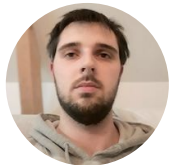


Super-acceleration with cyclical step-sizes

Baptiste Goujaud, Damien Scieur, Aymeric Dieuleveut, Adrien Taylor, Fabian Pedregosa



Google Research All-Russian seminar on optimization, 2022-03-16

Preprint: <https://arxiv.org/pdf/2106.09687.pdf>

HeavyBall

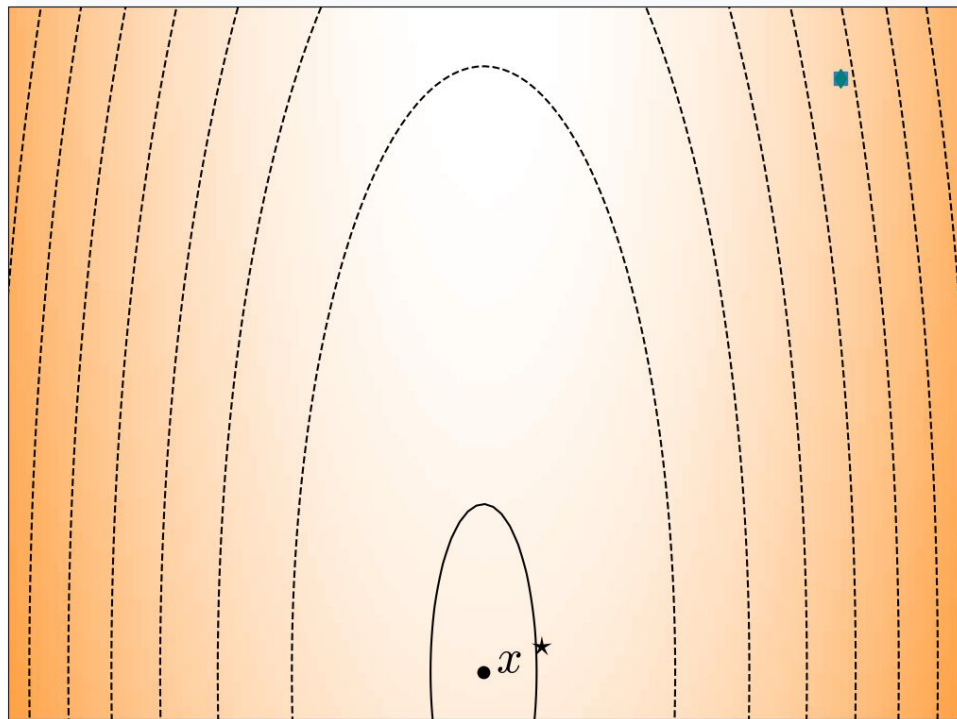
aka gradient descent with momentum

Two parameters; step-size $h > 0$ and momentum $m \in (0, 1)$

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

Optimal among gradient-based methods on quadratics.

Stochastic variant popular in deep learning.



- Gradient Descent
- ◆ HeavyBall

Cyclical HeavyBall

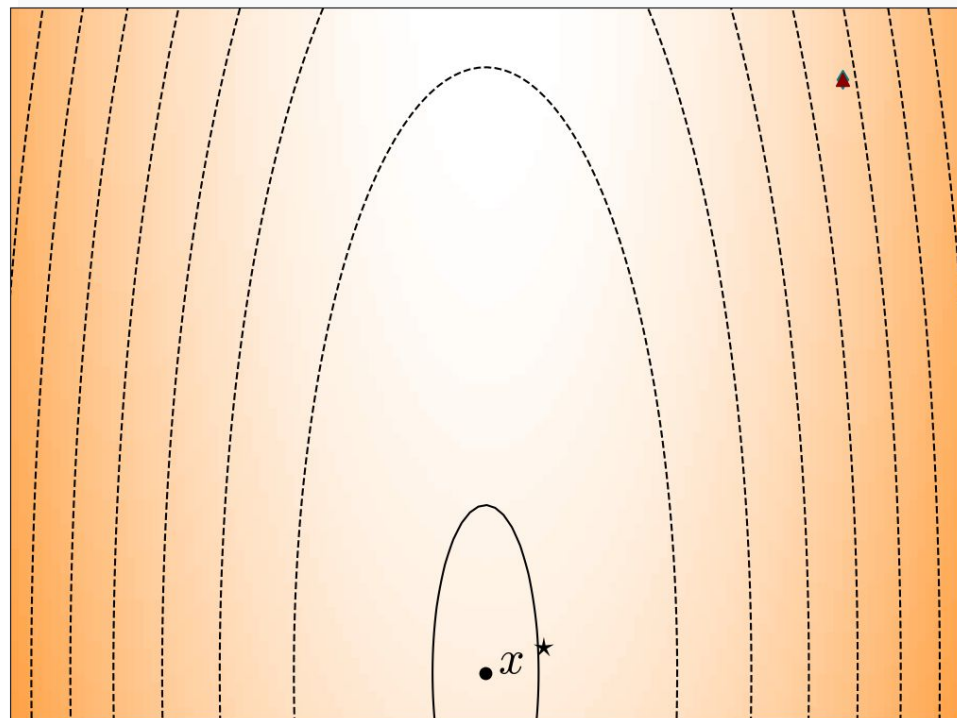
Alternates between two step-sizes h_0 and h_1

$$\begin{aligned} \text{Set } h_t &= h_0 \text{ if } t \text{ is odd and } h_t = h_1 \text{ otherwise} \\ \mathbf{x}_{t+1} &= \mathbf{x}_t - h_t \nabla f(\mathbf{x}_t) + m(\mathbf{x}_t - \mathbf{x}_{t-1}) \end{aligned}$$

Reported faster convergence ([Loshchilov and Hutter, 2017](#); [Smith, 2017](#))

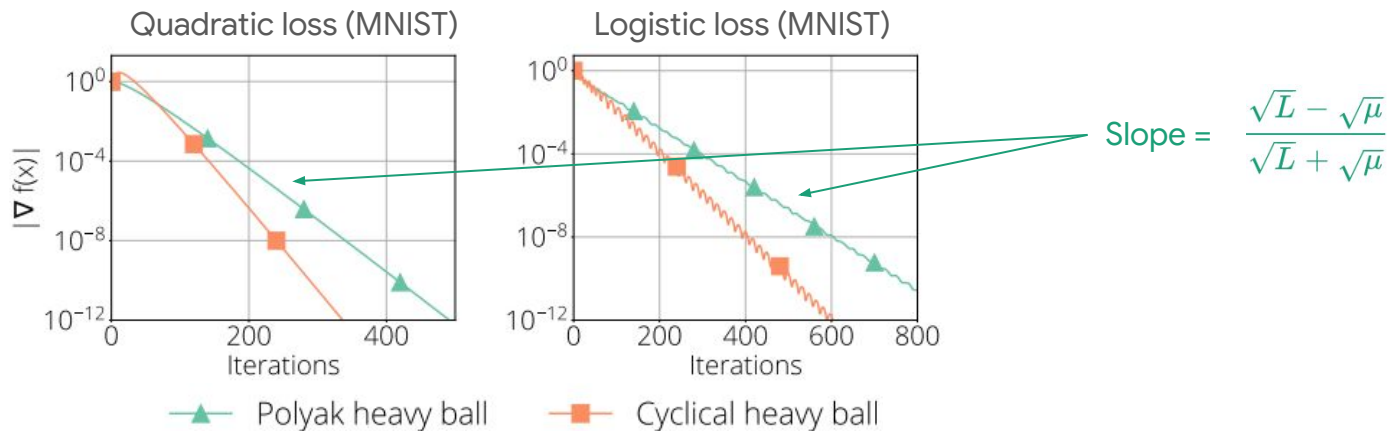
Pervasive (TF, PyTorch, optax, etc.)

No analysis that explains why/when it works.



—◆— HeavyBall
—▲— Cyclical HeavyBall

Benchmarks

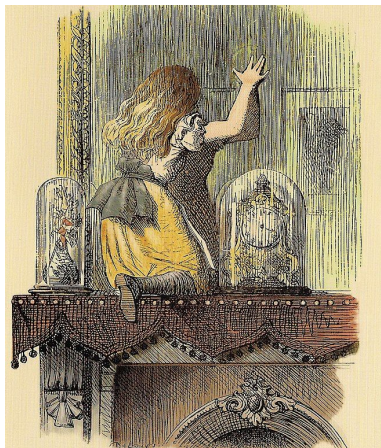


What is the slope of **cyclical heavy ball**?

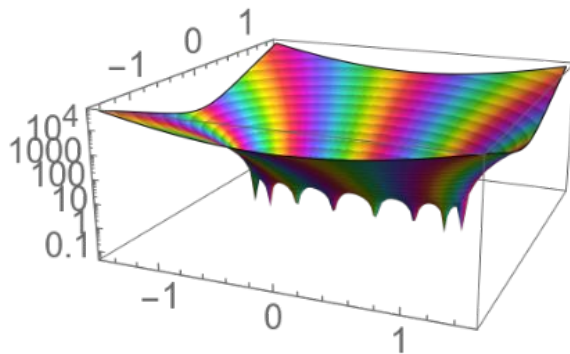
What are the optimal parameters?



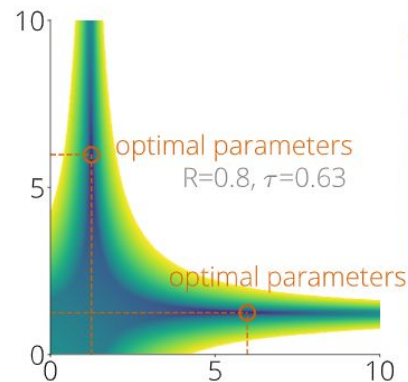
Today's topic



Optimization and
Polynomials



Cyclical
HeavyBall

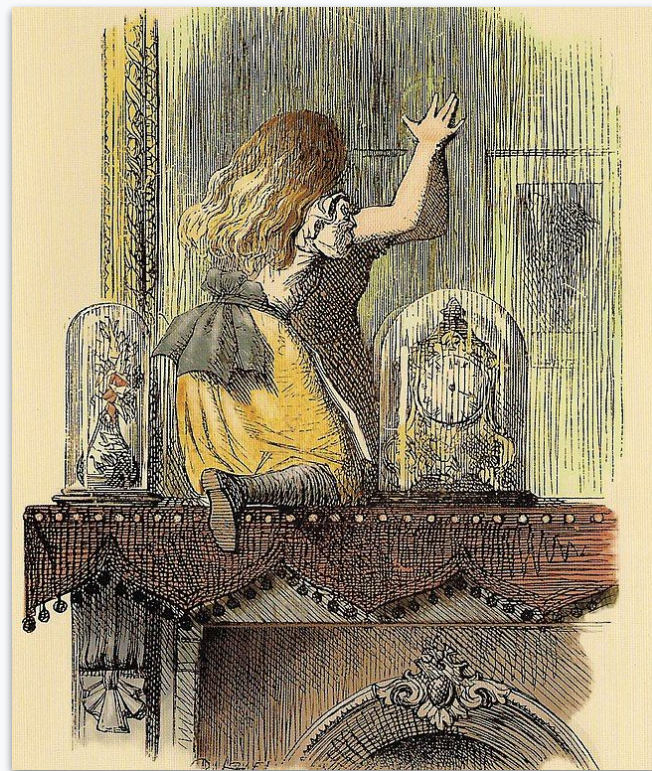


Simulations &
Open problems

Polynomials and optimization

Some problems can be posed in the space of polynomials.

Exploited in early numerical analysis [[Hestenes and Stiefel \(1952\)](#), [Rutishauser \(1959\)](#)]



Polynomials and Optimization

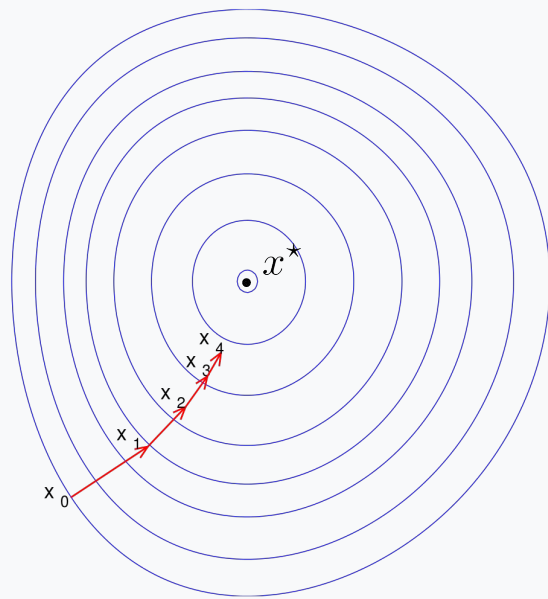
Consider Gradient Descent on

$$f(x) = \frac{1}{2}(x - x^*)H(x - x_*)$$

Then at iteration t we have

$$\begin{aligned}x_{t+1} - x_* &= x_t - x_* - \gamma H(x_t - x_*) \\&= (I - \gamma H)(x_t - x_*) \\&= \dots \\&= (I - \gamma H)^{t+1}(x_0 - x_*)\end{aligned}$$

Polynomial in H



Real-valued polynomials

Taking norms on the previous expression

Cauchy-Schwarz

$$\|\mathbf{x}_{t+1} - \mathbf{x}^*\|_2 \leq \left\| \left(\mathbf{I} - \frac{2}{L+\mu} \mathbf{H} \right)^{t+1} \right\|_2 \|\mathbf{x}_0 - \mathbf{x}^*\|_2$$

Matrix 2-norm

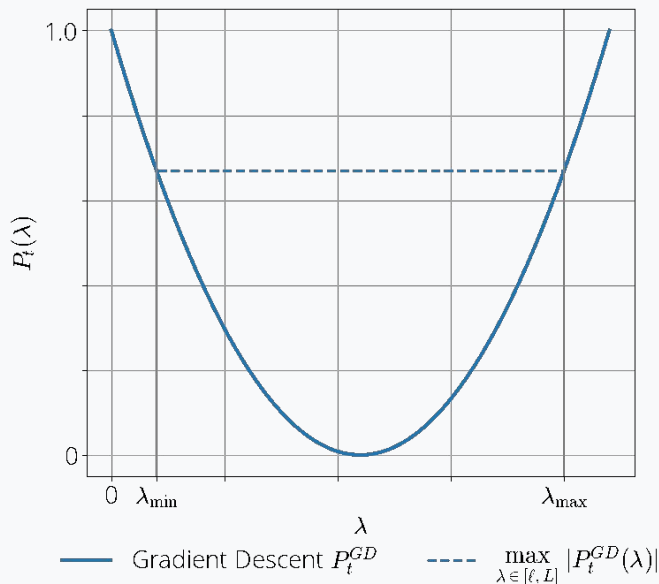
$$\leq \max_{\lambda \in [\mu, L]} \left| \left(1 - \frac{2}{L+\mu} \lambda \right)^{t+1} \right| \|\mathbf{x}_0 - \mathbf{x}^*\|_2$$

Convergence rate $\left(\frac{L - \mu}{L + \mu} \right)^t$



L, μ = largest and smallest eigenvalue of \mathbf{H}

The residual polynomial P_t^{GD} , with $t = 2$



Gradient-based Methods and Polynomials

Corollary (Convergence rate) *Let μ and L be the smallest and largest eigenvalue of \mathbf{H} respectively. Then for any gradient-based method with residual polynomial P_t , we have*

$$\|\mathbf{x}_t - \mathbf{x}^*\| \leq \underbrace{\max_{\lambda \in [\mu, L]}}_{\text{conditioning}} \underbrace{|P_t(\lambda)|}_{\text{algorithm}} \underbrace{\|\mathbf{x}_0 - \mathbf{x}^*\|}_{\text{initialization}}. \quad (17)$$

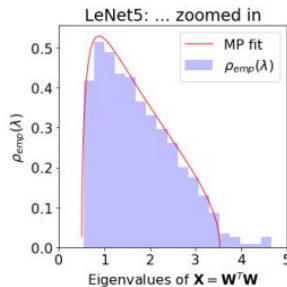
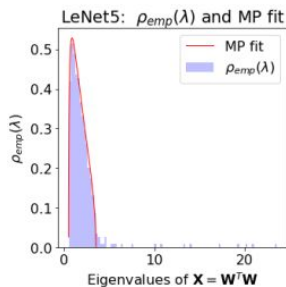
- **Problem difficulty** enters through $[\mu, L]$, interval that contains Hessian eigenvalues.
- **Algorithm** enters through polynomial P_t . This polynomial verifies $P_t(0)=1$

Parenthesis: average-case analysis

Same technique can be used to derive average-case analysis

$$\mathbb{E}\|x_t - x^*\|^2 = \overbrace{R^2}^{\text{initialization}} \int \underbrace{P_t^2}_{\text{algorithm}} \overbrace{d\mu}^{\text{problem}}.$$

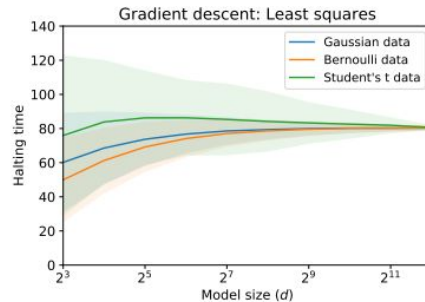
$d\mu$ = density of Hessian eigenvalues



[Martin and Mahoney, 2018](#)

Average-case acceleration [\[P. and Scieur, 2020\]](#)

Concentration [\[Paquette et al. 2022\]](#)



Google Research

HeavyBall

The HeavyBall update

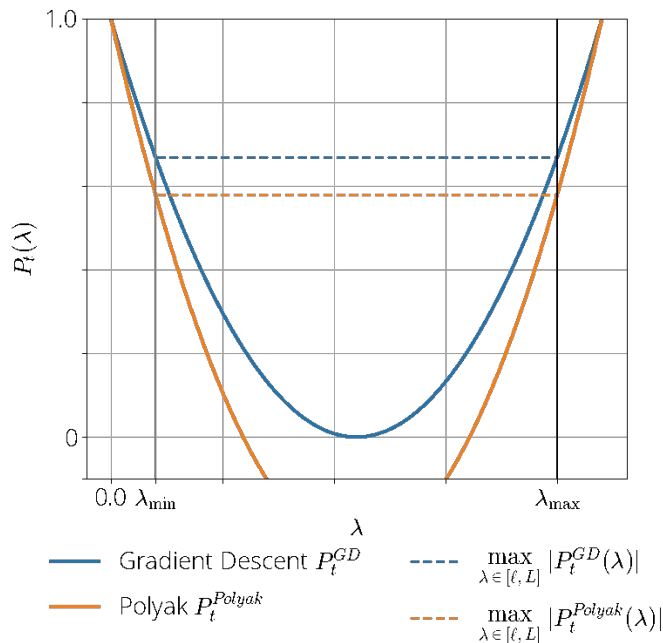
$$\mathbf{x}_{t+1} = \mathbf{x}_t + m(\mathbf{x}_t - \mathbf{x}_{t-1}) - h \nabla f(\mathbf{x}_t)$$

Gives the residual polynomial

$$P_t(\lambda) = m^{t/2} \left(\underbrace{\frac{2m}{1+m} T_t(\sigma(\lambda))}_{\text{Chebyshev 1st kind}} - \underbrace{\frac{m-1}{1+m} U_t(\sigma(\lambda))}_{\text{Chebyshev 2nd kind}} \right)$$

$$\text{with } \sigma(\lambda) = \frac{1}{2\sqrt{m}}(1 + m - h\lambda)$$

The residual polynomial P_t^{Polyak} , with $t = 2$



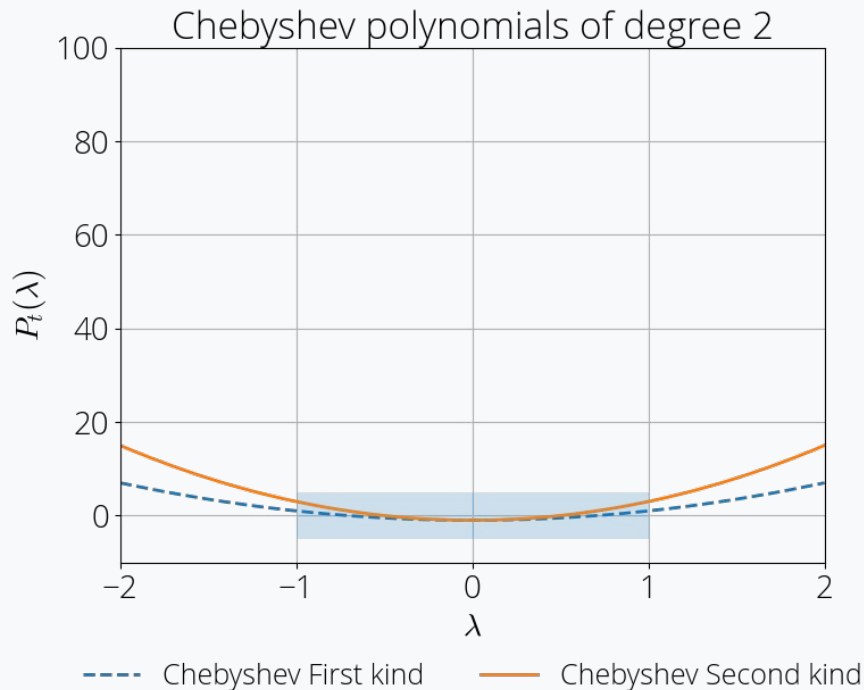
The two faces of Chebyshev polynomials

In the $[-1, 1]$ interval, Chebyshev polynomials are linearly bounded.

$$|T_t(\xi)| \leq 1 \quad \text{and} \quad |U_t(\xi)| \leq t + 1$$

Outside, they grow exponentially.

$$T_t(\xi) = \frac{1}{2} \left(\xi - \sqrt{\xi^2 - 1} \right)^t + \frac{1}{2} \left(\xi + \sqrt{\xi^2 - 1} \right)^t$$
$$U_t(\xi) = \frac{(\xi + \sqrt{\xi^2 - 1})^{t+1} - (\xi - \sqrt{\xi^2 - 1})^{t+1}}{2\sqrt{\xi^2 - 1}}.$$



Link function

$$\sigma(\lambda) = \frac{1}{2\sqrt{m}}(1 + m - h\lambda)$$

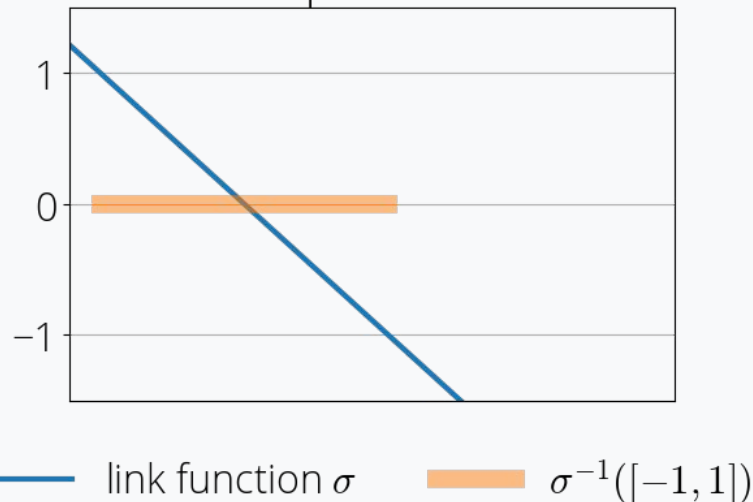
Pre-image is also an interval:

$$\sigma^{-1}([-1, 1]) = \left[\frac{(1 - \sqrt{m})^2}{h}, \frac{(1 + \sqrt{m})^2}{h} \right]$$

Robust region: Parameters for which

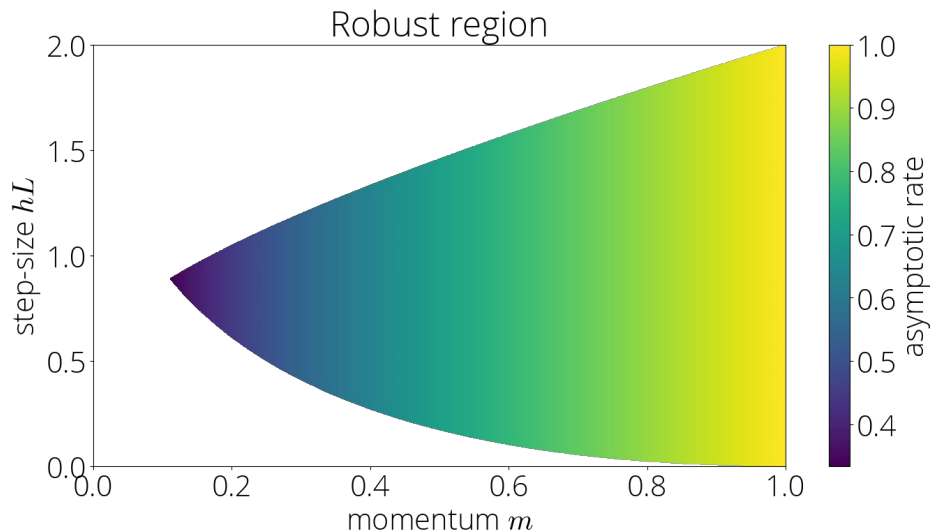
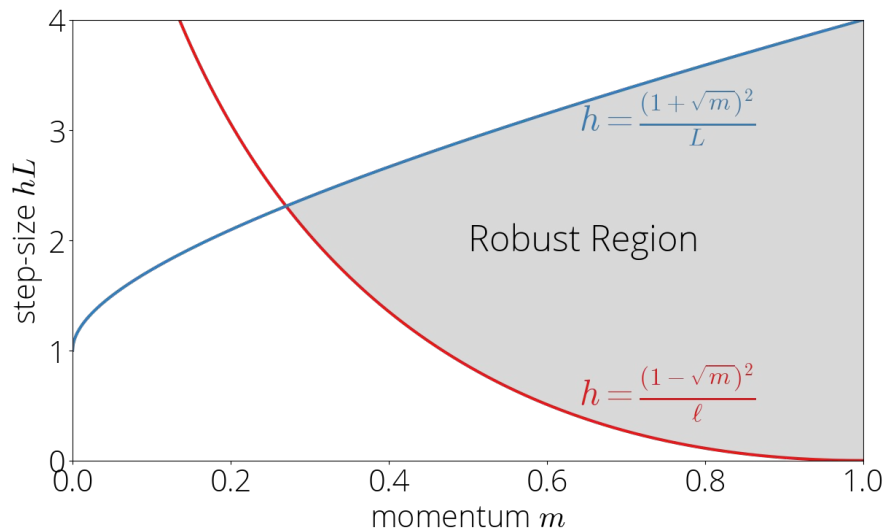
$$[\mu, L] \subseteq \sigma^{-1}([-1, 1])$$

Constant step-size link function



The asymptotic rate in the robust region is \sqrt{m} .

$$\equiv \|x_t - x_\star\| = \mathcal{O}(\sqrt{m}^t)$$



Optimal parameters (aka Polyak HeavyBall)

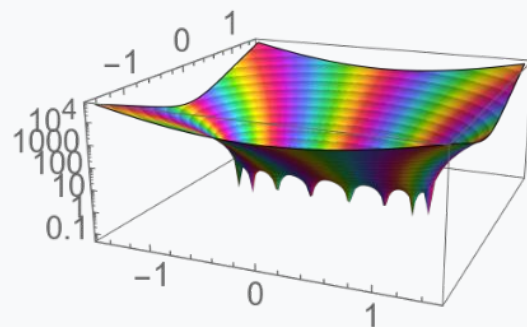
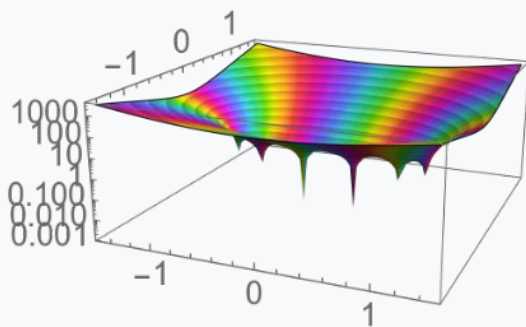
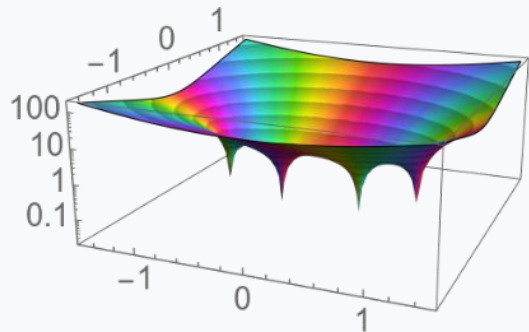
Minimizing m in the robust region results in (worst-case) optimal params

$$m = \left(\frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}} \right)^2$$

$$h = \left(\frac{2}{\sqrt{L} + \sqrt{\mu}} \right)^2$$

Asymptotic convergence rate: $\sqrt{m} = \frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}}$

2. Cyclical HeavyBall



Cyclical HeavyBall

Alternates between 2 step-sizes

Set $h_t = h_0$ if t is odd and $h_t = h_1$ otherwise
 $\mathbf{x}_{t+1} = \mathbf{x}_t - h_t \nabla f(\mathbf{x}_t) + m(\mathbf{x}_t - \mathbf{x}_{t-1})$

Analysis of Cyclical HeavyBall

- Coefficients in recurrence now depends on t

$$P_{2t+1}(\lambda) = (1 + m + h_1 \lambda) P_{2t}(\lambda) - m P_{2t-1}(\lambda)$$

$$P_{2t+2}(\lambda) = (1 + m + h_0 \lambda) P_{2t+1}(\lambda) - m P_{2t}(\lambda)$$

- Known in the OP field as "orthogonal polynomials with varying coefficients" [[Chihara \(1968\)](#), [Van Assche \(1985\)](#)]



TS Chihara

Better to
chain
iterations!

- Analyzed by chaining iterations:

$$P_{2t+2}(\lambda) = ((1 + m + h_0 \lambda)(1 + m + h_1 \lambda) - 2m) P_{2t}(\lambda) - m^2 P_{2t-2}(\lambda)$$

Cyclical step-sizes

The residual polynomial for the cyclical HeavyBall method at even iterations is

$$P_{2t}(\lambda) = m^t \left(\frac{2m}{1+m} T_{2t}(\zeta(\lambda)) + \frac{1-m}{1+m} U_{2t}(\zeta(\lambda)) \right), \quad (6)$$

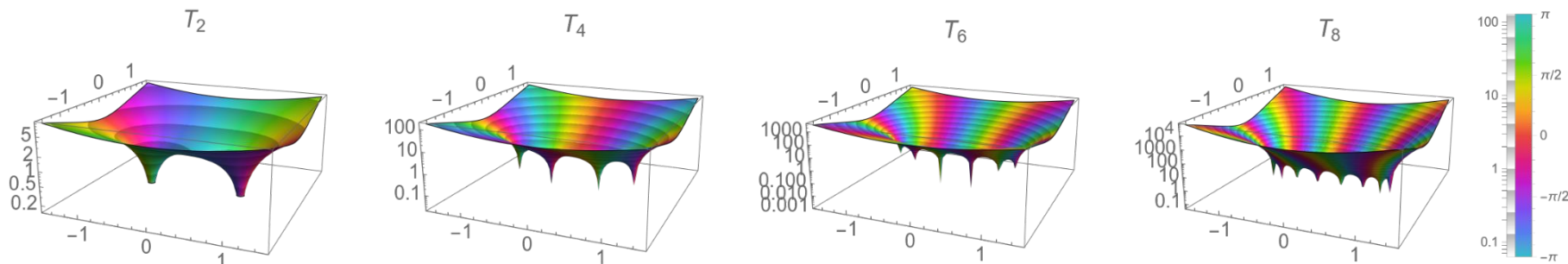
$$\text{with } \zeta(\lambda) = \frac{1+m}{2\sqrt{m}} \sqrt{\left(1 - \frac{h_0}{1+m} \lambda\right) \left(1 - \frac{h_1}{1+m} \lambda\right)}.$$



Same than HeavyBall except for link function ζ

Complex Chebyshev polynomials

Image of link function can now be real or imaginary



Chebyshev polynomials grow exponentially in $\mathbb{C} \setminus [-1, 1]$

Link function

Pre-image no longer interval.

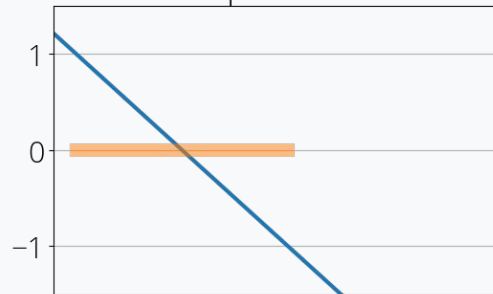
union of two intervals

Robust region: If

$$[\mu, L] \subseteq \sigma^{-1}([-1, 1])$$

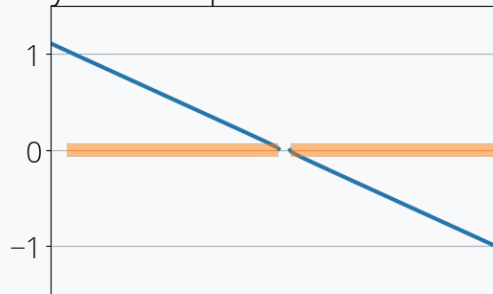
Then $\|x_t - x_\star\| = \mathcal{O}(\sqrt{m}^t)$

Constant step-size link function



— link function σ $\sigma^{-1}([-1, 1])$

Cyclical step-size link function



— link function ζ $\zeta^{-1}([-1, 1])$

A finer model for the Hessian eigenvalues

Consider eigenvalues in union of two disjoint intervals

$$\Lambda = [\mu_1, L_1] \cup [\mu_2, L_2], \quad \underbrace{L_1 - \mu_1 = L_2 - \mu_2}_{\text{same size}}$$

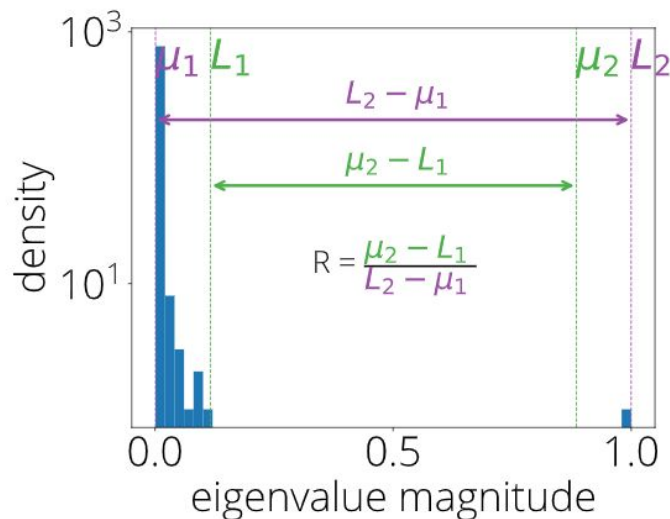
The ratio R will play an important role:

$$R \triangleq \frac{\mu_2 - L_1}{L_2 - \mu_1}$$

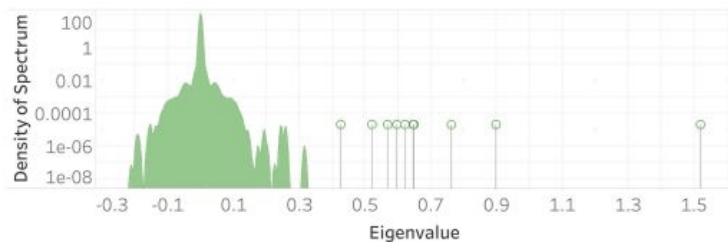
$R = 0$, one interval

$R=1$, all eigenvalues are at extremes.

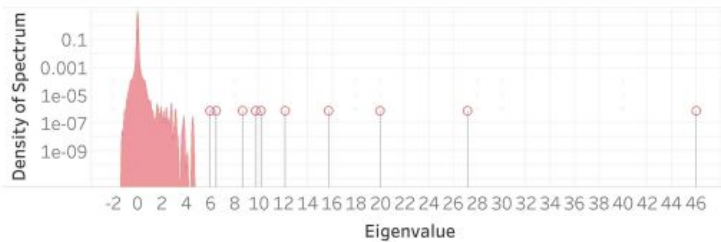
MNIST



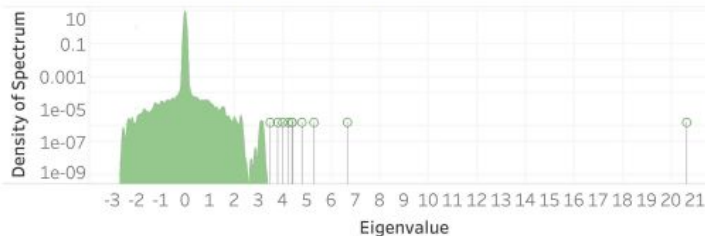
Eigengaps Everywhere



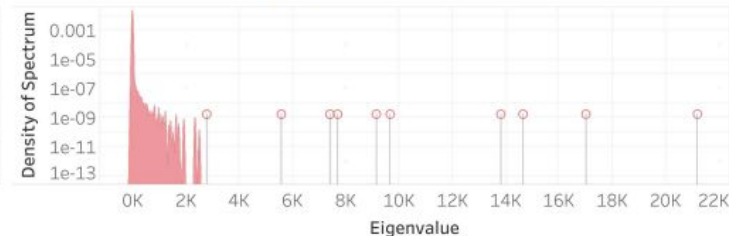
(a) MNIST, train



(b) MNIST, test



(e) CIFAR10, train



(f) CIFAR10, test

[\(Papayan 2020\)](#)

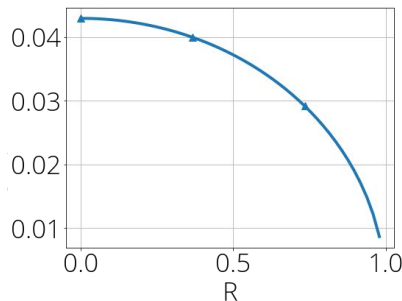
Optimal parameters

Minimize m s.t. $\underbrace{[\mu, L] \subseteq \sigma^{-1}([-1, 1])}_{\text{robust region}}$

Solution

$$m = \left(\frac{\sqrt{\rho^2 - R^2} - \sqrt{\rho^2 - 1}}{\sqrt{1 - R^2}} \right)^2 \quad \text{with} \quad \rho \stackrel{\text{def}}{=} \frac{L_2 + \mu_1}{L_2 - \mu_1}$$

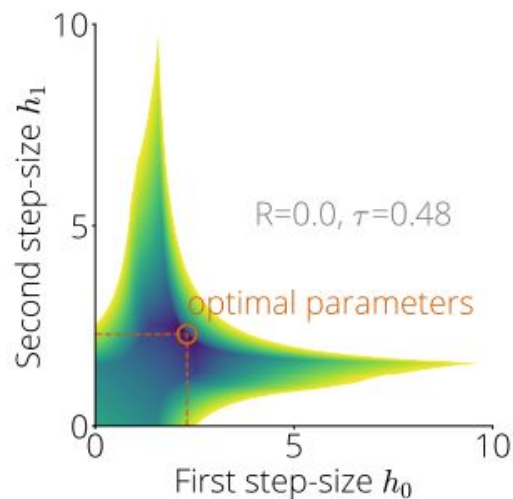
- $R=0$ we recover Polyak HeavyBall
- Decreasing in R



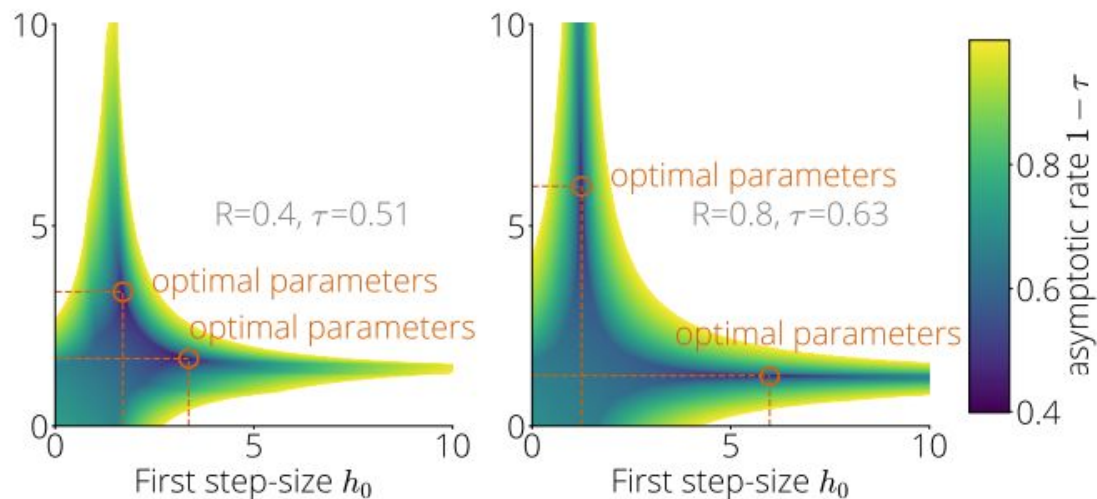
Optimal step-sizes

$$h_t = \frac{1+m}{L_1} \text{ if } t \text{ is odd and } h_t = \frac{1+m}{\mu_2} \text{ otherwise}$$

Gap is small



Gap is large

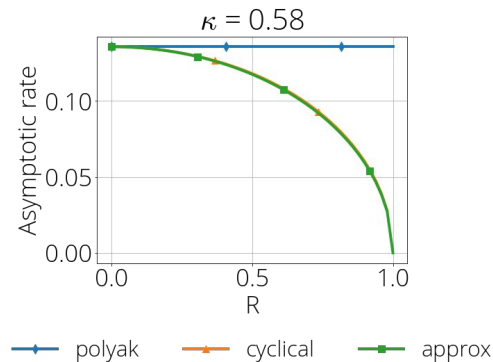


Convergence Rates

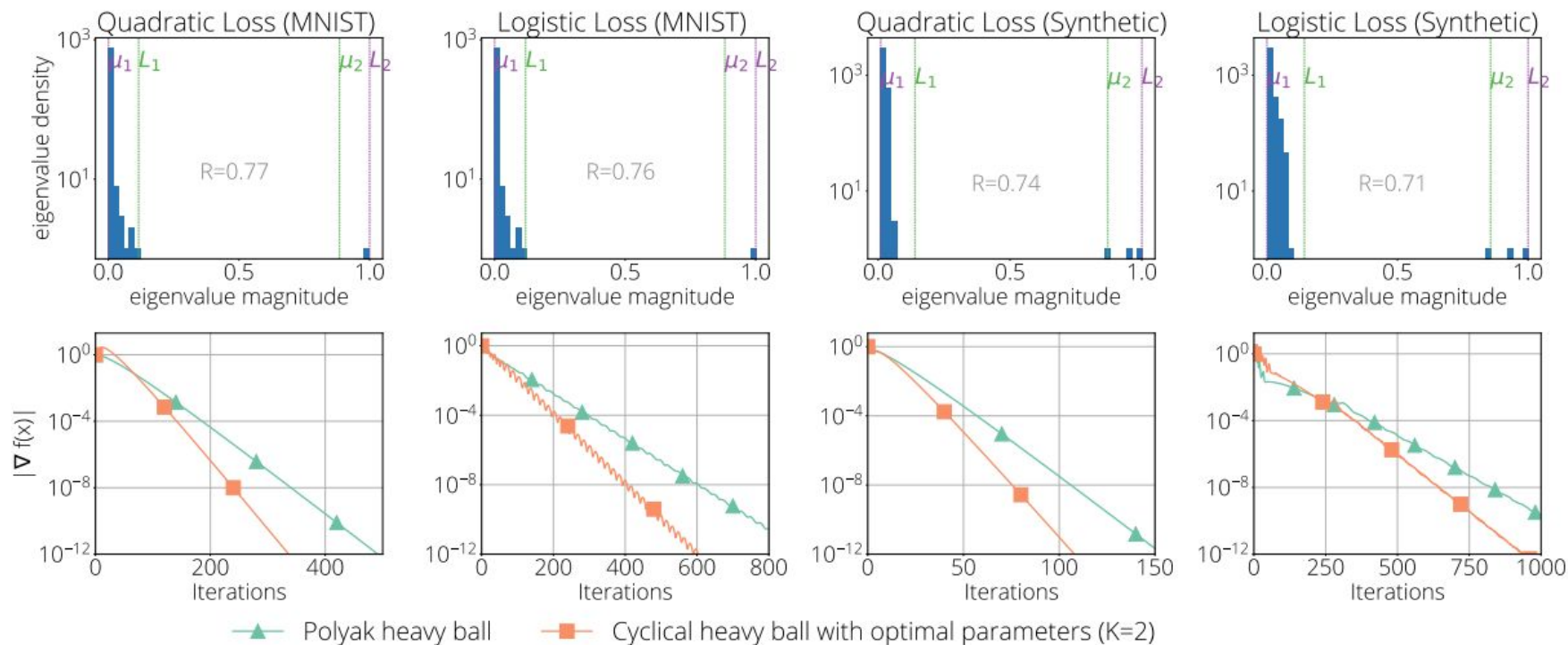
$$\text{Asymptotic rate} = \sqrt{m} = \frac{\sqrt{\rho^2 - R^2} - \sqrt{\rho^2 - 1}}{\sqrt{1 - R^2}}$$

For ill-conditioned problems ($\mu \ll L$),

$$\sqrt{m} \approx \sqrt{1 - R^2} r^{\text{Polyak}}$$

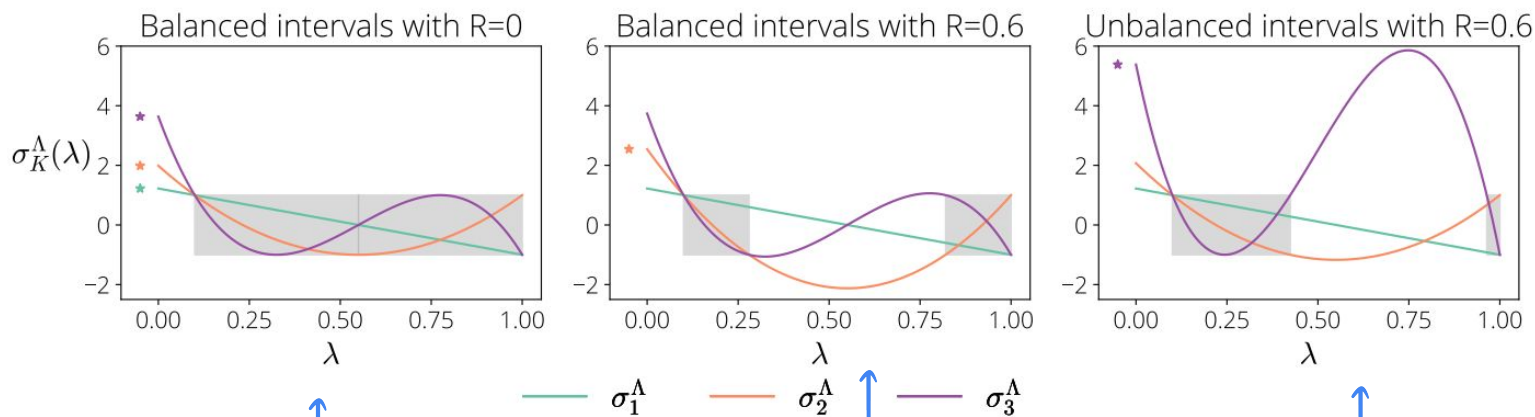


Benchmarks



Beyond cycles of length 2

Link functions are optimal if $2\zeta - 1$ hit ± 1 at edges and $\notin [-1, 1]$ outside



1 cycle is optimal, 2/3 cycles are optimal, 3 cycles is optimal

Conclusions

Cyclical Heavy Ball converges faster in the presence of spectral gap.

Assuming knowledge of this gap, converges at a rate $\approx \sqrt{1 - R^2} r^{\text{Polyak}}$

Speedup observed also on non-quadratic objectives.

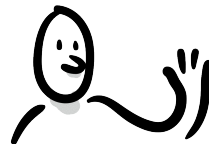
Open Problems

More complex Hessian support: closed form for larger cycles.

Interpolating step-sizes

How to estimate the eigen-gap?

Stochastic algorithm? Non-quadratic objectives?



Local convergence for non-quadratics

Theorem 5.1 (Local convergence). *Let $f : \mathbb{R}^d \mapsto \mathbb{R}$ be a (potentially non-quadratic) twice continuously differentiable function, x_* a local minimizer, and H be the Hessian of f at x_* with $\text{Sp}(H) \subseteq \Lambda$. Let x_t denote the result of running Algorithm 1 with parameters h_1, h_2, \dots, h_K, m , and let $1 - \tau$ be the linear convergence rate on the quadratic objective (OPT). Then we have*

$$\forall \varepsilon > 0, \exists \text{ open set } V_\varepsilon : x_0, x_* \in V_\varepsilon \implies \|x_t - x_*\| = O((1 - \tau + \varepsilon)^t) \|x_0 - x_*\|. \quad (24)$$

Potential extension: [A Modular Analysis of Provable Acceleration via Polyak's Momentum: Training a Wide ReLU Network and a Deep Linear Network](#), Jun-Kun Wang, Chi-Heng Lin, Jacob Abernethy

3. Acceleration & Simulations

