

Nesterov's estimate sequence: 1. What is it and how to construct one

Andersen Ang

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

Homepage angms.science

Version: July 28, 2023

First draft: Nov 21, 2021

Content

Nesterov's estimate sequence: $\{\phi_k(\mathbf{x}), \lambda_k\}_{k=0}^{\infty}$, $\lambda_k \geq 0$ that

$$\lambda_k \xrightarrow{k \rightarrow \infty} 0, \quad \phi_k(\mathbf{x}) \leq (1 - \lambda_k)f(\mathbf{x}) + \lambda_k\phi_0(\mathbf{x})$$

Why estimate sequence: $f(\mathbf{x}_k) - f^* \leq \lambda_k(\phi_0(\mathbf{x}^*) - f^*) \xrightarrow{k \rightarrow \infty} 0$.

How to construct an estimate sequence for str-cvx smooth f

Reference

Yurii Nesterov, Introductory lectures on convex optimization: a basic course, Kluwer Academic Publishers, 2003.

Yurii Nesterov, Lectures on convex optimization. Vol. 137. Berlin: 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 ($\mu = 0$)
- f is continuous
- f is continuously differentiable
- ∇f is globally L -Lipschitz, $L > 0$

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

$\operatorname{dom} f$ is a convex set and $\operatorname{epi} f$ is a convex set
 $f - \frac{\mu}{2} \|\mathbf{x}\|_2^2$ is convex

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

no jump
 $\nabla f(\mathbf{x})$ exists for all $\mathbf{x} \in \operatorname{dom} f$

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

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

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

solution set, assumed nonempty
minimizer
optimal function value

Nesterov's estimate sequence: the definition

- ▶ Also called Nesterov's estimating sequence¹

▶ **Definition 2.2.1** A sequences pair $\{\phi_k(\mathbf{x}), \lambda_k\}_{k=0}^{\infty}$ is estimate sequence of $f(\cdot)$ if

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

- ▶ At this stage

- ▶ We haven't specify what is λ_0
 - ▶ If $\lambda_0 > 1$ then Def2 is not convex combination but linear combination. That's why we put quote "convex combination"
- ▶ We haven't specify how we get λ_k
- ▶ We haven't specify what is ϕ_0
- ▶ We haven't specify what property ϕ_k has

- ▶ At this stage, from Definition 2.2.1, we only know $\{\lambda_k\}_{k \in \mathbb{N}}$ converges to 0. But we don't know how it converges to 0, we also don't know is $\{\lambda_k\}_{k \in \mathbb{N}}$ monotonically converges to 0.

- ▶ For example, the following oscillating sequence fulfills Def0 and Def1

$$\frac{\sin x + 1}{x + 0.1}, x \geq 0 : \{1.6, 0.9, 0.3, 0.05, 0.008, 0.11, 0.23, \dots \text{ for } x = \{1, 2, 3, 4, \dots\}\}$$

¹Nesterov used the term "estimate sequence" in his 2003 book and then used "estimating sequence" in his 2018 book.

Nesterov's estimate sequence: the λ_k

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

$$\begin{array}{llll} \text{(Def0)} & \lambda_k & \geq & 0 & (\forall k) \\ \text{(Def1)} & \lambda_k & \xrightarrow{k \rightarrow \infty} & 0 & (\forall k) \\ \text{(Def2)} & \phi_k(\mathbf{x}) & \leq & (1 - \lambda_k)f(\mathbf{x}) + \lambda_k \phi_0(\mathbf{x}) & (\forall k)(\forall \mathbf{x} \in \mathbb{R}^n) \end{array}$$

► **Lemma 2.2.2 (Partly)** Assume that

$$\text{(L2.2.2 A4a)} \quad \alpha_k \in]0, 1[\quad (\forall k) \quad \alpha_k \text{ strictly positive and strictly smaller than 1}$$

$$\text{(L2.2.2 A4b)} \quad \sum_{k=0}^{\infty} \alpha_k = +\infty \quad \{\alpha_k\} \text{ is not a summable sequence}$$

$$\text{(L2.2.2 A5)} \quad \lambda_0 := 1 \quad \text{we initialize } \lambda_0$$

$$\text{(L2.2.2 A6)} \quad \lambda_{k+1} = (1 - \alpha_k)\lambda_k \quad (\forall k) \quad \text{define how we update } \lambda_k$$

► With **Lemma 2.2.2 (Partly)**, now

► $\{\lambda_k\}_{k \in \mathbb{N}}$ is *monotonically decreasing*:

$$\lambda_{k+1} \stackrel{\text{L2.2.2A6}}{=} (1 - \alpha_k)\lambda_k \stackrel{\text{L2.2.2A4a}}{<} \lambda_k \stackrel{\text{L2.2.2A6}}{=} (1 - \alpha_{k-1})\lambda_{k-1} \stackrel{\text{L2.2.2A4a}}{<} \lambda_{k-1} < \dots < \lambda_0 := 1 \quad (\#)$$

Reading (#) from right to left also means that Def 0 is satisfied, i.e., all $\lambda_k \geq 0$

(L2.2.2 A4) to (L2.2.2 A6) imply (Def2) $\lambda_{k+1} \rightarrow 0$ is satisfied

Definition 2.2.1 Lemma 2.2.2 (Partly) Assume that

(L2.2.2 A4a) $\alpha_k \in]0, 1[$ ($\forall k$) α_k strictly positive and strictly smaller than 1

(L2.2.2 A4b) $\sum_{k=0}^{\infty} \alpha_k = +\infty$ $\{\alpha_k\}$ is not a summable sequence

(L2.2.2 A5) $\lambda_0 := 1$ we initialize λ_0

(L2.2.2 A6) $\lambda_{k+1} = (1 - \alpha_k)\lambda_k$ ($\forall k$) define how we update λ_k

► With **Lemma 2.2.2 (Partly)**,

$$\lambda_{k+1} \stackrel{L2.2.2A6}{=} (1 - \alpha_k)\lambda_k \stackrel{L2.2.2A6}{=} (1 - \alpha_k)(1 - \alpha_{k-1})\lambda_{k-1} \stackrel{L2.2.2A6}{=} \dots \stackrel{L2.2.2A6}{=} \prod_{i=1}^k (1 - \alpha_i)\lambda_0 \stackrel{L2.2.2A5}{=} \prod_{i=1}^k (1 - \alpha_i) \quad (!)$$

► Now we show that L2.2.2 A4 implies $\prod_{k=1}^{\infty} (1 - \alpha_k) = 0$.

Notice that L2.2.2 A4b is a sum but what we want to prove is product, this gives the hint that we should take log.

Let $S = \prod_{k=1}^{\infty} (1 - \alpha_k) = 0$, now consider

$$\begin{aligned} \log S &= \sum_{k=1}^{\infty} \log(1 - \alpha_k) \leq - \sum_{k=1}^{\infty} \alpha_k && \log(1 - x) \text{ is concave so it is under its 1st-order Taylor expansion} \\ &= -\infty && L2.2.2A4b \\ \iff S &= e^{-\infty} = 0 \end{aligned}$$

Therefore, by (!), we have $\lambda_{\infty} = S = 0$, i.e., $\lambda_k \xrightarrow{k \rightarrow +\infty} 0$.

Why study Nesterov's estimate sequence?

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

$$\begin{aligned} \text{(Def0)} \quad & \lambda_k \geq 0 && (\forall k) \\ \text{(Def1)} \quad & \lambda_k \xrightarrow{k \rightarrow \infty} 0 && (\forall k) \\ \text{(Def2)} \quad & \phi_k(\mathbf{x}) \leq (1 - \lambda_k)f(\mathbf{x}) + \lambda_k \phi_0(\mathbf{x}) && (\forall k)(\forall \mathbf{x} \in \mathbb{R}^n) \end{aligned}$$

► **Lemma 2.2.1** IF for a sequence $\{\mathbf{x}_k\}_{k \in \mathbb{N}}$ we have

$$f(\mathbf{x}_k) \leq \phi_k^* := \min_{\mathbf{x} \in \mathbb{R}^n} \phi_k(\mathbf{x}), \quad (2.2.3)$$

THEN

$$\underbrace{f(\mathbf{x}_k) - f^*}_{\text{a constant}} \leq \lambda_k \underbrace{(\phi_0(\mathbf{x}^*) - f^*)}_{\text{a constant}} \xrightarrow{\text{Def1}} 0. \quad (3)$$

- It forms a global upper bound the of the cost optimality gap $\underbrace{f(\mathbf{x}_k) - f^*}_{\text{a constant}}$
- This upper bound converges to 0 by Def1. (Note $\phi_0(\mathbf{x}^*) - f^*$ is a constant.)

⇒ the convergence rate of $\{f(\mathbf{x}_k) - f^*\}_{k \in \mathbb{N}}$ follows that of $\{\lambda_k\}_{k \in \mathbb{N}}$ the reason why we study estimate sequence

Proof

$$\begin{aligned} f(\mathbf{x}_k) & \stackrel{(2.2.3)}{\leq} \phi_k^* \stackrel{(2.2.3)}{:=} \min_{\mathbf{x} \in \mathbb{R}^n} \phi_k(\mathbf{x}) \stackrel{(\text{Def2})}{\leq} \min_{\mathbf{x} \in \mathbb{R}^n} (1 - \lambda_k)f(\mathbf{x}) + \lambda_k \phi_0(\mathbf{x}) \leq (1 - \lambda_k)f(\mathbf{x}^*) + \lambda_k \phi_0(\mathbf{x}^*) \\ \iff f(\mathbf{x}_k) - f^* & \leq \lambda_k (\phi_0(\mathbf{x}^*) - f^*) \xrightarrow{(\text{Def1})} 0. \quad \square \end{aligned}$$

Nesterov's estimate sequence

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

$$\text{(Def0)} \quad \lambda_k \geq 0 \quad (\forall k)$$

$$\text{(Def1)} \quad \lambda_k \xrightarrow{k \rightarrow \infty} 0 \quad (\forall k)$$

$$\text{(Def2)} \quad \phi_k(\mathbf{x}) \leq (1 - \lambda_k)f(\mathbf{x}) + \lambda_k\phi_0(\mathbf{x}) \quad (\forall k)(\forall \mathbf{x} \in \mathbb{R}^n)$$

► **Lemma 2.2.1** IF for a sequence $\{\mathbf{x}_k\}_{k \in \mathbb{N}}$ we have

$$f(\mathbf{x}_k) \leq \phi_k^* := \min_{\mathbf{x} \in \mathbb{R}^n} \phi_k(\mathbf{x}), \quad (2.2.3)$$

THEN

$$f(\mathbf{x}_k) - f^* \leq \lambda_k(\phi_0(\mathbf{x}^*) - f^*) \xrightarrow{\text{Def1}} 0. \quad (3)$$

► Now we know estimate sequence is **useful to derive convergence rate**

► The next questions is: how to construct an estimate sequence?

► how should we pick ϕ_0 ?

► how should we update λ_k and ϕ_k ?

A way to construct an estimate sequence for smooth convex f

► **Lemma 2.2.2 IF**

(L2.2.2 A1)	f is L -smooth μ -strongly convex	possibly $\mu = 0$ (convex not strongly convex)
(L2.2.2 A2)	$\phi_0(\cdot)$ is a convex function on \mathbb{R}^n	any arbitrary convex function
(L2.2.2 A3)	$\{\mathbf{y}_k\}_{k=0}^\infty$ is a sequence in \mathbb{R}^n	any arbitrary sequence
(L2.2.2 A4a)	$\alpha_k \in]0, 1[$	$(\forall k)$ α_k strictly positive and strictly smaller than 1
(L2.2.2 A4b)	$\sum_{k=0}^\infty \alpha_k = \infty$	$\{\alpha_k\}$ is not a summable sequence
(L2.2.2 A5)	$\lambda_0 := 1$	we initialize λ_0

THEN the sequence-pair $\{\phi_k(\mathbf{x}), \lambda_k\}_{k=0}^\infty$ defined as

(L2.2.2 A6) $\lambda_{k+1} = (1 - \alpha_k)\lambda_k$ $(\forall k)$ how we update λ_k

(L2.2.2 A7) $\phi_{k+1}(\mathbf{x}) = (1 - \alpha_k)\phi_k(\mathbf{x}) + \alpha_k \left(f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x} - \mathbf{y}_k \rangle + \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}_k\|_2^2 \right)$ $(\forall k)$ how we update ϕ_k

is an estimate sequence of $f(\mathbf{x})$.

► To prove $\{\phi_k(\mathbf{x}), \lambda_k\}_{k \in \mathbb{N}}$ is an estimate sequence of $f(\mathbf{x})$, we need to show

- P0 $\{\lambda_k\}_{k \in \mathbb{N}}$ defined in this way is nonnegative
- P1 $\{\lambda_k\}_{k \in \mathbb{N}}$ defined in this way converges to 0
- P2 $\{\phi_k(\mathbf{x})\}_{k \in \mathbb{N}}$ defined in this way satisfies $\phi_k \leq \phi_k(\mathbf{x}) \leq (1 - \lambda_k)f(\mathbf{x}) + \lambda_k\phi_0(\mathbf{x}) \quad \forall k$

► Showing P0 is simple: by (!) we have $\lambda_{k+1} = \prod_{i=1}^k (1 - \alpha_i) > 0$.

Now we have $P0 : \lambda_k \geq 0$.

Proof $P1$: showing $\lambda_k \rightarrow 0$

- ▶ **Proposition** By definition $\lambda_{k+1} = (1 - \alpha_k)\lambda_k$ with assumption $\alpha_k \in]0, 1[$, the sequence $\{\lambda_k\}_{k \in \mathbb{N}}$ is monotonically decreasing.

Proof by ratio test

$$\begin{aligned}\lambda_{k+1} = (1 - \alpha_k)\lambda_k &\iff \frac{\lambda_{k+1}}{\lambda_k} = 1 - \alpha_k \\ &\iff \frac{\lambda_{k+1}}{\lambda_k} < \underset{\alpha_k \in]0, 1[}{1}\end{aligned}$$

- ▶ By $P0 : \lambda_k \geq 0$, the sequence $\{\lambda_k\}_{k \in \mathbb{N}}$ is bounded below by 0.

- ▶ **Theorem (Real analysis 101)**

- ▶ Monotonic decreasing AND bounded below $\implies \{\lambda_k\}_{k \in \mathbb{N}}$ has a limit c

- ▶ What we need to do: show $c = 0$.

There are three ways to show $\lambda_k \rightarrow 0$

► Way 1: By **Monotone convergence theorem (Real analysis 101)**, $c = \inf\{\lambda_k\}_{k \in \mathbb{N}} = 0$.

► Way 2: By contradiction.

► By $\lambda_k \geq 0$, suppose the sequence $\{\lambda_k\}_{k \in \mathbb{N}}$ converges to a positive number $c > 0$.

► Now consider $\lambda_k - \lambda_{k+1} = \alpha_k - \lambda_k = \alpha_k \lambda_k$. It forms a telescoping sum, sum it from 0 to k gives

$$\lambda_0 - \lambda_{k+1} = \sum_{i=0}^k \alpha_i \lambda_i \geq \sum_{i=0}^k \alpha_i c = c \sum_{i=0}^k \alpha_i \quad (*)$$

where the \geq is based on the fact that we assume $\{\lambda_k\}_{k \in \mathbb{N}}$ converges (from above: all $\lambda_k \geq c$ for all k) to c .

► Now $\lambda_0 - \lambda_{k+1} \stackrel{(*)}{\geq} c \sum_{i=0}^k \alpha_i$. Take limit $k \rightarrow \infty$ gives $\lambda_0 - c \geq c \sum_{i=0}^{\infty} \alpha_i$. By $\sum_{i=0}^{\infty} \alpha_i = \infty$ so $\lambda_0 - c \geq +\infty$,

which is impossible (because $\lambda_0 := 1$), a contradiction, therefore $c = 0$.

► Way 3: By $\lambda_{\infty} = S = 0$ using L2.2.2 A4, A5, A6 and property of $\log(1-x)$ as what we did previously.

Proof part 2: on ϕ_k by induction

► Base case $k = 0$: $\phi_0(\mathbf{x}) \leq (1 - \lambda_0)f(\mathbf{x}) + \lambda_0\phi_0(\mathbf{x}) = \phi_0(\mathbf{x})$ by $\lambda_0 := 1$.

► Induction hypothesis $\phi_k(\mathbf{x}) \leq (1 - \lambda_k)f(\mathbf{x}) + \lambda_k\phi_0(\mathbf{x})$

► Case $k + 1$

$$\phi_{k+1}(\mathbf{x}) = (1 - \alpha_k)\phi_k(\mathbf{x}) + \alpha_k \left(f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x} - \mathbf{y}_k \rangle + \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}_k\|_2^2 \right) \quad \text{by A7 def of } \phi_{k+1}$$

$$\leq (1 - \alpha_k)\phi_k(\mathbf{x}) + \alpha_k f(\mathbf{x})$$

A1: f μ -str cvx, $\mu \geq 0$

$$\stackrel{\text{tricky}}{=} (1 - \alpha_k) \left(\phi_k(\mathbf{x}) + \underbrace{(1 - \lambda_k)f(\mathbf{x}) - (1 - \lambda_k)f(\mathbf{x})}_{=0} \right) + \alpha_k f(\mathbf{x})$$

$$= (1 - \alpha_k) \left(\underbrace{\phi_k(\mathbf{x}) - (1 - \lambda_k)f(\mathbf{x})}_{\leq \lambda_k \phi_0(\mathbf{x})} \right) + \underbrace{(1 - \alpha_k)(1 - \lambda_k)f(\mathbf{x})}_{= ((1 - \alpha_k) - (1 - \alpha_k)\lambda_k)f} + \alpha_k f(\mathbf{x})$$

$$\leq (1 - \alpha_k)\lambda_k\phi_0(\mathbf{x}) + \left(1 - (1 - \alpha_k)\lambda_k \right) f(\mathbf{x})$$

case k & $\alpha_k \stackrel{A4a}{<} 1$

$$= \lambda_{k+1}\phi_0(\mathbf{x}) + (1 - \lambda_{k+1})f(\mathbf{x}).$$

$\lambda_{k+1} \stackrel{A6}{=} (1 - \alpha_k)\lambda_k$

So case $k + 1$ is true. By induction, the proof ϕ_k is completed. □

The framework carries over to convex but not strongly convex f

- In the proof

$$\begin{aligned}\phi_{k+1}(\mathbf{x}) &= (1 - \alpha_k)\phi_k(\mathbf{x}) + \alpha_k \left(f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x} - \mathbf{y}_k \rangle + \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}_k\|_2^2 \right) && \text{by A7 def of } \phi_{k+1} \\ &\leq (1 - \alpha_k)\phi_k(\mathbf{x}) + \alpha_k f(\mathbf{x}) && \text{A1: } f \text{ } \mu\text{-str cvx, } \mu \geq 0\end{aligned}$$

The argument holds if f convex but not strongly convex

- In fact the whole framework assume $\mu \geq 0$, which includes the case $\mu = 0 \iff f$ is convex but not strongly convex
- When f is convex but not strongly convex, we construct ϕ_{k+1} as L2.2.2 A7 with $\mu = 0$, i.e.,

$$\begin{aligned}\phi_{k+1}(\mathbf{x}) &= (1 - \alpha_k)\phi_k(\mathbf{x}) + \alpha_k \left(f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x} - \mathbf{y}_k \rangle + \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}_k\|_2^2 \right) && \text{L2.2.2A7} \\ &= (1 - \alpha_k)\phi_k(\mathbf{x}) + \alpha_k \left(f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x} - \mathbf{y}_k \rangle \right) && \mu = 0\end{aligned}$$

The framework carries over to nondifferentiable convex f

- ▶ Note that in the whole proof we never explicitly make use of the assumption that f is L -smooth
- ▶ The only place we make use of f is differentiable is where we assume $\nabla f(\mathbf{x})$ exists at \mathbf{y}_k
- ▶ In the proof

$$\begin{aligned}\phi_{k+1}(\mathbf{x}) &= (1 - \alpha_k)\phi_k(\mathbf{x}) + \alpha_k \left(f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x} - \mathbf{y}_k \rangle + \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}_k\|_2^2 \right) && \text{by A7 def of } \phi_{k+1} \\ &\leq (1 - \alpha_k)\phi_k(\mathbf{x}) + \alpha_k f(\mathbf{x}) && \text{A1: } f \text{ } \mu\text{-str cvx, } \mu \geq 0\end{aligned}$$

The argument holds if f is convex but not differentiable

- ▶ When f is convex but not differentiable, we replace ∇f by subdifferential / subgradient

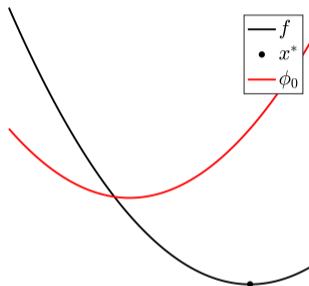
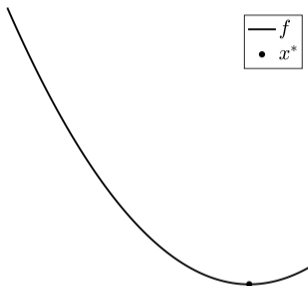
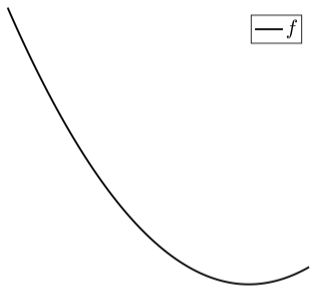
$$\phi_{k+1}(\mathbf{x}) = (1 - \alpha_k)\phi_k(\mathbf{x}) + \alpha_k \left(f(\mathbf{y}_k) + \langle \partial f(\mathbf{y}_k), \mathbf{x} - \mathbf{y}_k \rangle \right)$$

What are these $\phi_k, \alpha_k, \lambda_k$ actually?

$$\phi_{k+1}(\mathbf{x}) = (1 - \alpha_k)\phi_k(\mathbf{x}) + \alpha_k \underbrace{\left(f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x} - \mathbf{y}_k \rangle + \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}_k\|_2^2 \right)}_{\psi(\mathbf{x})} \quad (\text{L2.2.2 A7})$$

- ▶ λ_k is defined by α_k so you can treat them as the same thing under different expression
- ▶ By A4, A5, A6, we can think of λ_k as the coefficient of convex combination and thus think of ϕ_{k+1} as convex combination of ϕ_k and $\psi(\mathbf{x})$
 - ▶ What is ψ : an global support / global under-estimator of f at a point \mathbf{y}_k
 - ▶ Therefore $\phi_{k+1} = \text{cvx}(\phi_k, \psi) = \text{cvx}(\phi_k, \text{lower estimate of } f)$

Understand ϕ_k through pictures ... 1/2

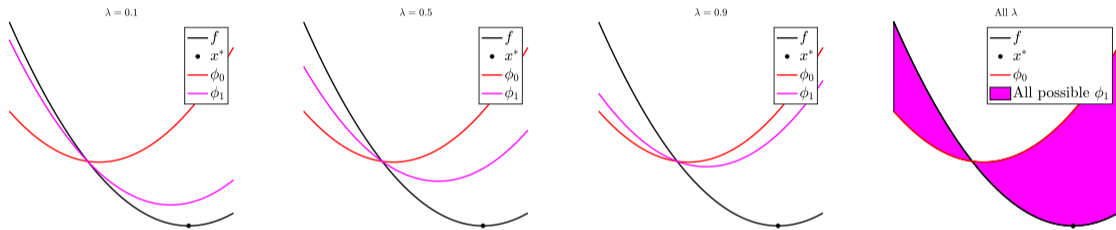


Recall

- ▶ We can pick any ϕ_0 as long as it is convex
- ▶ If f is strongly-convex, a simple ϕ_0 is a quadratic
- ▶ If f is convex, a simple ϕ_0 is a affine function (a line here)
- ▶ If f is convex and nondifferentiable, a simple ϕ_0 is a affine function (where we replace gradient by subgradient)

Understand ϕ_k through pictures ... 2/2

You build ϕ_1 using ϕ_0 and f via a convex combination with weights λ_1



- ▶ An observation: in all the cases, $\{\text{minimizer of } \phi_1\}$ is closer to x^* than $\{\text{minimizer of } \phi_0\}$ to x^*
- ▶ Therefore, if we can somehow find $\{\text{minimizer of } \phi_1\}$ and then use it to construct/update to ϕ_2 , we move closer to x^*
- ▶ By Lemma 2.2.1, the convergence speed of such process is bounded above by how fast λ_k approaches to 0

Small summary

► **[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)]**

$$\begin{array}{l} \text{A1 } f \text{ } L\text{-smooth } \mu\text{-strongly cvx} \\ \text{A2 } \phi_0(\cdot) \text{ a cvx function} \\ \text{A3 } \{\mathbf{y}_k\}_{k=0}^{\infty} \text{ is a sequence} \end{array} \quad \text{THEN } \{\phi_k(\mathbf{x}), \lambda_k\}_{k=0}^{\infty} \text{ defined by A6 A7 is an estimate sequence of } f$$

$$\text{IF } \begin{array}{l} \text{A4a } \alpha_k \in]0, 1[\quad \forall k \\ \text{A4b } \sum_{k=0}^{\infty} \alpha_k = \infty \\ \text{A5 } \lambda_0 := 1 \end{array} \quad \begin{array}{l} \text{A6 } \lambda_{k+1} = (1 - \alpha_k)\lambda_k \\ \text{A7 } \phi_{k+1}(\mathbf{x}) = (1 - \alpha_k)\phi_k(\mathbf{x}) + \alpha_k \left(f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x} - \mathbf{y}_k \rangle + \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}_k\|_2^2 \right) \end{array}$$

► **What now**

- What is ϕ_0 ?
- Well we can use any convex ϕ_0 by A2
- We can use a simple function, like a quadratic

A simple quadratic ϕ_0

- ▶ We can just define ϕ_0 as

$$\phi_0(\mathbf{x}) := \phi_0^* + \frac{\gamma_0}{2} \|\mathbf{x} - \mathbf{v}_0\|_2^2. \quad (\text{Phi-0})$$

- ▶ We now introduce three new things: ϕ , γ and \mathbf{v}
- ▶ ϕ_0, γ_0 are scalars and $\mathbf{v}_0 \in \mathbb{R}^n$ is a vector
- ▶ ϕ_0^* is a shifting parameter, shifting the parabola up and down
- ▶ γ_0 is a slope parameter
- ▶ \mathbf{v}_0 is a shifting parameter, shifting the parabola horizontally
- ▶ ϕ_k, γ_k and \mathbf{v}_k are all sequence that keep changing

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

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

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

Why ϕ_k, γ_k and \mathbf{v}_k are updated this way is not intuitive and can be considered as black magic by Nesterov.

- ▶ Here μ is the strong convexity parameter of f

Lemma on ϕ_k

$$\text{A7 : } \phi_{k+1}(\mathbf{x}) = (1 - \alpha_k)\phi_k(\mathbf{x}) + \alpha_k \left(f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x} - \mathbf{y}_k \rangle + \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}_k\|_2^2 \right)$$

► Lemma 2.2.3 If

$$\phi_0(\mathbf{x}) := \phi_0^* + \frac{\gamma_0}{2} \|\mathbf{x} - \mathbf{v}_0\|_2^2. \quad (\text{Phi-0})$$

Then $\{\phi_k(\mathbf{x})\}_{k=0}^{\infty}$ defined by A7 in Lemma 2.2.2 preserves the canonical form of $\{\phi_k(\mathbf{x})\}_{k \in \mathbb{N}}$

$$\phi_k(\mathbf{x}) = \phi_k^* + \frac{\gamma_k}{2} \|\mathbf{x} - \mathbf{v}_k\|_2^2. \quad (\text{Phi-k})$$

where ϕ_k, γ_k and \mathbf{v}_k are defined as

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

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

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

► Proof First by definition (Phi-0) gives

$$\nabla^2 \phi_0(\mathbf{x}) \stackrel{(\text{Phi-0})}{=} \gamma_0 \mathbf{I}_n. \quad (\dagger)$$

What next we show $\nabla^2 \phi_k(\mathbf{x})$ has the same form as $\nabla^2 \phi_0(\mathbf{x})$, i.e., we want to show

$$\nabla^2 \phi_k(\mathbf{x}) = \gamma_k \mathbf{I}_n.$$

We do so by induction.

Prove $\nabla^2 \phi_k(\mathbf{x}) = \gamma_k \mathbf{I}_n$.

$\nabla^2 \phi_0(\mathbf{x}) \stackrel{(\text{Phi-0})}{=} \gamma_0 \mathbf{I}_n$	(†)
$\phi_{k+1}(\mathbf{x}) = (1 - \alpha_k)\phi_k(\mathbf{x}) + \alpha_k \left(f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x} - \mathbf{y}_k \rangle + \frac{\mu}{2} \ \mathbf{x} - \mathbf{y}_k\ _2^2 \right)$	(A7)
$\gamma_{k+1} = (1 - \alpha_k)\gamma_k + \alpha_k \mu$	(i)

- ▶ Base case is proved by (†)
- ▶ Induction hypothesis: $\nabla^2 \phi_k(\mathbf{x}) = \gamma_k \mathbf{I}_n$
- ▶ Case $k + 1$

$$\begin{aligned} \phi_{k+1}(\mathbf{x}) &= (1 - \alpha_k)\phi_k(\mathbf{x}) + \alpha_k \left(f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x} - \mathbf{y}_k \rangle + \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}_k\|_2^2 \right) && \text{by def (A7)} \\ \nabla \phi_{k+1}(\mathbf{x}) &= (1 - \alpha_k)\nabla \phi_k(\mathbf{x}) + \alpha_k \left(\nabla f(\mathbf{y}_k) + \mu(\mathbf{x} - \mathbf{y}_k) \right) \\ \nabla^2 \phi_{k+1}(\mathbf{x}) &= (1 - \alpha_k)\nabla^2 \phi_k(\mathbf{x}) + \alpha_k \mu \mathbf{I} \\ &= (1 - \alpha_k)\gamma_k \mathbf{I}_n + \alpha_k \mu \mathbf{I} && \text{induction hypothesis} \\ &= \left((1 - \alpha_k)\gamma_k + \alpha_k \mu \right) \mathbf{I} \\ &= \gamma_{k+1} \mathbf{I} \end{aligned}$$

- ▶ Hence now we have showed $\nabla^2 \phi_k(\mathbf{x}) = \gamma_k \mathbf{I}_n$. This equation means that if we perform the antiderivative twice we get

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

for a scalar ϕ_k^* and a vector \mathbf{v}_k . Our remaining tasks are to

- ▶ shows \mathbf{v}_k satisfies (ii) in Lemma 2.2.3
- ▶ shows ϕ_k^* satisfies (iii) in Lemma 2.2.3

Proving \mathbf{v}_{k+1} .

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

$$\phi_{k+1}(\mathbf{x}) = (1 - \alpha_k)\phi_k(\mathbf{x}) + \alpha_k \left(f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x} - \mathbf{y}_k \rangle + \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}_k\|_2^2 \right) \quad (A7)$$

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

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

► First combine (A7) and ($\dagger\dagger$)

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

► What are we going to do now is to find the minimizer of ϕ_{k+1} and denote it as \mathbf{v}_{k+1} . I.e., find $\mathbf{v}_{k+1} = \operatorname{argmin} \phi_{k+1}$. This is basically the idea from the pictures of ϕ_k we previously seen.

► Take gradient

$$\nabla \phi_{k+1}(\mathbf{x}) = (1 - \alpha_k)\gamma_k(\mathbf{x} - \mathbf{v}_k) + \alpha_k \left(\nabla f(\mathbf{y}_k) + \mu(\mathbf{x} - \mathbf{y}_k) \right)$$

► Consider at minimizer \mathbf{v}_{k+1} that $\nabla \phi_{k+1}(\mathbf{v}_{k+1}) = \mathbf{0}$

$$(1 - \alpha_k)\gamma_k(\mathbf{v}_{k+1} - \mathbf{v}_k) + \alpha_k \nabla f(\mathbf{y}_k) + \alpha_k \mu(\mathbf{v}_{k+1} - \mathbf{y}_k) = \mathbf{0}$$

$$\iff ((1 - \alpha_k)\gamma_k + \alpha_k \mu)\mathbf{v}_{k+1} + \alpha_k \nabla f(\mathbf{y}_k) - (1 - \alpha_k)\gamma_k \mathbf{v}_k - \alpha_k \mu \mathbf{y}_k = \mathbf{0}$$

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

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

$$\iff (ii)$$

Proving ϕ_{k+1}^* .

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

$$\phi_{k+1}(\mathbf{x}) = (1 - \alpha_k)\phi_k(\mathbf{x}) + \alpha_k \left(f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x} - \mathbf{y}_k \rangle + \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}_k\|_2^2 \right) \quad (A7)$$

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

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

► $(\dagger\dagger) = (A7)$ at $k + 1$

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

► Put $\mathbf{x} = \mathbf{y}_k$

$$\begin{aligned} \phi_{k+1}^* + \frac{\gamma_{k+1}}{2} \|\mathbf{y}_k - \mathbf{v}_{k+1}\|_2^2 &= (1 - \alpha_k)\phi_k(\mathbf{y}_k) + \alpha_k f(\mathbf{y}_k) \\ &\stackrel{(\dagger\dagger)}{=} (1 - \alpha_k) \left(\phi_k^* + \frac{\gamma_k}{2} \|\mathbf{y}_k - \mathbf{v}_k\|_2^2 \right) + \alpha_k f(\mathbf{y}_k) \quad (\dagger\dagger\dagger) \end{aligned}$$

► By (ii)

$$\begin{aligned} \mathbf{v}_{k+1} - \mathbf{y}_k &= \frac{(1 - \alpha_k)\gamma_k \mathbf{v}_k + \alpha_k \mu \mathbf{y}_k - \alpha_k \nabla f(\mathbf{y}_k)}{\gamma_{k+1}} - \mathbf{y}_k \\ &= \frac{(1 - \alpha_k)\gamma_k \mathbf{v}_k + \alpha_k \mu \mathbf{y}_k - \gamma_{k+1} \mathbf{y}_k - \alpha_k \nabla f(\mathbf{y}_k)}{\gamma_{k+1}} \\ &\stackrel{(i)}{=} \frac{(1 - \alpha_k)\gamma_k (\mathbf{v}_k - \mathbf{y}_k) - \alpha_k \nabla f(\mathbf{y}_k)}{\gamma_{k+1}} \end{aligned}$$

$$\frac{\gamma_{k+1}}{2} \|\mathbf{v}_{k+1} - \mathbf{y}_k\|_2^2 = \frac{(1 - \alpha_k)^2 \gamma_k^2 \|\mathbf{v}_k - \mathbf{y}_k\|_2^2 - 2 \langle (1 - \alpha_k)\gamma_k (\mathbf{v}_k - \mathbf{y}_k), \alpha_k \nabla f(\mathbf{y}_k) \rangle + \alpha_k^2 \|\nabla f(\mathbf{y}_k)\|_2^2}{2\gamma_{k+1}} \quad (\dagger\dagger\dagger\dagger)$$

► Put $(\dagger\dagger\dagger\dagger)$ into $(\dagger\dagger\dagger)$ will give (iii) , trust me.



$$\phi_{k+1}^* + \frac{\gamma_{k+1}}{2} \|\mathbf{y}_k - \mathbf{v}_{k+1}\|_2^2 = (1 - \alpha_k) \left(\phi_k^* + \frac{\gamma_k}{2} \|\mathbf{y}_k - \mathbf{v}_k\|_2^2 \right) + \alpha_k f(\mathbf{y}_k) \quad (\dagger \dagger \dagger)$$

$$\frac{\gamma_{k+1}}{2} \|\mathbf{v}_{k+1} - \mathbf{y}_k\|_2^2 = \frac{(1 - \alpha_k)^2 \gamma_k^2 \|\mathbf{v}_k - \mathbf{y}_k\|_2^2 - 2 \langle (1 - \alpha_k) \gamma_k (\mathbf{v}_k - \mathbf{y}_k), \alpha_k \nabla f(\mathbf{y}_k) \rangle + \alpha_k^2 \|\nabla f(\mathbf{y}_k)\|_2^2}{2\gamma_{k+1}} \quad (\dagger \dagger \dagger \dagger)$$

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

What to do: show $\{(\dagger \dagger \dagger) \text{ and } (\dagger \dagger \dagger \dagger)\} - (\text{iii}) = 0$

$$\begin{aligned} & (1 - \alpha_k) \left(\phi_k^* + \frac{\gamma_k}{2} \|\mathbf{y}_k - \mathbf{v}_k\|_2^2 \right) + \alpha_k f(\mathbf{y}_k) - \frac{\gamma_{k+1}}{2} \|\mathbf{y}_k - \mathbf{v}_{k+1}\|_2^2 \\ & - \left\{ (1 - \alpha_k) \phi_k^* + \alpha_k f(\mathbf{y}_k) - \frac{\alpha_k^2}{2\gamma_{k+1}} \|\nabla f(\mathbf{y}_k)\|_2^2 + \frac{\alpha_k (1 - \alpha_k) \gamma_k}{\gamma_{k+1}} \left(\frac{\mu}{2} \|\mathbf{y}_k - \mathbf{v}_k\|_2^2 + \langle \nabla f(\mathbf{y}_k), \mathbf{v}_k - \mathbf{y}_k \rangle \right) \right\} \\ = & (1 - \alpha_k) \frac{\gamma_k}{2} \|\mathbf{y}_k - \mathbf{v}_k\|_2^2 - \frac{(1 - \alpha_k)^2 \gamma_k^2 \|\mathbf{v}_k - \mathbf{y}_k\|_2^2 - 2 \langle (1 - \alpha_k) \gamma_k (\mathbf{v}_k - \mathbf{y}_k), \alpha_k \nabla f(\mathbf{y}_k) \rangle + \alpha_k^2 \|\nabla f(\mathbf{y}_k)\|_2^2}{2\gamma_{k+1}} \\ & - \left\{ -\frac{\alpha_k^2}{2\gamma_{k+1}} \|\nabla f(\mathbf{y}_k)\|_2^2 + \frac{\alpha_k (1 - \alpha_k) \gamma_k}{\gamma_{k+1}} \left(\frac{\mu}{2} \|\mathbf{y}_k - \mathbf{v}_k\|_2^2 + \langle \nabla f(\mathbf{y}_k), \mathbf{v}_k - \mathbf{y}_k \rangle \right) \right\} \\ = & \frac{(1 - \alpha_k) \gamma_k}{2} \left[1 - \frac{(1 - \alpha_k) \gamma_k}{\gamma_{k+1}} \right] \|\mathbf{y}_k - \mathbf{v}_k\|_2^2 + \frac{(1 - \alpha_k) \alpha_k \gamma_k}{\gamma_{k+1}} \langle \mathbf{v}_k - \mathbf{y}_k, \nabla f(\mathbf{y}_k) \rangle \\ & - \left\{ \frac{\alpha_k (1 - \alpha_k) \gamma_k}{\gamma_{k+1}} \frac{\mu}{2} \|\mathbf{y}_k - \mathbf{v}_k\|_2^2 + \frac{\alpha_k (1 - \alpha_k) \gamma_k}{\gamma_{k+1}} \langle \nabla f(\mathbf{y}_k), \mathbf{v}_k - \mathbf{y}_k \rangle \right\} \\ = & 0 \text{ by (i) } \gamma_{k+1} = (1 - \alpha_k) \gamma_k + \alpha_k \mu \end{aligned}$$

Last page

- [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)]

$$\begin{array}{l} \text{A1 } f \text{ } L\text{-smooth } \mu\text{-strongly cvx} \\ \text{A2 } \phi_0(\cdot) \text{ a cvx function} \\ \text{A3 } \{\mathbf{y}_k\}_{k=0}^{\infty} \text{ is a sequence} \end{array} \quad \text{THEN } \{\phi_k(\mathbf{x}), \lambda_k\}_{k=0}^{\infty} \text{ defined by A6 A7 is an estimate sequence of } f$$

$$\begin{array}{l} \text{IF } \text{A4a } \alpha_k \in]0, 1[\quad \forall k \\ \text{A4b } \sum_{k=0}^{\infty} \alpha_k = \infty \\ \text{A5 } \lambda_0 := 1 \end{array} \quad \begin{array}{l} \text{A6 } \lambda_{k+1} = (1 - \alpha_k)\lambda_k \\ \text{A7 } \phi_{k+1}(\mathbf{x}) = (1 - \alpha_k)\phi_k(\mathbf{x}) + \alpha_k (f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x} - \mathbf{y}_k \rangle + \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}_k\|_2^2) \end{array}$$

- [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 $\{\phi_k(\mathbf{x})\}_{k \in \mathbb{N}}$

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

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

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

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