Super-acceleration with cyclical step-sizes

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









Google Research

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

HeavyBall

aka gradient descent with momentum

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

$$oldsymbol{x}_{t+1} = oldsymbol{x}_t + oldsymbol{m}(oldsymbol{x}_t - oldsymbol{x}_{t-1}) - oldsymbol{h}
abla f(oldsymbol{x}_t)$$

Optimal among gradient-based methods on quadratics.

Stochastic variant popular in deep learning.







Cyclical HeavyBall

Alternates between two step-sizes h_o and h₁

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

Reported faster convergence (<u>Loshchilov and</u> <u>Hutter, 2017; Smith, 2017</u>)

Pervasive (TF, PyTorch, optax, etc.)

No analysis that explains why/when it works.





Benchmarks



What is the slope of cyclical heavy ball? What are the optimal parameters?









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

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

Then at iteration t we have

$$x_{t+1} - x_{\star} = x_t - x_{\star} - \gamma H(x_t - x_{\star})$$
$$= (I - \gamma H)(x_t - x_{\star})$$
$$= \dots$$
$$= (I - \gamma H)^{t+1}(x_0 - x_{\star})$$

Polynomial in H



Real-valued polynomials

Taking norms on the previous expression

Cauchy-Schwarz

 $\|oldsymbol{x}_{t+1} - oldsymbol{x}^\star\|_2 \leq \|(oldsymbol{I} - rac{2}{L+\mu}oldsymbol{H})^{t+1}\|_2\|oldsymbol{x}_0 - oldsymbol{x}^\star\|_2$

Matrix 2-norm

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

 $0 \xrightarrow{\lambda}_{min} \xrightarrow{\lambda}_{max}$ Gradient Descent $P_t^{GD} \xrightarrow{----} \max_{\lambda \in [t, L]} |P_t^{GD}(\lambda)|$

Google Research

L, μ = largest and smallest eigenvalue of H

Gradient-based Methods and Polynomials

Corollary (Convergence rate) Let μ and L be the smallest and largest eigenvalue of \boldsymbol{H} respectively. Then for any gradient-based method with residual polynomial P_t , we have $\|\boldsymbol{x}_t - \boldsymbol{x}^{\star}\| \leq \max_{\lambda \in [\mu, L]} \|P_t(\lambda)\| \|\boldsymbol{x}_0 - \boldsymbol{x}^{\star}\|.$ (17)

conditioning algorithm initialization

 Algorithm enters through polynomial P_t. This polynomial verifies P_t(0)=1

HeavyBall

The HeavyBall update

$$oldsymbol{x}_{t+1} = oldsymbol{x}_t + oldsymbol{m}(oldsymbol{x}_t - oldsymbol{x}_{t-1}) - oldsymbol{h}
abla f(oldsymbol{x}_t)$$

Gives the residual polynomial

$$P_t(\lambda) = m^{t/2} \left(rac{2m}{1+m} \, T_t(\sigma(\lambda)) - rac{m-1}{1+m} \, U_t(\sigma(\lambda))
ight)$$

Chebyshev 1st kind

$$+m$$

Chebyshev 2nd kind

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



The two faces of Chebyshev polynomials

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

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

Outside, they grow exponentially.

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



Link function

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

Pre-image is also an interval:

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

Robust region: Parameters for which

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





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



A Hitchhiker's Guide to Momentum, http://fa.bianp.net/blog/2021/hitchhiker/

Optimal parameters (aka Polyak HeavyBall)

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

$$m = \Big(rac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}}\Big)^2 \qquad \qquad h = \Big(rac{2}{\sqrt{L} + \sqrt{\mu}}\Big)^2$$

Asymptotic convergence rate:

$$\sqrt{m} = \frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}}$$

2. Cyclical HeavyBall



Cyclical HeavyBall

Alternates between 2 step-sizes

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

Analysis of Cyclical HeavyBall

• Coefficients in recurrence now depends on t

 $egin{aligned} P_{2t+1}(\lambda) &= (1+m+h_1\lambda)P_{2t}(\lambda) - mP_{2t-1}(\lambda) \ P_{2t+2}(\lambda) &= (1+m+h_0\lambda)P_{2t+1}(\lambda) - mP_{2t}(\lambda) \end{aligned}$

- Known in the OP field as "orthogonal polynomials with varying coefficients" [Chihara (1968), Van Assche (1985)]
- Analyzed by chaining iterations:

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

Better to chain iterations!





Google Research

Cyclical step-sizes

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

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

Link function

Pre-image no longer interval.

union of two intervals

Robust region: If $[\mu, L] \subseteq \sigma^{-1}([-1, 1])$ Then $||x_t - x_\star|| = \mathcal{O}(\sqrt{m}^t)$



A finer model for the Hessian eigenvalues

Consider eigenvalues in union of two

disjoint intervals

$$\Lambda = [\boldsymbol{\mu}_1, \, \boldsymbol{L}_1] \cup [\boldsymbol{\mu}_2, \, \boldsymbol{L}_2], \underbrace{\boldsymbol{L}_1 - \boldsymbol{\mu}_1 = \boldsymbol{L}_2 - \boldsymbol{\mu}_2}_{\text{same size}}.$$

The ratio *R* will play an important role:

$$R riangleq rac{oldsymbol{\mu}_2 - oldsymbol{L}_1}{oldsymbol{L}_2 - oldsymbol{\mu}_1}$$

R = 0, one interval R=1, all eigenvalues are at extremes.

 $\begin{array}{c} \begin{array}{c} \mu_{1} \ L_{1} \\ \mu_{2} - \mu_{1} \\ \mu_{2} - L_{1} \\ R = \frac{\mu_{2} - L_{1}}{L_{2} - \mu_{1}} \\ 0.0 \\ 0.5 \end{array}$

 10^{3}

0.5 1.0 eigenvalue magnitude

Google Research

U2 L2

MNIST

Eigengaps Everywhere





(Papyan 2020)

Optimal parameters Minimize **m** s.t. $[\mu, L] \subseteq \sigma^{-1}([-1, 1])$ robust region $m = \left(\frac{\sqrt{\rho^2 - R^2} - \sqrt{\rho^2 - 1}}{\sqrt{1 - R^2}}\right)^2 \quad \text{with} \quad \rho \stackrel{\text{def}}{=} \frac{L_2 + \mu_1}{L_2 - \mu_1}$

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



Optimal step-sizes



Convergence Rates

Asymptotic rate =
$$\sqrt{m} = rac{\sqrt{
ho^2 - R^2} - \sqrt{
ho^2 - 1}}{\sqrt{1 - R^2}}$$



Benchmarks



Google Research

Beyond cycles of length 2

Link functions re optimal if 2ζ -1 hit \pm 1 at edges and \notin [-1, 1] outside



Conclusions

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

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

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?