

Nesterov's estimate sequence: 2. General scheme of optimal method

Andersen Ang

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

Version: August 8, 2023
First draft: July 31, 2023

Content

- 1 Initialize $\mathbf{x}_0 \in \mathbb{R}^n$, $\gamma_0 > 0$, $\mathbf{v}_0 = \mathbf{x}_0$
- 2 **for** $k = 1, 2, \dots$ **do**
- 3 Compute $\alpha_k \in]0, 1[$ [from $L\alpha_k^2 = (1 - \alpha_k)\gamma_k + \alpha_k\mu$
- 4 $\gamma_{k+1} = (1 - \alpha_k)\gamma_k + \alpha_k\mu$
- 5 $\mathbf{y}_k = \frac{\alpha_k\gamma_k\mathbf{v}_k + \gamma_{k+1}\mathbf{x}_k}{\gamma_k + \alpha_k\mu}$
- 6 Find \mathbf{x}_{k+1} s.t. $f(\mathbf{x}_{k+1}) \leq f(\mathbf{y}_k) - \frac{1}{2L}\|\nabla f(\mathbf{y}_k)\|_2^2$. E.g., a gradient step
- 7 $\mathbf{v}_k = \frac{(1 - \alpha_k)\gamma_k\mathbf{v}_k + \alpha_k\mu\mathbf{y}_k - \alpha_k\nabla f(\mathbf{y}_k)}{\gamma_{k+1}}$

$$f(\mathbf{x}_k) - f^* \leq \lambda_k \left[f(\mathbf{x}_0) - f^* + \frac{\gamma_0}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 \right], \quad \lambda_0 = 1, \quad \lambda_k = \prod_{i=0}^{k-1} (1 - \alpha_i)$$

$$\text{Convergence of } \{\lambda_k\}_{k \in \mathbb{N}}: \quad \lambda_k \leq \frac{\mu}{(\gamma_0 - \mu)} \frac{1}{\left[\sinh \left((k+1) \sqrt{\frac{\mu}{L}} \right) \right]^2}$$

$$f(\mathbf{x}_k) - f^* \leq \frac{f(\mathbf{x}_0) - f^* + \frac{\gamma_0}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{4(\gamma_0 - \mu)} \frac{L}{(k+1)^2}$$

Reference Yurii Nesterov, Lectures on convex optimization. Springer, 2018.

Problem setup: unconstrained convex smooth optimization

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

► $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is μ -strongly convex and L -smooth.

► f is convex

► f is μ -strongly convex, $\mu \geq 0$

► The assumption subsume the case for f is convex ($\mu = 0$)

► If a function ψ is convex, we have $\psi(\mathbf{a}) \geq \psi(\mathbf{b}) + \langle \nabla \psi(\mathbf{b}), \mathbf{a} - \mathbf{b} \rangle$

► f is continuous

► f is continuously differentiable

► ∇f is globally L -Lipschitz, $L > 0$

For the details of convexity, epigraph, smoothness, see [here](#).

► We also assume a solution $\mathbf{x}^* \in \mathcal{X}^*$ exists.

► $\mathcal{X}^* := \underset{\mathbf{x}}{\operatorname{argmin}} f(\mathbf{x})$

► $\mathbf{x}^* \in \mathcal{X}^*$

► $f^* := f(\mathbf{x}^*)$

$$f \in \mathcal{C}_L^{1,1}$$

$\operatorname{dom} f$ is a convex set and $\operatorname{epi} f$ is a convex set

$$f - \frac{\mu}{2} \|\mathbf{x}\|_2^2 \text{ is convex}$$

no jump

$\nabla f(\mathbf{x})$ exists for all $\mathbf{x} \in \operatorname{dom} f$

$$(\forall \mathbf{x} \forall \mathbf{y} \neq \mathbf{x}) \left(\frac{\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|}{\|\mathbf{x} - \mathbf{y}\|} \leq L \right)$$

solution set, assumed nonempty

minimizer

optimal function value

Things from part 1



► [Definition 2.2.1 (“what is” estimate sequence)] A sequences pair $\{\phi_k(\mathbf{x}), \lambda_k\}_{k=0}^{\infty}$ is estimate sequence of $f(\cdot)$ if

(Def0)	λ_k	\geq	0	(∀k)	$\{\lambda_k\}_{k \in \mathbb{N}}$ is nonnegative $\{\lambda_k\}_{k \in \mathbb{N}}$ converges to 0 $\{\phi_k\}_{k \in \mathbb{N}} \leq$ “convex combination” of f, ϕ_0
(Def1)	λ_k	$\xrightarrow{k \rightarrow \infty}$	0	(∀k)	
(Def2)	$\phi_k(\mathbf{x})$	\leq	$(1 - \lambda_k)f(\mathbf{x}) + \lambda_k \phi_0(\mathbf{x})$	(∀k)(∀ $\mathbf{x} \in \mathbb{R}^n$)	

► [Lemma 2.2.1 (“why of” estimate sequence)] Assume \mathbf{x}^* exists. For a sequence $\{\mathbf{x}_k\}_{k \in \mathbb{N}}$:

IF $f(\mathbf{x}_k) \leq \phi_k^* := \min_{\mathbf{x} \in \mathbb{R}^n} \phi_k(\mathbf{x})$ THEN $f(\mathbf{x}_k) - f^* \leq \lambda_k(\phi_0(\mathbf{x}^*) - f^*) \xrightarrow{\text{Def1}} 0$.

► [Lemma 2.2.2 (“how to” estimate sequence)]

A1	f L -smooth μ -strongly cvx	THEN	$\{\phi_k(\mathbf{x}), \lambda_k\}_{k=0}^{\infty}$ defined by A6 A7 is an estimate sequence of f
A2	$\phi_0(\cdot)$ a cvx function		
A3	$\{\mathbf{y}_k\}_{k=0}^{\infty}$ is a sequence		
IF A4a	$\alpha_k \in]0, 1[\quad \forall k$	A6	$\lambda_{k+1} = (1 - \alpha_k)\lambda_k$
A4b	$\sum_{k=0}^{\infty} \alpha_k = \infty$	A7	$\phi_{k+1}(\mathbf{x}) = (1 - \alpha_k)\phi_k(\mathbf{x}) + \alpha_k(f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x} - \mathbf{y}_k \rangle + \frac{\mu}{2} \ \mathbf{x} - \mathbf{y}_k\ _2^2)$
A5	$\lambda_0 := 1$		

► [Lemma 2.2.3 (a quadratic ϕ_0)] IF $\phi_0(\mathbf{x}) := \phi_0^* + \frac{\gamma_0}{2} \|\mathbf{x} - \mathbf{v}_0\|_2^2$ THEN $\{\phi_k(\mathbf{x})\}_{k=0}^{\infty}$ defined as A7 in Lemma 2.2.2 preserves the canonical form of $\{\phi_k(\mathbf{x})\}_{k \in \mathbb{N}}$

$$\phi_k(\mathbf{x}) = \phi_k^* + \frac{\gamma_0}{2} \|\mathbf{x} - \mathbf{v}_k\|_2^2.$$

where $\gamma_{k+1} = (1 - \alpha_k)\gamma_k + \alpha_k\mu$ (i)

$$\mathbf{v}_{k+1} = \frac{(1 - \alpha_k)\gamma_k \mathbf{v}_k + \alpha_k\mu \mathbf{y}_k - \alpha_k \nabla f(\mathbf{y}_k)}{\gamma_{k+1}}$$
 (ii)

$$\phi_{k+1}^* = (1 - \alpha_k)\phi_k^* + \alpha_k f(\mathbf{y}_k) - \frac{\alpha_k^2}{2\gamma_{k+1}} \|\nabla f(\mathbf{y}_k)\|_2^2 + \frac{\alpha_k(1 - \alpha_k)\gamma_k}{\gamma_{k+1}} \left(\frac{\mu}{2} \|\mathbf{y}_k - \mathbf{v}_k\|_2^2 + \langle \nabla f(\mathbf{y}_k), \mathbf{v}_k - \mathbf{y}_k \rangle \right)$$
 (iii)

Things to focus

► [Definition 2.2.1 (“what is” estimate sequence)] A sequences pair $\{\phi_k(\mathbf{x}), \lambda_k\}_{k=0}^{\infty}$ is estimate sequence of $f(\cdot)$ if

$$\begin{array}{l} \text{(Def0)} \quad \lambda_k \geq 0 \quad (\forall k) \\ \text{(Def1)} \quad \lambda_k \xrightarrow{k \rightarrow \infty} 0 \quad (\forall k) \\ \text{(Def2)} \quad \phi_k(\mathbf{x}) \leq (1 - \lambda_k)f(\mathbf{x}) + \lambda_k \phi_0(\mathbf{x}) \quad (\forall k)(\forall \mathbf{x} \in \mathbb{R}^n) \end{array} \left| \begin{array}{l} \{\lambda_k\}_{k \in \mathbb{N}} \text{ is nonnegative} \\ \{\lambda_k\}_{k \in \mathbb{N}} \text{ converges to } 0 \\ \{\phi_k\}_{k \in \mathbb{N}} \leq \text{“convex combination” of } f, \phi_0 \end{array} \right.$$

► [Lemma 2.2.1 (“why of” estimate sequence)] Assume \mathbf{x}^* exists. For a sequence $\{\mathbf{x}_k\}_{k \in \mathbb{N}}$:

$$\text{IF } f(\mathbf{x}_k) \leq \phi_k^* := \min_{\mathbf{x} \in \mathbb{R}^n} \phi_k(\mathbf{x}) \quad \text{THEN } f(\mathbf{x}_k) - f^* \leq \lambda_k (\phi_0(\mathbf{x}^*) - f^*) \xrightarrow{\text{Def1}} 0.$$

► [Lemma 2.2.2 (“how to” estimate sequence)]

$$\begin{array}{ll} \text{A1} & f \text{ } L\text{-smooth } \mu\text{-strongly cvx} \\ \text{A2} & \phi_0(\cdot) \text{ a cvx function} \\ \text{A3} & \{\mathbf{y}_k\}_{k=0}^{\infty} \text{ is a sequence} \\ \text{IF A4a} & \alpha_k \in]0, 1[\quad \forall k \\ \text{A4b} & \sum_{k=0}^{\infty} \alpha_k = \infty \\ \text{A5} & \lambda_0 := 1 \end{array} \quad \text{THEN } \{\phi_k(\mathbf{x}), \lambda_k\}_{k=0}^{\infty} \text{ defined by A6 A7 is an estimate sequence of } f$$

$$\begin{array}{ll} \text{A6} & \lambda_{k+1} = (1 - \alpha_k)\lambda_k \\ \text{A7} & \phi_{k+1}(\mathbf{x}) = (1 - \alpha_k)\phi_k(\mathbf{x}) + \alpha_k \left(f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x} - \mathbf{y}_k \rangle + \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}_k\|_2^2 \right) \end{array}$$

► [Lemma 2.2.3 (a quadratic ϕ_0)] IF $\phi_0(\mathbf{x}) := \phi_0^* + \frac{\gamma_0}{2} \|\mathbf{x} - \mathbf{v}_0\|_2^2$ THEN $\{\phi_k(\mathbf{x})\}_{k=0}^{\infty}$ defined as A7 in Lemma 2.2.2 preserves the canonical form of $\{\phi_k(\mathbf{x})\}_{k \in \mathbb{N}}$

$$\phi_k(\mathbf{x}) = \phi_k^* + \frac{\gamma_0}{2} \|\mathbf{x} - \mathbf{v}_k\|_2^2.$$

$$\text{where } \gamma_{k+1} = (1 - \alpha_k)\gamma_k + \alpha_k \mu \quad (i)$$

$$\mathbf{v}_{k+1} = \frac{(1 - \alpha_k)\gamma_k \mathbf{v}_k + \alpha_k \mu \mathbf{y}_k - \alpha_k \nabla f(\mathbf{y}_k)}{\gamma_{k+1}} \quad (ii)$$

$$\phi_{k+1}^* = (1 - \alpha_k)\phi_k^* + \alpha_k f(\mathbf{y}_k) - \frac{\alpha_k^2}{2\gamma_{k+1}} \|\nabla f(\mathbf{y}_k)\|_2^2 + \frac{\alpha_k(1 - \alpha_k)\gamma_k}{\gamma_{k+1}} \left(\frac{\mu}{2} \|\mathbf{y}_k - \mathbf{v}_k\|_2^2 + \langle \nabla f(\mathbf{y}_k), \mathbf{v}_k - \mathbf{y}_k \rangle \right) \quad (iii)$$

Showing $f(\mathbf{x}_k) \leq \phi_k^* \dots 1/3$

f is L -smooth and μ -strongly convex

(A1)

IF $f(\mathbf{x}_k) \leq \phi_k^* := \min_{\mathbf{x} \in \mathbb{R}^n} \phi_k(\mathbf{x})$ THEN $f(\mathbf{x}_k) - f^* \leq \lambda_k (\phi_0(\mathbf{x}^*) - f^*) \xrightarrow{\text{Def1}} 0$. (Lemma 2.2.1)

$$\phi_{k+1}^* = (1 - \alpha_k) \phi_k^* + \alpha_k f(\mathbf{y}_k) - \frac{\alpha_k^2}{2\gamma_{k+1}} \|\nabla f(\mathbf{y}_k)\|_2^2 + \frac{\alpha_k(1 - \alpha_k)\gamma_k}{\gamma_{k+1}} \left(\frac{\mu}{2} \|\mathbf{y}_k - \mathbf{v}_k\|_2^2 + \langle \nabla f(\mathbf{y}_k), \mathbf{v}_k - \mathbf{y}_k \rangle \right) \quad (iii)$$

► Idea: if ϕ_{k+1}^* defined by (iii) satisfies the IF condition in Lemma 2.2.1 then we have convergence result of $f(\mathbf{x}_k)$ in the THEN statement

► By (iii), assume $\phi_k^* \geq f(\mathbf{x}_k)$, then (iii) becomes

$$\phi_{k+1}^* \geq (1 - \alpha_k) f(\mathbf{x}_k) + \alpha_k f(\mathbf{y}_k) - \frac{\alpha_k^2}{2\gamma_{k+1}} \|\nabla f(\mathbf{y}_k)\|_2^2 + \frac{\alpha_k(1 - \alpha_k)\gamma_k}{\gamma_{k+1}} \left(\frac{\mu}{2} \|\mathbf{y}_k - \mathbf{v}_k\|_2^2 + \langle \nabla f(\mathbf{y}_k), \mathbf{v}_k - \mathbf{y}_k \rangle \right) \quad (iii)$$

► Since $\frac{\alpha_k(1 - \alpha_k)\gamma_k}{\gamma_{k+1}} \frac{\mu}{2} \|\mathbf{y}_k - \mathbf{v}_k\|_2^2 \geq 0$, then from (iii) we have

$$\phi_{k+1}^* \geq (1 - \alpha_k) f(\mathbf{x}_k) + \alpha_k f(\mathbf{y}_k) - \frac{\alpha_k^2}{2\gamma_{k+1}} \|\nabla f(\mathbf{y}_k)\|_2^2 + \frac{\alpha_k(1 - \alpha_k)\gamma_k}{\gamma_{k+1}} \langle \nabla f(\mathbf{y}_k), \mathbf{v}_k - \mathbf{y}_k \rangle$$

► By (A1) f is convex: $f(\mathbf{x}_k) \geq f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x}_k - \mathbf{y}_k \rangle$, so

$$\begin{aligned} \phi_{k+1}^* &\geq (1 - \alpha_k) \left(f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x}_k - \mathbf{y}_k \rangle \right) + \alpha_k f(\mathbf{y}_k) - \frac{\alpha_k^2}{2\gamma_{k+1}} \|\nabla f(\mathbf{y}_k)\|_2^2 + \frac{\alpha_k(1 - \alpha_k)\gamma_k}{\gamma_{k+1}} \langle \nabla f(\mathbf{y}_k), \mathbf{v}_k - \mathbf{y}_k \rangle \\ &= f(\mathbf{y}_k) + (1 - \alpha_k) \langle \nabla f(\mathbf{y}_k), \mathbf{x}_k - \mathbf{y}_k \rangle - \frac{\alpha_k^2}{2\gamma_{k+1}} \|\nabla f(\mathbf{y}_k)\|_2^2 + \frac{\alpha_k(1 - \alpha_k)\gamma_k}{\gamma_{k+1}} \langle \nabla f(\mathbf{y}_k), \mathbf{v}_k - \mathbf{y}_k \rangle \\ &= f(\mathbf{y}_k) - \frac{\alpha_k^2}{2\gamma_{k+1}} \|\nabla f(\mathbf{y}_k)\|_2^2 + (1 - \alpha_k) \langle \nabla f(\mathbf{y}_k), \mathbf{x}_k - \mathbf{y}_k + \frac{\alpha_k\gamma_k}{\gamma_{k+1}} (\mathbf{v}_k - \mathbf{y}_k) \rangle \end{aligned} \quad (†)$$

Proving $f(\mathbf{x}_k) \leq \phi_k^* \dots$ 2/3

f is L -smooth and μ -strongly convex

(A1)

IF $f(\mathbf{x}_k) \leq \phi_k^* := \min_{\mathbf{x} \in \mathbb{R}^n} \phi_k(\mathbf{x})$ THEN $f(\mathbf{x}_k) - f^* \leq \lambda_k (\phi_0(\mathbf{x}^*) - f^*) \xrightarrow{\text{Def1}} 0$. (Lemma 2.2.1)

$$\phi_{k+1}^* \geq f(\mathbf{y}_k) - \frac{\alpha_k^2}{2\gamma_{k+1}} \|\nabla f(\mathbf{y}_k)\|_2^2 + (1 - \alpha_k) \left\langle \nabla f(\mathbf{y}_k), \mathbf{x}_k - \mathbf{y}_k + \frac{\alpha_k \gamma_k}{\gamma_{k+1}} (\mathbf{v}_k - \mathbf{y}_k) \right\rangle \quad (\dagger)$$

$$\gamma_{k+1} = (1 - \alpha_k) \gamma_k + \alpha_k \mu \quad (i)$$

► What we want: make (\dagger) fulfill the IF in Lemma 2.2.1, i.e., $\phi_{k+1}^* \geq f(\mathbf{x}_{k+1})$

► Note that $f(\mathbf{y}_k)$ in (\dagger) is on \mathbf{y} while in the IF we need $f(\mathbf{x})$.

► Hence we need to get \mathbf{x} from \mathbf{y} , a way is gradient descent with constant stepsize

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

(GD) on $\underbrace{\text{smooth function } f}_{\text{A1}}$ guarantees the sufficient descent lemma (proof [here](#))

$$f(\mathbf{y}_k) - \frac{1}{2L} \|\nabla f(\mathbf{y}_k)\|_2^2 \geq f(\mathbf{x}_{k+1}). \quad (\text{SDL})$$

► To make use of (SDL) we need to make $\frac{1}{2L} = \frac{\alpha_k^2}{2\gamma_{k+1}}$. This can be done by defining α_k as a positive root of

$$L\alpha_k^2 = (1 - \alpha_k)\gamma_k + \alpha_k \mu \stackrel{(i)}{=} \gamma_{k+1} \iff L\alpha_k^2 - \gamma_{k+1} = 0 \implies \alpha_k = \frac{-0 \pm \sqrt{0^2 - 4(L)(-\gamma_{k+1})}}{2L} \iff \alpha_k^2 = \frac{\gamma_{k+1}}{L}$$

► Apply (SDL) on (\dagger) gives

$$\phi_{k+1}^* \geq f(\mathbf{x}_{k+1}) + (1 - \alpha_k) \left\langle \nabla f(\mathbf{y}_k), \mathbf{x}_k - \mathbf{y}_k + \frac{\alpha_k \gamma_k}{\gamma_{k+1}} (\mathbf{v}_k - \mathbf{y}_k) \right\rangle \quad (\dagger')$$

Proving $f(\mathbf{x}_k) \leq \phi_k^* \dots$ 3/3

f is L -smooth and μ -strongly convex

(A1)

$\{\mathbf{y}_k\}_{k=0}^\infty$ is a sequence

(A3)

IF $f(\mathbf{x}_k) \leq \phi_k^* := \min_{\mathbf{x} \in \mathbb{R}^n} \phi_k(\mathbf{x})$ **THEN** $f(\mathbf{x}_k) - f^* \leq \lambda_k (\phi_0(\mathbf{x}^*) - f^*) \xrightarrow{\text{Def1}} 0$. (Lemma 2.2.1)

$$\phi_{k+1}^* \geq f(\mathbf{y}_k) - \frac{\alpha_k^2}{2\gamma_{k+1}} \|\nabla f(\mathbf{y}_k)\|_2^2 + (1 - \alpha_k) \left\langle \nabla f(\mathbf{y}_k), \mathbf{x}_k - \mathbf{y}_k + \frac{\alpha_k \gamma_k}{\gamma_{k+1}} (\mathbf{v}_k - \mathbf{y}_k) \right\rangle \quad (\dagger)$$

$$\phi_{k+1}^* \geq f(\mathbf{x}_{k+1}) + (1 - \alpha_k) \left\langle \nabla f(\mathbf{y}_k), \mathbf{x}_k - \mathbf{y}_k + \frac{\alpha_k \gamma_k}{\gamma_{k+1}} (\mathbf{v}_k - \mathbf{y}_k) \right\rangle \quad (\dagger')$$

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

$$f(\mathbf{x}_{k+1}) \geq f(\mathbf{y}_k) - \frac{1}{2L} \|\nabla f(\mathbf{y}_k)\|_2^2 \quad (\text{SDL})$$

$$\gamma_{k+1} = (1 - \alpha_k) \gamma_k + \alpha_k \mu \quad (i)$$

1. We have defined \mathbf{x} from \mathbf{y}

- ▶ (GD) $\xrightarrow{A1}$ (SDL) so we have $(\dagger) \xRightarrow{(\text{SDL})} (\dagger')$
- ▶ (GD) is a way to define \mathbf{x} from \mathbf{y}

2. We still haven't specify how we get \mathbf{y} from \mathbf{x}

- ▶ By (A3) \mathbf{y} is free

3. Our goal is to make (\dagger') satisfy \square , i.e., $\phi_{k+1}^* \geq f(\mathbf{x}_{k+1})$

- ▶ We can consider $(1 - \alpha_k) \left\langle \nabla f(\mathbf{y}_k), \mathbf{x}_k - \mathbf{y}_k + \frac{\alpha_k \gamma_k}{\gamma_{k+1}} (\mathbf{v}_k - \mathbf{y}_k) \right\rangle \geq 0$
- ▶ The borderline case is $\mathbf{x}_k - \mathbf{y}_k + \frac{\alpha_k \gamma_k}{\gamma_{k+1}} (\mathbf{v}_k - \mathbf{y}_k) = \mathbf{0}$. Now using (i) we have

$$\mathbf{y}_k \stackrel{(i)}{=} \frac{\alpha_k \gamma_k \mathbf{v}_k + \gamma_{k+1} \mathbf{x}_k}{\gamma_k + \alpha_k \mu}$$

Algorithm: Nesterov's general scheme of optimal 1st-order method

1 Initialize $\mathbf{x}_0 \in \mathbb{R}^n$, $\gamma_0 > 0$, $\mathbf{v}_0 = \mathbf{x}_0$

2 **for** $k = 1, 2, \dots$ **do**

3 (a) Compute $\alpha_k \in]0, 1[$ from the equation

$$L\alpha_k^2 = (1 - \alpha_k)\gamma_k + \alpha_k\mu$$

 Update $\gamma_{k+1} = (1 - \alpha_k)\gamma_k + \alpha_k\mu$

4

5 (b) $\mathbf{y}_k = \frac{\alpha_k\gamma_k\mathbf{v}_k + \gamma_{k+1}\mathbf{x}_k}{\gamma_k + \alpha_k\mu}$

6 Compute $f(\mathbf{y}_k)$ and $\nabla f(\mathbf{y}_k)$

7

8 (c) Find \mathbf{x}_{k+1} such that

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{y}_k) - \frac{1}{2L} \|\nabla f(\mathbf{y}_k)\|_2^2$$

 (For example, a gradient descent step)

9

10 (d) $\mathbf{v}_k = \frac{(1 - \alpha_k)\gamma_k\mathbf{v}_k + \alpha_k\mu\mathbf{y}_k - \alpha_k\nabla f(\mathbf{y}_k)}{\gamma_{k+1}}$

Algorithm: Nesterov's general scheme of optimal 1st-order method (compact form)

- 1 Initialize $\mathbf{x}_0 \in \mathbb{R}^n$, $\gamma_0 > 0$, $\mathbf{v}_0 = \mathbf{x}_0$
- 2 **for** $k = 1, 2, \dots$ **do**
- 3 Compute $\alpha_k \in]0, 1[$ from $L\alpha_k^2 = (1 - \alpha_k)\gamma_k + \alpha_k\mu$
- 4 $\gamma_{k+1} = (1 - \alpha_k)\gamma_k + \alpha_k\mu$
- 5 $\mathbf{y}_k = \frac{\alpha_k\gamma_k\mathbf{v}_k + \gamma_{k+1}\mathbf{x}_k}{\gamma_k + \alpha_k\mu}$
- 6 Find \mathbf{x}_{k+1} such that $f(\mathbf{x}_{k+1}) \leq f(\mathbf{y}_k) - \frac{1}{2L}\|\nabla f(\mathbf{y}_k)\|_2^2$. E.g., a gradient descent step
- 7 $\mathbf{v}_k = \frac{(1 - \alpha_k)\gamma_k\mathbf{v}_k + \alpha_k\mu\mathbf{y}_k - \alpha_k\nabla f(\mathbf{y}_k)}{\gamma_{k+1}}$

► **Theorem** For the sequence $\{\mathbf{x}_k\}_{k \in \mathbb{N}}$ generated by optimal 1st-order method, we have

$$f(\mathbf{x}_k) - f^* \leq \lambda_k \left[f(\mathbf{x}_0) - f^* + \frac{\gamma_0}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 \right],$$

where $\lambda_0 = 1$ and $\lambda_k = \prod_{i=0}^{k-1} (1 - \alpha_i)$.

► **Proof**

- Let $\phi_0(\mathbf{x}) = f(\mathbf{x}_0) + \frac{\gamma_0}{2} \|\mathbf{x} - \mathbf{v}_0\|_2^2$, then $f(\mathbf{x}_0) = \phi_0^*$ satisfies $f(\mathbf{x}_k) \leq \phi_k^*$ at $k = 0$
- The scheme satisfies $f(\mathbf{x}_k) \leq \phi_k^*$ at all k (this is exactly what we spent on 3 slides to show)
- Then

$$\begin{aligned} f(\mathbf{x}_k) - f^* &\leq \lambda_k \left[\phi_0(\mathbf{x}^*) - f^* \right] && \text{Lemma 2.2.1} \\ &\leq \lambda_k \left[\left(f(\mathbf{x}_0) + \frac{\gamma_0}{2} \|\mathbf{x} - \mathbf{v}_0\|_2^2 \right) \Big|_{\mathbf{x}=\mathbf{x}^*} - f^* \right] && \phi_0(\mathbf{x}) = f(\mathbf{x}_0) + \frac{\gamma_0}{2} \|\mathbf{x} - \mathbf{v}_0\|_2^2 \\ &\leq \lambda_k \left[f(\mathbf{x}_0) - f^* + \frac{\gamma_0}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 \right] && \mathbf{v}_0 = \mathbf{x}_0 \quad \square \end{aligned}$$

Convergence rate of method follows the convergence of $\{\lambda_k\}_{k \in \mathbb{N}}$

Algorithm: Nesterov's general scheme of optimal 1st-order method (compact form)

```

1 Initialize  $\mathbf{x}_0 \in \mathbb{R}^n$ ,  $\gamma_0 > 0$ ,  $\mathbf{v}_0 = \mathbf{x}_0$ 
2 for  $k = 1, 2, \dots$  do
3     Compute  $\alpha_k \in ]0, 1[$  from  $L\alpha_k^2 = (1 - \alpha_k)\gamma_k + \alpha_k\mu$ 
4      $\gamma_{k+1} = (1 - \alpha_k)\gamma_k + \alpha_k\mu$ 
5      $\mathbf{y}_k = \frac{\alpha_k\gamma_k\mathbf{v}_k + \gamma_{k+1}\mathbf{x}_k}{\gamma_k + \alpha_k\mu}$ 
6     Find  $\mathbf{x}_{k+1}$  such that  $f(\mathbf{x}_{k+1}) \leq f(\mathbf{y}_k) - \frac{1}{2L}\|\nabla f(\mathbf{y}_k)\|_2^2$ . E.g., a gradient descent step
7      $\mathbf{v}_k = \frac{(1 - \alpha_k)\gamma_k\mathbf{v}_k + \alpha_k\mu\mathbf{y}_k - \alpha_k\nabla f(\mathbf{y}_k)}{\gamma_{k+1}}$ 

```

► **Theorem** For $\lambda_0 = 1$, $\lambda_k = \prod_{i=0}^{k-1} (1 - \alpha_i)$, we have

$$f(\mathbf{x}_k) - f^* \leq \lambda_k \left[f(\mathbf{x}_0) - f^* + \frac{\gamma_0}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 \right].$$

► The method converge: $f(\mathbf{x}_k) \rightarrow f^*$

because $\lambda_k \rightarrow 0$

► The convergence rate of method follows the convergence rate of $\{\lambda_k\}_{k \in \mathbb{N}}$

► **Theorem (Convergence rate of λ_k)** If $\gamma_0 \in]\mu, 3L + \mu]$, then

$$\lambda_k \leq \frac{\mu}{(\gamma_0 - \mu) \left[\sinh \left((k+1) \sqrt{\frac{\mu}{L}} \right) \right]^2} \quad \forall k \geq 0.$$

Proof of λ_k

γ_{k+1}	$= (1 - \alpha_k)\gamma_k + \alpha_k\mu$	(i)
λ_{k+1}	$= (1 - \alpha_k)\lambda_k$	(A6)
λ_0	$:= 1$	(A7)
$L\alpha_k^2$	$= (1 - \alpha_k)\gamma_k + \alpha_k\mu$	step 1 of algo

► First

$$\begin{aligned}
 \gamma_{k+1} - \mu &= (1 - \alpha_k)\gamma_k + \alpha_k\mu - \mu && (i) \\
 &= (1 - \alpha_k)(\gamma_k - \mu) \\
 &= (1 - \alpha_k)(1 - \alpha_{k-1})(\gamma_{k-1} - \mu) && \text{calling (i) on } k-1 \\
 &= \prod_{i=0}^k (1 - \alpha_i)(\gamma_0 - \mu) && \text{calling (i) multiple times} \\
 &= \lambda_{k+1}(\gamma_0 - \mu) && \text{by (A6, A7)}
 \end{aligned}$$

Hence

$$\gamma_{k+1} - \mu = \lambda_{k+1}(\gamma_0 - \mu) \implies \gamma_{k+1} = \mu + \lambda_{k+1}(\gamma_0 - \mu) \quad (\#)$$

► Next

$$L\alpha_k^2 = (1 - \alpha_k)\gamma_k + \alpha_k\mu \stackrel{(i)}{=} \gamma_{k+1} \iff L\alpha_k^2 - \gamma_{k+1} = 0 \implies \alpha_k = \frac{-0 \pm \sqrt{0^2 - 4(L)(-\gamma_{k+1})}}{2L} = \sqrt{\frac{\gamma_{k+1}}{L}}. \quad (\#\#)$$

Rearrange (A6)

$$1 - \frac{\lambda_{k+1}}{\lambda_k} \stackrel{(A6)}{=} \alpha_k \stackrel{(\#\#)}{=} \sqrt{\frac{\gamma_{k+1}}{L}} \stackrel{(\#)}{=} \sqrt{\frac{\mu + \lambda_{k+1}(\gamma_0 - \mu)}{L}} = \sqrt{\frac{\mu}{L} + \lambda_{k+1} \frac{\gamma_0 - \mu}{L}} \quad (\#\#\#)$$

Several facts on $\{\lambda_k\}_k$

α_k	$\in]0, 1[$	(A4a)
λ_{k+1}	$= (1 - \alpha_k)\lambda_k$	(A6)
λ_0	$:= 1$	(A7)
$1 - \frac{\lambda_{k+1}}{\lambda_k}$	$= \sqrt{\frac{\mu}{L} + \lambda_{k+1} \frac{\gamma_0 - \mu}{L}}$	(###)

► When we get (###) from (A6) by division of λ_k , we need to make sure $\lambda_k \neq 0$

► First $\lambda_{k+1} \stackrel{(A6)}{=} (1 - \alpha_k)\lambda_k \stackrel{(A4a)}{<} \lambda_k$ (Strict Decrease)

► By (A4a), (A6), (A7) that $\lambda_k = \prod_{i=0}^k (1 - \alpha_i) > 0$ so

- $\{\lambda_k\}_{k \in \mathbb{N}}$ is a strictly positive sequence that is upper bounded by 1
- we can take inverse of λ_k
- we can take square-root of λ_k

($\lambda_k > 0$)
(inverse ok)
(square-root ok)

► Let $q = \frac{\mu}{L}$. By (inverse ok) we can multiply (###) by $\frac{1}{\lambda_{k+1}}$ on both sides gives

$$\frac{1}{\lambda_{k+1}} - \frac{1}{\lambda_k} \stackrel{(###)}{=} \frac{1}{\lambda_{k+1}} \sqrt{q + \lambda_{k+1} \frac{\gamma_0 - \mu}{L}} \stackrel{(\text{square-root ok})}{=} \frac{1}{\sqrt{\lambda_{k+1}}} \sqrt{q \frac{1}{\lambda_{k+1}} + \frac{\gamma_0 - \mu}{L}} \quad (\heartsuit)$$

► By (inverse ok) and (square-root ok) we can apply a not-so-trivial decomposition $a - b = (\sqrt{a} + \sqrt{b})(\sqrt{a} - \sqrt{b})$

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

► By (inverse ok)

$$\lambda_{k+1} \stackrel{(\text{Strict Decrease})}{<} \lambda_k \quad \text{AND} \quad \{\lambda_k\}_k \stackrel{\lambda_k > 0}{>} 0 \quad \stackrel{(\text{inverse ok})}{\implies} \frac{1}{\lambda_k} < \frac{1}{\lambda_{k+1}} \quad \stackrel{\lambda > 0}{\implies}$$

$$\frac{1}{\sqrt{\lambda_k}} < \frac{1}{\sqrt{\lambda_{k+1}}}$$

Tricky step

$$\frac{1}{\lambda_{k+1}} - \frac{1}{\lambda_k} = \frac{1}{\sqrt{\lambda_{k+1}}} \sqrt{q \frac{1}{\lambda_{k+1}} + \frac{\gamma_0 - \mu}{L}} \quad (\heartsuit)$$

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

$$\frac{1}{\sqrt{\lambda_k}} < \frac{1}{\sqrt{\lambda_{k+1}}} \quad \blacksquare$$

$$\gamma_0 \in]\mu, 3L + \mu], \quad L > 0$$

- Now combine \blacksquare , (\heartsuit) and $(\heartsuit\heartsuit)$

$$\frac{1}{\sqrt{\lambda_{k+1}}} \sqrt{q \frac{1}{\lambda_{k+1}} + \frac{\gamma_0 - \mu}{L}} \stackrel{(\heartsuit)}{=} \frac{1}{\lambda_{k+1}} - \frac{1}{\lambda_k} \stackrel{(\heartsuit\heartsuit)}{=} \left(\frac{1}{\sqrt{\lambda_{k+1}}} + \frac{1}{\sqrt{\lambda_k}} \right) \left(\frac{1}{\sqrt{\lambda_{k+1}}} - \frac{1}{\sqrt{\lambda_k}} \right) \blacksquare < \frac{2}{\sqrt{\lambda_{k+1}}} \left(\frac{1}{\sqrt{\lambda_{k+1}}} - \frac{1}{\sqrt{\lambda_k}} \right).$$

Now we see that why are we doing these: we can cancel out $\frac{1}{\sqrt{\lambda_{k+1}}}$ to simplify

$$\sqrt{q \frac{1}{\lambda_{k+1}} + \frac{\gamma_0 - \mu}{L}} < 2 \left(\frac{1}{\sqrt{\lambda_{k+1}}} - \frac{1}{\sqrt{\lambda_k}} \right) \iff \frac{1}{\sqrt{\lambda_{k+1}}} - \frac{1}{\sqrt{\lambda_k}} > \frac{1}{2} \sqrt{q \frac{1}{\lambda_{k+1}} + \frac{\gamma_0 - \mu}{L}} \quad (\dagger)$$

- By assumption $\gamma_0 \in]\mu, 3L + \mu]$ we have $\gamma_0 > \mu$ so multiplying $\sqrt{\frac{L}{\gamma_0 - \mu}}$ to (\dagger) do not change the sign

$$\sqrt{\frac{1}{\lambda_{k+1}} \frac{L}{\gamma_0 - \mu}} - \sqrt{\frac{1}{\lambda_k} \frac{L}{\gamma_0 - \mu}} > \frac{1}{2} \sqrt{q \frac{1}{\lambda_{k+1}} \frac{L}{\gamma_0 - \mu} + 1}$$

- Denote $\xi_k := \sqrt{\frac{1}{\lambda_k} \frac{L}{\gamma_0 - \mu}}$, we have

$$\xi_{k+1} - \xi_k > \frac{1}{2} \sqrt{q \xi_{k+1}^2 + 1}$$

$$\begin{aligned}
 & L \geq \mu \geq 0 \\
 \xi_k & := \sqrt{\frac{L}{(\gamma_0 - \mu)\lambda_k}} \\
 \lambda_0 & := 1 \\
 \gamma_0 & \in]\mu, 3L + \mu]
 \end{aligned}$$

(Δ)

► Let $\delta = \frac{1}{2}\sqrt{q} = \frac{1}{2}\sqrt{\frac{\mu}{L}}$, we are going to use mathematical induction to show

$$\xi_k \geq \frac{e^{(k+1)\delta} - e^{-(k+1)\delta}}{4\delta} \quad \forall k \geq 0$$

► **Base case** ($k = 0$): we want to show

$$\sqrt{\frac{L}{\gamma_0 - \mu}} \geq \frac{e^\delta - e^{-\delta}}{4\delta}$$

► First, by $\gamma_0 \in]\mu, 3L + \mu]$, we have

$$0 < \gamma_0 - \mu \leq 3L \xrightarrow[\text{so division ok}]{\text{strict positive}} \frac{1}{\gamma_0 - \mu} \geq \frac{1}{3L} > 0 \xrightarrow[\text{ok to } \sqrt{\cdot}]{\text{strict positive}} \frac{1}{\sqrt{\gamma_0 - \mu}} \geq \frac{1}{\sqrt{3L}} \implies \sqrt{\frac{L}{\gamma_0 - \mu}} \geq \frac{1}{\sqrt{3}}$$

► By $L \geq \mu \geq 0$, we have

$$0 \leq q := \frac{\mu}{L} \leq 1 \implies 0 \leq \sqrt{q} \leq 1 \implies 0 \leq \delta \leq \frac{1}{2}$$

► A tricky step: the function $\frac{e^x - e^{-x}}{4x}$ for $0 \leq x \leq 0.5$ is bounded above by $\frac{e^x - e^{-x}}{4x} \Big|_{x=0.5}$. i.e.,

$$\frac{e^{1/2} - e^{-1/2}}{2} \geq \frac{e^\delta - e^{-\delta}}{4\delta}$$

$$0.5773\dots = \frac{1}{\sqrt{3}} > \frac{e^{1/2} - e^{-1/2}}{2} = 0.521095$$

► Lastly ■, ■ ■ \implies ■

$$\psi(t) := \frac{e^{(t+1)\delta} - e^{-(t+1)\delta}}{4\delta}$$

$$\psi'(t) = \frac{e^{(t+1)\delta} + e^{-(t+1)\delta}}{4}$$

$$\psi'(t) \rightarrow \infty \text{ so } \psi \text{ is a convex function}$$

$$\psi(a) \geq \psi(b) + \psi'(b)(a - b)$$

$$\psi'(t+1) := \frac{e^{(t+2)\delta} + e^{-(t+2)\delta}}{4}$$

$$\psi(t) \geq \psi(t+1) + \frac{e^{(t+2)\delta} + e^{-(t+2)\delta}}{4}$$

$$= \psi(t+1) + \frac{1}{2} \sqrt{\frac{\left(e^{(t+2)\delta} + e^{-(t+2)\delta}\right)^2}{4}}$$

$$= \psi(t+1) + \frac{1}{2} \sqrt{\frac{e^{2(t+2)\delta} + 2 + e^{-2(t+2)\delta}}{4}}$$

$$= \psi(t+1) + \frac{1}{2} \sqrt{\frac{e^{2(t+2)\delta} - 2 + e^{-2(t+2)\delta}}{4}} + \frac{4}{4}$$

$$= \psi(t+1) + \frac{1}{2} \sqrt{4\delta^2 \left(\frac{e^{(t+2)\delta} - e^{-(t+2)\delta}}{4\delta}\right)^2} + 1$$

$$> \psi(t+1) - \frac{1}{2} \sqrt{4\delta^2 \left(\frac{e^{(t+2)\delta} - e^{-(t+2)\delta}}{4\delta}\right)^2} + 1$$

thing we want to prove, with integer k generalized to real t

the derivative of ψ , note that δ in the denominator is gone

ψ is convex

$\psi(a) \geq \psi(b) + \psi'(b)(a - b)$ with $a = t, b = t + 1$

a WTF step

another WTF step: everything inside square-root is positive

because everything inside square-root is positive

After so many efforts, we now have

$$\psi(t) \geq \psi(t+1) - \frac{1}{2} \sqrt{4\delta^2 \left(\frac{e^{(t+2)\delta} - e^{-(t+2)\delta}}{4\delta}\right)^2} + 1$$

(WTF)

The last step: proving case $k + 1$ by contradiction

- Now we have

$$\begin{aligned} \psi(t) &:= \frac{e^{(t+1)\delta} - e^{-(t+1)\delta}}{4\delta} && \text{integer } k \text{ generalized to real } t \\ \xi_{k+1} - \xi_k &> \frac{1}{2} \sqrt{q\xi_{k+1}^2 + 1} && (\square) \\ \psi(t) &\geq \psi(t+1) - \frac{1}{2} \sqrt{4\delta^2 \left(\frac{e^{(t+2)\delta} - e^{-(t+2)\delta}}{4\delta} \right)^2 + 1} && (WTF) \\ \xi_k &\geq \frac{e^{(k+1)\delta} - e^{-(k+1)\delta}}{4\delta} && \text{induction hypothesis} \end{aligned}$$

- What are going to do: use these to show that at case $k + 1$

$$\xi_{k+1} \geq \frac{e^{(k+2)\delta} - e^{-(k+2)\delta}}{4\delta} \stackrel{t \equiv k}{=} \psi(t+1)$$

by contradiction: (\square) and induction hypothesis will contradict with (WTF) ← this is why we derive (WTF)

- A preparatory material

$$\xi_{k+1} - \frac{1}{2} \sqrt{q\xi_{k+1}^2 + 1} \stackrel{(\square)}{>} \xi_k \stackrel{\text{induction hypothesis}}{\geq} \frac{e^{(k+1)\delta} - e^{-(k+1)\delta}}{4\delta} = \psi(t)$$

which is

$$\psi(t) < \xi_{k+1} - \frac{1}{2} \sqrt{q\xi_{k+1}^2 + 1}$$

- For the purpose of contradiction, assume $\xi_{k+1} < \psi(t+1)$, now we have

$$\psi(t) < \psi(t+1) - \frac{1}{2} \sqrt{q\xi_{k+1}^2 + 1}$$

$$\text{Proof } \xi_{k+1} \geq \frac{e^{(k+2)\delta} - e^{-(k+2)\delta}}{4\delta} \stackrel{t \equiv k}{=} \psi(t+1)$$

$$\begin{aligned} \psi(t) &< \psi(t+1) - \frac{1}{2} \sqrt{q\xi_{k+1}^2 + 1} & (X) \\ \psi(t) &\geq \psi(t+1) - \frac{1}{2} \sqrt{4\delta^2 \left(\frac{e^{(t+2)\delta} - e^{-(t+2)\delta}}{4\delta} \right)^2 + 1} & (WTF) \\ \delta &:= \frac{1}{2} \sqrt{q} \\ \psi(t) &:= \frac{e^{(t+1)\delta} - e^{-(t+1)\delta}}{4\delta} \end{aligned}$$

$$\psi(t+1) - \frac{1}{2} \sqrt{4\delta^2 \left(\frac{e^{(t+2)\delta} - e^{-(t+2)\delta}}{4\delta} \right)^2 + 1} \stackrel{(WTF)}{\leq} \psi(t) \stackrel{(X)}{<} \psi(t+1) - \frac{1}{2} \sqrt{q\xi_{k+1}^2 + 1}$$

$$-\frac{1}{2} \sqrt{4\delta^2 \left(\frac{e^{(t+2)\delta} - e^{-(t+2)\delta}}{4\delta} \right)^2 + 1} < -\frac{1}{2} \sqrt{q\xi_{k+1}^2 + 1}$$

$$\sqrt{4\delta^2 \left(\frac{e^{(t+2)\delta} - e^{-(t+2)\delta}}{4\delta} \right)^2 + 1} > \sqrt{q\xi_{k+1}^2 + 1}$$

$$4\delta^2 \left(\frac{e^{(t+2)\delta} - e^{-(t+2)\delta}}{4\delta} \right)^2 > q\xi_{k+1}^2$$

$$2\delta \left(\frac{e^{(t+2)\delta} - e^{-(t+2)\delta}}{4\delta} \right) > \sqrt{q}\xi_{k+1}$$

$$\frac{e^{(t+2)\delta} - e^{-(t+2)\delta}}{4\delta} > \xi_{k+1}$$

$$\underbrace{\psi(t+1)} > \xi_{k+1}$$

contradicts with $\xi_{k+1} < \psi(t+1)$

What now: we just proved $\xi_k := \sqrt{\frac{L}{(\gamma_0 - \mu)\lambda_k}} \geq \frac{e^{(k+1)\delta} - e^{-(k+1)\delta}}{4\delta} \quad \forall k$

$$\begin{aligned} \lambda_k &\leq \frac{16\delta^2 L}{(\gamma_0 - \mu) \left[e^{(k+1)\delta} - e^{-(k+1)\delta} \right]^2} \quad \forall k && \text{rearrange} \\ &= \frac{4\mu}{(\gamma_0 - \mu) \left[\exp\left((k+1)\sqrt{\frac{\mu}{L}}\right) - \exp\left(- (k+1)\sqrt{\frac{\mu}{L}}\right) \right]^2} \quad \forall k \quad \delta := \frac{1}{2}\sqrt{q}, q := \frac{\mu}{L} \\ &= \frac{\mu}{(\gamma_0 - \mu)} \frac{1}{\left[\sinh\left((k+1)\sqrt{\frac{\mu}{L}}\right) \right]^2} \quad \forall k \end{aligned}$$

Thus

$$\lambda_k \leq \frac{\mu}{(\gamma_0 - \mu)} \frac{1}{\left[\sinh\left((k+1)\sqrt{\frac{\mu}{L}}\right) \right]^2} \quad \forall k \quad \gamma_0 \in]\mu, 3L + \mu]$$

The smallest γ_0 is $(1 + \epsilon)\mu$ and thus

$$\lambda_k \leq \frac{\mu}{\epsilon} \frac{1}{\left[\sinh\left((k+1)\sqrt{\frac{\mu}{L}}\right) \right]^2} \quad \forall k$$

The largest γ_0 is $3L + \mu$ and thus

$$\lambda_k \leq \frac{\mu}{3L} \frac{1}{\left[\sinh\left((k+1)\sqrt{\frac{\mu}{L}}\right) \right]^2} \quad \forall k$$

$$\lambda_k \leq \frac{\mu}{(\gamma_0 - \mu)} \frac{1}{\left[\sinh \left((k+1) \sqrt{\frac{\mu}{L}} \right) \right]^2} = \frac{4\mu}{(\gamma_0 - \mu) \left[\exp \left((k+1) \sqrt{\frac{\mu}{L}} \right) - \exp \left(- (k+1) \sqrt{\frac{\mu}{L}} \right) \right]^2}, \quad \gamma_0 \in]\mu, 3L + \mu]$$

Consider Taylor series of $\exp(\theta)$

$$\exp(\theta) = \frac{\theta^0}{0!} + \frac{\theta^1}{1!} + \frac{\theta^2}{2!} + \frac{\theta^3}{3!} + \dots$$

$$\exp \left((k+1) \sqrt{\frac{\mu}{L}} \right) = 1 + (k+1) \sqrt{\frac{\mu}{L}} + \frac{(k+1)^2 \frac{\mu}{L}}{2} + \frac{(k+1)^3 \left(\frac{\mu}{L} \right)^{\frac{3}{2}}}{6} + \dots$$

$$\exp \left(- (k+1) \sqrt{\frac{\mu}{L}} \right) = 1 - (k+1) \sqrt{\frac{\mu}{L}} + \frac{(k+1)^2 \frac{\mu}{L}}{2} - \frac{(k+1)^3 \left(\frac{\mu}{L} \right)^{\frac{3}{2}}}{6} + \dots$$

$$-\exp \left(- (k+1) \sqrt{\frac{\mu}{L}} \right) = -1 + (k+1) \sqrt{\frac{\mu}{L}} - \frac{(k+1)^2 \frac{\mu}{L}}{2} + \frac{(k+1)^3 \left(\frac{\mu}{L} \right)^{\frac{3}{2}}}{6} + \dots$$

$$\exp \left(- (k+1) \sqrt{\frac{\mu}{L}} \right) - \exp \left(- (k+1) \sqrt{\frac{\mu}{L}} \right) = 2 \left((k+1) \sqrt{\frac{\mu}{L}} + \frac{(k+1)^3 \left(\frac{\mu}{L} \right)^{\frac{3}{2}}}{6} + \frac{(k+1)^5 \left(\frac{\mu}{L} \right)^{\frac{5}{2}}}{120} + \dots \right)$$

$$\left[\exp \left(- (k+1) \sqrt{\frac{\mu}{L}} \right) - \exp \left(- (k+1) \sqrt{\frac{\mu}{L}} \right) \right]^2 = 4 \left[(k+1) \sqrt{\frac{\mu}{L}} + \underbrace{\frac{(k+1)^3 \left(\frac{\mu}{L} \right)^{\frac{3}{2}}}{6} + \frac{(k+1)^5 \left(\frac{\mu}{L} \right)^{\frac{5}{2}}}{120} + \dots}_{>0} \right]^2$$

$$> 4 \left[(k+1) \sqrt{\frac{\mu}{L}} \right]^2$$

$$= 4(k+1)^2 \frac{\mu}{L}$$

$$\lambda_k \leq \frac{\mu}{(\gamma_0 - \mu)} \frac{1}{\left[\sinh \left((k+1) \sqrt{\frac{\mu}{L}} \right) \right]^2} = \frac{4\mu}{(\gamma_0 - \mu) \left[\exp \left((k+1) \sqrt{\frac{\mu}{L}} \right) - \exp \left(- (k+1) \sqrt{\frac{\mu}{L}} \right) \right]^2}, \quad \gamma_0 \in]\mu, 3L + \mu]$$

$$\begin{aligned} \left[\exp \left(- (k+1) \sqrt{\frac{\mu}{L}} \right) - \exp \left(- (k+1) \sqrt{\frac{\mu}{L}} \right) \right]^2 &> 4(k+1)^2 \frac{\mu}{L} \\ \frac{1}{\left[\exp \left(- (k+1) \sqrt{\frac{\mu}{L}} \right) - \exp \left(- (k+1) \sqrt{\frac{\mu}{L}} \right) \right]^2} &< \frac{1}{4(k+1)^2 \frac{\mu}{L}} \end{aligned}$$

Theorem (Convergence rate) IF $\lambda_0 = 1$, $\lambda_k = \prod_{i=0}^{k-1} (1 - \alpha_i)$, $\gamma_0 \in]\mu, 3L + \mu]$, then for the sequence $\{\mathbf{x}_k\}_{k \in \mathbb{N}}$ generated by Nesterov's general scheme, we have that for the sequence $\{f(\mathbf{x}_k)\}_{k \in \mathbb{N}}$

$$f(\mathbf{x}_k) - f^* \leq \frac{L}{4(\gamma_0 - \mu)(k+1)^2} \left[f(\mathbf{x}_0) - f^* + \frac{\gamma_0}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 \right]$$

Proof

$$\begin{aligned} f(\mathbf{x}_k) - f^* &\leq \lambda_k \left[f(\mathbf{x}_0) - f^* + \frac{\gamma_0}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 \right] \\ &\leq \frac{\mu}{(\gamma_0 - \mu) \left[\sinh \left((k+1) \sqrt{\frac{\mu}{L}} \right) \right]^2} \left[f(\mathbf{x}_0) - f^* + \frac{\gamma_0}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 \right] \\ &< \frac{L}{4(\gamma_0 - \mu)(k+1)^2} \left[f(\mathbf{x}_0) - f^* + \frac{\gamma_0}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 \right] \end{aligned}$$