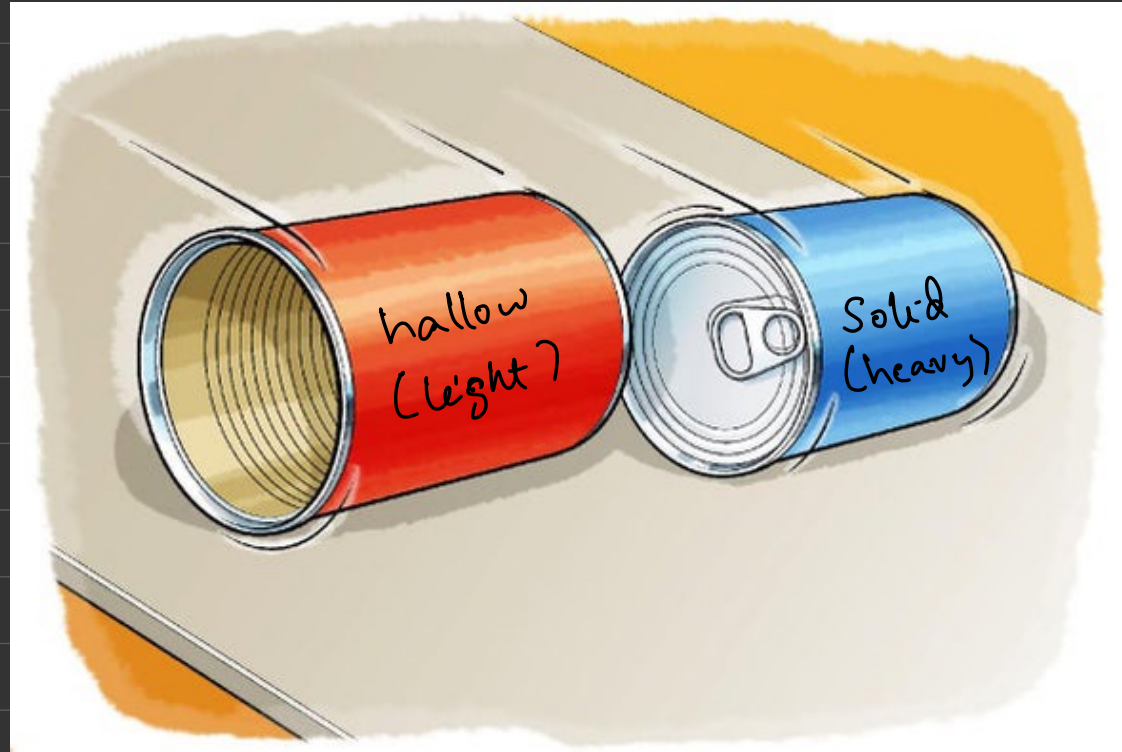


- Heavy-ball method
  - Quadratic function
  - iteration complexity  $O\left(\sqrt{k} \log\left(\frac{1}{\epsilon}\right)\right)$
- Nesterov's acceleration  
(just the idea today)



In which object goes down the ramp faster?

- added inertia acts as a smoother and an accelerates

# Heavy-ball method or Polyak momentum

B.T. Polyak (1964)

"Some methods of speeding up convergence of iteration methods."

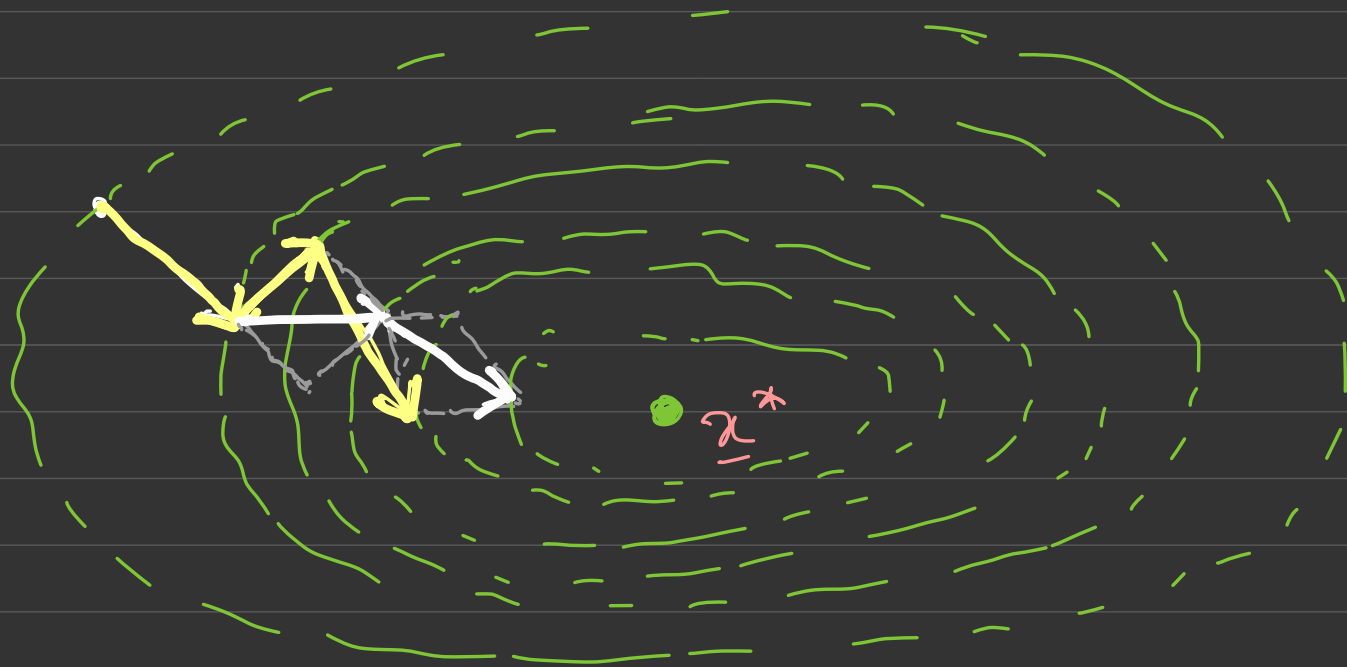
minimize  $f(\underline{x})$

$$\underline{x} \in \mathbb{R}^n$$

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

momentum term

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



→ gradient descent

→ heavy-ball method

Quadratic problem:

$$\underset{\underline{x}}{\text{minimize}} \quad \frac{1}{2} (\underline{x} - \underline{x}^*)^T Q (\underline{x} - \underline{x}^*)$$

Recall  $\nabla f(\underline{x}) = Q (\underline{x} - \underline{x}^*)$

Gradient descent:

$$\| \underline{x}_t - \underline{x}^* \|_2 \leq \left( \frac{\kappa - 1}{\kappa + 1} \right)^t \| \underline{x}_0 - \underline{x}^* \|_2$$

iteration complexity:  $O\left(\kappa \log\left(\frac{1}{\epsilon}\right)\right)$

Heavy-ball state-space model:

$$\begin{bmatrix} \bar{x}_{t+1} \\ \bar{x}_t \end{bmatrix} = \begin{bmatrix} (1 + \theta_t) \mathbf{I} & -\theta_t \mathbf{I} \\ \mathbf{I} & 0 \end{bmatrix} \begin{bmatrix} \bar{x}_t \\ \bar{x}_{t-1} \end{bmatrix} - \begin{bmatrix} \eta_t \nabla f(\bar{x}_t) \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} \bar{x}_{t+1} - \bar{x}^* \\ \bar{x}_t - \bar{x}^* \end{bmatrix} = \begin{bmatrix} (1 + \theta_t) \mathbf{I} & -\theta_t \mathbf{I} \\ \mathbf{I} & 0 \end{bmatrix} \begin{bmatrix} \bar{x}_t - \bar{x}^* \\ \bar{x}_{t-1} - \bar{x}^* \end{bmatrix} - \begin{bmatrix} \eta_t \nabla f(\bar{x}_t) \\ 0 \end{bmatrix}$$

$$= \begin{bmatrix} (1 + \theta_t) \mathbf{I} - \eta_t \mathbf{Q} & -\theta_t \mathbf{I} \\ \mathbf{I} & 0 \end{bmatrix} \begin{bmatrix} \bar{x}_t - \bar{x}^* \\ \bar{x}_{t-1} - \bar{x}^* \end{bmatrix}$$

Convergence depends on  
the spectrum of :  $H_t$

Find  $\eta_t$  and  $\theta_t$  to control the spectrum of  $H_t$

- Suppose  $\lambda_i$  is the  $i$ th eigenvalue of  $Q$
- Spectral radius of  $H_t$

$$\rho(H_t) = \rho \left( \begin{bmatrix} (1+\theta_t)\mathbb{I} - \eta_t \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} & -\theta_t \mathbb{I} \\ \mathbb{I} & 0 \end{bmatrix} \right)$$

$$\leq \max_{1 \leq i \leq n} \rho \left( \begin{bmatrix} 1 + \theta_t - \eta_t \lambda_i & -\theta_t \\ \mathbb{1} & 0 \end{bmatrix} \right)$$

two eigenvalues are

roots of  $s^2 - (1 + \theta_t - \eta_t \lambda_i) s + \theta_t = 0$

• If  $(1 + \theta_t - \eta_t \lambda_i)^2 - 4\theta_t \leq 0$ , then roots have the same magnitude  $\sqrt{\theta_t}$  (imaginary roots)

•  $(1 + \theta_t - \eta_t \lambda_i)^2 \leq 4\theta_t$  is satisfied if  $\theta_t \in [(1 - \sqrt{\eta_t \lambda_i})^2, (1 + \sqrt{\eta_t \lambda_i})^2]$

Suppose  $f$  is a  $L$ -smooth and  $\mu$  strongly convex quadratic function

$$0 \leq \mu I \leq \mathcal{Q} \leq L I$$

$$\theta_t = \max \left\{ (1 - \sqrt{\eta_t L})^2, (1 - \sqrt{\eta_t \mu})^2 \right\}$$

still yields  $(1 + \theta_t - \eta_t \lambda_i)^2 \leq 4\theta_t$



$$\Rightarrow e(H_t) \leq \sqrt{\theta_t}$$

Now,

$$1 - \sqrt{n_t L} = - (1 - \sqrt{n_t \mu}) \Rightarrow n_t = \frac{4}{(\sqrt{L} + \sqrt{\mu})^2}$$

$$\Rightarrow \theta_t = \max \left\{ \left( 1 - \frac{2\sqrt{L}}{\sqrt{L} + \sqrt{\mu}} \right)^2, \left( 1 - \frac{2\sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}} \right)^2 \right\}$$

$$= \left( \frac{\sqrt{k} - 1}{\sqrt{k} + 1} \right)^2$$

$$\Rightarrow e(H_t) \leq \frac{\sqrt{k} - 1}{\sqrt{k} + 1}$$

## Theorem:

Suppose  $f$  is a  $L$ -smooth and  $\mu$ -strongly convex quadratic function. Then  $\eta_t = \frac{4}{(\sqrt{L} + \sqrt{\mu})^2}$  and

$$\theta_t = \max \left\{ \left| 1 - \sqrt{\eta_t L} \right|, \left| 1 - \sqrt{\eta_t \mu} \right| \right\}^2 \quad \text{with } k = \frac{L}{\mu}$$

results in

$$\left\| \begin{bmatrix} \underline{x}_{t+1} - \underline{x}^* \\ \underline{x}_t - \underline{x}^* \end{bmatrix} \right\|_2 \leq \left( \frac{\sqrt{k} - 1}{\sqrt{k} + 1} \right)^t \left\| \begin{bmatrix} \underline{x}_1 - \underline{x}^* \\ \underline{x}_0 - \underline{x}^* \end{bmatrix} \right\|_2$$

Recall  $\exp(-x) \geq 1 - x$

$$\exp\left(-\frac{2}{\sqrt{k+1}}\right) \geq 1 - \frac{2}{\sqrt{k+1}}$$

$$\Rightarrow \left\| \begin{bmatrix} \bar{x}_{t+1} - \bar{x}^* \\ \bar{x}_t - \bar{x}^* \end{bmatrix} \right\|_2^2 \leq \exp\left(-\frac{4t}{\sqrt{k+1}}\right) \left\| \begin{bmatrix} \bar{x}_1 - \bar{x}^* \\ \bar{x}_0 - \bar{x}^* \end{bmatrix} \right\|_2^2$$

Iteration complexity:  $O\left(\sqrt{k} \log\left(\frac{1}{\epsilon}\right)\right)$

Suppose condition number is  $k = 100$

$$10 \log\left(\frac{1}{\epsilon}\right) \lesssim 100 \log\left(\frac{1}{\epsilon}\right)$$

How about more general convex functions?

# Nesterov's method :

Yuri Nesterov 1983

$$\underline{y}_{t+1} = \underline{x}_t - \eta \nabla f(\underline{x}_t)$$

$$\underline{x}_{t+1} = \underline{y}_{t+1} + \frac{t}{t+3} (\underline{y}_{t+1} - \underline{y}_t)$$

$$\underline{y}_{t+1} = \underline{x}_t - \eta_t \nabla f(\underline{x}_t) \quad : \text{normal}$$

$$\underline{z}_{t+1} = \underline{z}_t - \eta_t \cdot \frac{t+1}{2} \nabla f(\underline{x}_t) \quad : \text{aggressive}$$

$$\underline{x}_{t+1} = \frac{t+1}{t+3} \underline{y}_{t+1} + \frac{2}{t+3} \underline{z}_{t+1} \quad : \text{average}$$

lower weight