

Nesterov's accelerated gradient method

on m -strongly convex L -smooth function converges at $\mathcal{O}(\exp \frac{-k}{\sqrt{Q}})$

Andersen Ang

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

Homepage angms.science

Version: July 21, 2023
First draft: August 2, 2017

Content

Nesterov's estimate sequence

$$\Phi_0(\mathbf{x}) := f(\mathbf{x}_0) + \frac{m}{2} \|\mathbf{x} - \mathbf{x}_0\|_2^2$$

$$\Phi_{k+1}(\mathbf{x}) := \left(1 - \frac{1}{\sqrt{Q}}\right) \Phi_k(\mathbf{x}) + \frac{1}{\sqrt{Q}} \left(f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{m}{2} \|\mathbf{x} - \mathbf{x}_k\|_2^2 \right)$$

Lemma 1 $\Phi_k(\mathbf{x}) \leq f(\mathbf{x}) + \left(1 - \frac{1}{\sqrt{Q}}\right)^k \left(\Phi_0(\mathbf{x}) - f(\mathbf{x}) \right)$

Lemma 2 $\nabla^2 \Phi_k(\mathbf{x}) = m\mathbf{I}_n$

Lemma 3 $f(\mathbf{y}_k) \leq \Phi_k^* := \min_{\mathbf{x} \in \mathbb{R}^n} \Phi_k(\mathbf{x})$

Lemma 4 $\nu_k - \mathbf{x}_k = \sqrt{Q}(\mathbf{x}_k - \mathbf{y}_K)$

NAG convergence rate $f(\mathbf{y}_k) - f^* \leq \left(\frac{m+L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 \right) \exp \frac{-k}{\sqrt{Q}}$

Problem setup: unconstrained strongly convex smooth optimisation

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

► We consider Euclidean space

► $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is L -smooth

► f is continuously differentiable

► ∇f is globally L -Lipschitz

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

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

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

► $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is m -strongly convex

► f is convex

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

► f is m -strongly convex

► $f(\mathbf{x}) - \frac{m}{2} \|\mathbf{x}\|_2^2$ is convex

the global minima of \mathcal{P} is unique

all local minima of \mathcal{P} are global minima

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

Gradient Descent (GD)

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

- ▶ GD starts with initial point $\mathbf{x}_0 \in \mathbb{R}^n$, iterates

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

If stepsize is sufficiently small ($m_k < \frac{2}{L}$), then $\{\mathbf{x}_k\}_{k \in \mathbb{N}}$ converges to a stationary point of f .

- ▶ f convex \implies {all local minimizers are global}

$\{\mathbf{x}_k\}_{k \in \mathbb{N}} \rightarrow$ a global minimizer \mathbf{x}^* (if it exists)

- ▶ f strongly convex \implies {unique minimizer}

global minimizer \mathbf{x}^* is unique (if it exists)

- ▶ Notation $f^* := f(\mathbf{x}^*)$ and $Q = \frac{L}{m}$.

- ▶ If f is L -smooth and convex, $f_k - f^* \leq \mathcal{O}\left(\frac{1}{k}\right)$

convergence rate on $\{f_k\}_{k \in \mathbb{N}}$ is $\mathcal{O}\left(\frac{1}{k}\right)$

Details

- ▶ If f is L -smooth and m -strongly convex, $f_k - f^* \leq \mathcal{O}\left(\exp \frac{-k}{Q}\right)$

convergence rate on $\{f_k\}_{k \in \mathbb{N}}$ is $\mathcal{O}\left(\exp \frac{-k}{Q}\right)$

Details

Nesterov's accelerated gradient (NAG) method

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

If f is L -smooth and convex

Algorithm 1: NAG (for convex smooth f)

1 Initialize $\mathbf{x}_0 \in \mathbb{R}^n$, $\lambda_1 = 1$

2 **while** not converge **do**

3

$$\begin{aligned} \mathbf{y}_{k+1} &= \mathbf{x}_k - \frac{\nabla f(\mathbf{x}_k)}{L} \\ \mathbf{x}_{k+1} &= (1 - \gamma_k)\mathbf{y}_{k+1} + \gamma_k\mathbf{y}_k \\ \gamma_k &= \frac{1 - \lambda_k}{\lambda_{k+1}} \\ \lambda_k &= \frac{1 + \sqrt{1 + 4\lambda_{k-1}^2}}{2} \end{aligned}$$

Theorem The sequence $\{f(\mathbf{x}_k)\}_{k \in \mathbb{N}}$ produced by NAG on convex L -smooth function satisfies

$$f(\mathbf{y}_k) - f^* \leq \left(\frac{1}{k^2}\right).$$

Details

If f is L -smooth and m -strongly convex

Fix $\gamma_k = \frac{\sqrt{Q} - 1}{\sqrt{Q} + 1}$ where $Q = \frac{L}{m}$

Algorithm 2: NAG (for strongly convex smooth f)

1 Initialize $\mathbf{x}_0 \in \mathbb{R}^n$

2 **while** not converge **do**

3

$$\begin{aligned} \mathbf{y}_{k+1} &= \mathbf{x}_k - \frac{1}{L}\nabla f(\mathbf{x}_k) \\ \mathbf{x}_{k+1} &= \left(1 - \frac{\sqrt{Q} - 1}{\sqrt{Q} + 1}\right)\mathbf{y}_{k+1} + \frac{\sqrt{Q} - 1}{\sqrt{Q} + 1}\mathbf{y}_k \end{aligned}$$

Theorem The sequence $\{f(\mathbf{x}_k)\}_{k \in \mathbb{N}}$ produced by NAG on m -strongly convex L -smooth function satisfies

$$f(\mathbf{y}_k) - f^* \leq \frac{m + L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 \exp\left(\frac{-k}{\sqrt{Q}}\right).$$

This pdf: prove this.

Convergence rate of NAG - proof idea: Nesterov's estimate sequence

- ▶ There are a few ways to prove the convergence of NAG.
- ▶ A way is to use a non-trivial technique known as the Nesterov's estimate sequence.
- ▶ Consider a sequence of function $\{\Phi_k(\mathbf{x})\}_{k \in \mathbb{N}}$ that
 - ▶ $\Phi_k(\mathbf{x})$ has a general structure with "parameters" varies with iteration k .
 - ▶ $\Phi_k(\mathbf{x})$ is based on f
 - ▶ $\Phi_k(\mathbf{x})$ is m -strongly convex
- ▶ $\Phi_k(\mathbf{x})$ can be defined as

$$\begin{aligned}\Phi_0(\mathbf{x}) &:= f(\mathbf{x}_0) + \frac{m}{2} \|\mathbf{x} - \mathbf{x}_0\|_2^2 \\ \Phi_{k+1}(\mathbf{x}) &:= \left(1 - \frac{1}{\sqrt{Q}}\right) \Phi_k(\mathbf{x}) + \frac{1}{\sqrt{Q}} \left(f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{m}{2} \|\mathbf{x} - \mathbf{x}_k\|_2^2 \right)\end{aligned}$$

Details of the theory of Nesterov's estimating sequence.

Understanding Nesterov's estimate sequence

$$\Phi_0(\mathbf{x}) := f(\mathbf{x}_0) + \frac{m}{2} \|\mathbf{x} - \mathbf{x}_0\|_2^2$$

$$\Phi_{k+1}(\mathbf{x}) := \left(1 - \frac{1}{\sqrt{Q}}\right) \Phi_k(\mathbf{x}) + \frac{1}{\sqrt{Q}} \left(\underbrace{f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{m}{2} \|\mathbf{x} - \mathbf{x}_k\|_2^2}_{\text{2nd-order Taylor approximation of } f \text{ at } \mathbf{x}_k} \right)$$

- ▶ $\Phi_k(\mathbf{x})$ is based on f
- ▶ $\Phi_k(\mathbf{x})$ is m -strongly convex
- ▶ $\Phi_k(\mathbf{x})$ varies with iteration k .
- ▶ We can see Φ_{k+1} is in the form $\Phi_{k+1} = (1 - \lambda)a + \lambda b$.
 - ▶ Φ_{k+1} is a convex combination of Φ_k and the 2nd-order Taylor approximation of f at \mathbf{x}_k .
 - ▶ $Q = 1 \iff$ the level sets of f is circular: Φ_{k+1} is more like the Taylor approximation
In fact by definition of NAG, if $Q = 1$, there is no acceleration and NAG reduces to GD.
In this case GD should solve the optimization problem in 1 step. [Details](#).
 - ▶ $Q \gg 1 \iff$ the level sets of f is elliptic: Φ_{k+1} is more like previous Φ_k

The derivatives of Nesterov's estimating sequence

$$\Phi_0(\mathbf{x}) := f(\mathbf{x}_0) + \frac{m}{2} \|\mathbf{x} - \mathbf{x}_0\|_2^2$$

$$\Phi_{k+1}(\mathbf{x}) := \left(1 - \frac{1}{\sqrt{Q}}\right) \Phi_k(\mathbf{x}) + \frac{1}{\sqrt{Q}} \left(\underbrace{f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{m}{2} \|\mathbf{x} - \mathbf{x}_k\|_2^2}_{\text{Taylor approximation of } f \text{ at } \mathbf{x}_k} \right)$$

- ▶ With respect to \mathbf{x} , the gradient and Hessian are

$$\nabla \Phi_0(\mathbf{x}) = m(\mathbf{x} - \mathbf{x}_0) \tag{1}$$

$$\nabla^2 \Phi_0(\mathbf{x}) = m\mathbf{I}_n \tag{2}$$

$$\nabla \Phi_{k+1}(\mathbf{x}) = \left(1 - \frac{1}{\sqrt{Q}}\right) \nabla \Phi_k(\mathbf{x}) + \frac{1}{\sqrt{Q}} \left(\nabla f(\mathbf{x}_k) + m(\mathbf{x} - \mathbf{x}_k) \right) \tag{3}$$

$$\nabla^2 \Phi_{k+1}(\mathbf{x}) = \left(1 - \frac{1}{\sqrt{Q}}\right) \nabla^2 \Phi_k(\mathbf{x}) + \frac{1}{\sqrt{Q}} m\mathbf{I}_n \tag{4}$$

- ▶ In other words,

- ▶ Φ_{k+1} is a convex combination of Φ_k and the 2nd-order Taylor approximation of f at \mathbf{x}_k .

- ▶ $\nabla \Phi_{k+1}(\mathbf{x})$ is a convex combination of $\nabla \Phi_k(\mathbf{x})$ and $\nabla f(\mathbf{x}_k) + m(\mathbf{x} - \mathbf{x}_k)$.

- ▶ $\nabla^2 \Phi_{k+1}(\mathbf{x})$ is a convex combination of $\nabla^2 \Phi_k(\mathbf{x})$ and $m\mathbf{I}_n$.

In fact we are going to show $\nabla^2 \Phi_{k+1}(\mathbf{x}) = m\mathbf{I}_n$ in Lemma 2.

- ▶ In fact the derivatives of Φ_k plays an important role in the whole proof.

$\Phi_k(\mathbf{x})$ with $k = 0, 1$

$$\Phi_{k+1}(\mathbf{x}) := \left(1 - \frac{1}{\sqrt{Q}}\right)\Phi_k(\mathbf{x}) + \frac{1}{\sqrt{Q}}\left(f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{m}{2}\|\mathbf{x} - \mathbf{x}_k\|_2^2\right)$$

$$\underbrace{f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{m}{2}\|\mathbf{x} - \mathbf{x}_k\|_2^2}_{\leq f(\mathbf{x}) \forall \mathbf{x}, \mathbf{x}_k}$$

f is m -strongly cvx

$$\Phi_0(\mathbf{x}) := f(\mathbf{x}_0) + \frac{m}{2}\|\mathbf{x} - \mathbf{x}_0\|_2^2$$

$$\Phi_1(\mathbf{x}) := \left(1 - \frac{1}{\sqrt{Q}}\right)\Phi_0(\mathbf{x}) + \frac{1}{\sqrt{Q}}\left(f(\mathbf{x}_0) + \langle \nabla f(\mathbf{x}_0), \mathbf{x} - \mathbf{x}_0 \rangle + \frac{m}{2}\|\mathbf{x} - \mathbf{x}_0\|_2^2\right)$$

$$= f(\mathbf{x}) + \left(1 - \frac{1}{\sqrt{Q}}\right)\Phi_0(\mathbf{x}) + \frac{1}{\sqrt{Q}}\left(\underbrace{f(\mathbf{x}_0) + \langle \nabla f(\mathbf{x}_0), \mathbf{x} - \mathbf{x}_0 \rangle + \frac{m}{2}\|\mathbf{x} - \mathbf{x}_0\|_2^2}_{\leq f(\mathbf{x})} - f(\mathbf{x})\right)$$

$$\leq f(\mathbf{x}) + \left(1 - \frac{1}{\sqrt{Q}}\right)\Phi_0(\mathbf{x}) + \frac{1}{\sqrt{Q}}\underbrace{f(\mathbf{x}) - f(\mathbf{x})}_{=0}$$

$$= f(\mathbf{x}) + \left(1 - \frac{1}{\sqrt{Q}}\right)\Phi_0(\mathbf{x}) - \left(1 - \frac{1}{\sqrt{Q}}\right)f(\mathbf{x})$$

$$= f(\mathbf{x}) + \left(1 - \frac{1}{\sqrt{Q}}\right)\left(\Phi_0(\mathbf{x}) - f(\mathbf{x})\right)$$

$\Phi_k(\mathbf{x})$ with $k = 2$

$$\Phi_{k+1}(\mathbf{x}) := \left(1 - \frac{1}{\sqrt{Q}}\right)\Phi_k(\mathbf{x}) + \frac{1}{\sqrt{Q}}\left(f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{m}{2}\|\mathbf{x} - \mathbf{x}_k\|_2^2\right)$$

$$\underbrace{f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{m}{2}\|\mathbf{x} - \mathbf{x}_k\|_2^2}_{\leq f(\mathbf{x}) \quad \forall \mathbf{x}, \mathbf{x}_k}$$

f is m -strongly cvx

$$\Phi_0(\mathbf{x}) := f(\mathbf{x}_0) + \frac{m}{2}\|\mathbf{x} - \mathbf{x}_0\|_2^2$$

$$\Phi_1(\mathbf{x}) \leq f(\mathbf{x}) + \left(1 - \frac{1}{\sqrt{Q}}\right)\left(\Phi_0(\mathbf{x}) - f(\mathbf{x})\right)$$

$$\Phi_2(\mathbf{x}) := \left(1 - \frac{1}{\sqrt{Q}}\right)\Phi_1(\mathbf{x}) + \frac{1}{\sqrt{Q}}\left(f(\mathbf{x}_1) + \langle \nabla f(\mathbf{x}_1), \mathbf{x} - \mathbf{x}_1 \rangle + \frac{m}{2}\|\mathbf{x} - \mathbf{x}_1\|_2^2\right)$$

$$= f(\mathbf{x}) + \left(1 - \frac{1}{\sqrt{Q}}\right)\Phi_1(\mathbf{x}) + \frac{1}{\sqrt{Q}}\left(\underbrace{f(\mathbf{x}_1) + \langle \nabla f(\mathbf{x}_1), \mathbf{x} - \mathbf{x}_1 \rangle + \frac{m}{2}\|\mathbf{x} - \mathbf{x}_1\|_2^2}_{\leq f(\mathbf{x})}\right) - f(\mathbf{x})$$

$$\leq f(\mathbf{x}) + \left(1 - \frac{1}{\sqrt{Q}}\right)\Phi_1(\mathbf{x}) + \frac{1}{\sqrt{Q}}f(\mathbf{x}) - f(\mathbf{x})$$

$$= f(\mathbf{x}) + \left(1 - \frac{1}{\sqrt{Q}}\right)\Phi_1(\mathbf{x}) - \left(1 - \frac{1}{\sqrt{Q}}\right)f(\mathbf{x})$$

$$= f(\mathbf{x}) + \left(1 - \frac{1}{\sqrt{Q}}\right)\left(\Phi_1(\mathbf{x}) - f(\mathbf{x})\right)$$

$$\leq f(\mathbf{x}) + \left(1 - \frac{1}{\sqrt{Q}}\right)\left(1 - \frac{1}{\sqrt{Q}}\right)\left(\Phi_0(\mathbf{x}) - f(\mathbf{x})\right) = f(\mathbf{x}) + \left(1 - \frac{1}{\sqrt{Q}}\right)^2\left(\Phi_0(\mathbf{x}) - f(\mathbf{x})\right)$$

Lemma 1

$$\Phi_{k+1}(\mathbf{x}) := \left(1 - \frac{1}{\sqrt{Q}}\right) \Phi_k(\mathbf{x}) + \frac{1}{\sqrt{Q}} \left(f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{m}{2} \|\mathbf{x} - \mathbf{x}_k\|_2^2 \right)$$
$$\underbrace{f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{m}{2} \|\mathbf{x} - \mathbf{x}_k\|_2^2}_{\leq f(\mathbf{x}) \quad \forall \mathbf{x}, \mathbf{x}_k}$$

f is m -strongly cvx

$$\Phi_0(\mathbf{x}) := f(\mathbf{x}_0) + \frac{m}{2} \|\mathbf{x} - \mathbf{x}_0\|_2^2$$

$$\Phi_1(\mathbf{x}) \leq f(\mathbf{x}) + \left(1 - \frac{1}{\sqrt{Q}}\right) \left(\Phi_0(\mathbf{x}) - f(\mathbf{x}) \right)$$

$$\Phi_2(\mathbf{x}) \leq f(\mathbf{x}) + \left(1 - \frac{1}{\sqrt{Q}}\right)^2 \left(\Phi_0(\mathbf{x}) - f(\mathbf{x}) \right)$$

Lemma 1 For all $k \in \mathbb{N} = \{1, 2, \dots\}$,

$$\Phi_k(\mathbf{x}) \leq f(\mathbf{x}) + \left(1 - \frac{1}{\sqrt{Q}}\right)^k \left(\Phi_0(\mathbf{x}) - f(\mathbf{x}) \right).$$

Proof by induction

- ▶ Based case is already proved.
- ▶ For case $k + 1$, repeat the procedure on deriving Φ_2 and make use of the induction hypothesis.

Lemma 2 $\nabla^2\Phi_k(\mathbf{x}) = m\mathbf{I}_n$

$$\begin{aligned}\Phi_0(\mathbf{x}) &:= f(\mathbf{x}_0) + \frac{m}{2}\|\mathbf{x} - \mathbf{x}_0\|_2^2 \\ \Phi_{k+1}(\mathbf{x}) &:= \left(1 - \frac{1}{\sqrt{Q}}\right)\Phi_k(\mathbf{x}) + \frac{1}{\sqrt{Q}}\left(f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{m}{2}\|\mathbf{x} - \mathbf{x}_k\|_2^2\right)\end{aligned}$$

Proof by induction

► **Base case** $k = 0$

$\nabla^2\Phi_0(\mathbf{x}) = m\mathbf{I}_n$ by definition.

► **Induction Hypothesis** $\nabla^2\Phi_k(\mathbf{x}) = m\mathbf{I}_n$

► **Case** $k + 1$

$$\Phi_{k+1}(\mathbf{x}) = \left(1 - \frac{1}{\sqrt{Q}}\right)\Phi_k(\mathbf{x}) + \frac{1}{\sqrt{Q}}\left(f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{m}{2}\|\mathbf{x} - \mathbf{x}_k\|_2^2\right) \quad \text{by definition}$$

$$\nabla^2\Phi_{k+1}(\mathbf{x}) = \left(1 - \frac{1}{\sqrt{Q}}\right)\nabla^2\Phi_k(\mathbf{x}) + \frac{1}{\sqrt{Q}}m\mathbf{I}_n$$

$$= \left(1 - \frac{1}{\sqrt{Q}}\right)m\mathbf{I}_n + \frac{1}{\sqrt{Q}}m\mathbf{I}_n$$

$$= m\mathbf{I}_n \quad \square$$

induction hypothesis

Lemma 3 $f(\mathbf{y}_k) \leq \Phi_k^* := \min_{\mathbf{x} \in \mathbb{R}^n} \Phi_k(\mathbf{x}) \dots 1/7$

| | | |
|---------------------------------|--|---------------|
| $\Phi_0(\mathbf{x})$ | $:= f(\mathbf{x}_0) + \frac{m}{2} \ \mathbf{x} - \mathbf{x}_0\ _2^2$ | estimate seq. |
| \mathbf{x}_0 | $= \mathbf{y}_0$ | NAG def. |
| \mathbf{y}_{k+1} | $= \mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k)$ | NAG def. |
| $f(\mathbf{a}) - f(\mathbf{b})$ | $\leq \langle \nabla f(\mathbf{b}), \mathbf{a} - \mathbf{b} \rangle + \frac{L}{2} \ \mathbf{a} - \mathbf{b}\ _2^2$ | L-smooth |

Proof by induction

► **Base case** $k = 0$

$$\Phi_0^* = \min_{\mathbf{x} \in \mathbb{R}^n} \Phi_0(\mathbf{x}_0) = \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}_0) + \frac{m}{2} \|\mathbf{x} - \mathbf{x}_0\|_2^2 = f(\mathbf{x}_0) = f(\mathbf{y}_0)$$

► **Induction Hypothesis** $f(\mathbf{y}_k) \leq \Phi_k^*$

► **Case** $k + 1$ Consider $f(\mathbf{y}_{k+1})$ and L -smoothness of f

$$\begin{aligned} f(\mathbf{y}_{k+1}) &\leq f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{y}_{k+1} - \mathbf{x}_k \rangle + \frac{L}{2} \|\mathbf{y}_{k+1} - \mathbf{x}_k\|_2^2 \\ &= f(\mathbf{x}_k) + \left\langle \nabla f(\mathbf{x}_k), \frac{-\nabla f(\mathbf{x}_k)}{L} \right\rangle + \frac{L}{2} \left\| \frac{-\nabla f(\mathbf{x}_k)}{L} \right\|_2^2 && \text{NAG update} \\ &= f(\mathbf{x}_k) - \frac{1}{L} \|\nabla f(\mathbf{x}_k)\|_2^2 + \frac{1}{2L} \|\nabla f(\mathbf{x}_k)\|_2^2 \\ &= f(\mathbf{x}_k) - \frac{1}{2L} \|\nabla f(\mathbf{x}_k)\|_2^2 \end{aligned}$$

Now for shorthand notation we will let $g := \frac{1}{2L} \|\nabla f(\mathbf{x}_k)\|_2^2$, we have $f(\mathbf{y}_{k+1}) \leq f(\mathbf{x}_k) - g$.

Lemma 3 ... 2/7

| | |
|---|------------------------------------|
| $f(\mathbf{y}_k) \geq f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{y}_k - \mathbf{x}_k \rangle$ $f(\mathbf{y}_k) \leq \Phi_k^*$ | f convex Induction Hypothesis |
|---|------------------------------------|

► From $f(\mathbf{y}_{k+1}) \leq f(\mathbf{x}_k) - g$, two tricky steps to create $(1 - \frac{1}{\sqrt{Q}})$

$$\begin{aligned}
 f(\mathbf{y}_{k+1}) &\leq f(\mathbf{x}_k) - \frac{f(\mathbf{x}_k)}{\sqrt{Q}} + \frac{f(\mathbf{x}_k)}{\sqrt{Q}} - g \\
 &= \left(1 - \frac{1}{\sqrt{Q}}\right) f(\mathbf{x}_k) + \frac{f(\mathbf{x}_k)}{\sqrt{Q}} - g \\
 &= \left(1 - \frac{1}{\sqrt{Q}}\right) f(\mathbf{x}_k) - \left(1 - \frac{1}{\sqrt{Q}}\right) f(\mathbf{y}_k) + \left(1 - \frac{1}{\sqrt{Q}}\right) f(\mathbf{y}_k) + \frac{f(\mathbf{x}_k)}{\sqrt{Q}} - g \\
 &= \left(1 - \frac{1}{\sqrt{Q}}\right) (f(\mathbf{x}_k) - f(\mathbf{y}_k)) + \left(1 - \frac{1}{\sqrt{Q}}\right) f(\mathbf{y}_k) + \frac{f(\mathbf{x}_k)}{\sqrt{Q}} - g \\
 &\leq \left(1 - \frac{1}{\sqrt{Q}}\right) (f(\mathbf{x}_k) - f(\mathbf{y}_k)) + \left(1 - \frac{1}{\sqrt{Q}}\right) \Phi_k^* + \frac{1}{\sqrt{Q}} f(\mathbf{x}_k) - g && \text{induction hypothesis} \\
 &\leq \left(1 - \frac{1}{\sqrt{Q}}\right) \langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{y}_k \rangle + \left(1 - \frac{1}{\sqrt{Q}}\right) \Phi_k^* + \frac{1}{\sqrt{Q}} f(\mathbf{x}_k) - g && f \text{ convex} \\
 f(\mathbf{y}_{k+1}) &\leq \left(1 - \frac{1}{\sqrt{Q}}\right) \langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{y}_k \rangle + \left(1 - \frac{1}{\sqrt{Q}}\right) \Phi_k^* + \frac{1}{\sqrt{Q}} f(\mathbf{x}_k) - g && \text{(now we have)}
 \end{aligned}$$

► Recall our goal is to show $f(\mathbf{y}_{k+1}) \leq \Phi_{k+1}^*$, we can try to show

$$\left(1 - \frac{1}{\sqrt{Q}}\right) \langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{y}_k \rangle + \left(1 - \frac{1}{\sqrt{Q}}\right) \Phi_k^* + \frac{1}{\sqrt{Q}} f(\mathbf{x}_k) - g \leq \Phi_{k+1}^* \quad \text{(what we want to prove)}$$

Lemma 3 ... 3/7

► Now consider $\Phi_k(\mathbf{x})$. Lemma 2 $\nabla^2 \Phi_k(\mathbf{x}) = m\mathbf{I}_n$ implies $\Phi_k(\mathbf{x}) = \Phi_k^* + \frac{m}{2} \|\mathbf{x} - \boldsymbol{\nu}_k\|_2^2$ for some $\boldsymbol{\nu}_k \in \mathbb{R}^n$ implies

1. $\nabla \Phi_k(\mathbf{x}) = m(\mathbf{x} - \boldsymbol{\nu}_k)$
2. Φ_k is minimized at $\boldsymbol{\nu}_k$, which implies $\nabla \Phi_k(\boldsymbol{\nu}_k) = \mathbf{0}$
3. Points 1,2 work for all k , including $k+1$
4. From $\Phi_0(\mathbf{x}) = f(\mathbf{x}_0) + \frac{m}{2} \|\mathbf{x} - \mathbf{x}_0\|_2^2$, $\boldsymbol{\nu}_0 = \mathbf{x}_0$

► By definition of $\Phi_{k+1}(\mathbf{x})$ in Nesterov's estimate sequence

$$\begin{aligned}\Phi_{k+1}(\mathbf{x}) &= \left(1 - \frac{1}{\sqrt{Q}}\right) \Phi_k(\mathbf{x}) + \frac{1}{\sqrt{Q}} \left(f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{m}{2} \|\mathbf{x} - \mathbf{x}_k\|_2^2 \right) \\ \nabla \Phi_{k+1}(\mathbf{x}) &= \left(1 - \frac{1}{\sqrt{Q}}\right) \nabla \Phi_k(\mathbf{x}) + \frac{1}{\sqrt{Q}} \left(\nabla f(\mathbf{x}_k) + m(\mathbf{x} - \mathbf{x}_k) \right) \\ &= \left(1 - \frac{1}{\sqrt{Q}}\right) m(\mathbf{x} - \boldsymbol{\nu}_k) + \frac{1}{\sqrt{Q}} \left(\nabla f(\mathbf{x}_k) + m(\mathbf{x} - \mathbf{x}_k) \right) \\ \nabla \Phi_{k+1}(\boldsymbol{\nu}_{k+1}) &= \left(1 - \frac{1}{\sqrt{Q}}\right) m(\boldsymbol{\nu}_{k+1} - \boldsymbol{\nu}_k) + \frac{1}{\sqrt{Q}} \left(\nabla f(\mathbf{x}_k) + m(\boldsymbol{\nu}_{k+1} - \mathbf{x}_k) \right) \\ &= \mathbf{0} \qquad \qquad \qquad (2) \ \& \ (3) \ \text{gives} \ \nabla \Phi_{k+1}(\boldsymbol{\nu}_{k+1}) = \mathbf{0}\end{aligned}$$

Lemma 3 ... 4/7 (just some algebra)

$$\left(1 - \frac{1}{\sqrt{Q}}\right)m(\boldsymbol{\nu}_{k+1} - \boldsymbol{\nu}_k) + \frac{1}{\sqrt{Q}}\left(\nabla f(\mathbf{x}_k) + m(\boldsymbol{\nu}_{k+1} - \mathbf{x}_k)\right) = \mathbf{0}$$

$$\begin{aligned} & \left(1 - \frac{1}{\sqrt{Q}}\right)(\boldsymbol{\nu}_{k+1} - \boldsymbol{\nu}_k) + \frac{1}{\sqrt{Q}}\left(\frac{\nabla f(\mathbf{x}_k)}{m} + (\boldsymbol{\nu}_{k+1} - \mathbf{x}_k)\right) = \mathbf{0} \\ \Leftrightarrow & \left(1 - \frac{1}{\sqrt{Q}}\right)\boldsymbol{\nu}_{k+1} - \left(1 - \frac{1}{\sqrt{Q}}\right)\boldsymbol{\nu}_k + \frac{1}{\sqrt{Q}}\boldsymbol{\nu}_{k+1} + \frac{1}{\sqrt{Q}}\left(\frac{\nabla f(\mathbf{x}_k)}{m} - \mathbf{x}_k\right) = \mathbf{0} \end{aligned}$$

Now

$$\boldsymbol{\nu}_{k+1} = \left(1 - \frac{1}{\sqrt{Q}}\right)\boldsymbol{\nu}_k + \frac{1}{\sqrt{Q}}\left(\mathbf{x}_k - \frac{\nabla f(\mathbf{x}_k)}{m}\right) \quad (5)$$

$$\Leftrightarrow -\boldsymbol{\nu}_{k+1} = -\left(1 - \frac{1}{\sqrt{Q}}\right)\boldsymbol{\nu}_k - \frac{1}{\sqrt{Q}}\left(\mathbf{x}_k - \frac{\nabla f(\mathbf{x}_k)}{m}\right)$$

$$\Leftrightarrow \mathbf{x}_k - \boldsymbol{\nu}_{k+1} = \mathbf{x}_k - \left(1 - \frac{1}{\sqrt{Q}}\right)\boldsymbol{\nu}_k - \frac{1}{\sqrt{Q}}\mathbf{x}_k + \frac{1}{\sqrt{Q}}\frac{\nabla f(\mathbf{x}_k)}{m}$$

$$= \left(1 - \frac{1}{\sqrt{Q}}\right)(\mathbf{x}_k - \boldsymbol{\nu}_k) + \frac{\nabla f(\mathbf{x}_k)}{m\sqrt{Q}}$$

$$\Leftrightarrow \|\mathbf{x}_k - \boldsymbol{\nu}_{k+1}\|_2^2 = \left(1 - \frac{1}{\sqrt{Q}}\right)^2 \|\mathbf{x}_k - \boldsymbol{\nu}_k\|_2^2 + 2\left(1 - \frac{1}{\sqrt{Q}}\right)\frac{\langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \boldsymbol{\nu}_k \rangle}{m\sqrt{Q}} + \frac{\|\nabla f(\mathbf{x}_k)\|_2^2}{m^2Q}$$

Lemma 3 ... 5/7

$$\|\mathbf{x}_k - \boldsymbol{\nu}_{k+1}\|_2^2 = \left(1 - \frac{1}{\sqrt{Q}}\right)^2 \|\mathbf{x}_k - \boldsymbol{\nu}_k\|_2^2 + 2\left(1 - \frac{1}{\sqrt{Q}}\right) \frac{\langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \boldsymbol{\nu}_k \rangle}{m\sqrt{Q}} + \frac{\|\nabla f(\mathbf{x}_k)\|_2^2}{m^2 Q}$$

► Now consider $\Phi_{k+1}(\mathbf{x})$ evaluate at \mathbf{x}_k , from in slide 14 we have

$$\begin{aligned} \Phi_{k+1}(\mathbf{x}_k) &= \Phi_{k+1}^* + \frac{m}{2} \|\mathbf{x}_k - \boldsymbol{\nu}_{k+1}\|_2^2 \\ &= \Phi_{k+1}^* + \frac{m}{2} \left(1 - \frac{1}{\sqrt{Q}}\right)^2 \|\mathbf{x}_k - \boldsymbol{\nu}_k\|_2^2 + \left(1 - \frac{1}{\sqrt{Q}}\right) \frac{\langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \boldsymbol{\nu}_k \rangle}{\sqrt{Q}} + \frac{\|\nabla f(\mathbf{x}_k)\|_2^2}{2mQ} \\ &= \Phi_{k+1}^* + \frac{m}{2} \left(1 - \frac{1}{\sqrt{Q}}\right)^2 \|\mathbf{x}_k - \boldsymbol{\nu}_k\|_2^2 + \left(1 - \frac{1}{\sqrt{Q}}\right) \frac{\langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \boldsymbol{\nu}_k \rangle}{\sqrt{Q}} + g \quad (*) \end{aligned}$$

by using the fact $mQ = L$ and $g = \frac{1}{2L} \|\nabla f(\mathbf{x}_k)\|_2^2$.

► By definition of $\Phi_{k+1}(\mathbf{x})$ from page 5, $\Phi_{k+1}(\mathbf{x}_k)$ is

$$\begin{aligned} \Phi_{k+1}(\mathbf{x}_k) &= \left(1 - \frac{1}{\sqrt{Q}}\right) \Phi_k(\mathbf{x}_k) + \frac{1}{\sqrt{Q}} \left(f(\mathbf{x}_k) + \underbrace{\langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{x}_k \rangle}_{=0} + \frac{m}{2} \underbrace{\|\mathbf{x}_k - \mathbf{x}_k\|_2^2}_{=0} \right) \\ &= \left(1 - \frac{1}{\sqrt{Q}}\right) \Phi_k(\mathbf{x}_k) + \frac{1}{\sqrt{Q}} f(\mathbf{x}_k) \quad (**) \end{aligned}$$

► $(*) = (**)$ gives

$$\Phi_{k+1}^* + \frac{m}{2} \left(1 - \frac{1}{\sqrt{Q}}\right)^2 \|\mathbf{x}_k - \boldsymbol{\nu}_k\|_2^2 + \left(1 - \frac{1}{\sqrt{Q}}\right) \frac{\langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \boldsymbol{\nu}_k \rangle}{\sqrt{Q}} + g = \left(1 - \frac{1}{\sqrt{Q}}\right) \Phi_k(\mathbf{x}_k) + \frac{1}{\sqrt{Q}} f(\mathbf{x}_k)$$

Lemma 3 ... 6/7

$$\Phi_{k+1}^* = -\frac{m}{2} \left(1 - \frac{1}{\sqrt{Q}}\right)^2 \|\mathbf{x}_k - \boldsymbol{\nu}_k\|_2^2 - \left(1 - \frac{1}{\sqrt{Q}}\right) \frac{\langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \boldsymbol{\nu}_k \rangle}{\sqrt{Q}} - g + \left(1 - \frac{1}{\sqrt{Q}}\right) \Phi_k(\mathbf{x}_k) + \frac{1}{\sqrt{Q}} f(\mathbf{x}_k)$$

By $\Phi_k(\mathbf{x}) = \Phi_k^* + \frac{m}{2} \|\mathbf{x} - \boldsymbol{\nu}_k\|_2^2$ (slide 14)

$$\left(1 - \frac{1}{\sqrt{Q}}\right) \Phi_k(\mathbf{x}_k) = \left(1 - \frac{1}{\sqrt{Q}}\right) \Phi_k^* + \left(1 - \frac{1}{\sqrt{Q}}\right) \frac{m}{2} \|\mathbf{x}_k - \boldsymbol{\nu}_k\|_2^2$$

Hence

$$\begin{aligned} \Phi_{k+1}^* = & \underbrace{-\frac{m}{2} \left(1 - \frac{1}{\sqrt{Q}}\right)^2 \|\mathbf{x}_k - \boldsymbol{\nu}_k\|_2^2}_{\text{wavy}} - \left(1 - \frac{1}{\sqrt{Q}}\right) \frac{\langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \boldsymbol{\nu}_k \rangle}{\sqrt{Q}} - g \\ & + \left(1 - \frac{1}{\sqrt{Q}}\right) \Phi_k^* + \underbrace{\left(1 - \frac{1}{\sqrt{Q}}\right) \frac{m}{2} \|\mathbf{x}_k - \boldsymbol{\nu}_k\|_2^2}_{\text{wavy}} + \frac{f(\mathbf{x}_k)}{\sqrt{Q}} \end{aligned}$$

Simplify the term

$$\Phi_{k+1}^* = \underbrace{\frac{m}{2\sqrt{Q}} \left(1 - \frac{1}{\sqrt{Q}}\right)^2 \|\mathbf{x}_k - \boldsymbol{\nu}_k\|_2^2}_{\text{wavy}} - \left(1 - \frac{1}{\sqrt{Q}}\right) \frac{\langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \boldsymbol{\nu}_k \rangle}{\sqrt{Q}} + \left(1 - \frac{1}{\sqrt{Q}}\right) \Phi_k^* - g + \frac{f(\mathbf{x}_k)}{\sqrt{Q}}$$

To proceed, we need lemma 4.

Lemma 4 $\boldsymbol{\nu}_k - \mathbf{x}_k = \sqrt{Q}(\mathbf{x}_k - \mathbf{y}_k)$

| | | |
|--------------------|--|-------------|
| Q | $= \frac{L}{m}$ | def of Q |
| \mathbf{y}_{k+1} | $= \mathbf{x}_k - \frac{\nabla f(\mathbf{x}_k)}{L}$ | NAG def (1) |
| \mathbf{x}_{k+1} | $= \left(1 + \frac{\sqrt{Q}-1}{\sqrt{Q}+1}\right) \mathbf{y}_{k+1} - \frac{\sqrt{Q}-1}{\sqrt{Q}+1} \mathbf{y}_k$ | NAG def (2) |

Proof by induction

- ▶ **Base case** $k = 0$ is true by $\mathbf{x}_0 = \mathbf{y}_0$ hence $\boldsymbol{\nu}_0 = \mathbf{x}_0$.
- ▶ **Induction hypothesis** $\boldsymbol{\nu}_k - \mathbf{x}_k = \sqrt{Q}(\mathbf{x}_k - \mathbf{y}_k)$
- ▶ **Case** $k + 1$

$$\begin{aligned}
 \boldsymbol{\nu}_{k+1} &\stackrel{(5)}{=} \left(1 - \frac{1}{\sqrt{Q}}\right) \boldsymbol{\nu}_k + \frac{1}{\sqrt{Q}} \left(\mathbf{x}_k - \frac{\nabla f(\mathbf{x}_k)}{m}\right) \\
 &= \left(1 - \frac{1}{\sqrt{Q}}\right) \boldsymbol{\nu}_k + \frac{1}{\sqrt{Q}} \left(\mathbf{x}_k - \frac{Q \nabla f(\mathbf{x}_k)}{L}\right) && \text{def of } Q \\
 \underbrace{\boldsymbol{\nu}_{k+1} - \mathbf{x}_{k+1}} &= \left(1 - \frac{1}{\sqrt{Q}}\right) \boldsymbol{\nu}_k + \frac{1}{\sqrt{Q}} \left(\mathbf{x}_k - \frac{Q \nabla f(\mathbf{x}_k)}{L}\right) - \underbrace{\mathbf{x}_{k+1}} \\
 &= \left(1 - \frac{1}{\sqrt{Q}}\right) (\mathbf{x}_k + \sqrt{Q}(\mathbf{x}_k - \mathbf{y}_k)) + \frac{1}{\sqrt{Q}} \mathbf{x}_k - \sqrt{Q} \frac{\nabla f(\mathbf{x}_k)}{L} - \mathbf{x}_{k+1} && \text{induction hypothesis} \\
 &= \sqrt{Q} \left(\mathbf{x}_k - \frac{\nabla f(\mathbf{x}_k)}{L}\right) - (\sqrt{Q} - 1) \mathbf{y}_k - \mathbf{x}_{k+1} \\
 &= \sqrt{Q} \mathbf{y}_{k+1} + (\sqrt{Q} + 1) \mathbf{x}_{k+1} - 2\sqrt{Q} \mathbf{y}_{k+1} - \mathbf{x}_{k+1} && \text{NAG def (1) NAG def (2)} \\
 &= \sqrt{Q}(\mathbf{x}_{k+1} - \mathbf{y}_{k+1}) \quad \square
 \end{aligned}$$

Lemma 3 ... 7/7

Lemma 4 $\nu_k - \mathbf{x}_k = \sqrt{Q}(\mathbf{x}_k - \mathbf{y}_k)$

The proof of Lemma 3 stops at

$$\Phi_{k+1}^* = \frac{m}{2\sqrt{Q}} \left(1 - \frac{1}{\sqrt{Q}}\right)^2 \|\mathbf{x}_k - \nu_k\|_2^2 - \left(1 - \frac{1}{\sqrt{Q}}\right) \frac{\langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \nu_k \rangle}{\sqrt{Q}} + \left(1 - \frac{1}{\sqrt{Q}}\right) \Phi_k^* - g + \frac{f(\mathbf{x}_k)}{\sqrt{Q}}$$

By lemma 4 we have

$$\begin{aligned} \Phi_{k+1}^* &= \frac{m}{2\sqrt{Q}} \left(1 - \frac{1}{\sqrt{Q}}\right)^2 \|\mathbf{x}_k - \nu_k\|_2^2 - \left(1 - \frac{1}{\sqrt{Q}}\right) \frac{\langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \nu_k \rangle}{\sqrt{Q}} + \left(1 - \frac{1}{\sqrt{Q}}\right) \Phi_k^* - g + \frac{f(\mathbf{x}_k)}{\sqrt{Q}} \\ &= \frac{m\sqrt{Q}}{2} \left(1 - \frac{1}{\sqrt{Q}}\right)^2 \|\mathbf{x}_k - \mathbf{y}_k\|_2^2 + \left(1 - \frac{1}{\sqrt{Q}}\right) \langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{y}_k \rangle + \left(1 - \frac{1}{\sqrt{Q}}\right) \Phi_k^* - g + \frac{f(\mathbf{x}_k)}{\sqrt{Q}} \end{aligned}$$

Recall (slide 13)

$$f(\mathbf{y}_{k+1}) \leq \left(1 - \frac{1}{\sqrt{Q}}\right) \langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{y}_k \rangle + \left(1 - \frac{1}{\sqrt{Q}}\right) \Phi_k^* + \frac{1}{\sqrt{Q}} f(\mathbf{x}_k) - g \quad (\text{now we have})$$

By $a = \Phi_{k+1}^* = \underbrace{\quad}_{\geq 0} + \underbrace{\quad}_{\geq 0} \geq \underbrace{\quad}_{\text{now we have}} \geq f(\mathbf{y}_{k+1})$, we have proved for the case $k + 1$ that $f(\mathbf{y}_{k+1}) \leq \Phi_{k+1}^*$.

' By induction, Lemma 3 is now proved. \square

Proving NAG convergence rate

$$\text{Lemma 1 } \Phi_{k+1}(\mathbf{x}) \leq f(\mathbf{x}) + \left(1 - \frac{1}{\sqrt{Q}}\right)^k \left(\Phi_0(\mathbf{x}) - f(\mathbf{x})\right) \quad \forall k$$

$$\text{Lemma 3 } f(\mathbf{y}_k) \leq \Phi_k^* \quad \forall k$$

$$f \text{ L-smooth } f(\mathbf{a}) - f(\mathbf{b}) \leq \langle \nabla f(\mathbf{b}), \mathbf{a} - \mathbf{b} \rangle + \frac{L}{2} \|\mathbf{a} - \mathbf{b}\|_2^2$$

$$\nabla f(\mathbf{x}^*) = \mathbf{0}$$

► **Theorem** $f(\mathbf{y}_k) - f^* \leq \frac{m+L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 e^{\frac{-k}{\sqrt{Q}}}$

► **Proof**

$$\begin{aligned}
 f(\mathbf{y}_k) - f^* &\leq \Phi_k(\mathbf{x}^*) - f^* && \text{lemma 3} \\
 &\leq f(\mathbf{x}^*) + \left(1 - \frac{1}{\sqrt{Q}}\right)^k \left(\Phi_0(\mathbf{x}^*) - f(\mathbf{x}^*)\right) - f^* && \text{lemma 1} \\
 &= \left(\Phi_0(\mathbf{x}^*) - f^*\right) \left(1 - \frac{1}{\sqrt{Q}}\right)^k && f(\mathbf{x}^*) = f^* \\
 &= \left(f(\mathbf{x}_0) - f^* + \frac{m}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2\right) \left(1 - \frac{1}{\sqrt{Q}}\right)^k && \text{Def. of } \Phi_0(\mathbf{x}) \\
 &\leq \left(\langle \nabla f(\mathbf{x}^*), \mathbf{x}_0 - \mathbf{x}^* \rangle + \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 + \frac{m}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2\right) \left(1 - \frac{1}{\sqrt{Q}}\right)^k && f \text{ L-smooth} \\
 &\leq \frac{m+L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 \left(1 + \left(-\frac{1}{\sqrt{Q}}\right)\right)^k && \nabla f(\mathbf{x}^*) = \mathbf{0} \\
 &\leq \frac{m+L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 \left(\exp\left(-\frac{1}{\sqrt{Q}}\right)\right)^k && 1 + x \leq e^x \\
 &= \frac{m+L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 \exp\left(\frac{-k}{\sqrt{Q}}\right)
 \end{aligned}$$

Discussion

- ▶ If we stop the algorithm when ϵ -accuracy is achieved

$$\frac{m+L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 \exp\left(\frac{-k}{\sqrt{Q}}\right) \leq \epsilon.$$

Re-arrange

$$k \geq \sqrt{Q} \ln \frac{1}{\epsilon} + \text{constant}.$$

I.e. it takes $\mathcal{O}\left(\sqrt{Q} \ln \frac{1}{\epsilon}\right)$ steps for NAG to converges.

- ▶ Compared to GD with rate $\mathcal{O}\left(Q \ln \frac{1}{\epsilon}\right)$, the improvement $Q \rightarrow \sqrt{Q}$ is significant as m can be viewed as regularization parameter in various machine learning model (norm regularized) and $\frac{1}{m}$ can be as large as sample size. Here the number of step reduced from sample size to $\sqrt{\text{sample size}}$.

Last page - summary

- For unconstrained smooth strongly-convex problem $\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$, with $f : \mathbb{R}^n \rightarrow \mathbb{R}$ being L -smooth and m -strongly convex, the NAG algorithm iterates the following :

$$\mathbf{y}_{k+1} = \mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k), \quad \mathbf{x}_{k+1} = \left(1 - \frac{\sqrt{Q} - 1}{\sqrt{Q} + 1}\right) \mathbf{y}_{k+1} + \frac{\sqrt{Q} - 1}{\sqrt{Q} + 1} \mathbf{y}_k, \quad Q = \frac{L}{m}$$

with initial point $\mathbf{x}_0 = \mathbf{y}_0 \in \mathbb{R}^n$, will produce a sequences $\{f(\mathbf{y}_k)\}_{k \in \mathbb{N}}$ that

$$f(\mathbf{y}_k) - f^* \leq \frac{m + L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 \exp\left(\frac{-k}{\sqrt{Q}}\right).$$

End of document