorname{argmin}_{\mathbf{x}} \mathbb{E}[f(\mathbf{x}, \xi)]$$

Warning: Minimizing a different loss function is problematic for a ML point of view

Unified Analysis of Stochastic Variants

We propose a generalized algorithm with a clipped stepsize, covering methods like SPS_{\max} , SPS_{+} , etc.

- 1: **for** $t = 1, \dots, T$ **do**
- 2: Sample ξ_t
- 3: Get subgradient $\mathbf{g}_t \in \partial h(\mathbf{x}_t, \xi_t)$
- 4: $\eta_t = \min \left(\frac{1}{\|\mathbf{g}_t\|_2^2}, \frac{\gamma}{h(\mathbf{x}_t, \xi_t)} \right)$
- 5: $\mathbf{x}_{t+1} = \mathbf{x}_t - \eta_t h(\mathbf{x}_t, \xi_t) \mathbf{g}_t$
- 6: **end for**

General Convergence Guarantees

Theorem

- Let $H(\mathbf{x}) = \mathbb{E}_{\xi \sim D}[h(\mathbf{x}, \xi)]$. If $h(\cdot, \xi_t)$ is L -self bounded, we have

$$\min\left(\frac{1}{2L}, \gamma\right) \frac{1}{T} \sum_{t=1}^T \mathbb{E}[H(\mathbf{x}_t)] \leq \frac{\|\mathbf{x}_1 - \mathbf{x}^*\|^2}{T} + 2\gamma H(\mathbf{x}^*)$$

General Convergence Guarantees

Theorem

- Let $H(\mathbf{x}) = \mathbb{E}_{\xi \sim D}[h(\mathbf{x}, \xi)]$. If $h(\cdot, \xi_t)$ is L -self bounded, we have

$$\min\left(\frac{1}{2L}, \gamma\right) \frac{1}{T} \sum_{t=1}^T \mathbb{E}[H(\mathbf{x}_t)] \leq \frac{\|\mathbf{x}_1 - \mathbf{x}^*\|^2}{T} + 2\gamma H(\mathbf{x}^*)$$

- If $h(\cdot, \xi_t)$ is G -Lipschitz, then we have

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}[H(\mathbf{x}_t)] \leq \frac{\|\mathbf{x}_1 - \mathbf{x}^*\|^2}{\gamma T} + 2H(\mathbf{x}^*) + \frac{G\|\mathbf{x}_1 - \mathbf{x}^*\|}{\sqrt{T}} + G\sqrt{2\gamma H(\mathbf{x}^*)}$$

General Convergence Guarantees

Theorem

- Let $H(\mathbf{x}) = \mathbb{E}_{\xi \sim D}[h(\mathbf{x}, \xi)]$. If $h(\cdot, \xi_t)$ is L -self bounded, we have

$$\min\left(\frac{1}{2L}, \gamma\right) \frac{1}{T} \sum_{t=1}^T \mathbb{E}[H(\mathbf{x}_t)] \leq \frac{\|\mathbf{x}_1 - \mathbf{x}^*\|^2}{T} + 2\gamma H(\mathbf{x}^*)$$

- If $h(\cdot, \xi_t)$ is G -Lipschitz, then we have

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}[H(\mathbf{x}_t)] \leq \frac{\|\mathbf{x}_1 - \mathbf{x}^*\|^2}{\gamma T} + 2H(\mathbf{x}^*) + \frac{G\|\mathbf{x}_1 - \mathbf{x}^*\|}{\sqrt{T}} + G\sqrt{2\gamma H(\mathbf{x}^*)}$$

- If $h(\cdot, \xi)$ is L -self-bounded and $H(\mathbf{x})$ has μ -quadratic growth, then

$$\mathbb{E}[\|\mathbf{x}_{T+1} - \mathbf{x}^*\|^2] \leq \mathbb{E}[\|\mathbf{x}_1 - \mathbf{x}^*\|^2] a^{T+1} + b \frac{1 - a^{T+1}}{1 - a} H(\mathbf{x}^*),$$

where $a = \frac{\mu}{2} \min\left(\frac{1}{2L}, \gamma\right)$ and $b = 2\gamma - \min\left(\frac{1}{2L}, \gamma\right)$

When Things Go Wrong: $h(\mathbf{x}^\star) > 0$

What happens when the minimum of our surrogate-generator h is strictly positive?

- **Deterministic case:** We underestimate f^\star (e.g., $h(\mathbf{x}) = f(\mathbf{x}) - c$ with $c < f^\star$)
- The correct stepsize for the underlying surrogate $\psi = \frac{1}{2}(h - h^\star)^2$ is $\frac{1}{\lambda_t}$, but the algorithm uses a stepsize for $\frac{1}{2}h^2$:

$$\eta'_t = \left(\frac{h(\mathbf{x}_t)}{h(\mathbf{x}_t) - h^\star} \right) \frac{1}{\lambda_t}$$

When Things Go Wrong: $h(\mathbf{x}^\star) > 0$

What happens when the minimum of our surrogate-generator h is strictly positive?

- **Deterministic case:** We underestimate f^\star (e.g., $h(\mathbf{x}) = f(\mathbf{x}) - c$ with $c < f^\star$)
- The correct stepsize for the underlying surrogate $\psi = \frac{1}{2}(h - h^\star)^2$ is $\frac{1}{\lambda_t}$, but the algorithm uses a stepsize for $\frac{1}{2}h^2$:

$$\eta'_t = \left(\frac{h(\mathbf{x}_t)}{h(\mathbf{x}_t) - h^\star} \right) \frac{1}{\lambda_t}$$

Problem

As $\mathbf{x}_t \rightarrow \mathbf{x}^\star$, we have $h(\mathbf{x}_t) \rightarrow h^\star > 0$. The term $\frac{h(\mathbf{x}_t)}{h(\mathbf{x}_t) - h^\star}$ **blows up to $+\infty$!** The stepsize becomes enormous near the minimum, causing instability. Moreover, clipping will not fix this issue.

Unstable Fixed Points

This intuition is formalized in following Proposition

Proposition (Unstable Fixed Point)

For a wide class of functions h (e.g., self-bounded with quadratic growth), if $h(\mathbf{x}^*) > 0$, then the minimizer \mathbf{x}^* is an **unstable fixed point**.

There exists a neighborhood around \mathbf{x}^* where if you enter, the next step will take you **further away** from \mathbf{x}^* .

Unstable Fixed Points

This intuition is formalized in following Proposition

Proposition (Unstable Fixed Point)

For a wide class of functions h (e.g., self-bounded with quadratic growth), if $h(\mathbf{x}^*) > 0$, then the minimizer \mathbf{x}^* is an **unstable fixed point**.

There exists a neighborhood around \mathbf{x}^* where if you enter, the next step will take you **further away** from \mathbf{x}^* .

- This confirms that the neighborhood of non-convergence is not just an artifact of the analysis

Cycling of the Iterates

This instability can lead to more than just a failure to converge; it can lead to cycles

Proposition (Cycling)

Consider the simple 1D function $h(x) = x^2 + 1$. Here $h^* = 1 > 0$. There exists an initial point x_1 such that the update rule

$$x_{t+1} = x_t - \frac{h(x_t)}{\|\nabla h(x_t)\|_2^2} \nabla h(x_t)$$

cycles on points different than $x^* = 0$. Moreover, the suboptimality on the average of the iterates also fails to converge.

The update is

$$x_{t+1} = x_t - \frac{h(x_t)}{\|\nabla h(x_t)\|_2^2} \nabla h(x_t) = x_t - \frac{x_t^2 + 1}{2x_t} = \frac{x_t^2 - 1}{2x_t}$$

The update is

$$x_{t+1} = x_t - \frac{h(x_t)}{\|\nabla h(x_t)\|_2^2} \nabla h(x_t) = x_t - \frac{x_t^2 + 1}{2x_t} = \frac{x_t^2 - 1}{2x_t}$$

If we start at $x_1 = \cot(\theta)$, from the identity of $\cot(2x) = \frac{\cot^2 x - 1}{2 \cot x}$

The update is

$$x_{t+1} = x_t - \frac{h(x_t)}{\|\nabla h(x_t)\|_2^2} \nabla h(x_t) = x_t - \frac{x_t^2 + 1}{2x_t} = \frac{x_t^2 - 1}{2x_t}$$

If we start at $x_1 = \cot(\theta)$, from the identity of $\cot(2x) = \frac{\cot^2 \theta - 1}{2 \cot \theta}$

Hence, set $x_1 = \cot \pi/7$, to have

$$x_1 = \cot(\pi/7) \rightarrow x_2 = \cot(2\pi/7) \rightarrow x_3 = \cot(4\pi/7) \rightarrow x_4 = x_1$$

The Set of Good Initial Point has Measure Zero

- In the previous Proposition we found a very specific initial point such that the algorithm cycles
- Was it just a very unlucky initial point?

The Set of Good Initial Point has Measure Zero

- In the previous Proposition we found a very specific initial point such that the algorithm cycles
- Was it just a very unlucky initial point?

Proposition

For $h(x) = x^2/2 + a$, the set of initial points where the update $x_{t+1} = x_t - \frac{h(x_t)}{\|\nabla h(x_t)\|_2^2} \nabla h(x_t)$ converges to the minimum has **measure zero**

Proposition

There exist f_1 and f_2 quadratic 1-d functions and a starting point x_1 such that SPS on $F(x) = 0.5(f_1(x) + f_2(x))$ satisfies

$$\mathbb{E}[F(x_t)] - \min_x F(x) \geq 2/3, \forall t$$

Summary and Takeaways

The Good: A Unifying Perspective

- The Polyak stepsize is equivalent to GD on a surrogate $\phi(\mathbf{x}) = \frac{1}{2}(f - f^*)^2$
- The adaptivity comes from the fact that ϕ is always locally “smooth” with a known curvature constant $\|\mathbf{g}\|_2^2$
- This framework simplifies and unifies the analysis of many Polyak-like methods

Summary and Takeaways

The Good: A Unifying Perspective

- The Polyak stepsize is equivalent to GD on a surrogate $\phi(\mathbf{x}) = \frac{1}{2}(f - f^*)^2$
- The adaptivity comes from the fact that ϕ is always locally “smooth” with a known curvature constant $\|\mathbf{g}\|_2^2$
- This framework simplifies and unifies the analysis of many Polyak-like methods

The Bad: Fundamental Instability

- When the surrogate’s minimum value is positive (e.g., f^* underestimated or no interpolation), the dynamics change drastically
- The algorithm becomes unstable near the optimum, leading to cycles and non-convergence
- This neighborhood of convergence is not an analysis artifact but a fundamental property of the method

“New Perspectives on the Polyak Stepsize: Surrogate Functions and Negative Results” Francesco Orabona, Ryan D’Orazio, NeurIPS’25

P.S. We have multiple PhD/Post-Doc/Research Scientist positions, and 3000 GH200 GPUs

