

Nesterov's Accelerated Gradient Descent on L -smooth convex function

Proof approach 1

Andersen Ang

ECS, Uni. Southampton, UK
andersen.ang@soton.ac.uk
Homepage angms.science

Version: November 4, 2023
First draft: August 2, 2017

Content

Problem setup: smooth unconstrained convex optimisation
Nesterov's accelerated gradient descent (NAGD)
Proving NAGD converges rate $\mathcal{O}\left(\frac{1}{k^2}\right)$
Summary

Problem setup: smooth unconstrained convex optimisation

$$(\mathcal{P}) : \operatorname{argmin}_{\mathbf{x}} f(\mathbf{x}).$$

► We consider Euclidean space

► $f : \mathbb{R}^n \rightarrow \mathbb{R}$

► f is L -smooth

► f is continuously differentiable

$f \in \mathcal{C}^1$, i.e., $\nabla f(\mathbf{x})$ exists for all $\mathbf{x} \in \operatorname{dom} f$

► ∇f is L -Lipschitz

$L > 0$ is the least upper bound in $\frac{\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|}{\|\mathbf{x} - \mathbf{y}\|} \leq L$

$$\forall \mathbf{a}, \mathbf{b} \in \operatorname{dom} f : f(\mathbf{a}) - f(\mathbf{b}) \leq \langle \nabla f(\mathbf{b}), \mathbf{a} - \mathbf{b} \rangle + \frac{L}{2} \|\mathbf{a} - \mathbf{b}\|_2^2.$$

► f is convex all local minima of \mathcal{P} are global minima

$$(\forall \mathbf{x} \in \operatorname{dom} f)(\forall \mathbf{y} \in \operatorname{dom} f) \left\{ f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \right\}$$

► Details of convexity, L -smoothness, see [here](#)

Gradient Descent (GD)

► Notation

$$f_k := f(\mathbf{x}_k)$$

$$f^* := f(\mathbf{x}^*)$$

► GD: start with initial point $\mathbf{x}_0 \in \mathbb{R}^n$, iterates

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k).$$

For sufficiently small stepsize ($\alpha_k < \frac{2}{L}$), the sequence $\{\mathbf{x}_k\}_{k \in \mathbb{N}}$ converges to a stationary point of f .

As f is convex, the sequence converges to the global minimizer \mathbf{x}^* (if exists).

► GD convergence as $f_k - f^* \leq \mathcal{O}\left(\frac{1}{k}\right)$

Nesterov's Accelerated Gradient Descent (NAGD)

$$(\mathcal{P}) : \min_{\mathbf{x}} f(\mathbf{x})$$

- Start with initial point $\mathbf{y}_0 = \mathbf{x}_0 \in \mathbb{R}^n$ and $\lambda_0 = 0$, iterates

$$\text{Gradient update} \quad \mathbf{y}_{k+1} = \mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k) \quad (1)$$

$$\text{Extrapolation} \quad \mathbf{x}_{k+1} = (1 - \gamma_k) \mathbf{y}_{k+1} + \gamma_k \mathbf{y}_k \quad (2)$$

$$\text{Extrapolation weight} \quad \gamma_k = \frac{1 - \lambda_k}{\lambda_{k+1}} \quad (3)$$

$$\text{Extrapolation weight} \quad \lambda_k = \frac{1 + \sqrt{1 + 4\lambda_{k-1}^2}}{2} \quad (4)$$

Note that here fix stepsize is used: $\alpha_k = \frac{1}{L} \forall k$.

- **Theorem.** If $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is L -smooth and convex, the sequences $\{f(\mathbf{y}_k)\}_k$ produced by NAGD converges to the optimal value f^* at the rate $\mathcal{O}\left(\frac{1}{k^2}\right)$ as

$$f(\mathbf{y}_k) - f^* \leq \frac{2L \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{k^2}.$$

- The convergence rate $\mathcal{O}\left(\frac{1}{k^2}\right)$ is optimal. I.e., no 1st-order algo. can perform better than NAGD in terms of convergence rate. All 1st-order algorithm can only be at most as good as NAGD. [Proof here](#).
- If f is nonconvex, the sequence $\{f(\mathbf{y}_k)\}_k$ produced by NAGD converges to the closest stationary point with the same convergence rate.

NAGD converges rate $\mathcal{O}\left(\frac{1}{k^2}\right)$ proof 1/6 **Stage 1: make use of convexity & smoothness**

- f cvx: $(\forall \mathbf{x} \forall \mathbf{y}) \left\{ f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \right\}$ gives

$$-f(\mathbf{y}) \leq -f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x} - \mathbf{y} \rangle \quad (5)$$

- f L -smooth $(\forall \mathbf{a} \forall \mathbf{b}) \left\{ f(\mathbf{a}) - f(\mathbf{b}) \leq \langle \nabla f(\mathbf{b}), \mathbf{a} - \mathbf{b} \rangle + \frac{L}{2} \|\mathbf{a} - \mathbf{b}\|_2^2 \right\}$, with $\mathbf{a} = \mathbf{x} - \frac{1}{L} \nabla f(\mathbf{x})$, $\mathbf{b} = \mathbf{x}$,

$$f\left(\mathbf{x} - \frac{1}{L} \nabla f(\mathbf{x})\right) - f(\mathbf{x}) \leq -\frac{1}{L} \|\nabla f(\mathbf{x})\|_2^2 + \frac{1}{2L} \|\nabla f(\mathbf{x})\|_2^2 = \frac{-1}{2L} \|\nabla f(\mathbf{x})\|_2^2. \quad (6)$$

- (5) + (6) will cancel $-f(\mathbf{x})$ and give

$$f\left(\mathbf{x} - \frac{1}{L} \nabla f(\mathbf{x})\right) - f(\mathbf{y}) \leq \frac{-1}{2L} \|\nabla f(\mathbf{x})\|_2^2 + \langle \nabla f(\mathbf{x}), \mathbf{x} - \mathbf{y} \rangle. \quad (7)$$

- Put $\mathbf{x} = \mathbf{x}_k$, $\mathbf{y} = \mathbf{x}^*$ in (7)

$$f\left(\mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k)\right) - f^* \leq \frac{-1}{2L} \|\nabla f(\mathbf{x}_k)\|_2^2 + \langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{x}^* \rangle. \quad (8)$$

- put $\mathbf{x} = \mathbf{x}_k$, $\mathbf{y} = \mathbf{y}_k$ in (7)

$$f\left(\mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k)\right) - f(\mathbf{y}_k) \leq \frac{-1}{2L} \|\nabla f(\mathbf{x}_k)\|_2^2 + \langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{y}_k \rangle. \quad (9)$$

- Proof overview: (8), (9) link $f(\mathbf{y}_{k+1})$, $f(\mathbf{y}_k)$ and f^* . We see $\nabla f(\mathbf{x}_k)$ appear in (8), (9) but not in the convergence result, so we eliminate $\nabla f(\mathbf{x}_k)$ in (8), (9).

Proof 2/6 Stage 2: eliminate gradient

$$\mathbf{y}_{k+1} = \mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k) \quad (1)$$

$$f\left(\mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k)\right) - f^* \leq \frac{-1}{2L} \|\nabla f(\mathbf{x}_k)\|_2^2 + \langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{x}^* \rangle. \quad (8)$$

$$f\left(\mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k)\right) - f(\mathbf{y}_k) \leq \frac{-1}{2L} \|\nabla f(\mathbf{x}_k)\|_2^2 + \langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{y}_k \rangle. \quad (9)$$

► Simplify notation, let $\delta_k := f(\mathbf{y}_k) - f^*$, then

$$f\left(\mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k)\right) \stackrel{(1)}{=} f(\mathbf{y}_{k+1}) \quad (10)$$

$$f\left(\mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k)\right) - f^* \stackrel{(10), \delta_k}{=} \delta_{k+1} \quad (11)$$

$$\begin{aligned} f\left(\mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k)\right) - f(\mathbf{y}_k) &= f\left(\mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k)\right) - f^* - (f(\mathbf{y}_k) - f^*) \\ &= \delta_{k+1} - \delta_k \end{aligned} \quad (12)$$

$$\nabla f(\mathbf{x}_k) \stackrel{(1)}{=} -L(\mathbf{y}_{k+1} - \mathbf{x}_k) \quad (13)$$

$$\|\nabla f(\mathbf{x}_k)\|_2^2 \stackrel{(13)}{=} L^2 \|\mathbf{y}_{k+1} - \mathbf{x}_k\|_2^2 \quad (14)$$

► Put (11,13,14) into (8)

$$\delta_{k+1} \leq -\frac{L}{2} \|\mathbf{y}_{k+1} - \mathbf{x}_k\|_2^2 - L \langle \mathbf{y}_{k+1} - \mathbf{x}_k, \mathbf{x}_k - \mathbf{x}^* \rangle. \quad (15)$$

► Put (12,13,14) into (9)

$$\delta_{k+1} - \delta_k \leq -\frac{L}{2} \|\mathbf{y}_{k+1} - \mathbf{x}_k\|_2^2 - L \langle \mathbf{y}_{k+1} - \mathbf{x}_k, \mathbf{x}_k - \mathbf{y}_k \rangle. \quad (16)$$

Proof 3/6 Stage 3: form telescoping sum

$$\begin{aligned}\lambda_k &= \frac{1}{2} \left(1 + \sqrt{1 + 4\lambda_{k-1}^2} \right) & (4) \\ \delta_{k+1} &\leq -\frac{L}{2} \|\mathbf{y}_{k+1} - \mathbf{x}_k\|_2^2 - L \langle \mathbf{y}_{k+1} - \mathbf{x}_k, \mathbf{x}_k - \mathbf{x}^* \rangle & (15) \\ \delta_{k+1} - \delta_k &\leq -\frac{L}{2} \|\mathbf{y}_{k+1} - \mathbf{x}_k\|_2^2 - L \langle \mathbf{y}_{k+1} - \mathbf{x}_k, \mathbf{x}_k - \mathbf{y}_k \rangle & (16)\end{aligned}$$

► **Tricky step:** consider (15) + $(\lambda_k - 1)(16)$.

$$\text{Left-hand side of (15) + } (\lambda_k - 1)(16) = \delta_{k+1} + (\lambda_k - 1)(\delta_{k+1} - \delta_k) = \lambda_k \delta_{k+1} - (\lambda_k - 1)\delta_k.$$

► Right-hand side of (15) + $(\lambda_k - 1)(16)$

$$\begin{aligned}& -\frac{L}{2} \|\mathbf{y}_{k+1} - \mathbf{x}_k\|_2^2 - L \langle \mathbf{y}_{k+1} - \mathbf{x}_k, \mathbf{x}_k - \mathbf{x}^* \rangle + (\lambda_k - 1) \left(-\frac{L}{2} \|\mathbf{y}_{k+1} - \mathbf{x}_k\|_2^2 - L \langle \mathbf{y}_{k+1} - \mathbf{x}_k, \mathbf{x}_k - \mathbf{y}_k \rangle \right) \\&= -\frac{\lambda_k L}{2} \|\mathbf{y}_{k+1} - \mathbf{x}_k\|_2^2 - L \langle \mathbf{y}_{k+1} - \mathbf{x}_k, \mathbf{x}_k - \mathbf{x}^* + (\lambda_k - 1)(\mathbf{x}_k - \mathbf{y}_k) \rangle \\&= -\frac{\lambda_k L}{2} \|\mathbf{y}_{k+1} - \mathbf{x}_k\|_2^2 - L \langle \mathbf{y}_{k+1} - \mathbf{x}_k, \lambda_k \mathbf{x}_k - (\lambda_k - 1)\mathbf{y}_k - \mathbf{x}^* \rangle\end{aligned}$$

► By LHS = RHS $\lambda_k \delta_{k+1} - (\lambda_k - 1)\delta_k \leq -\frac{\lambda_k L}{2} \|\mathbf{y}_{k+1} - \mathbf{x}_k\|_2^2 - L \langle \mathbf{y}_{k+1} - \mathbf{x}_k, \lambda_k \mathbf{x}_k - (\lambda_k - 1)\mathbf{y}_k - \mathbf{x}^* \rangle$.

Multiply the inequality with λ_k :

$$\begin{aligned}\lambda_k^2 \delta_{k+1} - \lambda_k (\lambda_k - 1)\delta_k &\leq -\frac{\lambda_k^2 L}{2} \|\mathbf{y}_{k+1} - \mathbf{x}_k\|_2^2 - \lambda_k L \langle \mathbf{y}_{k+1} - \mathbf{x}_k, \lambda_k \mathbf{x}_k - (\lambda_k - 1)\mathbf{y}_k - \mathbf{x}^* \rangle \\&= -\frac{L}{2} \left(\lambda_k^2 \|\mathbf{y}_{k+1} - \mathbf{x}_k\|_2^2 + 2\lambda_k \langle \mathbf{y}_{k+1} - \mathbf{x}_k, \lambda_k \mathbf{x}_k - (\lambda_k - 1)\mathbf{y}_k - \mathbf{x}^* \rangle \right). \quad (\#)\end{aligned}$$

► (4) gives $(2\lambda_k - 1)^2 = 1 + 4\lambda_{k-1}^2 \iff 4\lambda_k^2 - 4\lambda_k + 1 = 1 + 4\lambda_{k-1}^2 \iff \lambda_{k-1}^2 = \lambda_k(\lambda_k - 1)$, put this into (#) gives

$$\lambda_k^2 \delta_{k+1} - \lambda_{k-1}^2 \delta_k \leq -\frac{L}{2} \left(\lambda_k^2 \|\mathbf{y}_{k+1} - \mathbf{x}_k\|_2^2 + 2\lambda_k \langle \mathbf{y}_{k+1} - \mathbf{x}_k, \lambda_k \mathbf{x}_k - (\lambda_k - 1)\mathbf{y}_k - \mathbf{x}^* \rangle \right) \quad (17)$$

Proof 4/6

$$\lambda_k = \frac{1}{2} \left(1 + \sqrt{1 + 4\lambda_{k-1}^2} \right) \quad (4)$$

$$\lambda_k^2 \delta_{k+1} - \lambda_k (\lambda_k - 1) \delta_k \leq -\frac{L}{2} \left(\lambda_k^2 \|\mathbf{y}_{k+1} - \mathbf{x}_k\|_2^2 + 2\lambda_k \langle \mathbf{y}_{k+1} - \mathbf{x}_k, \lambda_k \mathbf{x}_k - (\lambda_k - 1) \mathbf{y}_k - \mathbf{x}^* \rangle \right) \quad (17)$$

► Inspecting the inner product in (17) we see that it is completing squares (Thanks to Tony Silveti-Falls for figuring it out, 2023 Nov 3).

$$\|\lambda \mathbf{a} + \mathbf{b}\|_2^2 = \lambda^2 \|\mathbf{a}\|_2^2 + 2\lambda \langle \mathbf{a}, \mathbf{b} \rangle + \|\mathbf{b}\|_2^2 \iff \lambda^2 \|\mathbf{a}\|_2^2 + 2\lambda \langle \mathbf{a}, \mathbf{b} \rangle = \|\lambda \mathbf{a} + \mathbf{b}\|_2^2 - \|\mathbf{b}\|_2^2.$$

$$\begin{aligned} & \lambda_k^2 \|\mathbf{y}_{k+1} - \mathbf{x}_k\|_2^2 + 2\lambda_k \langle \mathbf{y}_{k+1} - \mathbf{x}_k, \lambda_k \mathbf{x}_k - (\lambda_k - 1) \mathbf{y}_k - \mathbf{x}^* \rangle \\ &= \|\lambda(\mathbf{y}_{k+1} - \mathbf{x}_k) + \lambda_k \mathbf{x}_k - (\lambda_k - 1) \mathbf{y}_k - \mathbf{x}^*\|_2^2 - \|\lambda_k \mathbf{x}_k - (\lambda_k - 1) \mathbf{y}_k - \mathbf{x}^*\|_2^2 \\ &= \|\lambda_k \mathbf{y}_{k+1} - (\lambda_k - 1) \mathbf{y}_k - \mathbf{x}^*\|_2^2 - \|\lambda_k \mathbf{x}_k - (\lambda_k - 1) \mathbf{y}_k - \mathbf{x}^*\|_2^2. \end{aligned}$$

► Using this (17) becomes

$$\lambda_k^2 \delta_{k+1} - \lambda_{k-1}^2 \delta_k \leq -\frac{L}{2} \left(\|\lambda_k \mathbf{y}_{k+1} - (\lambda_k - 1) \mathbf{y}_k - \mathbf{x}^*\|_2^2 - \|\lambda_k \mathbf{x}_k - (\lambda_k - 1) \mathbf{y}_k - \mathbf{x}^*\|_2^2 \right). \quad (18)$$

► We have $\lambda_k \mathbf{x}_k - (\lambda_k - 1) \mathbf{y}_k = (1 - \lambda_{k-1}) \mathbf{y}_{k-1} + \lambda_{k-1} \mathbf{y}_k$.

Proof: $\gamma_k \stackrel{(3)}{=} \frac{1 - \lambda_k}{\lambda_{k+1}} \iff \gamma_k \lambda_{k+1} = 1 - \lambda_k.$

By (2) $\mathbf{x}_{k+1} = (1 - \gamma_k) \mathbf{y}_{k+1} + \gamma_k \mathbf{y}_k$ gives $\mathbf{x}_{k+1} = \mathbf{y}_{k+1} + \gamma_k (\mathbf{y}_k - \mathbf{y}_{k+1})$, multiply with λ_{k+1} gives $\lambda_{k+1} \mathbf{x}_{k+1} = \lambda_{k+1} \mathbf{y}_{k+1} + \lambda_{k+1} \gamma_k (\mathbf{y}_k - \mathbf{y}_{k+1}) = \lambda_{k+1} \mathbf{y}_{k+1} + (1 - \lambda_k) (\mathbf{y}_k - \mathbf{y}_{k+1})$, rearrange gives $\lambda_{k+1} \mathbf{x}_{k+1} - \lambda_{k+1} \mathbf{y}_{k+1} = (1 - \lambda_k) (\mathbf{y}_k - \mathbf{y}_{k+1})$, add \mathbf{y}_{k+1} on both side gives $\lambda_{k+1} \mathbf{x}_{k+1} - (\lambda_{k+1} - 1) \mathbf{y}_{k+1} = (1 - \lambda_k) \mathbf{y}_k + \lambda_k \mathbf{y}_{k+1}$. Move counter k by -1 gives the result.

So (18) becomes

$$\lambda_k^2 \delta_{k+1} - \lambda_{k-1}^2 \delta_k \leq -\frac{L}{2} \left(\|\lambda_k \mathbf{y}_{k+1} - (\lambda_k - 1) \mathbf{y}_k - \mathbf{x}^*\|_2^2 - \|(1 - \lambda_{k-1}) \mathbf{y}_{k-1} + \lambda_{k-1} \mathbf{y}_k - \mathbf{x}^*\|_2^2 \right).$$

Proof ... 5/6

We have $\lambda_k^2 \delta_{k+1} - \lambda_{k-1}^2 \delta_k \leq -\frac{L}{2} \left(\|\lambda_k \mathbf{y}_{k+1} - (\lambda_k - 1) \mathbf{y}_k - \mathbf{x}^*\|_2^2 - \|(1 - \lambda_{k-1}) \mathbf{y}_{k-1} + \lambda_{k-1} \mathbf{y}_k - \mathbf{x}^*\|_2^2 \right).$

Rearrange the second term to make the terms in right-hand side have similar form

$$\lambda_k^2 \delta_{k+1} - \lambda_{k-1}^2 \delta_k \leq -\frac{L}{2} \left(\|\lambda_k \mathbf{y}_{k+1} - (\lambda_k - 1) \mathbf{y}_k - \mathbf{x}^*\|_2^2 - \|\lambda_{k-1} \mathbf{y}_k - (\lambda_{k-1} - 1) \mathbf{y}_{k-1} - \mathbf{x}^*\|_2^2 \right). \quad (19)$$

Let $\mathbf{u}_k = \lambda_k \mathbf{y}_{k+1} - (\lambda_k - 1) \mathbf{y}_k - \mathbf{x}^*$ so $\lambda_{k-1} \mathbf{y}_k - (\lambda_{k-1} - 1) \mathbf{y}_{k-1} - \mathbf{x}^* = \mathbf{u}_{k-1}$ and (19) becomes

$$\begin{aligned} \lambda_k^2 \delta_{k+1} - \lambda_{k-1}^2 \delta_k &\leq -\frac{L}{2} \left(\|\mathbf{u}_k\|_2^2 - \|\mathbf{u}_{k-1}\|_2^2 \right) \\ \lambda_1^2 \delta_2 - \lambda_0^2 \delta_1 &\leq -\frac{L}{2} \left(\|\mathbf{u}_1\|_2^2 - \|\mathbf{u}_0\|_2^2 \right) && \text{case } k = 1 \\ \lambda_2^2 \delta_3 - \lambda_1^2 \delta_2 &\leq -\frac{L}{2} \left(\|\mathbf{u}_2\|_2^2 - \|\mathbf{u}_1\|_2^2 \right) && \text{case } k = 2 \\ &\vdots \\ \lambda_{K-1}^2 \delta_K - \lambda_{K-2}^2 \delta_{K-1} &\leq -\frac{L}{2} \left(\|\mathbf{u}_{K-1}\|_2^2 - \|\mathbf{u}_{K-2}\|_2^2 \right) && \text{case } k = K - 1 \\ \lambda_{K-1}^2 \delta_K - \lambda_0^2 \delta_1 &\leq -\frac{L}{2} \left(\|\mathbf{u}_{K-1}\|_2^2 - \|\mathbf{u}_0\|_2^2 \right) && \text{sum } k = 1 \text{ to } k = K - 1 \\ &= \frac{L}{2} \left(\|\mathbf{u}_0\|_2^2 - \|\mathbf{u}_{K-1}\|_2^2 \right) \\ &\leq \frac{L}{2} \|\mathbf{u}_0\|_2^2 && \|\mathbf{u}_{K-1}\|_2^2 \geq 0 \end{aligned}$$

By definition, $\lambda_0 = 0$, $\mathbf{y}_0 = \mathbf{x}_0$, $\mathbf{u}_0 = \lambda_0 \mathbf{y}_1 - (\lambda_0 - 1) \mathbf{y}_0 - \mathbf{x}^* \stackrel{\lambda_0=0}{=} \mathbf{y}_0 - \mathbf{x}^* \stackrel{\mathbf{y}_0=\mathbf{x}_0}{=} \mathbf{x}_0 - \mathbf{x}^*$, thus

$$\lambda_{K-1}^2 \delta_K \leq \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 \implies \delta_K \leq \frac{L \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{2\lambda_{K-1}^2}.$$

Proof ... 6/6

Lemma. $\lambda_{k-1} \geq \frac{k}{2}$.

Proof (by induction)

► Case $k = 0$ and $\lambda_0 = 0$. It is trivial $0 \geq 0/2$.

► Case $k = 1$. By definition,

$$\lambda_k = \frac{1 + \sqrt{1 + 4\lambda_{k-1}^2}}{2} = \frac{1 + \sqrt{1 + 4 \cdot 0^2}}{2} = 1 > \frac{1}{2} = \frac{k}{2} \Big|_{k=1}$$

► Induction hypothesis: assume $\lambda_{n-1} \geq \frac{n}{2}$.

► Case $k = n$

$$\begin{aligned} \lambda_n &= \frac{1 + \sqrt{1 + 4\lambda_{n-1}^2}}{2} \\ &\geq \frac{1 + \sqrt{1 + 4\left(\frac{n}{2}\right)^2}}{2} && \text{[Induction hypothesis]} \\ &= \frac{1 + \sqrt{1 + n^2}}{2} \\ &> \frac{1 + \sqrt{n^2}}{2} \\ &= \frac{1 + n}{2}. \quad \square \end{aligned}$$

With $\lambda_{k-1} \geq \frac{k}{2}$, so

$$\frac{1}{\lambda_{k-1}^2} \leq \frac{4}{k^2}.$$

Therefore $\delta_K \leq \frac{L\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{2\lambda_{K-1}^2}$ becomes

$$f(\mathbf{y}_K) - f^* \leq \frac{2L\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{K^2}.$$

where $f(\mathbf{y}_K) - f^* =: \delta_K$. \square

Rename K as k gives

$$f(\mathbf{y}_k) - f^* \leq \frac{2L\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{k^2}.$$

This $\begin{cases} \text{complicated} \\ \text{highly-involved} \\ \text{non-intuitive} \end{cases}$ proof is now completed.

Last page - summary

For unconstrained convex smooth problem

$$(\mathcal{P}) : \underset{\mathbf{x}}{\operatorname{argmin}} f(\mathbf{x})$$

with $f : \mathbb{R}^n \rightarrow \mathbb{R}$ being convex, L -smooth, the NAGD algorithm starts with initial point $\mathbf{x}_0 = \mathbf{y}_0 \in \mathbb{R}^n$ and $\lambda_0 = 0$ and iterates the following:

$$\begin{array}{lll} \text{Gradient update} & \mathbf{y}_{k+1} & = \mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k) \\ \text{Extrapolation} & \mathbf{x}_{k+1} & = (1 - \gamma_k) \mathbf{y}_{k+1} + \gamma_k \mathbf{y}_k \\ \text{Extrapolation weight} & \gamma_k & = \frac{1 - \lambda_k}{\lambda_{k+1}} \\ \text{Extrapolation weight} & \lambda_k & = \frac{1 + \sqrt{1 + 4\lambda_{k-1}^2}}{2} \end{array}$$

the sequences $\{f(\mathbf{y}_k)\}_{k \in \mathbb{N}}$ produced will converges to the optimal f^* at order of $\mathcal{O}\left(\frac{1}{k^2}\right)$ as

$$f(\mathbf{y}_k) - f^* \leq \frac{2L \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{k^2}.$$

The proof can be used for proximal gradient descent.

End of document