

# **Accelerated gradient methods**

**Ye Shi**

ShanghaiTech University

# Outline

---

- Heavy-ball methods
- Nesterov's accelerated gradient methods
- Accelerated proximal gradient methods (FISTA)
- Convergence analysis
- Lower bounds

# (Proximal) gradient methods

---

Iteration complexities of (proximal) gradient methods

- strongly convex and smooth problems

$$O\left(\kappa \log \frac{1}{\varepsilon}\right)$$

- convex and smooth problems

$$O\left(\frac{1}{\varepsilon}\right)$$

Can one still hope to further accelerate convergence?

# Issues and possible solutions

---

## Issues:

- GD focuses on improving the cost per iteration, which might sometimes be too “short-sighted”
- GD might sometimes zigzag or experience abrupt changes

## Solutions:

- exploit information from the history (i.e. past iterates)
- add buffers (like momentum) to yield smoother trajectory

# **Heavy-ball methods**

— Polyak '64

# Heavy-ball method

---



$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^n} \quad f(\mathbf{x})$$

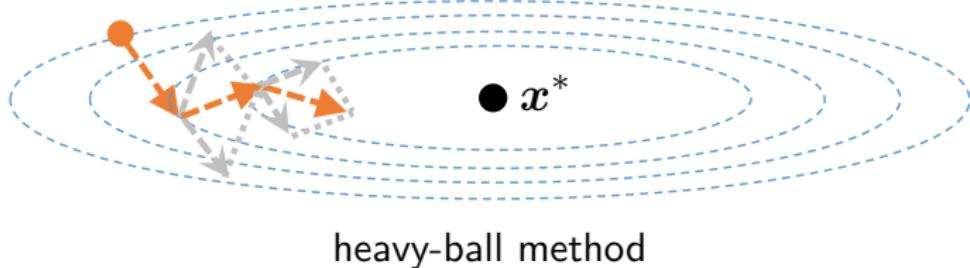
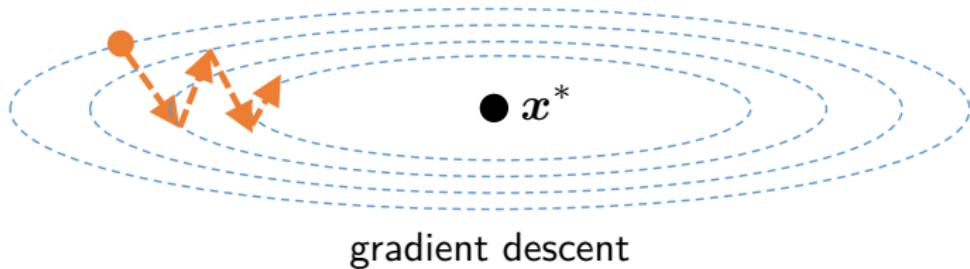
$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \nabla f(\mathbf{x}^t) + \underbrace{\theta_t (\mathbf{x}^t - \mathbf{x}^{t-1})}_{\text{momentum term}}$$

B. Polyak

- add inertia to the “ball” (i.e. include a momentum term) to mitigate zigzagging

# Heavy-ball method

---



# State-space models

---

$$\text{minimize}_{\boldsymbol{x}} \quad \frac{1}{2}(\boldsymbol{x} - \boldsymbol{x}^*)^\top \boldsymbol{Q}(\boldsymbol{x} - \boldsymbol{x}^*)$$

where  $\boldsymbol{Q} \succ \mathbf{0}$  has a condition number  $\kappa$

One can understand heavy-ball methods through dynamical systems

# State-space models

---

Consider the following dynamical system

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

or equivalently,

$$\underbrace{\begin{bmatrix} \mathbf{x}^{t+1} - \mathbf{x}^* \\ \mathbf{x}^t - \mathbf{x}^* \end{bmatrix}}_{\text{state}} = \begin{bmatrix} (1 + \theta_t)\mathbf{I} & -\theta_t\mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{x}^t - \mathbf{x}^* \\ \mathbf{x}^{t-1} - \mathbf{x}^* \end{bmatrix} - \begin{bmatrix} \eta_t \nabla f(\mathbf{x}^t) \\ \mathbf{0} \end{bmatrix}$$
$$= \underbrace{\begin{bmatrix} (1 + \theta_t)\mathbf{I} - \eta_t \mathbf{Q} & -\theta_t\mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{bmatrix}}_{\text{system matrix}} \begin{bmatrix} \mathbf{x}^t - \mathbf{x}^* \\ \mathbf{x}^{t-1} - \mathbf{x}^* \end{bmatrix}$$

# System matrix

---

$$\begin{bmatrix} \mathbf{x}^{t+1} - \mathbf{x}^* \\ \mathbf{x}^t - \mathbf{x}^* \end{bmatrix} = \underbrace{\begin{bmatrix} (1 + \theta_t)\mathbf{I} - \eta_t \mathbf{Q} & -\theta_t \mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{bmatrix}}_{=: \mathbf{H}_t \text{ (system matrix)}} \begin{bmatrix} \mathbf{x}^t - \mathbf{x}^* \\ \mathbf{x}^{t-1} - \mathbf{x}^* \end{bmatrix} \quad (7.1)$$

**implication:** convergence of heavy-ball methods depends on the spectrum of the system matrix  $\mathbf{H}_t$

**key idea:** find appropriate stepsizes  $\eta_t$  and momentum coefficients  $\theta_t$  to control the spectrum of  $\mathbf{H}_t$

# Convergence for quadratic problems

## Theorem 7.1 (Convergence of heavy-ball methods for quadratic functions)

Suppose  $f$  is a  $L$ -smooth and  $\mu$ -strongly convex quadratic function. Set  $\eta_t \equiv 4/(\sqrt{L} + \sqrt{\mu})^2$ ,  $\theta_t \equiv \max \{|1 - \sqrt{\eta_t L}|, |1 - \sqrt{\eta_t \mu}|\}^2$ , and  $\kappa = L/\mu$ . Then

$$\left\| \begin{bmatrix} \mathbf{x}^{t+1} - \mathbf{x}^* \\ \mathbf{x}^t - \mathbf{x}^* \end{bmatrix} \right\|_2 \lesssim \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^t \left\| \begin{bmatrix} \mathbf{x}^1 - \mathbf{x}^* \\ \mathbf{x}^0 - \mathbf{x}^* \end{bmatrix} \right\|_2$$

- iteration complexity:  $O(\sqrt{\kappa} \log \frac{1}{\varepsilon})$
- significant improvement over GD:  $O(\sqrt{\kappa} \log \frac{1}{\varepsilon})$  vs.  $O(\kappa \log \frac{1}{\varepsilon})$
- relies on knowledge of both  $L$  and  $\mu$

## Proof of Theorem 7.1

---

In view of (7.1), it suffices to control the spectrum of  $\mathbf{H}_t$  (which is time-invariant). Let  $\lambda_i$  be the  $i$ th eigenvalue of  $\mathbf{Q}$  and set

$\Lambda := \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix}$ , then the spectral radius (denoted by  $\rho(\cdot)$ ) of  $\mathbf{H}_t$  obeys

$$\begin{aligned}\rho(\mathbf{H}_t) &= \rho\left(\begin{bmatrix} (1 + \theta_t)\mathbf{I} - \eta_t \Lambda & -\theta_t \mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{bmatrix}\right) \\ &\leq \max_{1 \leq i \leq n} \rho\left(\begin{bmatrix} 1 + \theta_t - \eta_t \lambda_i & -\theta_t \\ 1 & 0 \end{bmatrix}\right)\end{aligned}$$

To finish the proof, it suffices to show

$$\max_i \rho\left(\begin{bmatrix} 1 + \theta_t - \eta_t \lambda_i & -\theta_t \\ 1 & 0 \end{bmatrix}\right) \leq \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \quad (7.2)$$

## Proof of Theorem 7.1

---

To show (7.2), note that the two eigenvalues of  $\begin{bmatrix} 1 + \theta_t - \eta_t \lambda_i & -\theta_t \\ 1 & 0 \end{bmatrix}$  are the roots of

$$z^2 - (1 + \theta_t - \eta_t \lambda_i)z + \theta_t = 0 \quad (7.3)$$

If  $(1 + \theta_t - \eta_t \lambda_i)^2 \leq 4\theta_t$ , then the roots of this equation have the same magnitudes  $\sqrt{\theta_t}$  (as they are either both imaginary or there is only one root).

In addition, one can easily check that  $(1 + \theta_t - \eta_t \lambda_i)^2 \leq 4\theta_t$  is satisfied if

$$\theta_t \in [(1 - \sqrt{\eta_t \lambda_i})^2, (1 + \sqrt{\eta_t \lambda_i})^2], \quad (7.4)$$

which would hold if one picks  $\theta_t = \max \{(1 - \sqrt{\eta_t L})^2, (1 - \sqrt{\eta_t \mu})^2\}$

## Proof of Theorem 7.1

---

With this choice of  $\theta_t$ , we have (from (7.3) and the fact that two eigenvalues have identical magnitudes)

$$\rho(\mathbf{H}_t) \leq \sqrt{\theta_t}$$

Finally, setting  $\eta_t = \frac{4}{(\sqrt{L} + \sqrt{\mu})^2}$  ensures  $1 - \sqrt{\eta_t L} = -(1 - \sqrt{\eta_t \mu})$ , which yields

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

This in turn establishes

$$\rho(\mathbf{H}_t) \leq \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}$$

## **Nesterov's accelerated gradient methods**

## Convex case

---

$$\text{minimize}_{\boldsymbol{x} \in \mathbb{R}^n} \quad f(\boldsymbol{x})$$

For a positive definite quadratic function  $f$ , including momentum terms allows to improve the iteration complexity from  $O(\kappa \log \frac{1}{\varepsilon})$  to  $O(\sqrt{\kappa} \log \frac{1}{\varepsilon})$

Can we obtain improvement for more general convex cases as well?

# Nesterov's idea

---



— Nesterov '83

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

$$\mathbf{y}^{t+1} = \mathbf{x}^{t+1} + \frac{t}{t+3} (\mathbf{x}^{t+1} - \mathbf{x}^t)$$

Y. Nesterov

- alternates between gradient updates and *proper* extrapolation
- each iteration takes nearly the same cost as GD
- not a descent method (i.e. we may not have  $f(\mathbf{x}^{t+1}) \leq f(\mathbf{x}^t)$ )
- one of the most *beautiful* and *mysterious* results in optimization

...

# Convergence of Nesterov's accelerated gradient method

---

Suppose  $f$  is convex and  $L$ -smooth. If  $\eta_t \equiv \eta = 1/L$ , then

$$f(\mathbf{x}^t) - f^{\text{opt}} \leq \frac{2L\|\mathbf{x}^0 - \mathbf{x}^*\|_2^2}{(t + 1)^2}$$

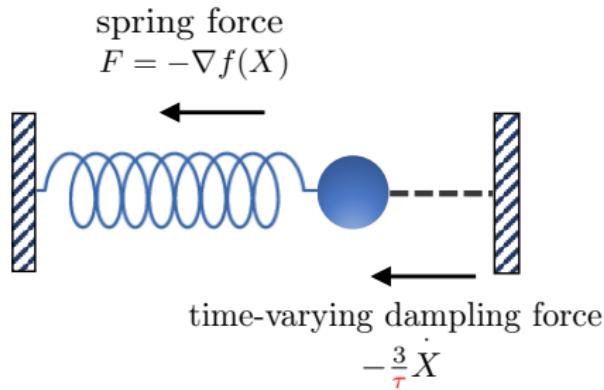
- iteration complexity:  $O\left(\frac{1}{\sqrt{\varepsilon}}\right)$
- much faster than gradient methods
- we'll provide proof for the (more general) proximal version later

## Interpretation using differential equations

---

Nesterov's momentum coefficient  $\frac{t}{t+3} = 1 - \frac{3}{t}$  is particularly mysterious

# Interpretation using differential equations



To develop insight into why Nesterov's method works so well, it's helpful to look at its continuous limits ( $\eta_t \rightarrow 0$ ), which is given by **second-order** ordinary differential equations (ODE)

$$\ddot{\mathbf{X}}(\tau) + \underbrace{\frac{3}{\tau}}_{\text{damping coefficient}} \dot{\mathbf{X}}(\tau) + \underbrace{\nabla f(\mathbf{X}(\tau))}_{\text{potential}} = \mathbf{0}$$

— Su, Boyd, Candes '14

## Heuristic derivation of ODE

---

To begin with, Nesterov's update rule is equivalent to

$$\frac{\mathbf{x}^{t+1} - \mathbf{x}^t}{\sqrt{\eta}} = \frac{t-1}{t+2} \frac{\mathbf{x}^t - \mathbf{x}^{t-1}}{\sqrt{\eta}} - \sqrt{\eta} \nabla f(\mathbf{y}^t) \quad (7.5)$$

Let  $t = \frac{\tau}{\sqrt{\eta}}$ . Set  $\mathbf{X}(\tau) \approx \mathbf{x}^{\tau/\sqrt{\eta}} = \mathbf{x}^t$  and  $\mathbf{X}(\tau + \sqrt{\eta}) \approx \mathbf{x}^{t+1}$ . Then the Taylor expansion gives

$$\begin{aligned}\frac{\mathbf{x}^{t+1} - \mathbf{x}^t}{\sqrt{\eta}} &\approx \dot{\mathbf{X}}(\tau) + \frac{1}{2} \ddot{\mathbf{X}}(\tau) \sqrt{\eta} \\ \frac{\mathbf{x}^t - \mathbf{x}^{t-1}}{\sqrt{\eta}} &\approx \dot{\mathbf{X}}(\tau) - \frac{1}{2} \ddot{\mathbf{X}}(\tau) \sqrt{\eta}\end{aligned}$$

which combined with (7.5) yields

$$\begin{aligned}\dot{\mathbf{X}}(\tau) + \frac{1}{2} \ddot{\mathbf{X}}(\tau) \sqrt{\eta} &\approx \left(1 - \frac{3\sqrt{\eta}}{\tau}\right) \left(\dot{\mathbf{X}}(\tau) - \frac{1}{2} \ddot{\mathbf{X}}(\tau) \sqrt{\eta}\right) - \sqrt{\eta} \nabla f(\mathbf{X}(\tau)) \\ \implies \dot{\mathbf{X}}(\tau) + \frac{3}{\tau} \dot{\mathbf{X}}(\tau) + \nabla f(\mathbf{X}(\tau)) &\approx \mathbf{0}\end{aligned}$$

# Convergence rate of ODE

---

$$\ddot{\mathbf{X}} + \frac{3}{\tau} \dot{\mathbf{X}} + \nabla f(\mathbf{X}) = \mathbf{0} \quad (7.6)$$

Standard ODE theory reveals that

$$f(\mathbf{X}(\tau)) - f^{\text{opt}} \leq O\left(\frac{1}{\tau^2}\right) \quad (7.7)$$

which somehow explains Nesterov's  $O(1/t^2)$  convergence

## Proof of (7.7)

---

Define  $\mathcal{E}(\tau) := \underbrace{\tau^2(f(\mathbf{X}) - f^{\text{opt}}) + 2\left\|\mathbf{X} + \frac{\tau}{2}\dot{\mathbf{X}} - \mathbf{X}^*\right\|_2^2}_{\text{Lyapunov function / energy function}}$ . This obeys

$$\begin{aligned}\dot{\mathcal{E}} &= 2\tau(f(\mathbf{X}) - f^{\text{opt}}) + \tau^2\langle \nabla f(\mathbf{X}), \dot{\mathbf{X}} \rangle + 4\left\langle \mathbf{X} + \frac{\tau}{2}\dot{\mathbf{X}} - \mathbf{X}^*, \frac{3}{2}\dot{\mathbf{X}} + \frac{\tau}{2}\ddot{\mathbf{X}} \right\rangle \\ &\stackrel{(i)}{=} 2\tau(f(\mathbf{X}) - f^{\text{opt}}) - 2\tau\langle \mathbf{X} - \mathbf{X}^*, \nabla f(\mathbf{X}) \rangle \stackrel{\text{(by convexity)}}{\leq} 0\end{aligned}$$

where (i) follows by replacing  $\tau\ddot{\mathbf{X}} + 3\dot{\mathbf{X}}$  with  $-\tau\nabla f(\mathbf{X})$

This means  $\mathcal{E}$  is non-decreasing in  $\tau$ , and hence

$$f(\mathbf{X}(\tau)) - f^{\text{opt}} \stackrel{\text{(defn of } \mathcal{E})}{\leq} \frac{\mathcal{E}(\tau)}{\tau^2} \leq \frac{\mathcal{E}(0)}{\tau^2} = O\left(\frac{1}{\tau^2}\right)$$

# Magic number 3

---

$$\ddot{\mathbf{X}} + \frac{3}{\tau} \dot{\mathbf{X}} + \nabla f(\mathbf{X}) = \mathbf{0}$$

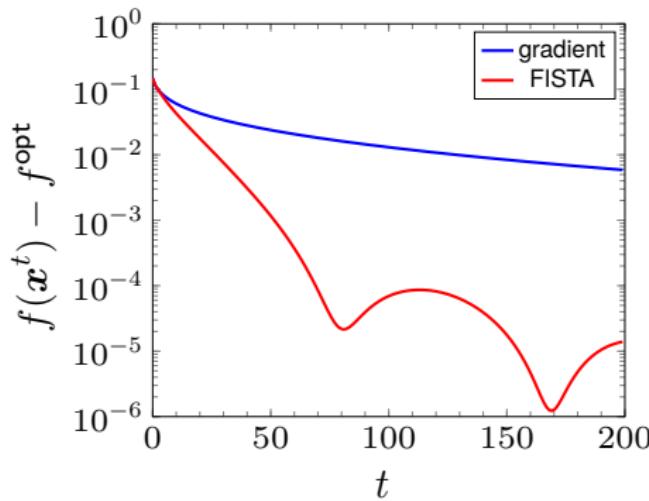
- 3 is the smallest constant that guarantees  $O(1/\tau^2)$  convergence, and can be replaced by any other  $\alpha \geq 3$
- in some sense, 3 minimizes the pre-constant in the convergence bound  $O(1/\tau^2)$  (see Su, Boyd, Candes '14)

# Numerical example

taken from UCLA EE236C

$$\underset{\mathbf{x}}{\text{minimize}} \quad \log \left( \sum_{i=1}^m \exp(\mathbf{a}_i^\top \mathbf{x} + b_i) \right)$$

with randomly generated problems and  $m = 2000$ ,  $n = 1000$



## Extension to composite models

---

$$\begin{aligned} & \text{minimize}_{\boldsymbol{x}} \quad F(\boldsymbol{x}) := f(\boldsymbol{x}) + h(\boldsymbol{x}) \\ & \text{subject to} \quad \boldsymbol{x} \in \mathbb{R}^n \end{aligned}$$

- $f$ : convex and smooth
- $h$ : convex (may not be differentiable)

let  $F^{\text{opt}} := \min_{\boldsymbol{x}} F(\boldsymbol{x})$  be the optimal cost

# FISTA (Beck & Teboulle '09)

---

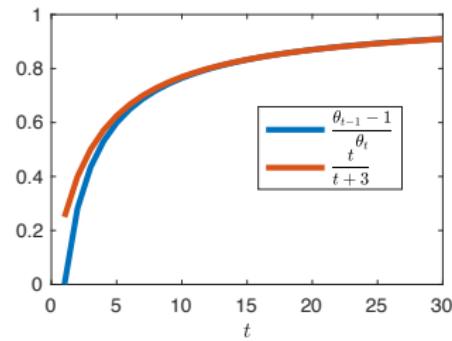
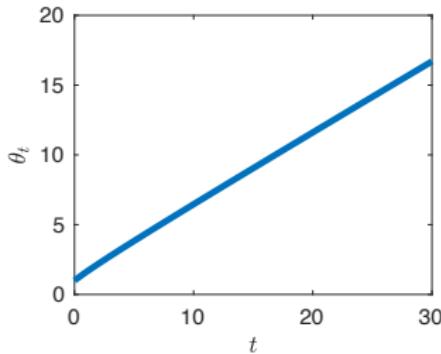
## Fast iterative shrinkage-thresholding algorithm

$$\begin{aligned}\mathbf{x}^{t+1} &= \text{prox}_{\eta_t h}(\mathbf{y}^t - \eta_t \nabla f(\mathbf{y}^t)) \\ \mathbf{y}^{t+1} &= \mathbf{x}^{t+1} + \frac{\theta_t - 1}{\theta_{t+1}} (\mathbf{x}^{t+1} - \mathbf{x}^t)\end{aligned}$$

where  $\mathbf{y}^0 = \mathbf{x}^0$ ,  $\theta_0 = 1$  and  $\theta_{t+1} = \frac{1+\sqrt{1+4\theta_t^2}}{2}$

- adopt the momentum coefficients originally proposed by Nesterov '83

# Momentum coefficient

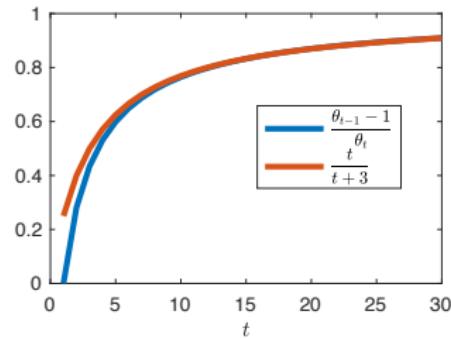
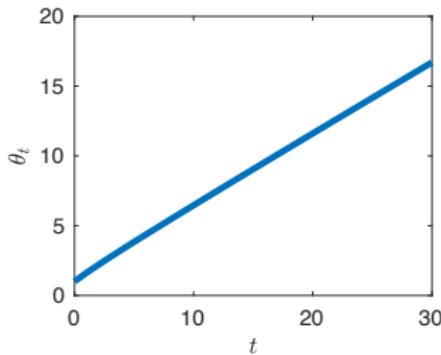


$$\theta_{t+1} = \frac{1 + \sqrt{1 + 4\theta_t^2}}{2} \quad \text{with } \theta_0 = 1$$

coefficient  $\frac{\theta_{t-1}}{\theta_{t+1}} = 1 - \frac{3}{t} + o\left(\frac{1}{t}\right)$  (homework)

- asymptotically equivalent to  $\frac{t}{t+3}$

# Momentum coefficient



$$\theta_{t+1} = \frac{1 + \sqrt{1 + 4\theta_t^2}}{2} \quad \text{with } \theta_0 = 1$$

## Fact 7.2

For all  $t \geq 1$ , one has  $\theta_t \geq \frac{t+2}{2}$  (homework)

## **Convergence analysis**

# Convergence for convex problems

---

## Theorem 7.3 (Convergence of accelerated proximal gradient methods for convex problems)

Suppose  $f$  is convex and  $L$ -smooth. If  $\eta_t \equiv 1/L$ , then

$$F(\mathbf{x}^t) - F^{\text{opt}} \leq \frac{2L\|\mathbf{x}^0 - \mathbf{x}^*\|_2^2}{(t + 1)^2}$$

- improved iteration complexity (i.e.  $O(1/\sqrt{\varepsilon})$ ) than proximal gradient method (i.e.  $O(1/\varepsilon)$ )
- fast if prox can be efficiently implemented

# Recap: the fundamental inequality for proximal method

---

Recall the following fundamental inequality shown in the last lecture:

## Lemma 7.4

Let  $\mathbf{y}^+ = \text{prox}_{\frac{1}{L}h}(\mathbf{y} - \frac{1}{L}\nabla f(\mathbf{y}))$ , then

$$F(\mathbf{y}^+) - F(\mathbf{x}) \leq \frac{L}{2}\|\mathbf{x} - \mathbf{y}\|_2^2 - \frac{L}{2}\|\mathbf{x} - \mathbf{y}^+\|_2^2$$

## Proof of Theorem 7.6

---

1. build a discrete-time version of “Lyapunov function”
2. ***magic happens!***
  - “Lyapunov function” is non-increasing when Nesterov’s momentum coefficients are adopted

# Proof of Theorem 7.6

**Key lemma:** monotonicity of a certain “Lyapunov function”

## Lemma 7.5

Let  $\mathbf{u}^t = \underbrace{\theta_{t-1}\mathbf{x}^t - (\mathbf{x}^* + (\theta_{t-1} - 1)\mathbf{x}^{t-1})}_{\text{or } \theta_{t-1}(\mathbf{x}^t - \mathbf{x}^*) - (\theta_{t-1} - 1)(\mathbf{x}^{t-1} - \mathbf{x}^*)}$ . Then

$$\|\mathbf{u}^{t+1}\|_2^2 + \frac{2}{L}\theta_t^2(F(\mathbf{x}^{t+1}) - F^{\text{opt}}) \leq \|\mathbf{u}^t\|_2^2 + \frac{2}{L}\theta_{t-1}^2(F(\mathbf{x}^t) - F^{\text{opt}})$$

- quite similar to  $2\|\dot{\mathbf{X}} + \frac{\tau}{2}\dot{\mathbf{X}} - \mathbf{X}^*\|_2^2 + \tau^2(f(\mathbf{X}) - f^{\text{opt}})$  (Lyapunov function) as discussed before (think about  $\theta_t \approx t/2$ )

## Proof of Theorem 7.6

---

With Lemma 7.5 in place, one has

$$\begin{aligned}\frac{2}{L} \theta_{t-1}^2 (F(\mathbf{x}^t) - F^{\text{opt}}) &\leq \|\mathbf{u}^1\|_2^2 + \frac{2}{L} \theta_0^2 (F(\mathbf{x}^1) - F^{\text{opt}}) \\ &= \|\mathbf{x}^1 - \mathbf{x}^*\|_2^2 + \frac{2}{L} (F(\mathbf{x}^1) - F^{\text{opt}})\end{aligned}$$

To bound the RHS of this inequality, we use Lemma 7.4 and  $\mathbf{y}^0 = \mathbf{x}^0$  to get

$$\begin{aligned}\frac{2}{L} (F(\mathbf{x}^1) - F^{\text{opt}}) &\leq \|\mathbf{y}^0 - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^1 - \mathbf{x}^*\|_2^2 = \|\mathbf{x}^0 - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^1 - \mathbf{x}^*\|_2^2 \\ \iff \quad \|\mathbf{x}^1 - \mathbf{x}^*\|_2^2 + \frac{2}{L} (F(\mathbf{x}^1) - F^{\text{opt}}) &\leq \|\mathbf{x}^0 - \mathbf{x}^*\|_2^2\end{aligned}$$

As a result,

$$\begin{aligned}\frac{2}{L} \theta_{t-1}^2 (F(\mathbf{x}^t) - F^{\text{opt}}) &\leq \|\mathbf{x}^1 - \mathbf{x}^*\|_2^2 + \frac{2}{L} (F(\mathbf{x}^1) - F^{\text{opt}}) \leq \|\mathbf{x}^0 - \mathbf{x}^*\|_2^2, \\ \implies F(\mathbf{x}^t) - F^{\text{opt}} &\leq \frac{L \|\mathbf{x}^0 - \mathbf{x}^*\|_2^2}{2 \theta_{t-1}^2} \stackrel{\text{(Fact 7.2)}}{\leq} \frac{2L \|\mathbf{x}^0 - \mathbf{x}^*\|_2^2}{(t+1)^2}\end{aligned}$$

## Proof of Lemma 7.5

---

Take  $\mathbf{x} = \frac{1}{\theta_t} \mathbf{x}^* + (1 - \frac{1}{\theta_t}) \mathbf{x}^t$  and  $\mathbf{y} = \mathbf{y}^t$  in Lemma 7.4 to get

$$F(\mathbf{x}^{t+1}) - F\left(\theta_t^{-1} \mathbf{x}^* + (1 - \theta_t^{-1}) \mathbf{x}^t\right) \quad (7.8)$$

$$\begin{aligned} &\leq \frac{L}{2} \|\theta_t^{-1} \mathbf{x}^* + (1 - \theta_t^{-1}) \mathbf{x}^t - \mathbf{y}^t\|_2^2 - \frac{L}{2} \|\theta_t^{-1} \mathbf{x}^* + (1 - \theta_t^{-1}) \mathbf{x}^t - \mathbf{x}^{t+1}\|_2^2 \\ &= \frac{L}{2\theta_t^2} \|\mathbf{x}^* + (\theta_t - 1) \mathbf{x}^t - \theta_t \mathbf{y}^t\|_2^2 - \frac{L}{2\theta_t^2} \underbrace{\|\mathbf{x}^* + (\theta_t - 1) \mathbf{x}^t - \theta_t \mathbf{x}^{t+1}\|_2^2}_{= -\mathbf{u}^{t+1}} \end{aligned}$$

$$\stackrel{(i)}{=} \frac{L}{2\theta_t^2} (\|\mathbf{u}^t\|_2^2 - \|\mathbf{u}^{t+1}\|_2^2), \quad (7.9)$$

where (i) follows from the definition of  $\mathbf{u}^t$  and  $\mathbf{y}^t = \mathbf{x}^t + \frac{\theta_{t-1}-1}{\theta_t} (\mathbf{x}^t - \mathbf{x}^{t-1})$

## Proof of Lemma 7.5 (cont.)

---

We will also lower bound (7.8). By convexity of  $F$ ,

$$\begin{aligned} F\left(\theta_t^{-1}\mathbf{x}^* + (1 - \theta_t^{-1})\mathbf{x}^t\right) &\leq \theta_t^{-1}F(\mathbf{x}^*) + (1 - \theta_t^{-1})F(\mathbf{x}^t) \\ &= \theta_t^{-1}F^{\text{opt}} + (1 - \theta_t^{-1})F(\mathbf{x}^t) \\ \iff F\left(\theta_t^{-1}\mathbf{x}^* + (1 - \theta_t^{-1})\mathbf{x}^t\right) - F(\mathbf{x}^{t+1}) &\leq (1 - \theta_t^{-1})(F(\mathbf{x}^t) - F^{\text{opt}}) - (F(\mathbf{x}^{t+1}) - F^{\text{opt}}) \end{aligned}$$

Combining this with (7.9) and  $\theta_t^2 - \theta_t = \theta_{t-1}^2$  yields

$$\begin{aligned} \frac{L}{2}(\|\mathbf{u}^t\|_2^2 - \|\mathbf{u}^{t+1}\|_2^2) &\geq \theta_t^2(F(\mathbf{x}^{t+1}) - F^{\text{opt}}) - (\theta_t^2 - \theta_t)(F(\mathbf{x}^t) - F^{\text{opt}}) \\ &= \theta_t^2(F(\mathbf{x}^{t+1}) - F^{\text{opt}}) - \theta_{t-1}^2(F(\mathbf{x}^t) - F^{\text{opt}}), \end{aligned}$$

thus finishing the proof

# Convergence for strongly convex problems

---

$$\mathbf{x}^{t+1} = \text{prox}_{\eta_t h}(\mathbf{y}^t - \eta_t \nabla f(\mathbf{y}^t))$$

$$\mathbf{y}^{t+1} = \mathbf{x}^{t+1} + \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} (\mathbf{x}^{t+1} - \mathbf{x}^t)$$

**Theorem 7.6 (Convergence of accelerated proximal gradient methods for strongly convex case)**

Suppose  $f$  is  $\mu$ -strongly convex and  $L$ -smooth. If  $\eta_t \equiv 1/L$ , then

$$F(\mathbf{x}^t) - F^{\text{opt}} \leq \left(1 - \frac{1}{\sqrt{\kappa}}\right)^t \left( F(\mathbf{x}^0) - F^{\text{opt}} + \frac{\mu \|\mathbf{x}^0 - \mathbf{x}^*\|_2^2}{2} \right)$$

## A practical issue

---

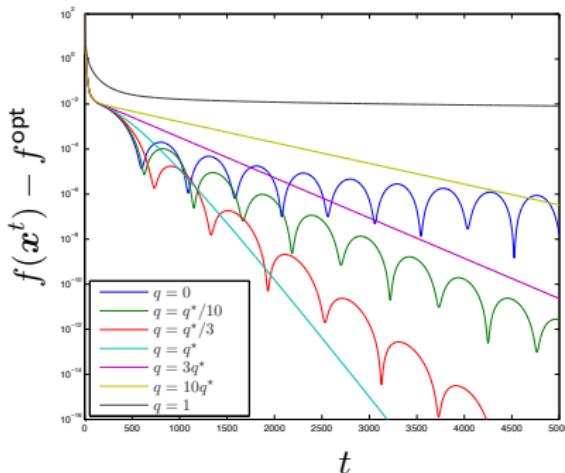
Fast convergence requires knowledge of  $\kappa = L/\mu$

- in practice, estimating  $\mu$  is typically very challenging

A common observation: ripples / bumps in the traces of cost values

# Rippling behavior

Numerical example: take  $\mathbf{y}^{t+1} = \mathbf{x}^{t+1} + \frac{1-\sqrt{q}}{1+\sqrt{q}}(\mathbf{x}^{t+1} - \mathbf{x}^t)$ ;  $q^* = 1/\kappa$



period of ripples is often proportional to  $\sqrt{L/\mu}$

O'Donoghue, Candes '12

- when  $q > q^*$ : we underestimate momentum  $\rightarrow$  slower convergence
- when  $q < q^*$ : we overestimate momentum ( $\frac{1-\sqrt{q}}{1+\sqrt{q}}$  is large)  
 $\rightarrow$  overshooting / rippling behavior

## Adaptive restart (O'Donoghue, Candes '12)

---

When a certain criterion is met, restart running FISTA with

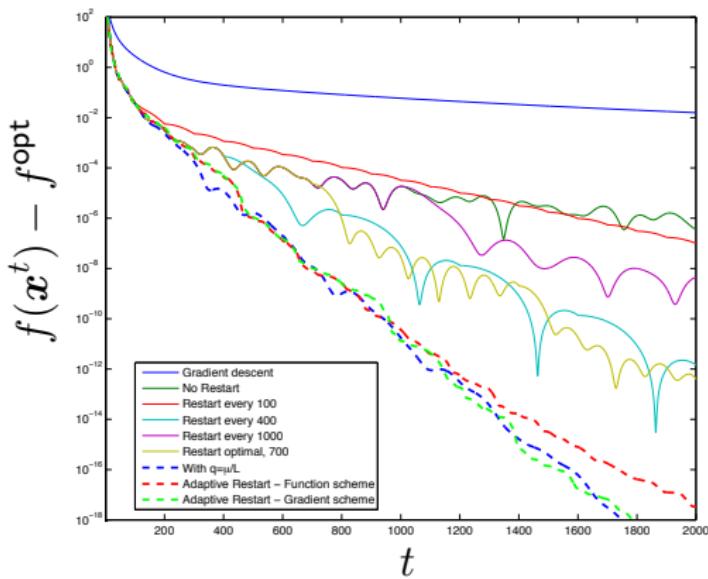
$$\mathbf{x}^0 \leftarrow \mathbf{x}^t$$

$$\mathbf{y}^0 \leftarrow \mathbf{x}^t$$

$$\theta_0 = 1$$

- take the current iterate as a new starting point
- erase all memory of previous iterates and reset the momentum back to zero

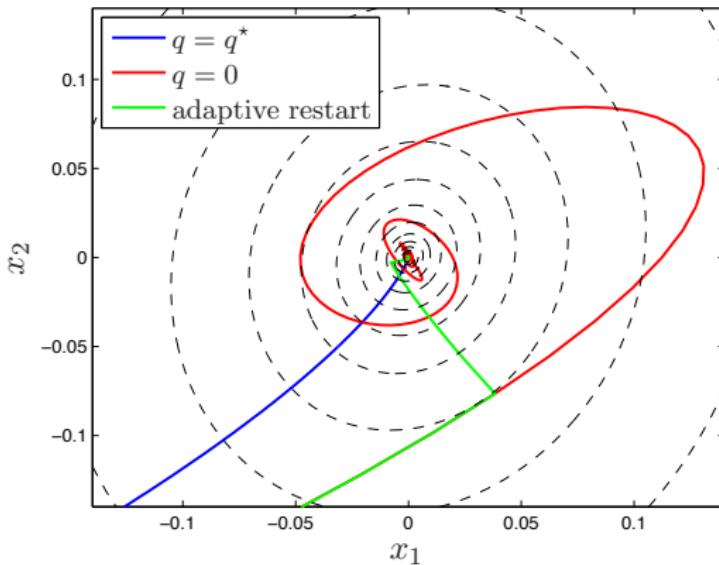
# Numerical comparisons of adaptive restart schemes



- function scheme: restart when  $f(\mathbf{x}^t) > f(\mathbf{x}^{t-1})$
- gradient scheme: restart when  $\underbrace{\langle \nabla f(\mathbf{y}^{t-1}), \mathbf{x}^t - \mathbf{x}^{t-1} \rangle}_{\text{restart when momentum lead us towards a bad direction}} > 0$

# Illustration

---



- with overestimated momentum (e.g.  $q = 0$ ), one sees spiralling trajectory
- adaptive restart helps mitigate this issue

## **Lower bounds**

# Optimality of Nesterov's method

---

Interestingly, no first-order methods can improve upon Nesterov's results in general

More precisely,  $\exists$  convex and  $L$ -smooth function  $f$  s.t.

$$f(\mathbf{x}^t) - f^{\text{opt}} \geq \frac{3L\|\mathbf{x}^0 - \mathbf{x}^*\|_2^2}{32(t+1)^2}$$

as long as  $\underbrace{\mathbf{x}^k \in \mathbf{x}^0 + \text{span}\{\nabla f(\mathbf{x}^0), \dots, \nabla f(\mathbf{x}^{k-1})\}}_{\text{definition of first-order methods}}$  for all  $1 \leq k \leq t$

— Nemirovski, Yudin '83

# Example

---

$$\underset{\mathbf{x} \in \mathbb{R}^{(2n+1)}}{\text{minimize}} \quad f(\mathbf{x}) = \frac{L}{4} \left( \frac{1}{2} \mathbf{x}^\top \mathbf{A} \mathbf{x} - \mathbf{e}_1^\top \mathbf{x} \right)$$

where  $\mathbf{A} = \begin{bmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & 2 \end{bmatrix} \in \mathbb{R}^{(2n+1) \times (2n+1)}$

- $f$  is convex and  $L$ -smooth
- the optimizer  $\mathbf{x}^*$  is given by  $x_i^* = 1 - \frac{i}{2n+2}$  ( $1 \leq i \leq n$ ) obeying

$$f^{\text{opt}} = \frac{L}{8} \left( \frac{1}{2n+2} - 1 \right) \quad \text{and} \quad \|\mathbf{x}^*\|_2^2 \leq \frac{2n+2}{3}$$

# Example

---

$$\underset{\mathbf{x} \in \mathbb{R}^{(2n+1)}}{\text{minimize}} \quad f(\mathbf{x}) = \frac{L}{4} \left( \frac{1}{2} \mathbf{x}^\top \mathbf{A} \mathbf{x} - \mathbf{e}_1^\top \mathbf{x} \right)$$

where  $\mathbf{A} = \begin{bmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & 2 \end{bmatrix} \in \mathbb{R}^{(2n+1) \times (2n+1)}$

- $\nabla f(\mathbf{x}) = \frac{L}{4} \mathbf{A} \mathbf{x} - \frac{L}{4} \mathbf{e}_1$
- $\underbrace{\text{span}\{\nabla f(\mathbf{x}^0), \dots, \nabla f(\mathbf{x}^{k-1})\}}_{=: \mathcal{K}_k} = \text{span}\{\mathbf{e}_1, \dots, \mathbf{e}_k\}$  if  $\mathbf{x}^0 = \mathbf{0}$ 
  - every iteration of first-order methods expands the search space by *at most* one dimension

## Example (cont.)

---

If we start with  $\mathbf{x}^0 = \mathbf{0}$ , then

$$\begin{aligned} f(\mathbf{x}^n) &\geq \inf_{\mathbf{x} \in \mathcal{K}_n} f(\mathbf{x}) = \frac{L}{8} \left( \frac{1}{n+1} - 1 \right) \\ \implies \frac{f(\mathbf{x}^n) - f^{\text{opt}}}{\|\mathbf{x}^0 - \mathbf{x}^*\|_2^2} &\geq \frac{\frac{L}{8} \left( \frac{1}{n+1} - \frac{1}{2n+2} \right)}{\frac{1}{3}(2n+2)} = \frac{3L}{32(n+1)^2} \end{aligned}$$

# Summary: accelerated proximal gradient

---

	stepsize rule	convergence rate	iteration complexity
convex & smooth problems	$\eta_t = \frac{1}{L}$	$O\left(\frac{1}{t^2}\right)$	$O\left(\frac{1}{\sqrt{\varepsilon}}\right)$
strongly convex & smooth problems	$\eta_t = \frac{1}{L}$	$O\left(\left(1 - \frac{1}{\sqrt{\kappa}}\right)^t\right)$	$O\left(\sqrt{\kappa} \log \frac{1}{\varepsilon}\right)$

# Reference

---

- [1] "*Some methods of speeding up the convergence of iteration methods,*" B. Polyak, *USSR Computational Mathematics and Mathematical Physics*, 1964
- [2] "*A method of solving a convex programming problem with convergence rate  $O(1/k^2)$ ,*" Y. Nesterov, *Soviet Mathematics Doklady*, 1983.
- [3] "*A fast iterative shrinkage-thresholding algorithm for linear inverse problems,*" A. Beck and M. Teboulle, *SIAM journal on imaging sciences*, 2009.
- [4] "*First-order methods in optimization,*" A. Beck, Vol. 25, SIAM, 2017.
- [5] "*Mathematical optimization, MATH301 lecture notes,*" E. Candes, Stanford.
- [6] "*Gradient methods for minimizing composite functions,*", Y. Nesterov, *Technical Report*, 2007.

# Reference

---

- [7] "*Large-scale numerical optimization, MS&E318 lecture notes,*" M. Saunders, Stanford.
- [8] "*Analysis and design of optimization algorithms via integral quadratic constraints,*" L. Lessard, B. Recht, A. Packard, *SIAM Journal on Optimization*, 2016.
- [9] "*Proximal algorithms,*" N. Parikh and S. Boyd, *Foundations and Trends in Optimization*, 2013.
- [10] "*Convex optimization: algorithms and complexity,*" S. Bubeck, *Foundations and trends in machine learning*, 2015.
- [11] "*Optimization methods for large-scale systems, EE236C lecture notes,*" L. Vandenberghe, UCLA.
- [12] "*Problem complexity and method efficiency in optimization,*" A. Nemirovski, D. Yudin, Wiley, 1983.

# Reference

---

- [13] "*Introductory lectures on convex optimization: a basic course,*" Y. Nesterov, 2004
- [14] "*A differential equation for modeling Nesterov's accelerated gradient method,*" W. Su, S. Boyd, E. Candes, NIPS, 2014.
- [15] "*On accelerated proximal gradient methods for convex-concave optimization,*" P. Tseng, 2008
- [16] "*Adaptive restart for accelerated gradient schemes,*" B. O'donoghue, and E. Candes, *Foundations of computational mathematics*, 2012