

SD

$$\frac{1}{\sqrt{k}}$$

GD

$$\frac{1}{k}$$

$$\Theta \left(\frac{2-1}{2+1} \right)$$

$$\mu = 0$$

$$\mu > 0$$

Discover **acceleration** of gradient descent

Daniil Merkulov

Optimization methods. MIPT

$$\frac{1}{k^2}$$

$$\Theta \left(\frac{\sqrt{2}-1}{2\sqrt{2}+1} \right)$$

Neural
Network

Coordinate shift

Consider the following quadratic optimization problem:

$$\min_{x \in \mathbb{R}^d} f(x) = \min_{x \in \mathbb{R}^d} \frac{1}{2} x^\top A x - b^\top x + c, \text{ where } A \in \mathbb{S}_{++}^d.$$

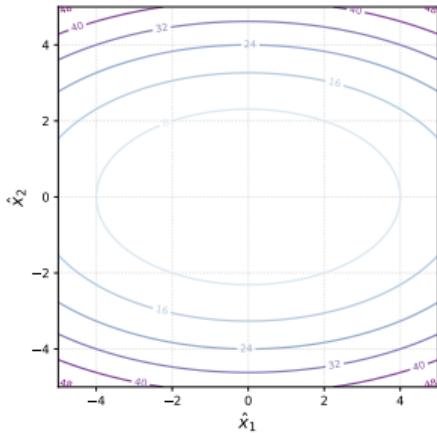
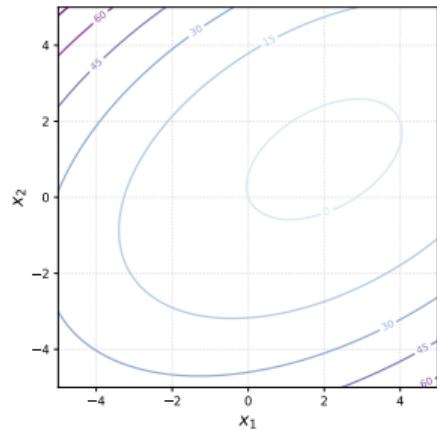
$$x^* = A^{-1} b$$

Coordinate shift

Consider the following quadratic optimization problem:

$$\min_{x \in \mathbb{R}^d} f(x) = \min_{x \in \mathbb{R}^d} \frac{1}{2} x^\top A x - b^\top x + c, \text{ where } A \in \mathbb{S}_{++}^d.$$

- Firstly, without loss of generality we can set $c = 0$, which will not affect optimization process.



Coordinate shift

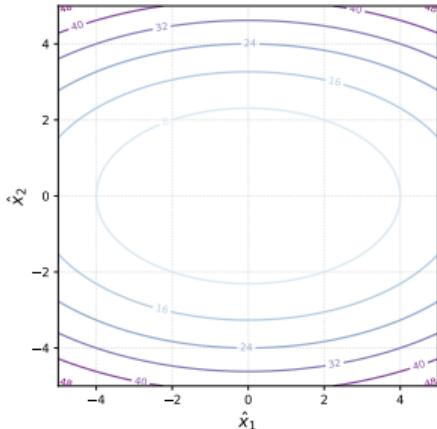
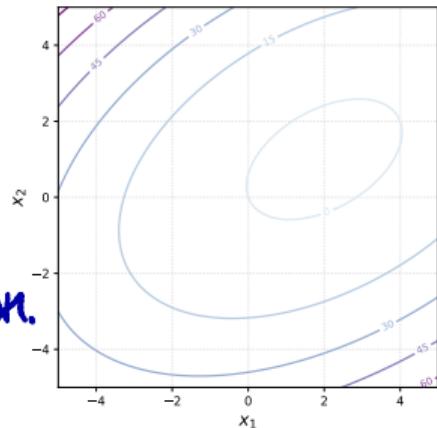
Consider the following quadratic optimization problem:

$$\min_{x \in \mathbb{R}^d} f(x) = \min_{x \in \mathbb{R}^d} \frac{1}{2} x^\top A x - b^\top x + c, \text{ where } A \in \mathbb{S}_{++}^d.$$

- Firstly, without loss of generality we can set $c = 0$, which will not affect optimization process.
- Secondly, we have a spectral decomposition of the matrix A :

$$A = Q \Lambda Q^T$$

$$\begin{aligned} Q^T Q &= I && \text{option.} \\ Q Q^T &= I \end{aligned}$$



Coordinate shift

Consider the following quadratic optimization problem:

$$\min_{x \in \mathbb{R}^d} f(x) = \min_{x \in \mathbb{R}^d} \frac{1}{2} x^\top A x - b^\top x + c, \text{ where } A \in \mathbb{S}_{++}^d.$$

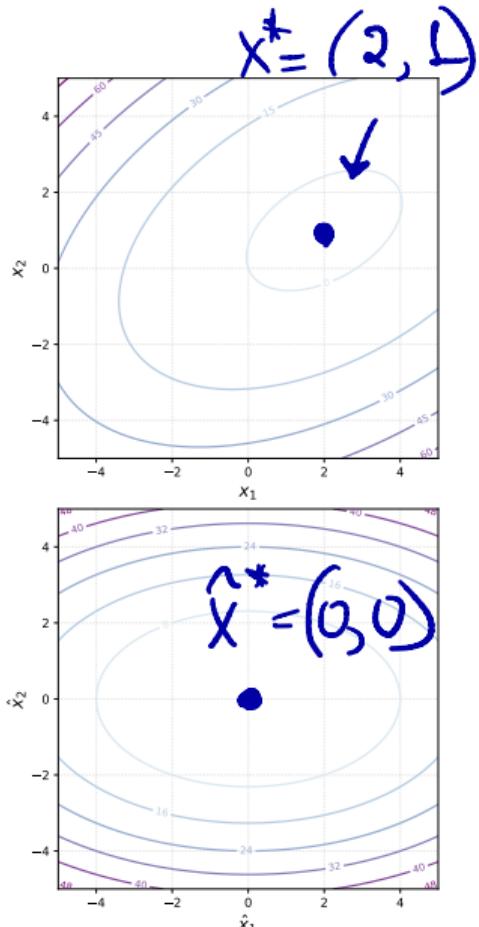
- Firstly, without loss of generality we can set $c = 0$, which will not affect optimization process.
- Secondly, we have a spectral decomposition of the matrix A :

$$A = Q \Lambda Q^T$$

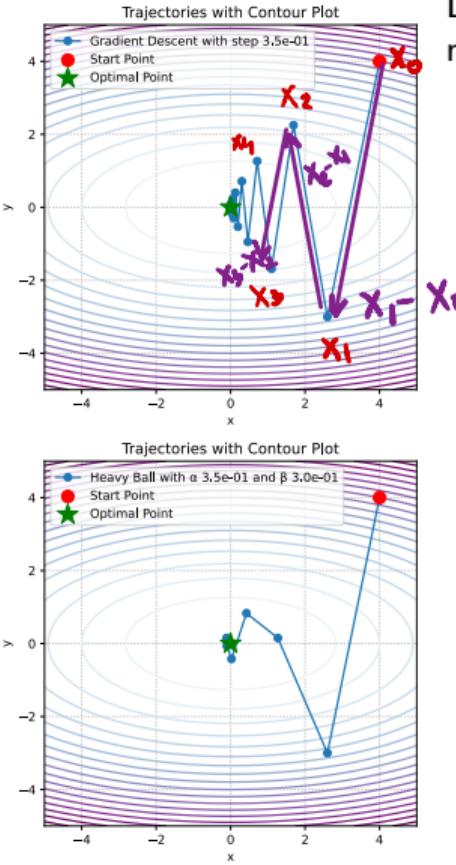
noboot + ebur

- Let's show, that we can switch coordinates in order to make an analysis a little bit easier. Let $\hat{x} = Q^T(x - x^*)$, where x^* is the minimum point of initial function, defined by $Ax^* = b$. At the same time $x = Q\hat{x} + x^*$.

$$\begin{aligned} f(\hat{x}) &= \frac{1}{2}(Q\hat{x} + x^*)^\top A(Q\hat{x} + x^*) - b^\top(Q\hat{x} + x^*) \\ &= \frac{1}{2}\hat{x}^T Q^T A Q \hat{x} + (x^*)^T A Q \hat{x} + \frac{1}{2}(x^*)^T A (x^*)^T - b^T Q \hat{x} - b^T x^* \\ &= \frac{1}{2}\hat{x}^T \Lambda \hat{x} \end{aligned}$$



Polyak Heavy ball method



Let's introduce the idea of momentum, proposed by Polyak in 1964. Recall that the momentum update is

размер шага

коэффициент (моментум)

$$\beta < 1$$

$$x_{k+1} = x_k - \alpha \nabla f(x_k) + \beta(x_k - x_{k-1}). \quad (=)$$

$$x_k = x_{k-1} - \alpha \nabla f(x_{k-1}) + \beta(x_{k-1} - x_{k-2})$$

$$= x_k - \alpha \nabla f(x_k) + \beta \cdot (-\alpha \nabla f(x_{k-1}) + \beta(x_{k-1} - x_{k-2}))$$

$$= x_k - \alpha \nabla f(x_k) - \alpha \cdot \beta \nabla f(x_{k-1}) + \beta^2(x_{k-1} - x_{k-2}) =$$

$$= x_k - \alpha [\nabla f(x_k) + \beta \cdot \nabla f(x_{k-1}) + \beta^2 \nabla f(x_{k-2}) + \beta^3 \dots]$$

Polyak Heavy ball method

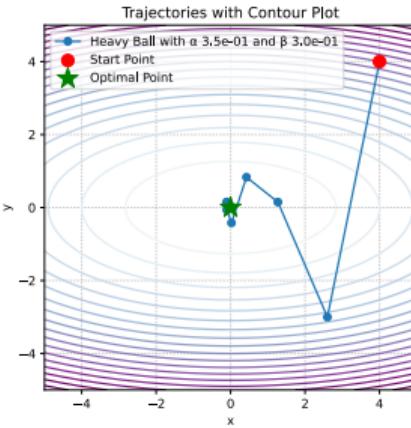
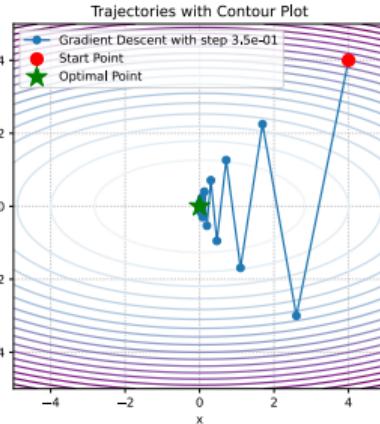
Let's introduce the idea of momentum, proposed by Polyak in 1964. Recall that the momentum update is

$$\nabla f(x_k) = \nabla f(\hat{x}_k) = \Delta x_k$$

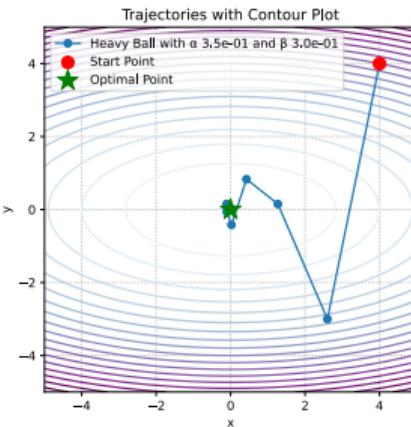
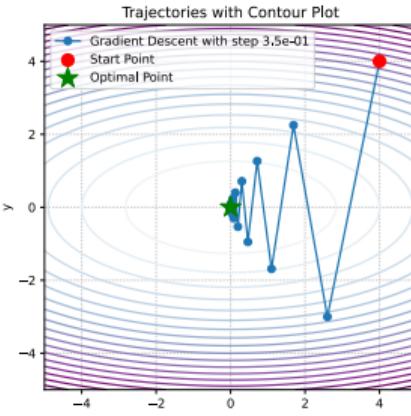
$$x_{k+1} = x_k - \alpha \nabla f(x_k) + \beta(x_k - x_{k-1}).$$

Which is in our case is

$$\hat{x}_{k+1} = \hat{x}_k - \alpha \Lambda \hat{x}_k + \beta(\hat{x}_k - \hat{x}_{k-1}) = (I - \alpha \Lambda + \beta I)\hat{x}_k - \beta \hat{x}_{k-1}$$



Polyak Heavy ball method



Let's introduce the idea of momentum, proposed by Polyak in 1964. Recall that the momentum update is

$$x_{k+1} = x_k - \alpha \nabla f(x_k) + \beta(x_k - x_{k-1}).$$

Which is in our case is

$$\hat{x}_{k+1} = \hat{x}_k - \alpha \Lambda \hat{x}_k + \beta(\hat{x}_k - \hat{x}_{k-1}) = (I - \alpha \Lambda + \beta I)\hat{x}_k - \beta \hat{x}_{k-1}$$

This can be rewritten as follows

$$\begin{cases} \hat{x}_{k+1} = (I - \alpha \Lambda + \beta I)\hat{x}_k - \beta \hat{x}_{k-1}, \\ \hat{x}_k = \hat{x}_k. \end{cases}$$

$\stackrel{K}{\overbrace{M \cdot \hat{z}_0}}$

Let's use the following notation $\hat{z}_k = \begin{bmatrix} \hat{x}_{k+1} \\ \hat{x}_k \end{bmatrix}$. Therefore $\hat{z}_{k+1} = M\hat{z}_k$, where the iteration matrix M is:

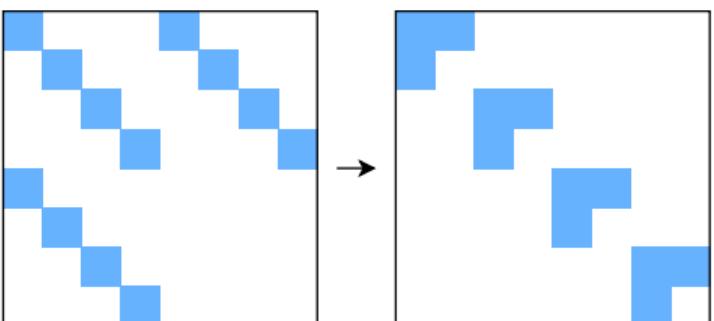
$$M = \begin{bmatrix} I - \alpha \Lambda + \beta I & -\beta I \\ I & 0_d \end{bmatrix} \cdot \hat{z}_{k-1} = \begin{bmatrix} \hat{x}_k \\ \hat{x}_{k-1} \end{bmatrix}$$

Reduction to a scalar case

Note, that M is $2d \times 2d$ matrix with 4 block-diagonal matrices of size $d \times d$ inside. It means, that we can rearrange the order of coordinates to make M block-diagonal in the following form. Note that in the equation below, the matrix M denotes the same as in the notation above, except for the described permutation of rows and columns. We use this slight abuse of notation for the sake of clarity.

Reduction to a scalar case

Note, that M is $2d \times 2d$ matrix with 4 block-diagonal matrices of size $d \times d$ inside. It means, that we can rearrange the order of coordinates to make M block-diagonal in the following form. Note that in the equation below, the matrix M denotes the same as in the notation above, except for the described permutation of rows and columns. We use this slight abuse of notation for the sake of clarity.



$$\begin{bmatrix} \hat{x}_k^{(1)} \\ \vdots \\ \hat{x}_k^{(d)} \\ \hat{x}_{k-1}^{(1)} \\ \vdots \\ \hat{x}_{k-1}^{(d)} \end{bmatrix} \rightarrow \begin{bmatrix} \hat{x}_k^{(1)} \\ \hat{x}_{k-1}^{(1)} \\ \vdots \\ \hat{x}_k^{(d)} \\ \hat{x}_{k-1}^{(d)} \end{bmatrix} \quad M = \begin{bmatrix} M_1 & & & \\ & M_2 & & \\ & & \ddots & \\ & & & M_d \end{bmatrix}$$

Figure 1: Illustration of matrix M rearrangement

where $\hat{x}_k^{(i)}$ is i -th coordinate of vector $\hat{x}_k \in \mathbb{R}^d$ and M_i stands for 2×2 matrix. This rearrangement allows us to study the dynamics of the method independently for each dimension. One may observe, that the asymptotic convergence rate of the $2d$ -dimensional vector sequence of \hat{z}_k is defined by the worst convergence rate among its block of coordinates. Thus, it is enough to study the optimization in a one-dimensional case.

Reduction to a scalar case

Скорост барота: $\max_{i \in \{s-d\}} |\lambda_i| = \rho(M)$

Тенеп

гостило үтб

былбасы i-ын

координаты

For i -th coordinate with λ_i as an i -th eigenvalue of matrix W we have:

$$\rho(M) = \max_i |\lambda_i|$$

$$M_i = \begin{bmatrix} 1 - \alpha \lambda_i + \beta & -\beta \\ 1 & 0 \end{bmatrix}.$$

$$\begin{pmatrix} \overset{\wedge}{x_k^i} \\ \overset{\wedge}{x_{k+1}^i} \\ \vdots \\ \overset{\wedge}{x_k^i} \\ \overset{\wedge}{x_{k-1}^i} \end{pmatrix} = M_i \begin{pmatrix} \overset{\wedge}{x_k^i} \\ \overset{\wedge}{x_{k+1}^i} \\ \vdots \\ \overset{\wedge}{x_k^i} \\ \overset{\wedge}{x_{k-1}^i} \end{pmatrix}$$

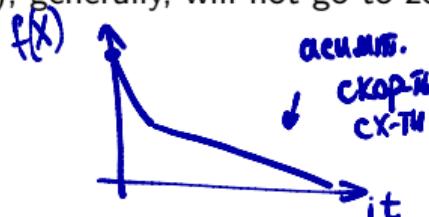
The method will be convergent if $\rho(M) < 1$, and the optimal parameters can be computed by optimizing the spectral radius

$$\alpha^*, \beta^* = \arg \min_{\alpha, \beta} \max_{\lambda \in [\mu, L]} |\lambda|$$

$$\alpha^* = \frac{4}{(\sqrt{L} + \sqrt{\mu})^2};$$

$$\beta^* = \left(\frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}} \right)^2.$$

It can be shown, that for such parameters the matrix M has complex eigenvalues, which forms a conjugate pair, so the distance to the optimum (in this case, $\|z_k\|$), generally, will not go to zero monotonically.



Heavy ball quadratic convergence

$$\det(M - \lambda I) = 0$$

$$\downarrow \quad \lambda^*, \beta^*$$

We can explicitly calculate the eigenvalues of M_i :

$$\lambda_1^M, \lambda_2^M = \lambda \left(\begin{bmatrix} 1 - \alpha\lambda_i + \beta & -\beta \\ 1 & 0 \end{bmatrix} \right) = \frac{1 + \beta - \alpha\lambda_i \pm \sqrt{(1 + \beta - \alpha\lambda_i)^2 - 4\beta}}{2}.$$

When α and β are optimal (α^*, β^*), the eigenvalues are complex-conjugated pair $(1 + \beta - \alpha\lambda_i)^2 - 4\beta \leq 0$, i.e. $\beta \geq (1 - \sqrt{\alpha\lambda_i})^2$.

$$\text{Re}(\lambda_1^M) = \frac{L + \mu - 2\lambda_i}{(\sqrt{L} + \sqrt{\mu})^2}; \quad \text{Im}(\lambda_1^M) = \frac{\pm 2\sqrt{(L - \lambda_i)(\lambda_i - \mu)}}{(\sqrt{L} + \sqrt{\mu})^2}; \quad |\lambda_1^M| = \frac{L - \mu}{(\sqrt{L} + \sqrt{\mu})^2}.$$

And the convergence rate does not depend on the stepsize and equals to $\sqrt{\beta^*}$.

$$= \sqrt{\beta^*} = \sqrt{\frac{L^2 - \mu^2}{(L + \mu)^2}} = \sqrt{\frac{L - \mu}{L + \mu}}$$

Heavy ball method summary

- Ensures accelerated convergence for strongly convex quadratic problems



Heavy ball method summary

- Ensures accelerated convergence for strongly convex quadratic problems
- Local accelerated convergence was proved in the original paper.

Heavy ball method summary

- Ensures accelerated convergence for strongly convex quadratic problems
- Local accelerated convergence was proved in the original paper.
- Recently was proved, that there is no global accelerated convergence for the method.

Heavy ball method summary



- Ensures accelerated convergence for strongly convex quadratic problems
- Local accelerated convergence was proved in the original paper.
- Recently was proved, that there is no global accelerated convergence for the method.
- Method was not extremely popular until the ML boom

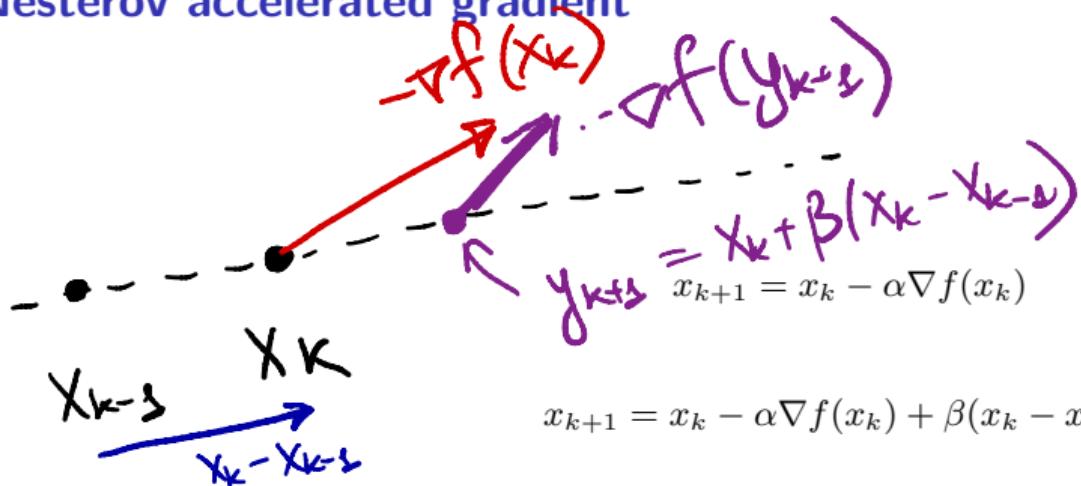
Heavy ball method summary

2 exellent

HB , NAG

- Ensures accelerated convergence for strongly convex quadratic problems
- Local accelerated convergence was proved in the original paper.
- Recently was proved, that there is no global accelerated convergence for the method.
- Method was not extremely popular until the ML boom
- Nowadays, it is de-facto standard for practical acceleration of gradient methods, even for the non-convex problems (neural network training)

Nesterov accelerated gradient



$$x_{k+1} = x_k - \alpha \nabla f(x_k) + \beta(x_k - x_{k-1})$$

(HB)

$$\begin{cases} y_{k+1} = x_k + \beta(x_k - x_{k-1}) \\ x_{k+1} = y_{k+1} - \alpha \nabla f(y_{k+1}) \end{cases}$$

$$x_{k+1} = x_k + \beta(x_k - x_{k-1}) - \alpha \cdot \nabla f(x_k + \beta(x_k - x_{k-1})) + \frac{\beta}{k^2}$$

(1981)

$$\exp\left(\frac{\sqrt{\alpha}-1}{\sqrt{\alpha}+1}\right)$$

Nesterov accelerated gradient

$$x_{k+1} = x_k - \alpha \nabla f(x_k)$$

$$y_{k+1} = x_k - x_{k-1} \quad (\text{GD})$$

$$x_{k+1} = x_k - \alpha \nabla f(x_k) + \beta(x_k - x_{k-1}) \quad | \quad x_{k+1} = x_k - \alpha \nabla f(x_k) + \beta y_{k+1} \quad (\text{HB})$$

$$\begin{cases} y_{k+1} = x_k + \beta(x_k - x_{k-1}) \\ x_{k+1} = y_{k+1} - \alpha \nabla f(y_{k+1}) \end{cases} \quad (\text{NAG})$$

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

Proof approach 1

Andersen Ang

ECS, Uni. Southampton, UK

andersen.ang@soton.ac.uk

Homepage angms.science

Version: November 4, 2023

First draft: August 2, 2017

Content

Problem setup: smooth unconstrained convex optimisation
Nesterov's accelerated gradient descent (NAGD)

Proving NAGD converges rate $\mathcal{O}\left(\frac{1}{k^2}\right)$

Summary

Problem setup: smooth unconstrained convex optimisation

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

- We consider Euclidean space
- $f : \mathbb{R}^n \rightarrow \mathbb{R}$
- f is L -smooth

- f is continuously differentiable
- $f \in \mathcal{C}^1$, i.e., $\nabla f(\mathbf{x})$ exists for all $\mathbf{x} \in \text{dom } f$
- ∇f is L -Lipschitz
 - $L > 0$ is the least upper bound in $\frac{\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|}{\|\mathbf{x} - \mathbf{y}\|} \leq L$

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

- f is convex
 - all local minima of \mathcal{P} are global minima
- $$(\forall \mathbf{x} \in \text{dom } f)(\forall \mathbf{y} \in \text{dom } f) \left\{ f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \right\}$$
- Details of convexity, L-smoothness, see [here](#)

Gradient Descent (GD)

- Notation

$$\begin{aligned} f_k &:= f(\mathbf{x}_k) \\ f^* &:= f(\mathbf{x}^*) \end{aligned}$$

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

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

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

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

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

Nesterov's Accelerated Gradient Descent (NAGD)

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

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

Gradient update

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

Extrapolation

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

Extrapolation weight

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

Extrapolation weight

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

Proof 2/6 Stage 2: eliminate gradient

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

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

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

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

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

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

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

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

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

► Put (11,13,14) into (8)

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

► Put (12,13,14) into (9)

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

Proof 3/6 Stage 3: form telescoping sum

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

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

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

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

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

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

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

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

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

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

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

Proof 4/6

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

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

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

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

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

► Using this (17) becomes

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

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

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

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

So (18) becomes

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

Proof ... 5/6

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

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

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

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

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

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

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

Proof ... 6/6

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

Proof (by induction)

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

► Case $k = 1$. By definition,

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

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

► Case $k = n$

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

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

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

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

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

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

Rename K as k gives

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

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

Last page - summary

For unconstrained convex smooth problem

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

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

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

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

Extrapolation weight $\gamma_k = \frac{1 - \lambda_k}{\lambda_{k+1}}$

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

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

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

The proof can be used for proximal gradient descent.

End of document

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

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

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

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

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

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

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

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

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

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

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

$L > 0$ 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

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

- 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}
- f strongly convex \implies {unique minimizer}

$\{\mathbf{x}_k\}_{k \in \mathbb{N}} \rightarrow$ a global minimizer \mathbf{x}^* (if it exists)
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**, $f_k - f^* \leq \mathcal{O}\left(\frac{1}{k}\right)$
[Details](#)

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

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

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

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

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

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

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

- ▶ $\Phi_k(\mathbf{x})$ 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 $\underbrace{\text{2nd-order Taylor approximation of } f \text{ at } \mathbf{x}_k}_{\text{f}(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{m}{2} \|\mathbf{x} - \mathbf{x}_k\|_2^2}$.
 - ▶ $Q = 1 \iff$ 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

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

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

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

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

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

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

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

- 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 $k = 0, 1$

$$\boxed{\begin{aligned}\Phi_{k+1}(\mathbf{x}) &:= \left(1 - \frac{1}{\sqrt{Q}}\right)\Phi_k(\mathbf{x}) + \frac{1}{\sqrt{Q}}\left(f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{m}{2}\|\mathbf{x} - \mathbf{x}_k\|_2^2\right) \\ f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{m}{2}\|\mathbf{x} - \mathbf{x}_k\|_2^2 &\leq f(\mathbf{x}) \quad \forall \mathbf{x}, \mathbf{x}_k\end{aligned}\quad f \text{ is } m\text{-strongly cvx}}$$

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

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

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

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

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

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

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

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

f is m -strongly cvx

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

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

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

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

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

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

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

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

Lemma 1

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

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

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

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

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** $k + 1$

$$\begin{aligned}\Phi_{k+1}(\mathbf{x}) &= \left(1 - \frac{1}{\sqrt{Q}}\right) \Phi_k(\mathbf{x}) + \frac{1}{\sqrt{Q}} \left(f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{m}{2} \|\mathbf{x} - \mathbf{x}_k\|_2^2 \right) \quad \text{by definition} \\ \nabla^2 \Phi_{k+1}(\mathbf{x}) &= \left(1 - \frac{1}{\sqrt{Q}}\right) \nabla^2 \Phi_k(\mathbf{x}) + \frac{1}{\sqrt{Q}} m\mathbf{I}_n \\ &= \left(1 - \frac{1}{\sqrt{Q}}\right) m\mathbf{I}_n + \frac{1}{\sqrt{Q}} m\mathbf{I}_n \quad \text{induction hypothesis} \\ &= m\mathbf{I}_n \quad \square\end{aligned}$$

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

$\Phi_0(\mathbf{x})$	\coloneqq	$f(\mathbf{x}_0) + \frac{m}{2} \ \mathbf{x} - \mathbf{x}_0\ _2^2$	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** $k = 0$

$$\Phi_0^* = \min_{x \in \mathbb{R}^n} \Phi_0(x) = \min_{x \in \mathbb{R}^n} f(x_0) + \frac{m}{2} \|x - x_0\|_2^2 = f(x_0) = f(y_0)$$

► **Induction Hypothesis** $f(y_k) \leq \Phi_k^*$

► **Case $k + 1$** Consider $f(y_{k+1})$ and L -smoothness of f

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

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

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 convex
$f(\mathbf{y}_k)$	\leq	Φ_k^*	Induction Hypothesis

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

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

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

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

This is what we are going to do in the next 4 - 5 slides.

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 $\nabla \Phi_k(\boldsymbol{\nu}_k) = 0$
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

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

(2) & (3) gives $\nabla \Phi_{k+1}(\boldsymbol{\nu}_{k+1}) = \mathbf{0}$

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

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

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

Now

$$\begin{aligned} \boldsymbol{\nu}_{k+1} &= \left(1 - \frac{1}{\sqrt{Q}}\right)\boldsymbol{\nu}_k + \frac{1}{\sqrt{Q}}\left(\mathbf{x}_k - \frac{\nabla f(\mathbf{x}_k)}{m}\right) \tag{5} \\ \iff -\boldsymbol{\nu}_{k+1} &= -\left(1 - \frac{1}{\sqrt{Q}}\right)\boldsymbol{\nu}_k - \frac{1}{\sqrt{Q}}\left(\mathbf{x}_k - \frac{\nabla f(\mathbf{x}_k)}{m}\right) \\ \iff \mathbf{x}_k - \boldsymbol{\nu}_{k+1} &= \mathbf{x}_k - \left(1 - \frac{1}{\sqrt{Q}}\right)\boldsymbol{\nu}_k - \frac{1}{\sqrt{Q}}\mathbf{x}_k + \frac{1}{\sqrt{Q}}\frac{\nabla f(\mathbf{x}_k)}{m} \\ &= \left(1 - \frac{1}{\sqrt{Q}}\right)(\mathbf{x}_k - \boldsymbol{\nu}_k) + \frac{\nabla f(\mathbf{x}_k)}{m\sqrt{Q}} \\ \iff \|\mathbf{x}_k - \boldsymbol{\nu}_{k+1}\|_2^2 &= \left(1 - \frac{1}{\sqrt{Q}}\right)^2 \|\mathbf{x}_k - \boldsymbol{\nu}_k\|_2^2 + 2\left(1 - \frac{1}{\sqrt{Q}}\right) \frac{\langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \boldsymbol{\nu}_k \rangle}{m\sqrt{Q}} + \frac{\|\nabla f(\mathbf{x}_k)\|_2^2}{m^2 Q} \end{aligned}$$

Lemma 3 ... 5/7

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

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

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

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

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

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

► $(*) = (**)$ gives

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

Lemma 3 ... 6/7

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

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

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

Hence

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

Simplify the term $\sim\sim$

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

To proceed, we need lemma 4.

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

$Q = \frac{L}{m}$	def of Q
$\mathbf{y}_{k+1} = \mathbf{x}_k - \frac{\nabla f(\mathbf{x}_k)}{\frac{L}{m}}$	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 $\nu_0 = \mathbf{x}_0$.
- **Induction hypothesis** $\nu_k - \mathbf{x}_k = \sqrt{Q}(\mathbf{x}_k - \mathbf{y}_k)$
- **Case** $k + 1$

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

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 $a = \Phi_{k+1}^* = \underbrace{\textcolor{yellow}{\dots}}_{\geq 0} + \textcolor{pink}{\dots} \geq \textcolor{pink}{\dots} \stackrel{\text{now we have}}{\geq} f(\mathbf{y}_{k+1})$, 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

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

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

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

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

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

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

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 $\min_{x \in \mathbb{R}^n} f(x)$, 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