



# 262B-Lecture 7

Date created: 2021.02.09  
N. of Pages: 12

-  $\min f(x)$  where  $mI \preceq \nabla^2 f(x) \preceq MI \quad \forall x$

$\Rightarrow$  Linear convergence

-  $\min f(x)$  if  $m=0 \Rightarrow$  sublinear convergence

More precisely:  $O(1/k)$

Assumptions:  $\nabla f(x^{(k)})^T \Delta x^{(k)} \leq -c \|\nabla f(x^{(k)})\|^2$  ①

$\left( x^{(k+1)} = x^{(k)} + \underbrace{\alpha^{(k)}}_{\in [\varepsilon, (2-\varepsilon)\alpha^{(k)}]} \Delta x^{(k)} \right)$

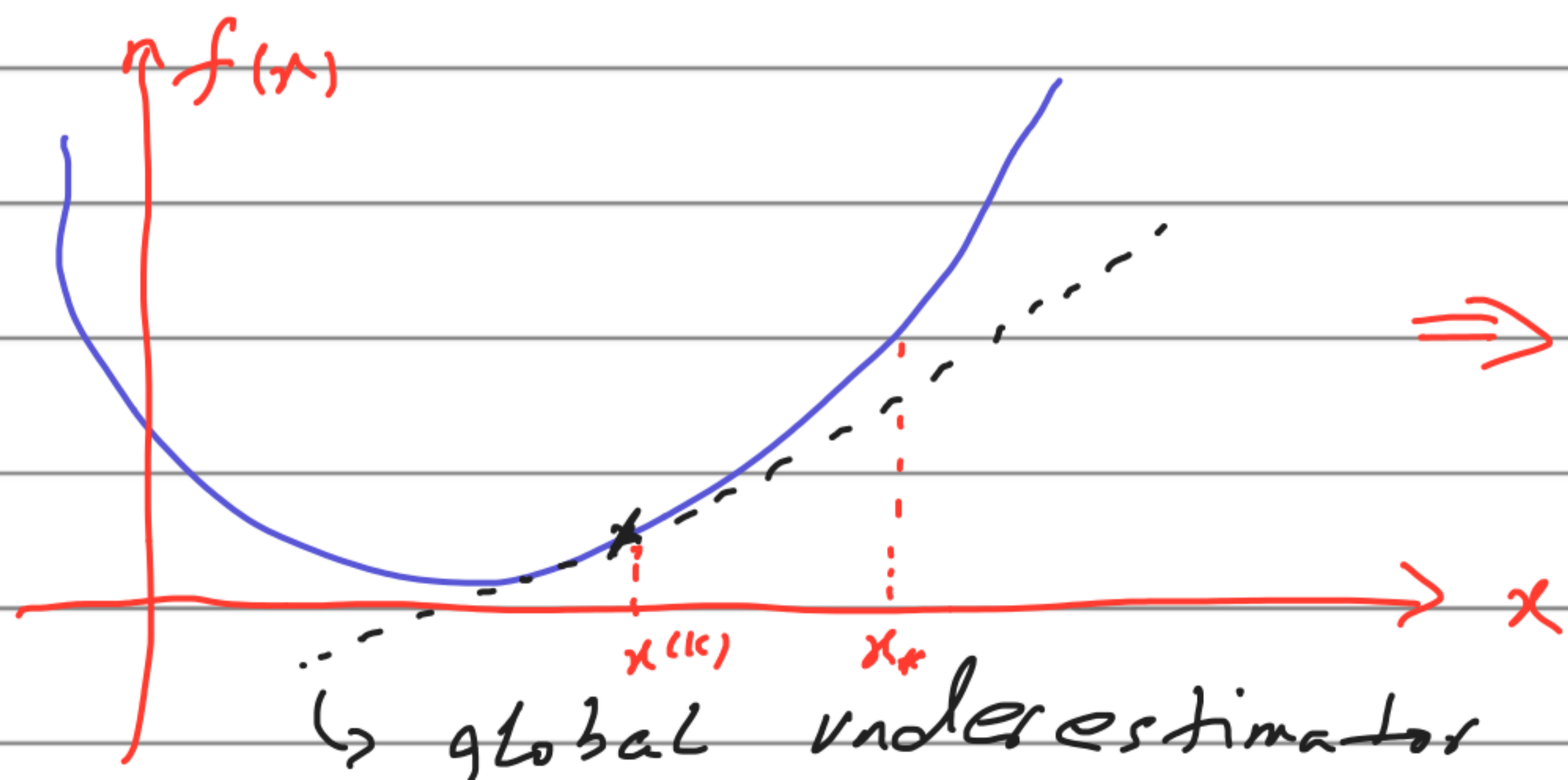
Define:  $X_*$  = set of all global minima

distance:  $d(x^{(k)}, X_*) = \min_{x_* \in X_*} \|x^{(k)} - x_*\|$

Proof: Based on a previous theorem, we know

that  $\|\nabla f(x^{(k)})\| \rightarrow 0$  as  $k \rightarrow \infty$

Now:  $\nabla^2 f(x) \succeq 0 \quad \forall x \Rightarrow$  Convex



$\Rightarrow f(x^{(k)}) - f(x_*) \leq \nabla f(x^{(k)})^T (x^{(k)} - x_*) \leq \|\nabla f(x^{(k)})\| \|x^{(k)} - x_*\|$  ②

minimize both sides of  $\textcircled{*}$  with respect to

$$x_* : f_* = \min f(x)$$

$$\Rightarrow f(x^{(k)}) - f_* \leq \|\nabla f(x^{(k)})\| \times d(x^{(k)}, x_*) \quad \textcircled{2}$$

$\textcircled{1}, \textcircled{2} \Rightarrow$  if error:  $f(x) - f_*$   
 $\textcircled{3}$

$$\text{then, } e(x^{(k+1)}) \leq e(x^{(k)}) - \frac{c\varepsilon^2 e(x^{(k)})^2}{2d(x^{(k)}, x_*)} \quad \forall k$$

proof for constant stepsize:

$$f(x^{(k+1)}) - f(x^{(k)}) \leq -\frac{\varepsilon^2}{2} |\nabla f(x^{(k)})^T \Delta x^{(k)}| \quad \textcircled{3}$$

$$\Rightarrow e(x^{(k+1)}) \leq e(x^{(k)}) \left( 1 - \underbrace{\frac{c\varepsilon^2 e(x^{(k)})}{2d(x^{(k)}, x_*)^2}}_{a^{(k)}} \right)$$

since  $e(x) \geq 0 \Rightarrow 1 - a^{(k)} \geq 0$

$$(1 - a^{(k)})^{-1} = 1 + a^{(k)} + (a^{(k)})^2 + (a^{(k)})^3 + \dots \geq 1 + a^{(k)}$$

$$\begin{aligned} \Rightarrow e(x^{(k+1)})^{-1} &\geq e(x^{(k)})^{-1} (1 - a^{(k)})^{-1} \\ &\geq e(x^{(k)})^{-1} (1 + a^{(k)}) \\ &= e(x^{(k)})^{-1} + \frac{c\varepsilon^2}{2d(x^{(k)}, X_*)^2} \end{aligned}$$

$\Rightarrow$  add them up, then take the inverse:

$$e(x^{(k)}) \leq \left( \frac{1}{e(x^{(0)})} + \frac{c\varepsilon^2}{2} \sum_{i=0}^{k-1} \frac{1}{d(x^{(i)}, X_*)^2} \right)^{-1}$$

$$\Rightarrow k e(x^{(k)}) \leq \left( \frac{1}{k e(x^{(0)})} + \frac{c\varepsilon^2}{2k} \sum_{i=0}^{k-1} \frac{1}{d(x^{(i)}, X_*)^2} \right)^{-1} \quad A$$

$$k \rightarrow \infty : e(x^{(k)}) \rightarrow 0 \Rightarrow d(x^{(i)}, X_*) \rightarrow 0$$

$$\Rightarrow \frac{1}{d(x^{(i)}, X_*)^2} \rightarrow \infty \Rightarrow \frac{1}{k} \sum_{i=0}^{k-1} \frac{1}{d(x^{(i)}, X_*)^2} \rightarrow \infty$$

$$\Rightarrow A = (\infty)^{-1} = 0$$

$$\Rightarrow \lim_{k \rightarrow \infty} k e(x^{(k)}) = 0 \Rightarrow e(x^{(k)}) = o\left(\frac{1}{k}\right)$$

Summary :	( $m > 0$ ) Strongly Convex	( $m \geq 0$ ) Convex
	$O\left(\log\left(\frac{1}{\epsilon}\right)\right)$	$O\left(\frac{1}{\epsilon}\right)$
	heavy ball	
	Nesterov	Can we improve it?

Acceleration :

~~$$x^{(k+1)} = x^{(k)} - \alpha^{(k)} \nabla f(x^{(k)})$$~~



$$x^{(k+1)} = x^{(k)} - \alpha^{(k)} \nabla f(x^{(k)}) + \beta^{(k)} \underbrace{(x^{(k)} - x^{(k-1)})}_{\text{momentum}}$$

If  $\alpha^{(k)} = \alpha$ ,  $\beta^{(k)} = \beta$

⇒ Heavy-ball method

Case study :  $\min_x \frac{1}{2} x^T Q x$  ( $m > 0$ )

Gradient alg :  $\frac{e(x^{(k+1)})}{e(x^{(k)})} \leq \frac{\text{c.d.}(\alpha) - 1}{\text{c.d.}(\alpha) + 1}$

accelerated gradient :  $\frac{e(x^{(k+1)})}{e(x^{(k)})} \leq \frac{\sqrt{\text{c.d.}(\alpha)} - 1}{\sqrt{\text{c.d.}(\alpha)} + 1}$  for some choices of  $\alpha$  and  $\beta$

Proof: Homework 2

Ex:  $\text{c.d.}(Q) = 10^4 \rightarrow$  ill-conditioned

$\Rightarrow$  extremely slow convergence for gradient alg.

But  $\sqrt{\text{c.d.}(Q)} = 10^2 \rightarrow$  much faster

$\Rightarrow$  This sort of acceleration is good

when  $\nabla^2 f(x_*) > 0$ .

previously:  $\left( \frac{\text{c.d.}(Q) - 1}{\text{c.d.}(Q) + 1} \right)^k \leq \varepsilon$

Now:  $\left( \frac{\sqrt{\text{c.d.}} - 1}{\sqrt{\text{c.d.}} + 1} \right)^k \leq \varepsilon$

$\Rightarrow$  # of iterations =  $O\left(\log \frac{1}{\varepsilon}\right)$

acceleration improves the constant of  $O$ .

Ex:  $\underbrace{1000} \times \log \frac{1}{\varepsilon} \Rightarrow \underbrace{20} \times \log \frac{1}{\varepsilon}$

improvement

Nesterov's acceleration method:

$$\cancel{x^{(k+1)} = x^{(k)} - \alpha^{(k)} \nabla f(x^{(k)})}$$



$$\begin{cases} y^{(k)} = x^{(k)} + \beta^{(k)} (x^{(k)} - x^{(k-1)}) & : \text{intermediate parameter} \\ x^{(k+1)} = y^{(k)} - \alpha \nabla f(y^{(k)}) & : \text{apply gradient alg. to } y^{(k)} \end{cases}$$

Assume:  $\beta^{(k)} \rightarrow 1$  as  $k \rightarrow \infty$

original proof:  $\beta^{(k)} = \frac{k-1}{k+2}$

special:  $f(x^{(k+1)}) \not\leq f(x^{(k)})$

⇒ Not a descent algorithm

Assume:  $0 \preceq \nabla^2 f(x) \preceq M I \Rightarrow m=0$

If  $\alpha \leq \frac{1}{m} \Rightarrow$

$$f(x^{(k)}) - f_* \leq \frac{2d(x^{(0)}, X_*)^2}{\alpha(k+1)^2}$$

⇒ Convergence rate =  $O(\frac{1}{k^2})$

If we want  $f(x^{(k)}) - f_* \leq \epsilon$

$$\Rightarrow \frac{2d(x^{(0)}, X_*)^2}{\alpha(k+1)^2} \leq \epsilon \Rightarrow \cdot$$

$$\Rightarrow \# \text{ of iterations} = O\left(\frac{1}{\sqrt{\epsilon}}\right)$$

Strongly convex  
( $m > 0$ )

Convex  
( $m \geq 0$ )

Gradient:  $O\left(\log \frac{1}{\epsilon}\right)$

↓ heavy ball

$O\left(\log \frac{1}{\epsilon}\right)$

but with a  
better constant

Gradient  $O\left(\frac{1}{\epsilon}\right)$

↓ Nesterov

$O\left(\frac{1}{\sqrt{\epsilon}}\right)$

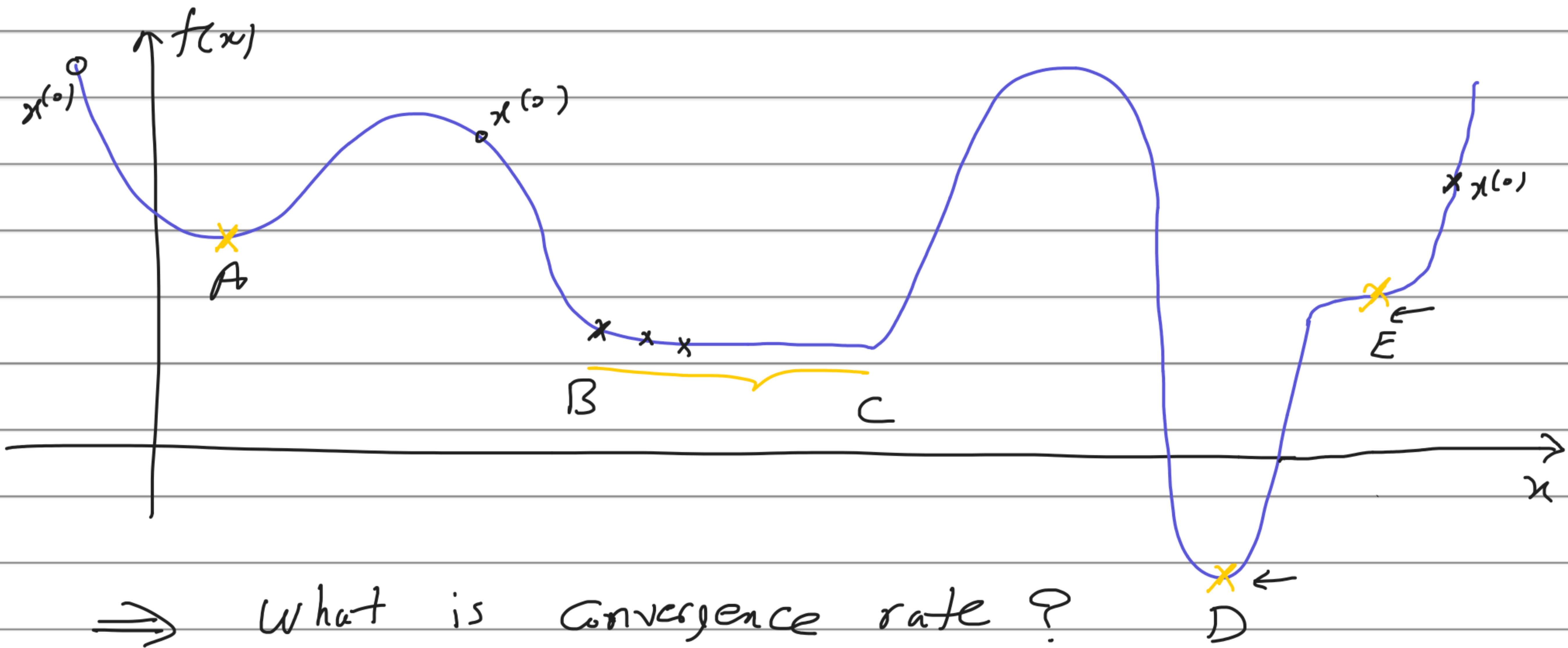
Ex:  $\epsilon = 10^{-4} \Rightarrow$  Compare  $10^4 = \frac{1}{\epsilon}$  to  $10^2 = \frac{1}{\sqrt{\epsilon}}$



Most of our analysis :  $\nabla^2 f(x) \succeq \epsilon I$   
 $\forall x$   
 $\Downarrow$   
Conver case

Generalize to the non-conver case using  
the idea of sub-level set if  $x^{(0)}$  is  
close to a solution  $x_*$ .

Question : what if  $f(x)$  is highly non-con  
 $\forall \epsilon x$   
and  $x^{(0)}$  is random (far away from  $x_*$ ) ?



$$\min f(x) \longrightarrow \nabla f(x) = 0 \implies x_* : \text{FOC}$$

Stationary point

Convergence: get to a neighborhood of  
a stationary point?

$$\text{Gradient alg: } x^{(k+1)} = x^{(k)} - \alpha \nabla f(x^{(k)})$$

$$\text{Assume: Constant stepsize } \alpha \in \left[ \frac{\bar{\epsilon}}{L}, \frac{2-\bar{\epsilon}}{L} \right]$$

$\bar{\epsilon} > 0$

$$\implies \text{stopping criterion: } \|\nabla f(x^{(k)})\| \leq \epsilon$$

previous theorem:  $\{x^{(k)}\}$ : limit points are stationary points.

$\Downarrow$  Proof

$$f(x^{(k+1)}) - f(x^{(k)}) \leq -\frac{\bar{\epsilon}^2}{2} \|\nabla f(x^{(k)})\|^2$$

$\Downarrow$  add this up over different values of  $k$ .

$$\implies \left( f(x^{(k+1)}) - f(x^{(0)}) \right) \leq -\frac{\bar{\epsilon}^2}{2} \sum_{i=0}^k \|\nabla f(x^{(i)})\|^2$$

$\geq \underbrace{\min f(x)}_{\text{global min}} - f(x^{(0)})$

①

$$\frac{\bar{\epsilon}^2}{2} \sum_{i=0}^k \|\nabla f(x^{(i)})\|^2 \geq \frac{\bar{\epsilon}^2}{2} x(k+1)$$

$$\times \min_{0 \leq i \leq k} \|\nabla f(x^{(i)})\|^2 \quad (2)$$

If we want to make sure that

$$\min_{0 \leq i \leq k} \|\nabla f(x^{(i)})\|^2 \leq \epsilon^2 \quad (3)$$

(1), (2), (3)  $\Rightarrow$

$$\min_{0 \leq i \leq k} \|\nabla f(x^{(i)})\|^2 \leq \frac{f(x^{(0)}) - \underbrace{(\min f(x))}_{\text{global min}}}{\frac{\bar{\epsilon}^2 (k+1)}{2}} \leq \epsilon^2$$

$\Rightarrow$  enough to have:

$$\# \text{ of iterations : } k = O\left(\frac{1}{\epsilon^2}\right)$$

$$\frac{f(x^{(0)}) - \min f(x)}{\bar{\epsilon}^2 (k+1) \times \frac{1}{2}} \leq \epsilon^2 \Rightarrow k+1 \sim \frac{1}{\epsilon^2}$$

strongly  
convex ( $m > 0$ )

convex  
( $m \geq 0$ )

non-convex

Gradient:

$$O\left(\log \frac{1}{\epsilon}\right)$$

heavy ball

$$O\left(\log \frac{1}{\epsilon}\right)$$

better constant

Gradient:

$$O\left(\frac{1}{\epsilon}\right)$$

nestrov

$$O\left(\frac{1}{\sqrt{\epsilon}}\right)$$

Gradient:

$$O\left(\frac{1}{\epsilon^2}\right)$$

?

(state-of-the-art)

first-order methods: (we use gradient to generate the points)

$$\text{Gradient: } x^{(k+1)} = x^{(k)} - \alpha^{(k)} \nabla f(x^{(k)})$$

↓

scaled

$$\text{gradient: } x^{(k+1)} = x^{(k)} - \alpha^{(k)} \underbrace{D^{(k)}}_{\approx} \nabla f(x^{(k)})$$

Assume  $D^{(k)} = D$ ,  $f(x) = \frac{1}{2} x^T Q x$

$$\Rightarrow x^{(k+1)} = x^{(k)} - \alpha^{(k)} D Q x^{(k)} \quad \times D^{-1/2}$$

Define:  $y^{(k)} = D^{-1/2} x^{(k)}$

$$\Rightarrow y^{(k+1)} = y^{(k)} - \alpha^{(k)} D^{-1/2} D Q D^{1/2} y^{(k)}$$

$$\Rightarrow y^{(k+1)} = y^{(k)} - \alpha^{(k)} (D^{1/2} Q D^{1/2}) y^{(k)}$$

$\Downarrow$

gradient alg for  $\min \frac{1}{2} y^T (D^{1/2} Q D^{1/2}) y$

$$\Rightarrow \frac{\|y^{(k+1)}\|}{\|y^{(k)}\|} \leq \frac{\text{c.d.}(D^{1/2} Q D^{1/2}) - 1}{\text{c.d.}(D^{1/2} Q D^{1/2}) + 1} \quad (\text{optimal step size})$$

$$\Rightarrow \frac{\sqrt{(x^{(k+1)})^T D^{-1} x^{(k+1)}}}{\sqrt{(x^{(k)})^T D x^{(k)}}} \leq \frac{\text{c.d.}(D^{1/2} Q D^{1/2}) - 1}{\text{c.d.}(D^{1/2} Q D^{1/2}) + 1}$$

error in weighted norm