

Lecture 8. Accelerated Gradient Methods

Bao Wang

Department of Mathematics

Scientific Computing and Imaging Institute

University of Utah

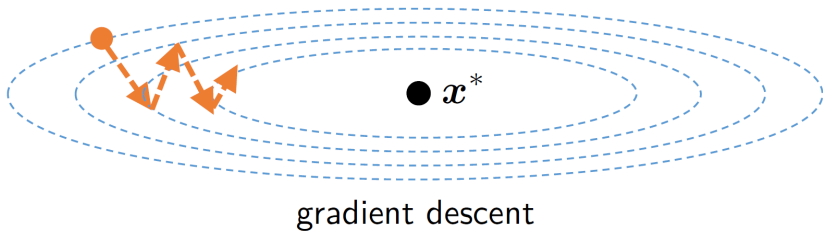
Math 5750/6880, Fall 2021

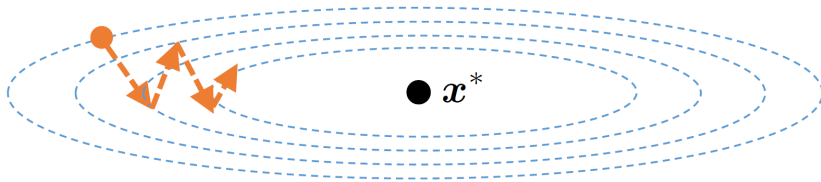
The iteration complexity of (proximal) gradient methods
strongly convex and smooth problems $O(\kappa \log \frac{1}{\epsilon})$.
convex and smooth problems $O(\frac{1}{\epsilon})$.

Can one still hope to further accelerate convergence?

Issues:

- 1) GD focuses on improving the cost per iteration, which might sometimes be too "short-sighted";
- 2) GD might sometimes zigzag or experience abrupt changes.





gradient descent

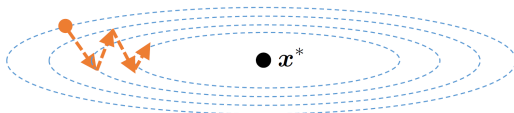
Solutions:

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

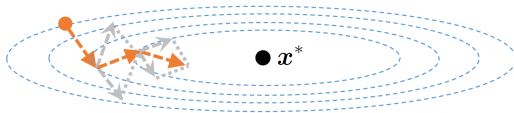
Heavy-ball method

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

where we add inertia to the "ball" (i.e. include a momentum term) to mitigate zigzagging.



gradient descent



heavy-ball method

State-space method

Consider

$$\min_{\mathbf{x}} \frac{1}{2} (\mathbf{x} - \mathbf{x}^*)^\top \mathbf{Q} (\mathbf{x} - \mathbf{x}^*),$$

where $\mathbf{Q} \succ 0$ has a condition number κ . One can understand heavy-ball methods through dynamical systems.

Consider the following dynamical system

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

State-space method

or equivalently

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

The convergence of heavy-ball methods depends on the spectrum of the system matrix \mathbf{H}_t . We need to find appropriate stepsizes η_t and momentum coefficients θ_t to control the spectrum of \mathbf{H}_t .

Convergence of heavy-ball methods for quadratic functions

Theorem 1. [Convergence of heavy-ball methods for quadratic functions] Suppose f is a L -smooth and μ -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

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

Note that the iteration complexity is $O(\sqrt{\kappa} \log \frac{1}{\epsilon})$ (vs. $O(\kappa \log \frac{1}{\epsilon})$ of GD), the convergence rate relies on knowledge of both L and μ .

Proof of Theorem 1. In view of (1), it suffices to control the spectrum of \mathbf{H}_t (which is time-invariant). Let λ_i be the i th eigenvalue of \mathbf{Q} and let $\Lambda := \text{diag}(\lambda_1, \dots, \lambda_n)$, then there exists an orthogonal matrix \mathbf{U} such that $\mathbf{Q} = \mathbf{U}\Lambda\mathbf{U}^\top$ and we have

$$\left\| \begin{bmatrix} \mathbf{U} & \mathbf{0} \\ \mathbf{0} & \mathbf{U} \end{bmatrix} \begin{bmatrix} (1 + \theta_t)\mathbf{I} - \eta\mathbf{Q} & -\theta_t\mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{U} & \mathbf{0} \\ \mathbf{0} & \mathbf{U} \end{bmatrix}^\top \right\|_2 = \left\| \begin{bmatrix} (1 + \theta_t)\mathbf{I} - \eta\Lambda & -\theta_t\mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{bmatrix} \right\|_2.$$

Further, note that the characteristic polynomial of the right-hand side matrix satisfies

$$\begin{aligned} \text{ch} \left(\begin{bmatrix} (1 + \theta_t)\mathbf{I} - \eta\Lambda & -\theta_t\mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{bmatrix} \right) &= \det \begin{bmatrix} \lambda\mathbf{I} - ((1 + \theta_t)\mathbf{I} - \eta\Lambda) & \theta_t\mathbf{I} \\ -\mathbf{I} & \lambda\mathbf{I} \end{bmatrix} \\ &= \det \left(\lambda^2\mathbf{I} - \lambda((1 + \theta_t)\mathbf{I} - \eta\Lambda) + \theta_t\mathbf{I} \right). \end{aligned}$$

The matrix $\lambda^2 \mathbf{I} - \lambda((1 + \theta_t)\mathbf{I} - \eta_t \Lambda) + \theta_t \mathbf{I}$ is diagonal with each diagonal entry be the characteristic polynomial of the following 2×2 matrix

$$\begin{bmatrix} 1 + \theta_t - \eta_t \lambda_i & -\theta_t \\ 1 & 0 \end{bmatrix}.$$

Then the spectral radius (denoted by $\rho(\cdot)$)¹ of \mathbf{H}_t obeys

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

To finish the proof, it suffices to show

$$\max_i \rho\left(\begin{bmatrix} 1 + \theta_t - \eta_t \lambda_i & -\theta_t \\ 1 & 0 \end{bmatrix}\right) \leq \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}. \quad (2)$$

¹The largest absolute value of its eigenvalues

To show (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

$$z^2 - (1 + \theta_t - \eta_t \lambda_i)z + \theta_t = 0. \quad (3)$$

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

$$\theta_t \in [(1 - \sqrt{\eta_t \lambda_i})^2, (1 + \sqrt{\eta_t \lambda_i})^2], \quad (4)$$

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

With this choice of θ_t , we have (from (3) and the fact that two eigenvalues have identical magnitudes) [Vieta's formula: $z_1 z_2 = \theta_t$]

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

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

$$\mathbf{x}^{t+1} = \mathbf{y}^t - \eta_t \nabla f(\mathbf{y}^t); \quad \mathbf{y}^{t+1} = \mathbf{x}^{t+1} + \frac{t}{t+3}(\mathbf{x}^{t+1} - \mathbf{x}^t).$$

- > 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 $f(\mathbf{x}^{t+1}) \leq f(\mathbf{x}^t)$), we will see it later.

Theorem 2. [Convergence of Nesterov's accelerated gradient method] Suppose f is convex and L -smooth. if $\eta_t \equiv \eta = 1/L$, then

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

Remark. The iteration complexity is $O(\frac{1}{\sqrt{\epsilon}})$, which is much faster than gradient methods.

ODE analogy of Nesterov's accelerated gradient

To develop insight into why Nesterov's method works so well, it's helpful to look at its continuous limits ($\eta_t \rightarrow 0$). To begin with, Nesterov's update rule is equivalent to

$$\frac{\mathbf{x}^{t+1} - \mathbf{x}^t}{\sqrt{\eta}} = \frac{t-1}{t+2} \frac{\mathbf{x}^t - \mathbf{x}^{t-1}}{\sqrt{\eta}} - \sqrt{\eta} \nabla f(\mathbf{y}^t). \quad (5)$$

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

$$\frac{\mathbf{x}^{t+1} - \mathbf{x}^t}{\sqrt{\eta}} \approx \dot{\mathbf{X}}(\tau) + \frac{1}{2} \ddot{\mathbf{X}}(\tau) \sqrt{\eta}; \quad \frac{\mathbf{x}^t - \mathbf{x}^{t-1}}{\sqrt{\eta}} \approx \dot{\mathbf{X}}(\tau) - \frac{1}{2} \ddot{\mathbf{X}}(\tau) \sqrt{\eta},$$

ODE analogy of Nesterov's accelerated gradient

which combined with (5) yields

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

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

What is the ODE limit of the heavy-ball method?

Heavy-ball method

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta \nabla f(\mathbf{x}^t) + \theta(\mathbf{x}^t - \mathbf{x}^{t-1}),$$

Let $\mathbf{m}^t := (\mathbf{x}^{t+1} - \mathbf{x}^t)/\sqrt{\eta}$ and let $\theta := 1 - \gamma\sqrt{\eta}$, where $\gamma \geq 0$ is another hyperparameter. Then we can rewrite the heavy-ball method as

$$\mathbf{m}^{t+1} = (1 - \gamma\sqrt{\eta})\mathbf{m}^t - \sqrt{\eta}\nabla f(\mathbf{x}^t); \quad \mathbf{x}^{t+1} = \mathbf{x}^t + \sqrt{s}\mathbf{m}^{t+1}.$$

Let $s \rightarrow 0$; we obtain the following system of first-order ODEs

$$\frac{d\mathbf{X}(t)}{dt} = \mathbf{M}(t); \quad \frac{d\mathbf{M}(t)}{dt} = -\gamma\mathbf{M}(t) - \nabla f(\mathbf{X}(t)),$$

which can be further written as

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

Heavy-ball method vs. Nesterov's acceleration

Heavy-ball method:

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta \nabla f(\mathbf{x}^t) + \theta(\mathbf{x}^t - \mathbf{x}^{t-1}),$$

and the ODE-limit is

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

Nesterov's acceleration:

$$\mathbf{x}^{t+1} = \mathbf{y}^t - \eta \nabla f(\mathbf{y}^t); \quad \mathbf{y}^{t+1} = \mathbf{x}^{t+1} + \frac{t}{t+3}(\mathbf{x}^{t+1} - \mathbf{x}^t),$$

and the ODE-limit is

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

Convergence rate inspired by the ODE analysis

By the standard ODE theory, we can show that

$$f(\mathbf{X}(\tau)) - f^{opt} \leq O\left(\frac{1}{\tau^2}\right), \quad (6)$$

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

Convergence rate inspired by the ODE analysis

Proof. Define $E(\tau) := \tau^2(f(\mathbf{X}) - f^{opt}) + 2\|\mathbf{X} + \frac{\tau}{2}\dot{\mathbf{X}} - \mathbf{X}^*\|_2^2$ (Lyapunov function).

This obeys

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

where (i) follows by replacing $\tau\ddot{\mathbf{X}} + 3\dot{\mathbf{X}}$ with $-\tau\nabla f(\mathbf{X})$. This means E is non-decreasing in τ , and hence

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

Extend Nesterov's acceleration to composite models

$$\min_{\mathbf{x}} F(\mathbf{x}) := f(\mathbf{x}) + h(\mathbf{x}) \quad \text{s.t. } \mathbf{x} \in \mathbb{R}^n,$$

where f is convex and smooth and h is convex (may not be differentiable). Let $F^{opt} := \min_{\mathbf{x}} F(\mathbf{x})$ be the optimal cost.

FISTA (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 $\theta_{t+1} = \frac{1 + \sqrt{1 + 4\theta_t^2}}{2}$.

We can show that $\frac{\theta_t - 1}{\theta_{t+1}} = 1 - \frac{3}{t} + o(\frac{1}{t})$ [Homework..](#)

We can also show that $\theta_t \geq \frac{t+2}{2}$. (Math induction.)

Theorem 3. [Convergence of accelerated proximal gradient methods for convex problems] Suppose f is convex and L -smooth. If $\eta_t \equiv 1/L$, then

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

Remark. The algorithm is fast if prox can be efficiently implemented.

Remark. The algorithm is particularly useful for ℓ_1 -regularization problem in e.g. image processing (total variation in the wavelet space) and compressed sensing.

Lemma 1. [Fundamental inequality for proximal method] Let

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

then

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

To proof Theorem 3, we follow: 1) build a discrete-time version of "Lyapunov function"; 2) "Lyapunov function" is non-increasing when Nesterov's momentum coefficients are adopted.

Proof of Lemma 1. More precisely, we have

$$F(\mathbf{y}^+) - F(\mathbf{x}) \leq \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 - \frac{L}{2} \|\mathbf{x} - \mathbf{y}^+\|_2^2 - \underbrace{g(\mathbf{x}, \mathbf{y})}_{\geq 0 \text{ by convexity}}$$

where $g(\mathbf{x}, \mathbf{y}) := f(\mathbf{x}) - f(\mathbf{y}) - \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle$.

Define $\phi(\mathbf{z}) = f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{z} - \mathbf{y} \rangle + \frac{L}{2} \|\mathbf{z} - \mathbf{y}\|_2^2 + h(\mathbf{z})$. It is easily seen that $\mathbf{y}^+ = \arg \min_{\mathbf{z}} \phi(\mathbf{z})$. Two important properties:

1. Since $\phi(\mathbf{z})$ is L -strongly convex, one has

$$\phi(\mathbf{x}) \geq \phi(\mathbf{y}^+) + \frac{L}{2} \|\mathbf{x} - \mathbf{y}^+\|_2^2.$$

2. From smoothness,

$$\phi(\mathbf{y}^+) = \underbrace{f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{y}^+ - \mathbf{y} \rangle + \frac{L}{2} \|\mathbf{y}^+ - \mathbf{y}\|_2^2}_{\text{upper bound on } f(\mathbf{y}^+) \text{ (L-smoothness)}} + h(\mathbf{y}^+) \geq f(\mathbf{y}^+) + h(\mathbf{y}^+) = F(\mathbf{y}^+).$$

Taken collectively, these yield

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

which together with the definition of $\phi(\mathbf{x})$ gives

$$\underbrace{f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + h(\mathbf{x})}_{=f(\mathbf{x})+h(\mathbf{x})-g(\mathbf{x},\mathbf{y})=F(\mathbf{x})-g(\mathbf{x},\mathbf{y})} + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 \geq F(\mathbf{y}^+) + \frac{L}{2} \|\mathbf{x} - \mathbf{y}^+\|_2^2$$

which finishes the proof.

Lemma 2. [Monotonicity of certain "Lyapunov function"]

Let

$$\mathbf{u}^t = \theta_{t-1} \mathbf{x}^t - \left(\mathbf{x}^* + (\theta_{t-1} - 1) \mathbf{x}^{t-1} \right).$$

Then

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

Remark. Note that this is quite similar to $2\|\mathbf{X} + \frac{\tau}{2}\dot{\mathbf{X}} - \mathbf{X}^*\|_2^2 + \tau^2(f(\mathbf{X}) - f^{opt})$, think about $\theta_t \approx t/2$.

Proof of Lemma 2. Take $\mathbf{x} = \frac{1}{\theta_t}\mathbf{x}^* + (1 - \frac{1}{\theta_t})\mathbf{x}^t$ and $\mathbf{y} = \mathbf{y}^t$ (based on FISTA $\mathbf{x}^{t+1} = \text{prox}_{\frac{1}{L}h}(\mathbf{y}^t - \frac{1}{L}\nabla f(\mathbf{y}^t))$, we have $\mathbf{x}^{t+1} = \text{prox}_{\frac{1}{L}h}(\mathbf{y} - \frac{1}{L}\nabla f(\mathbf{y}))$) in Lemma 1 to get

$$\begin{aligned}
 & F(\mathbf{x}^{t+1}) - F(\theta_t^{-1}\mathbf{x}^* + (1 - \theta_t^{-1})\mathbf{x}^t) \\
 & \leq \frac{L}{2}\|\theta_t^{-1}\mathbf{x}^* + (1 - \theta_t^{-1})\mathbf{x}^t - \mathbf{y}^t\|_2^2 - \frac{L}{2}\|\theta_t^{-1}\mathbf{x}^* + (1 - \theta_t^{-1})\mathbf{x}^t - \mathbf{x}^{t+1}\|_2^2 \\
 & = \frac{L}{2\theta_t^2}\|\mathbf{x}^* + (\theta_t - 1)\mathbf{x}^t - \theta_t\mathbf{y}^t\|_2^2 - \frac{L}{2\theta_t^2}\underbrace{\|\mathbf{x}^* + (\theta_t - 1)\mathbf{x}^t - \theta_t\mathbf{x}^{t+1}\|_2^2}_{=-\mathbf{u}^{t+1}} \quad (7) \\
 & \underbrace{=}_{(i)} \frac{L}{2\theta_t^2}(\|\mathbf{u}^t\|_2^2 - \|\mathbf{u}^{t+1}\|_2^2),
 \end{aligned}$$

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

We will also lower bound (7). By convexity of F ,

$$F\left(\theta_t^{-1}\mathbf{x}^* + (1 - \theta_t^{-1})\mathbf{x}^t\right) \leq \theta_t^{-1}F(\mathbf{x}^*) + (1 - \theta_t^{-1})F(\mathbf{x}^t) = \theta_t^{-1}F^{opt} + (1 - \theta_t^{-1})F(\mathbf{x}^t)$$

$$\Rightarrow F\left(\theta_t^{-1}\mathbf{x}^* + (1 - \theta_t^{-1})\mathbf{x}^t\right) - F(\mathbf{x}^{t+1}) \leq (1 - \theta_t^{-1})(F(\mathbf{x}^t) - F^{opt}) - (F(\mathbf{x}^{t+1}) - F^{opt})$$

Combining this with (7) (last equation) and $\theta_t^2 - \theta_t = \theta_{t-1}^2$ yields

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

thus finishing the proof.

Proof of Theorem 3. With Lemma 2, one has

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

To bound the RHS of this inequality, we use Lemma 1 and $\mathbf{y}^0 = \mathbf{x}^0$ ($\mathbf{y}^+ = \mathbf{x}^1$) and take $\mathbf{x} = \mathbf{x}^*$ to get

$$\begin{aligned} \frac{2}{L}(F(\mathbf{x}^1) - F^{opt}) &\leq \|\mathbf{y}^0 - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^1 - \mathbf{x}^*\|_2^2 = \|\mathbf{x}^0 - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^1 - \mathbf{x}^*\|_2^2 \\ &\Leftrightarrow \|\mathbf{x}^1 - \mathbf{x}^*\|_1^2 + \frac{2}{L}(F(\mathbf{x}^*) - F^{opt}) \leq \|\mathbf{x}^0 - \mathbf{x}^*\|_2^2 \end{aligned}$$

As a result,

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

Hence,

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

Lower bounds

Interestingly, no first-order methods can improve upon Nesterov's results in general. More precisely, \exists convex and L -smooth function f s.t.

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

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

Example

Consider $\min_{\mathbf{x} \in \mathbb{R}^{(2n+1)}} \frac{L}{4} \left(\frac{1}{2} \mathbf{x}^\top \mathbf{A} \mathbf{x} - \mathbf{e}_1^\top \mathbf{x} \right)$ where

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

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

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

Example

Also, $\nabla f(\mathbf{x}) = \frac{L}{4}\mathbf{A}\mathbf{x} - \frac{L}{4}\mathbf{e}_1$ and $\underbrace{\text{span}\{\nabla f(\mathbf{x}^0), \dots, \nabla f(\mathbf{x}^{k-1})\}}_{:=\mathcal{K}_k} = \text{span}\{\mathbf{e}_1, \dots, \mathbf{e}_k\}$ if

$\mathbf{x}^0 = 0$. That is, every iteration of first-order methods expands the search space by at most one dimension.

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

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

Numerical example

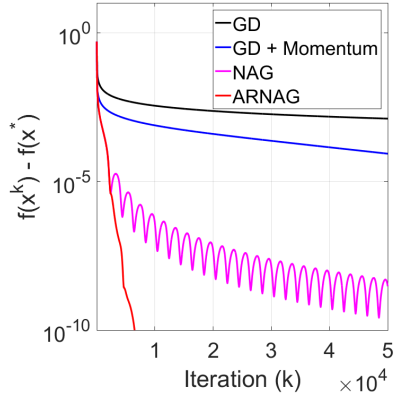
Consider

$$\min_{\mathbf{w}} f(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{L} \mathbf{w} - \mathbf{w}^T \mathbf{e}_1,$$

where

$$\mathbf{L} = \begin{pmatrix} 2 & -1 & 0 & \cdots & 0 & -1 \\ -1 & 2 & -1 & 0 & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ 0 & \cdots & 0 & -1 & 2 & -1 \\ -1 & 0 & \cdots & 0 & -1 & 2 \end{pmatrix}_{1000 \times 1000},$$

and \mathbf{e}_1 is a 1000-dim vector whose first entry is 1 and all the other entries are 0.



Nesterov's acceleration is not a monotonic method! We can further accelerated it via restart, resulting in linear convergence with some further assumption. See V. Roulet, A. d'Aspremont, "Sharpness, Restart and Acceleration", NeurIPS 2017.

Nesterov's method for strongly convex problems

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

Theorem 4. [Convergence of accelerated proximal gradient methods for strongly convex case] Suppose f is μ -strongly convex and L -smooth. If $\eta_t \equiv 1/L$, then

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