



# 262B-Lecture 2

Date created: 2021.01.21  
N. of Pages: 15

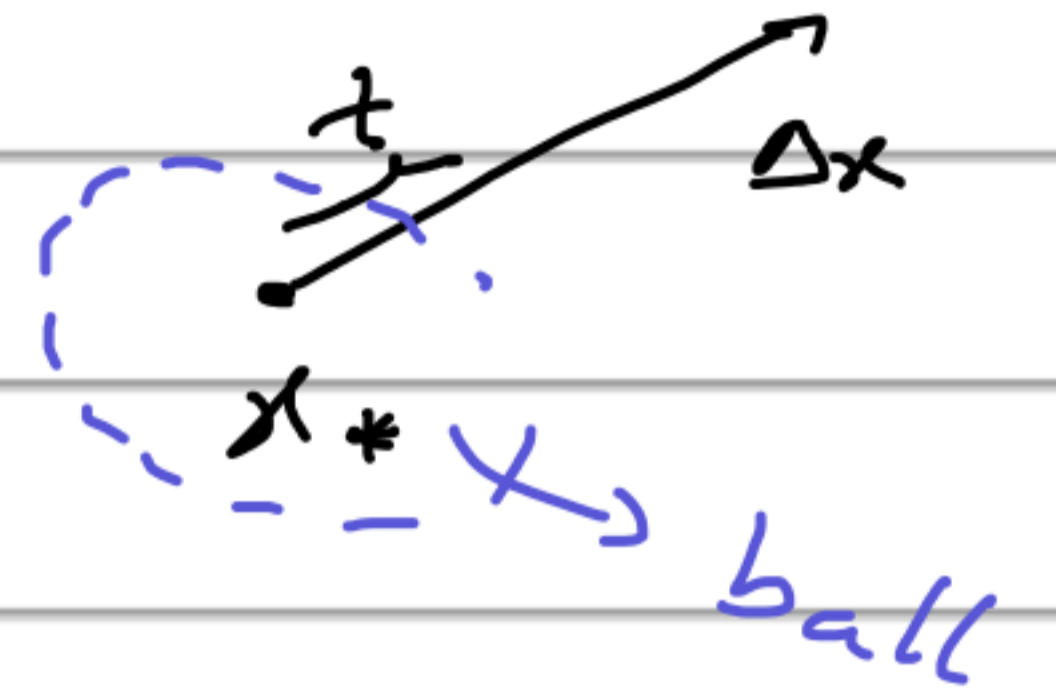
$$\min f(x) \longrightarrow x_* : \text{local min} \longrightarrow \nabla f(x_*) = 0$$

FOC

$$\text{Proof } \underline{1} \longrightarrow \text{Proof } \underline{2} \longrightarrow \text{Proof } \underline{3}$$

$$g(t) = f(x_* + t \Delta x)$$

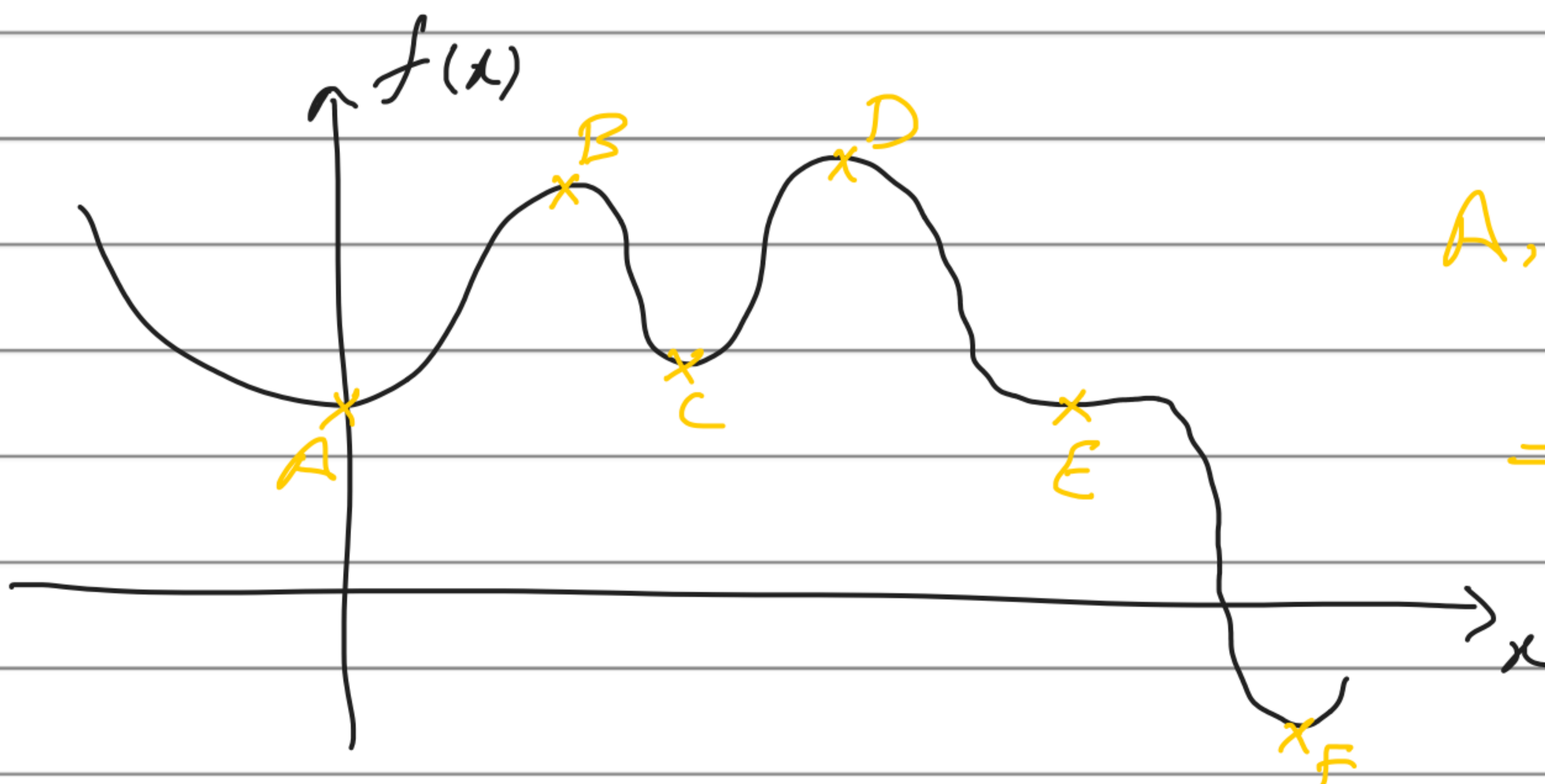
$$g'(0) = \lim_{\substack{t \rightarrow 0 \\ t \downarrow 0}} \frac{f(x_* + t \Delta x) - f(x_*)}{t} \geq 0$$



$$\Rightarrow 0 \leq g'(0) = \nabla f(x_*)^T \Delta x \quad \forall \Delta x$$

$$\Delta x = -\nabla f(x_*) \Rightarrow 0 \leq -\|\nabla f(x_*)\|^2$$

$$\Rightarrow \|\nabla f(x_*)\| = 0 \Rightarrow \nabla f(x_*) = 0$$



A, B, C, D, E, F

$$\Rightarrow \nabla f(x_*) = 0$$

SOC (necessary): second-order necessary condition

Thm: If  $x_*$  is a local min  $\Rightarrow \nabla^2 f(x_*) \succeq_0$   
↙  
PSD

$A \succeq_0$  PSD (positive semi-definite)

$A \succ_0$  PD

$A \preceq_0$  NSD (eigs!)

$A \prec_0$  ND

proof: (by contradiction)

If  $\nabla^2 f(x_*) \not\succeq_0 \Rightarrow$

$\exists \Delta x : \Delta x^T \nabla^2 f(x_*) \Delta x < 0$

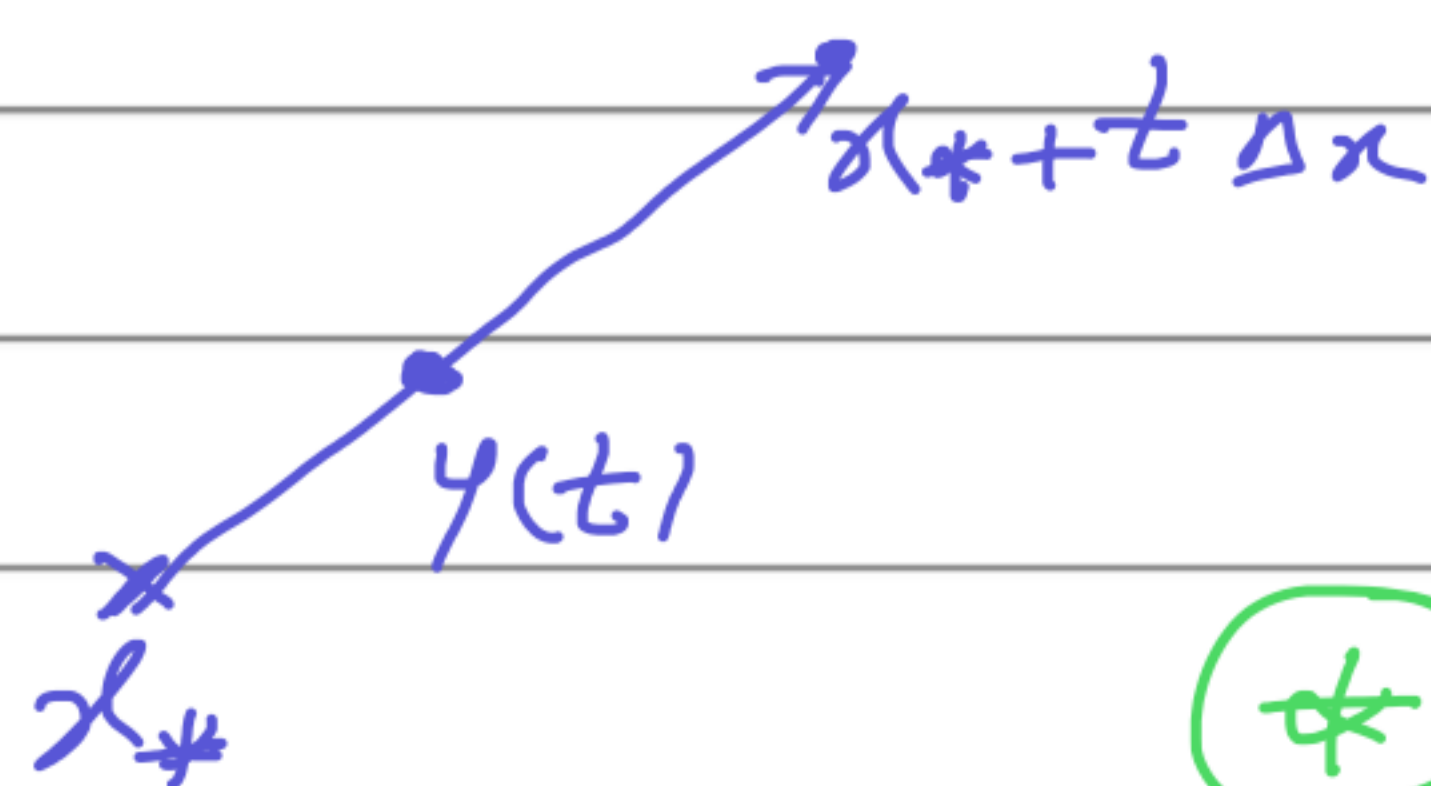
$$f(\underbrace{x_*}_{\text{nominal}} + \underbrace{t \Delta x}_{\text{perturbation}}) = f(x_*) + t \nabla f(x_*)^T \Delta x$$

$$+ \frac{1}{2} \Delta x^T \nabla^2 f(y(t)) \Delta x$$

new point

depends on  $\underline{t}$

$t \rightarrow 0 : y(t) \rightarrow x_* \quad \text{so} \quad \nabla^2 f(y(t)) \rightarrow \nabla^2 f(x_*)$



$\circledast$

$\Rightarrow$

$$f(x_* + t \Delta x) < f(x_*) < 0$$

new point

local min

$\times$

✓  
 $\min x^2$

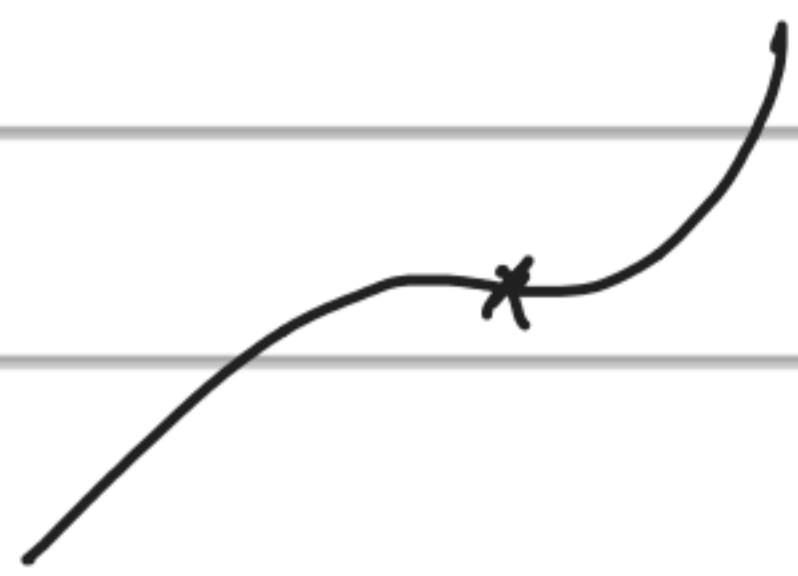


$$f'(0) = 0$$

$$f''(0) \geq 0$$

$$x_* = 0$$

✓  
 $\min x^3$

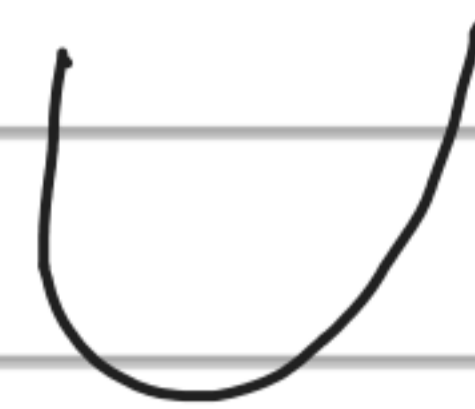


$$f'(0) = 0$$

$$f''(0) \geq 0$$

satisfies

✓  
 $\min x^4$



$$f'(0) = 0$$

$$f''(0) \geq 0$$

FOC, SOC

SOC (sufficient) : second-order sufficient condition

Thm: If  $x_*$  s.t.  $\nabla f(x_*) = 0$ ,  $\nabla^2 f(x_*) > 0$

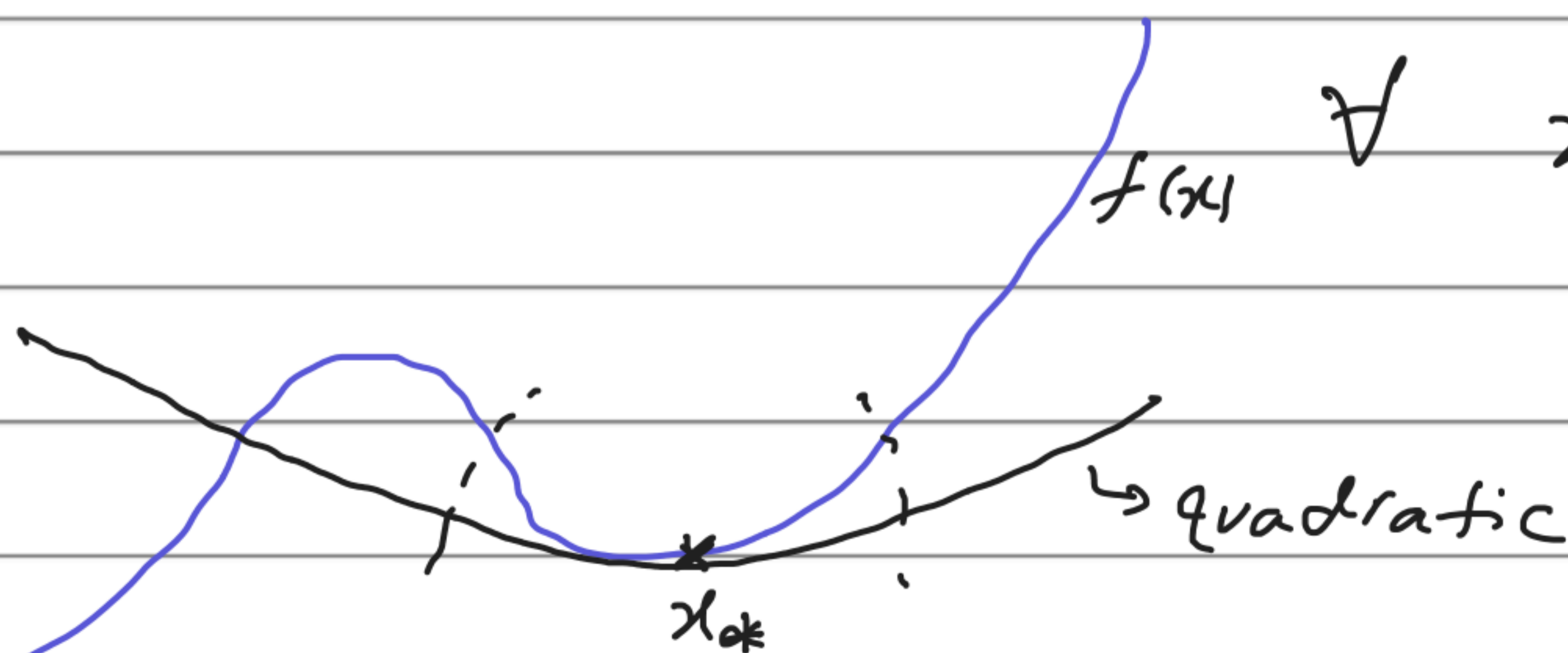
then: 1 -  $x_*$  is a strict local min

2 -  $\exists \epsilon > 0, \mu > 0$  s.t.

$$f(x) \geq f(x_*) + \mu \|x - x_*\|^2$$

quadratic in  $x$

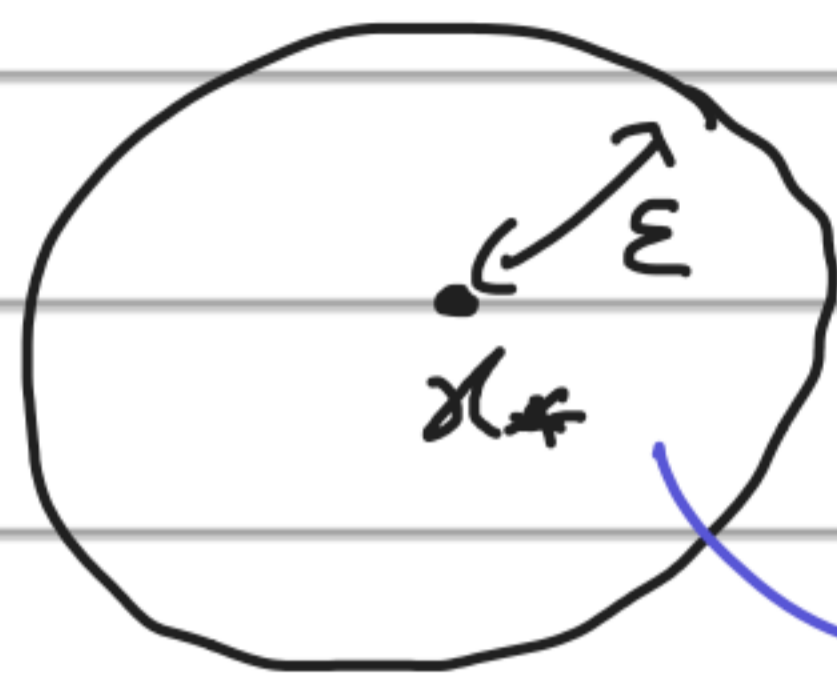
$\forall x \in B$  centered at  $x_*$  of radius  $\epsilon$



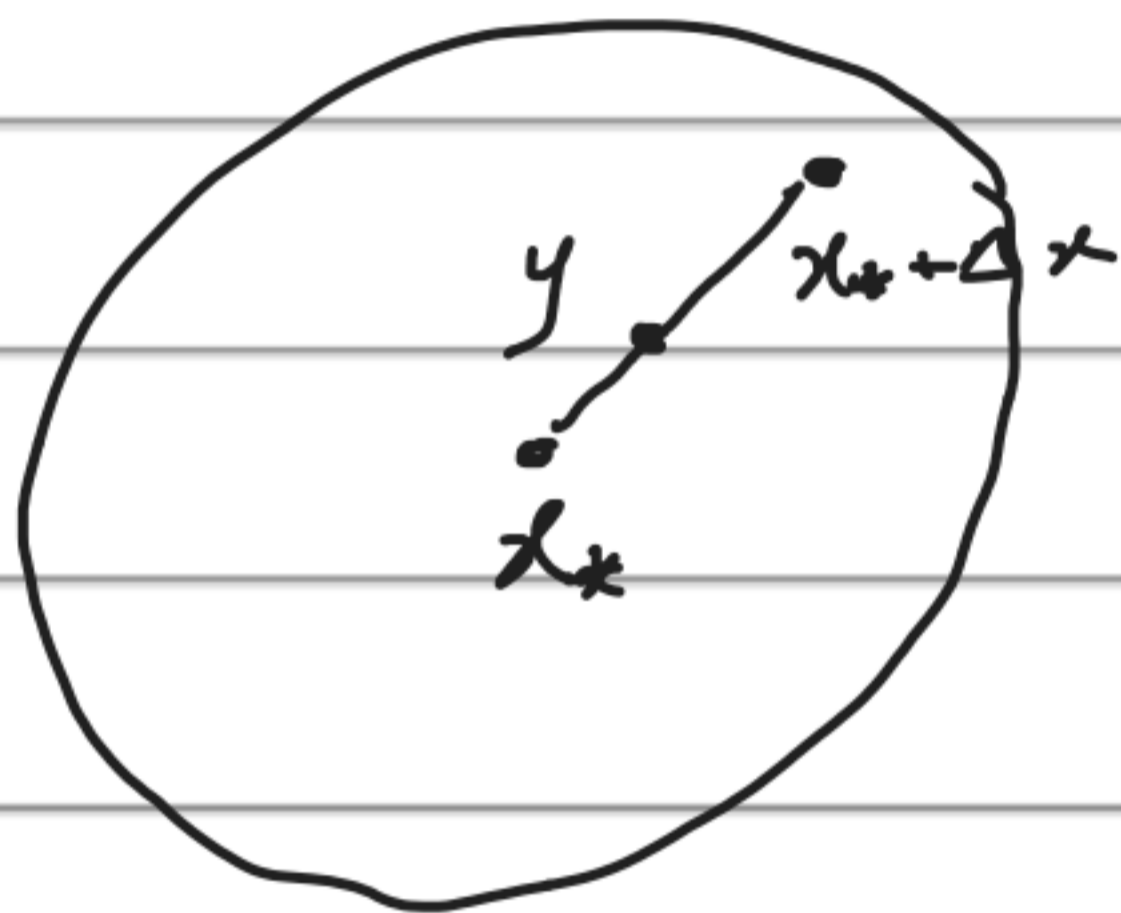
Proof:  $\nabla^2 f(x_*) > 0 \Rightarrow \exists \epsilon, \lambda > 0$

s.t.  $\min \text{eig}(\nabla^2 f(x)) \geq \lambda \quad \forall x: \|x - x_*\| \leq \epsilon$

$x_*$   
 ↙  
 Hessian is PD



Hessian at each point inside the ball is PD.



$$f(x_*) + \nabla f(x_*) \Delta x + \frac{1}{2} \Delta x^T \nabla^2 f(y) \Delta x \geq f(x_*) + \frac{\lambda}{2} \|\Delta x\|^2$$

$$\Delta x^T A \Delta x \geq \min \text{eig}(A) \|\Delta x\|^2$$

$$\Rightarrow f(x) \geq f(x_*) + \frac{\lambda}{2} \|x - x_*\|^2$$

$$\min f(x) \rightarrow \underbrace{\nabla f(x_*) = 0}_{\text{FOC}}, \quad \underbrace{\nabla^2 f(x_*) \succeq 0}_{\text{SOC}}$$

$\Rightarrow$  How to find such point  $x_*$  ?

Descent methods: Can't find  $x_*$  in one shot.

$\Rightarrow$  Find it gradually

$$x^{(0)} \rightarrow x^{(1)} \rightarrow x^{(2)} \rightarrow \dots \rightarrow x^{(k)} \rightarrow x^{(k+1)} \rightarrow \dots$$

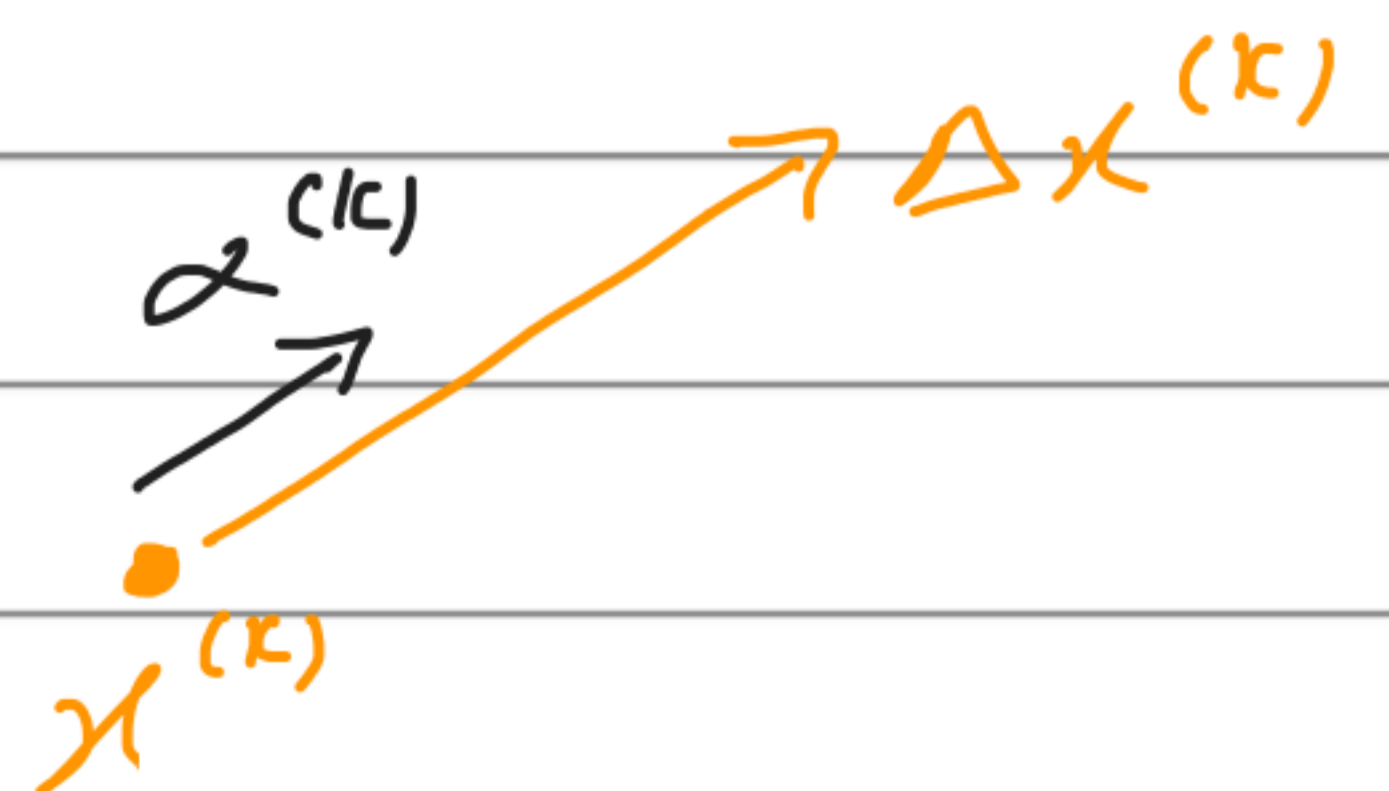
(generate a sequence of point)

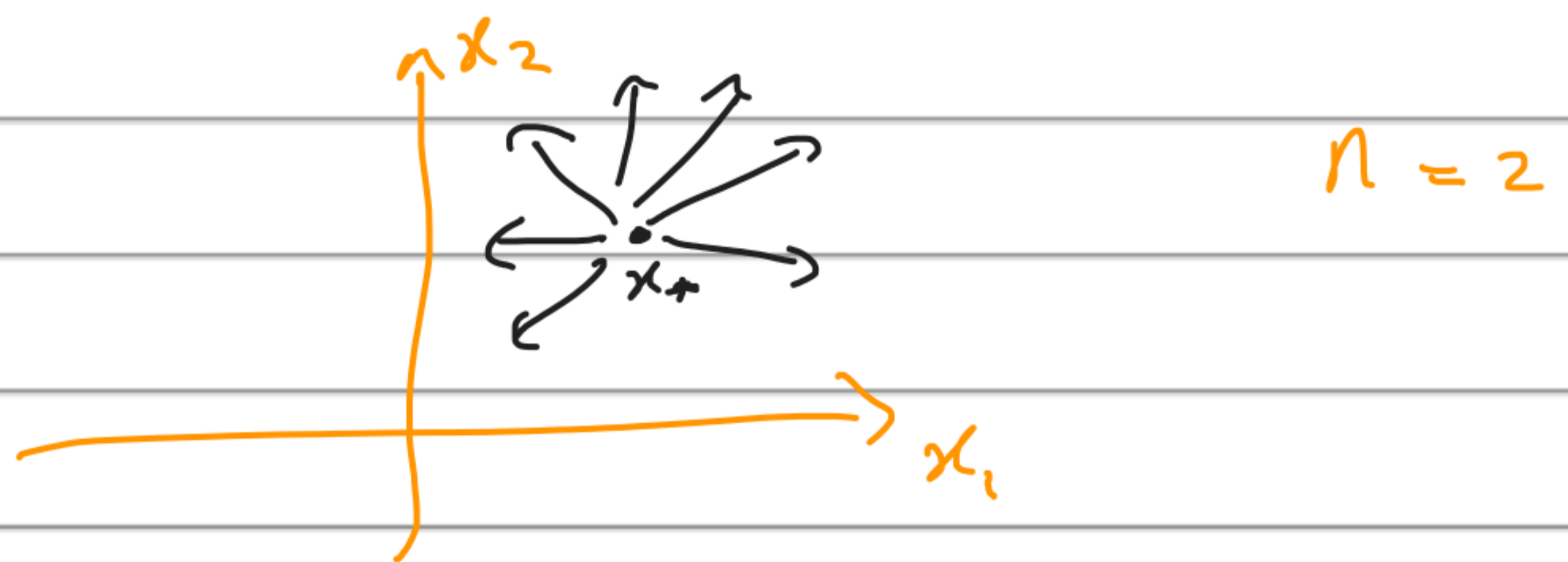
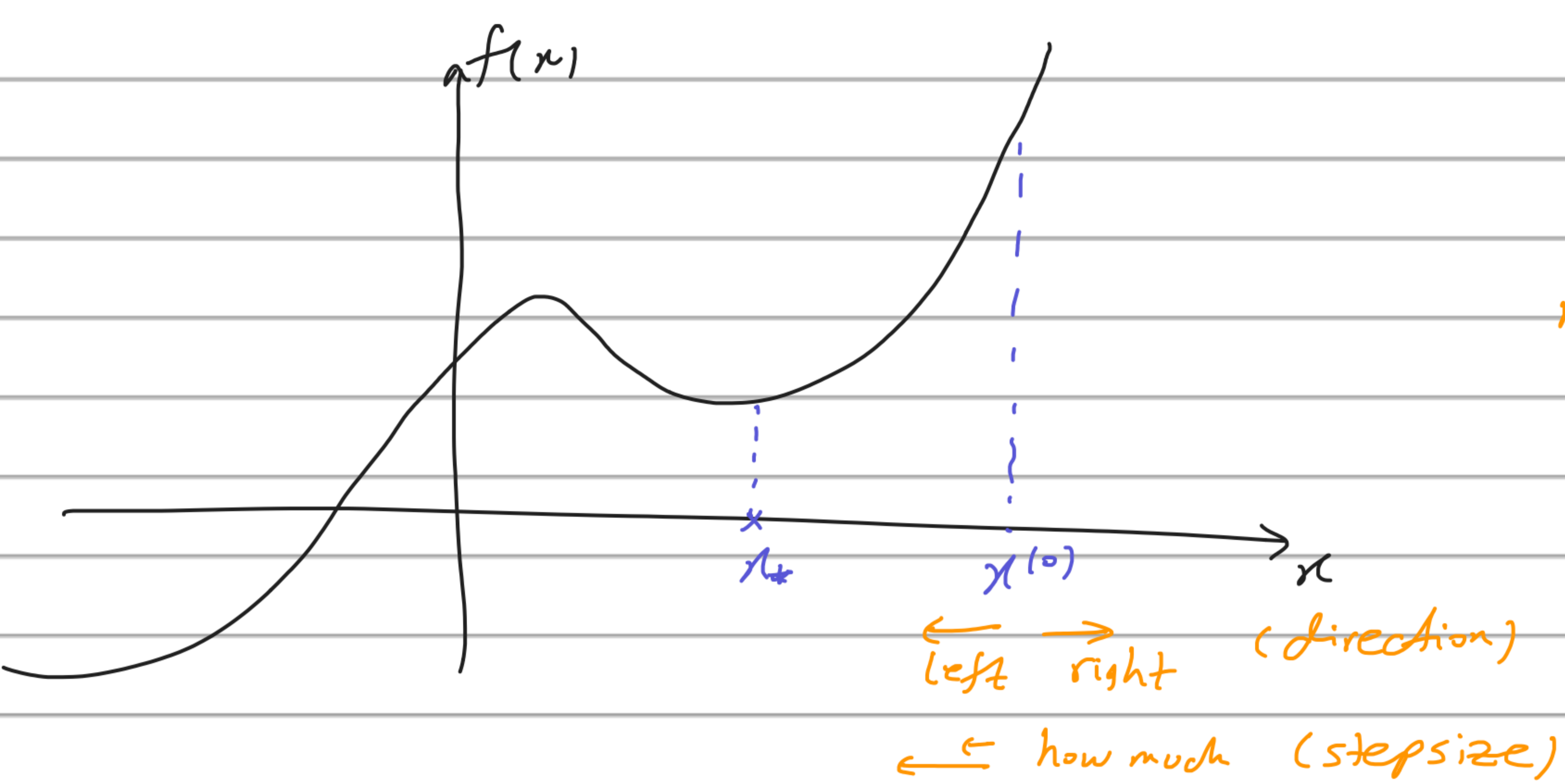
improvement at each iteration:

$$f(x^{(0)}) > f(x^{(1)}) > f(x^{(2)}) > \dots$$

$$x^{(k)} \rightarrow x^{(k+1)} \quad \text{updating rule:}$$

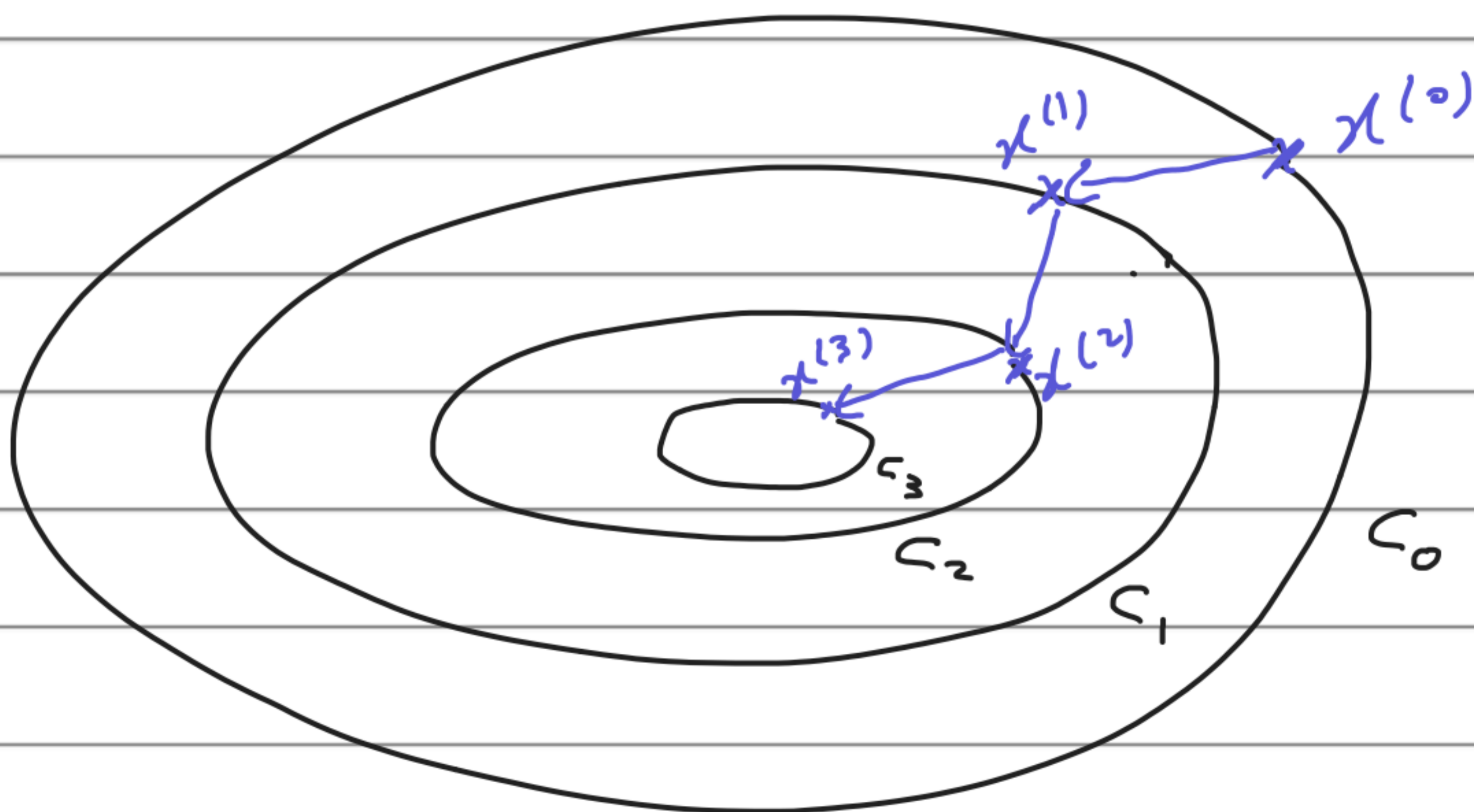
$$\underbrace{x^{(k+1)}}_{\text{new point}} = \underbrace{x^{(k)}}_{\text{old point}} + \underbrace{\alpha^{(k)}}_{\text{stepsize}} \underbrace{\Delta x^{(k)}}_{\text{direction}}$$





Level set  $(c) = \{ x \mid f(x) = c \}$   
 $\downarrow$   
 $c \in \mathbb{R}$

$n = 2$



$$\underbrace{x^{(k+1)}}_{\text{new}} = \underbrace{x^{(k)}}_{\text{old}} + \underbrace{\alpha^{(k)} \Delta x^{(k)}}_{\text{perturbation}} \alpha$$

$$f(x^{(k+1)}) = f(x^{(k)}) + \alpha^{(k)} \nabla f(x^{(k)})^T \Delta x^{(k)} + \dots$$

$\underbrace{\hspace{10em}}_{\text{first term}} \xrightarrow{\text{dominates}} \underbrace{\hspace{2em}}_{\text{rest}}$

improvement:  $f(x^{(k+1)}) < f(x^{(k)})$

Thm: 1 - If  $\nabla f(x^{(k)})^T \Delta x^{(k)} > 0$ , then

$$\exists \tau > 0 \text{ s.t.}$$

$$f(x^{(k+1)}) > f(x^{(k)}) \quad \forall \alpha^{(k)} \in [0, \tau]$$

2 - If  $\nabla f(x^{(k)})^T \Delta x^{(k)} < 0$ , then

$$\exists \tau > 0 \text{ s.t.}$$

$$f(x^{(k+1)}) < f(x^{(k)}) \quad \forall \alpha^{(k)} \in [0, \underline{\tau}]$$

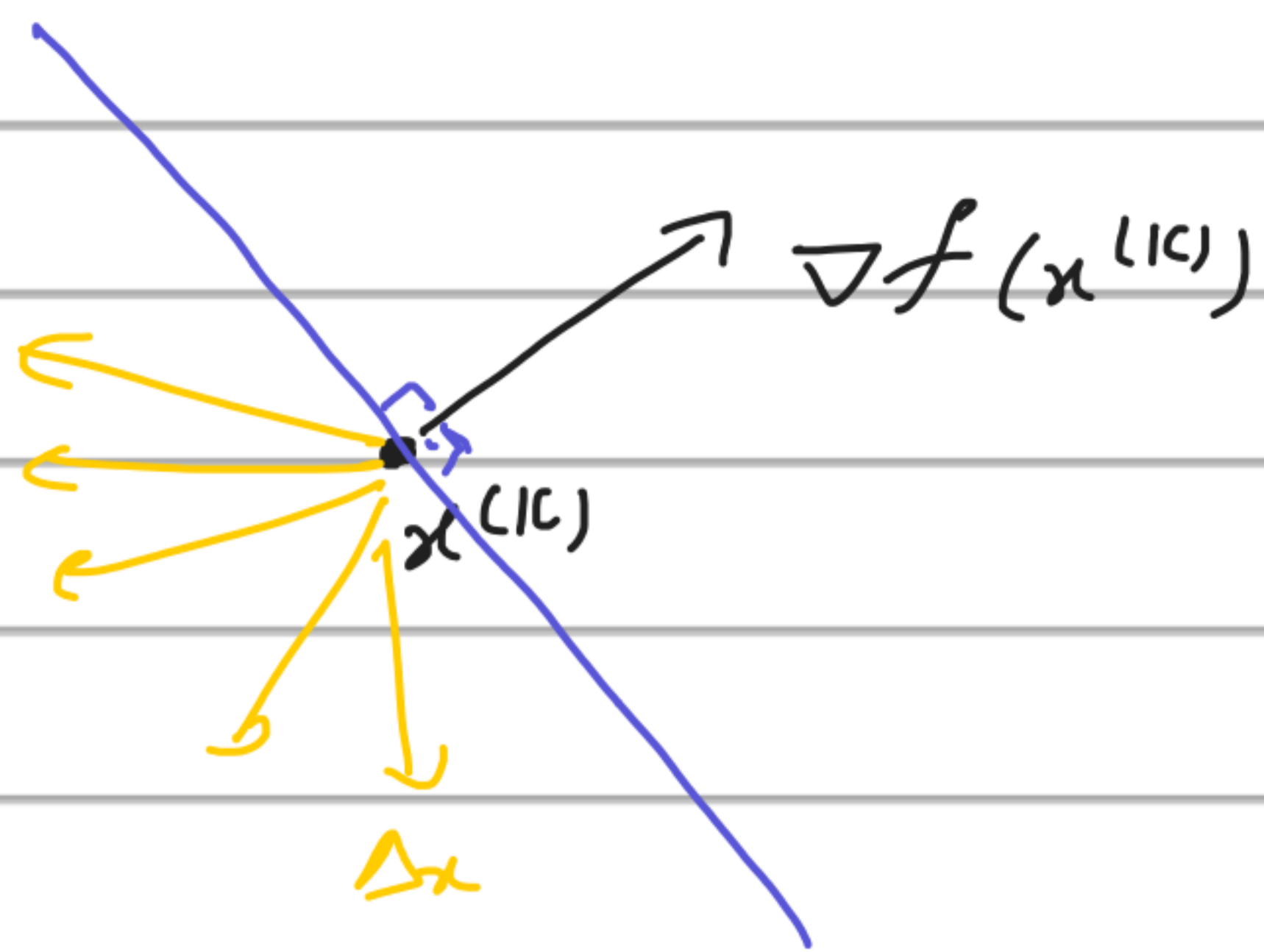
3 - If  $\nabla f(x^{(k)})^T \Delta x^{(k)} = 0 \Rightarrow$

$$\text{look at sign } (\Delta x^{(k)})^T \nabla^2 f(x^{(k)}) \Delta x^{(k)}$$

Descent direction:  $\Delta x$  is called descent

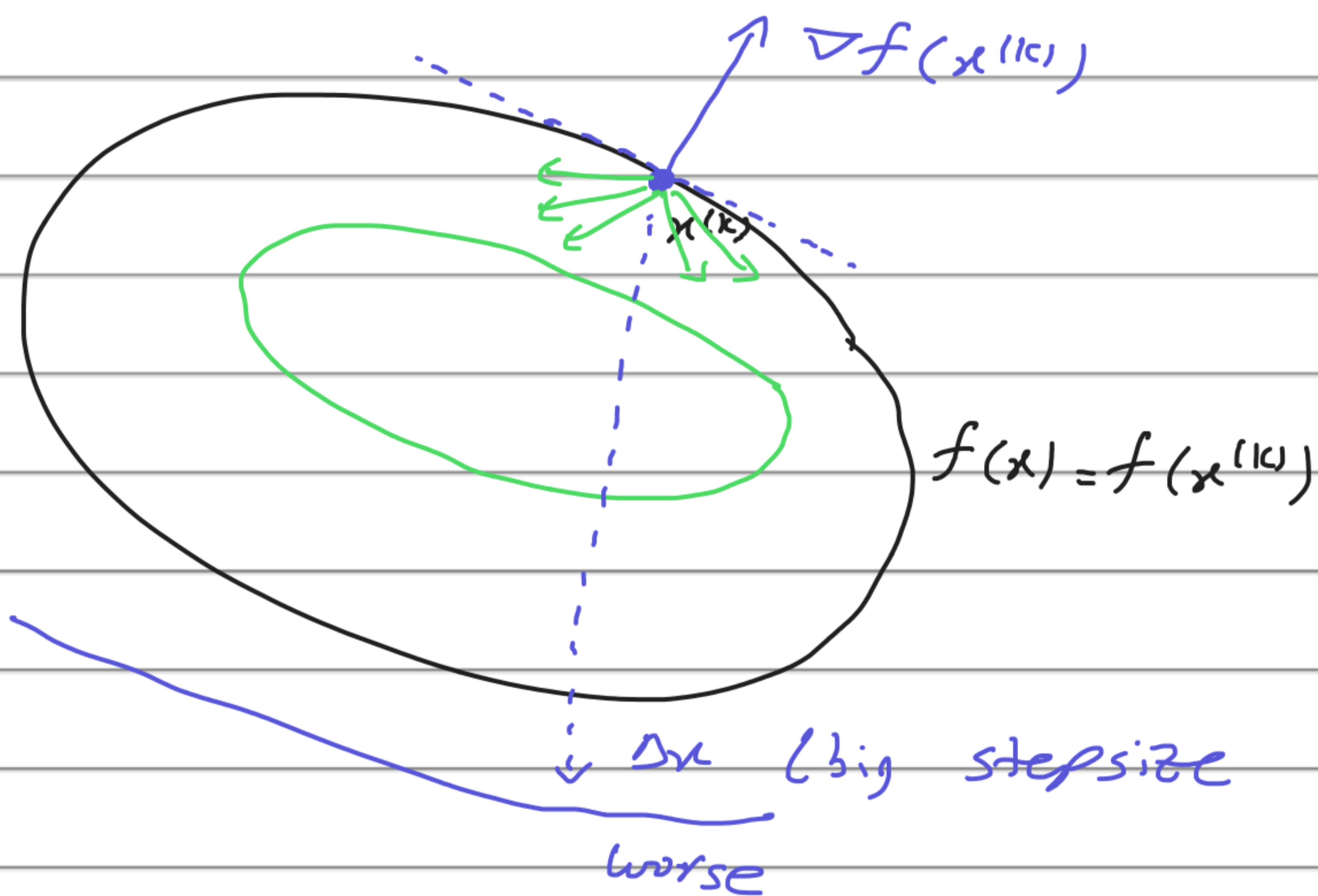
if  $\nabla f(x^{(k)})^T \Delta x < 0$





angle between  
direction & gradient

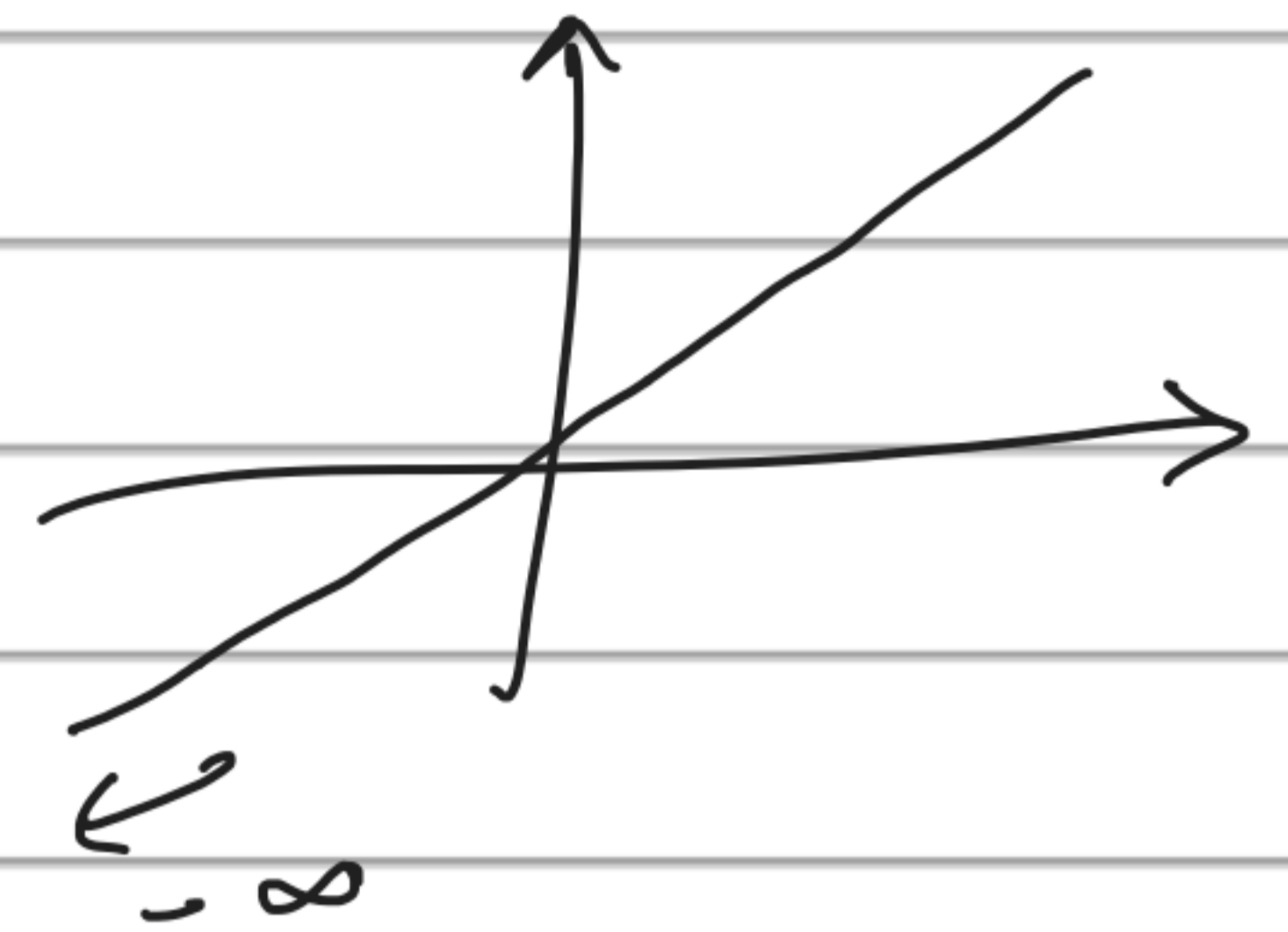
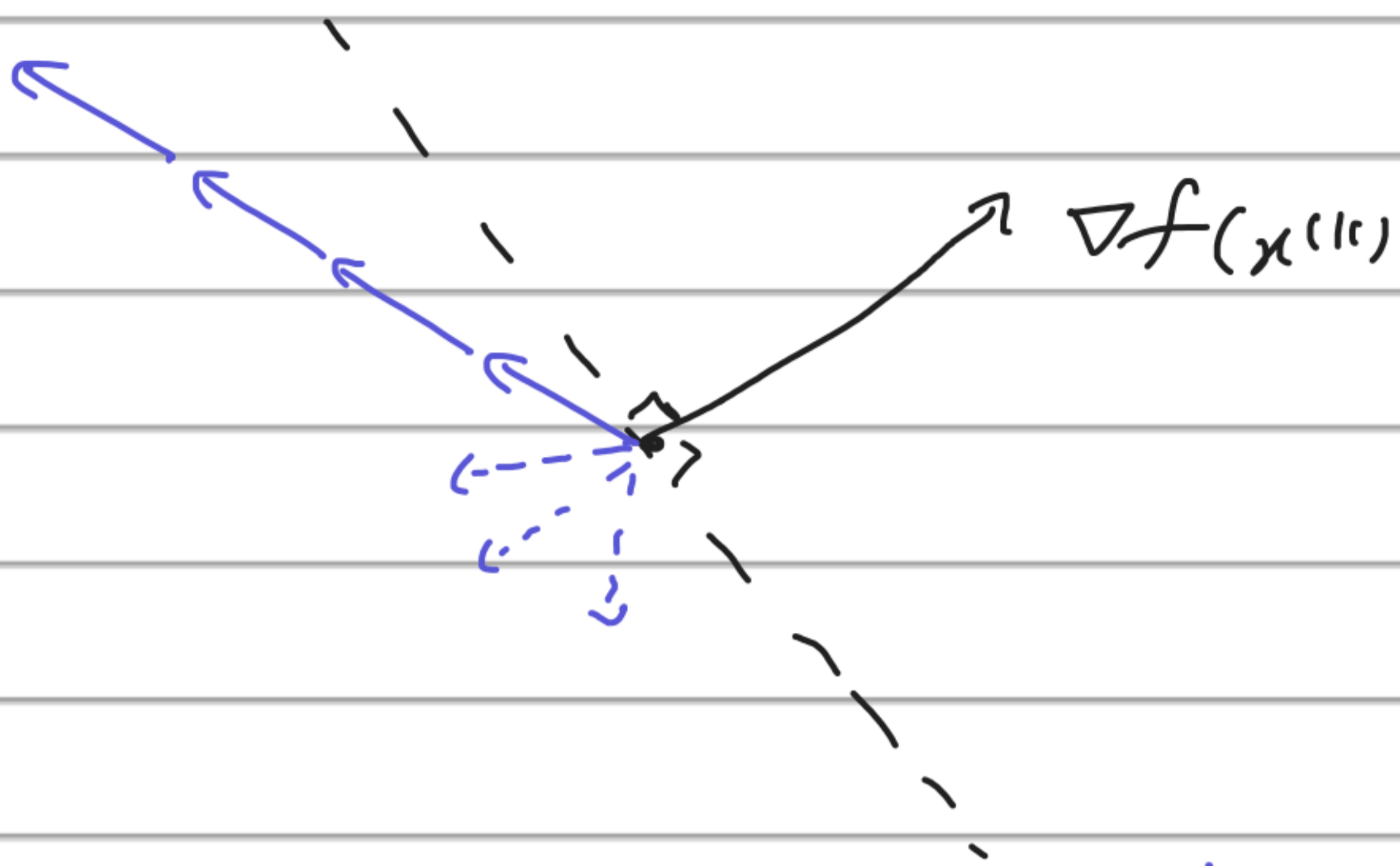
$\Rightarrow$  There are infinitely many descent directions



Design the best direction?

$$\underbrace{f(x^{(k+1)})}_{\text{reduce}} = f(x^{(k)}) + \underbrace{\alpha^{(k)} \nabla f(x^{(k)})^T \Delta x^{(k)}}_{\text{reduce}} + \dots$$

Best direction :  $\min_{\Delta x} \underbrace{\nabla f(x^{(k)})^T \Delta x}_{\text{Linear}}$



shouldn't confuse direction

with stepsize

$$\Rightarrow \min_{\Delta x} \underbrace{\nabla f(x^{(k)})^T \Delta x}_{\text{with stepsize}}$$

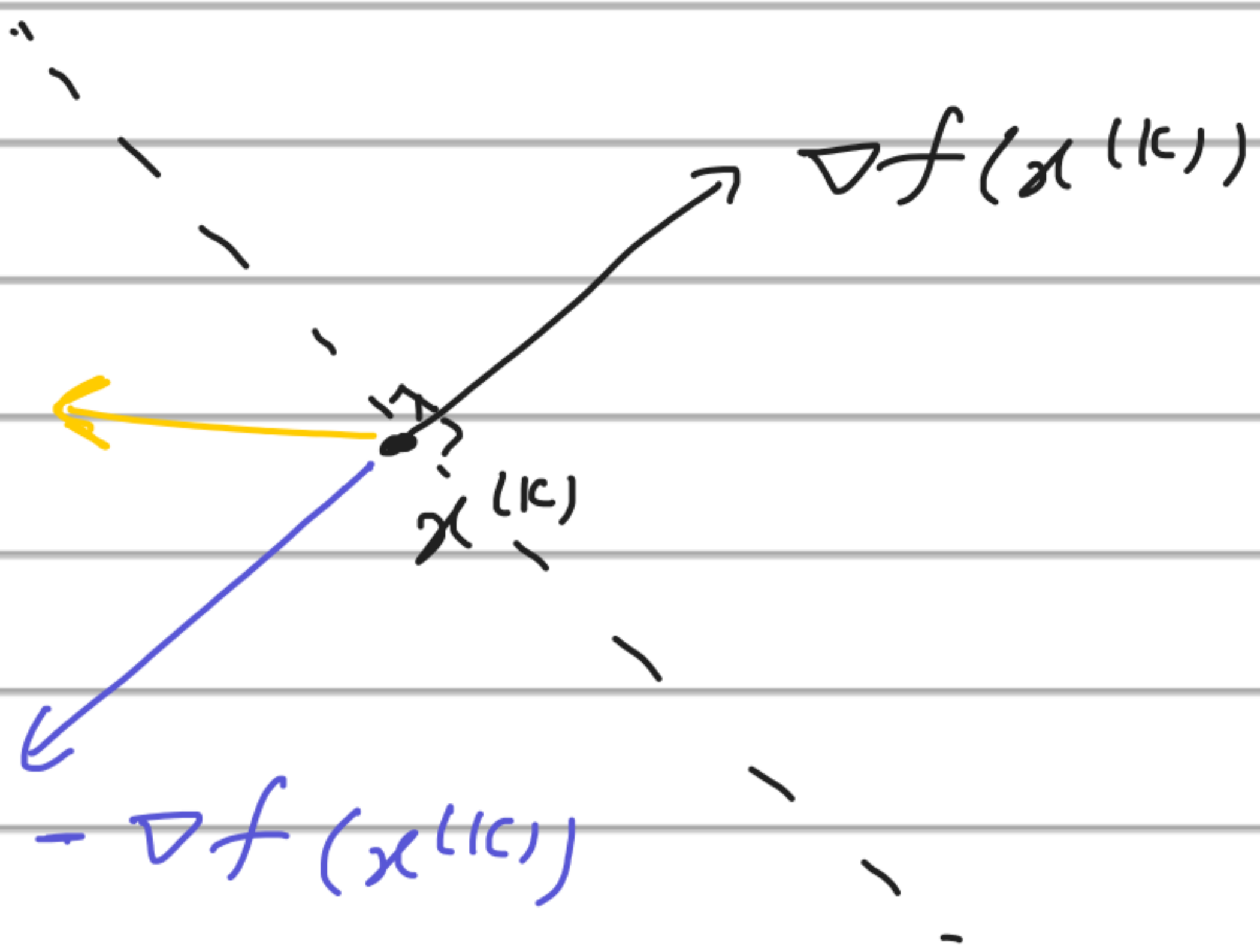
s.t.  $\|\Delta x\| \leq 1$

$$\|x\|_1 = \sum |x_i|, \quad \|x\|_2 = \left( \sum x_i^2 \right)^{1/2}$$

$$\|x\|_\infty = \max_i (|x_i|), \dots$$

Standard norm : length,  $\|\Delta x\|_2$

$$\Rightarrow \Delta x = - \frac{\nabla f(x^{(k)})}{\|\nabla f(x^{(k)})\|} \xrightarrow{\text{proportional}} - \nabla f(x^{(k)})$$



steepest descent (w.r.t.  $\|\cdot\|_2$ )

$$\Delta x^{(k)} = -\nabla f(x^{(k)}) \quad (\text{Gradient algorithm})$$

generalize

$$\Delta x^{(k)} = -\underbrace{D^{(k)}}_{\text{PD}} \nabla f(x^{(k)})$$

Thm: If  $D^{(k)} >_0 \Rightarrow \Delta x^{(k)}$ : descent direction

$$\nabla f(x^{(k)})^T \Delta x^{(k)} = - \underbrace{\nabla f(x^{(k)})^T}_{\text{vector}} \underbrace{D^{(k)}}_{\substack{\text{matrix} \\ >_0}} \underbrace{\nabla f(x^{(k)})}_{\text{vector}}$$

$$\leq 0$$

$$D^{(k)} = \nabla^2 f(x^{(k)})^{-1} \quad \text{if } \nabla^2 f(x^{(k)}) > 0$$

$\Rightarrow$  Newton's method

old  $\checkmark$   $x^{(k)}$   $\rightarrow$  new?  $x^{(k+1)} = x^{(k)} + \alpha^{(k)} \Delta x^{(k)}$   $\rightarrow$  direction?

$$f(x^{(k+1)}) = f(x^{(k)}) + \nabla f(x^{(k)}) \Delta x^{(k)} + \frac{1}{2} (\Delta x^{(k)})^T \nabla^2 f(x^{(k)}) \Delta x^{(k)} + \dots$$

$$f(x^{(k+1)}) < f(x^{(k)})$$

minimize

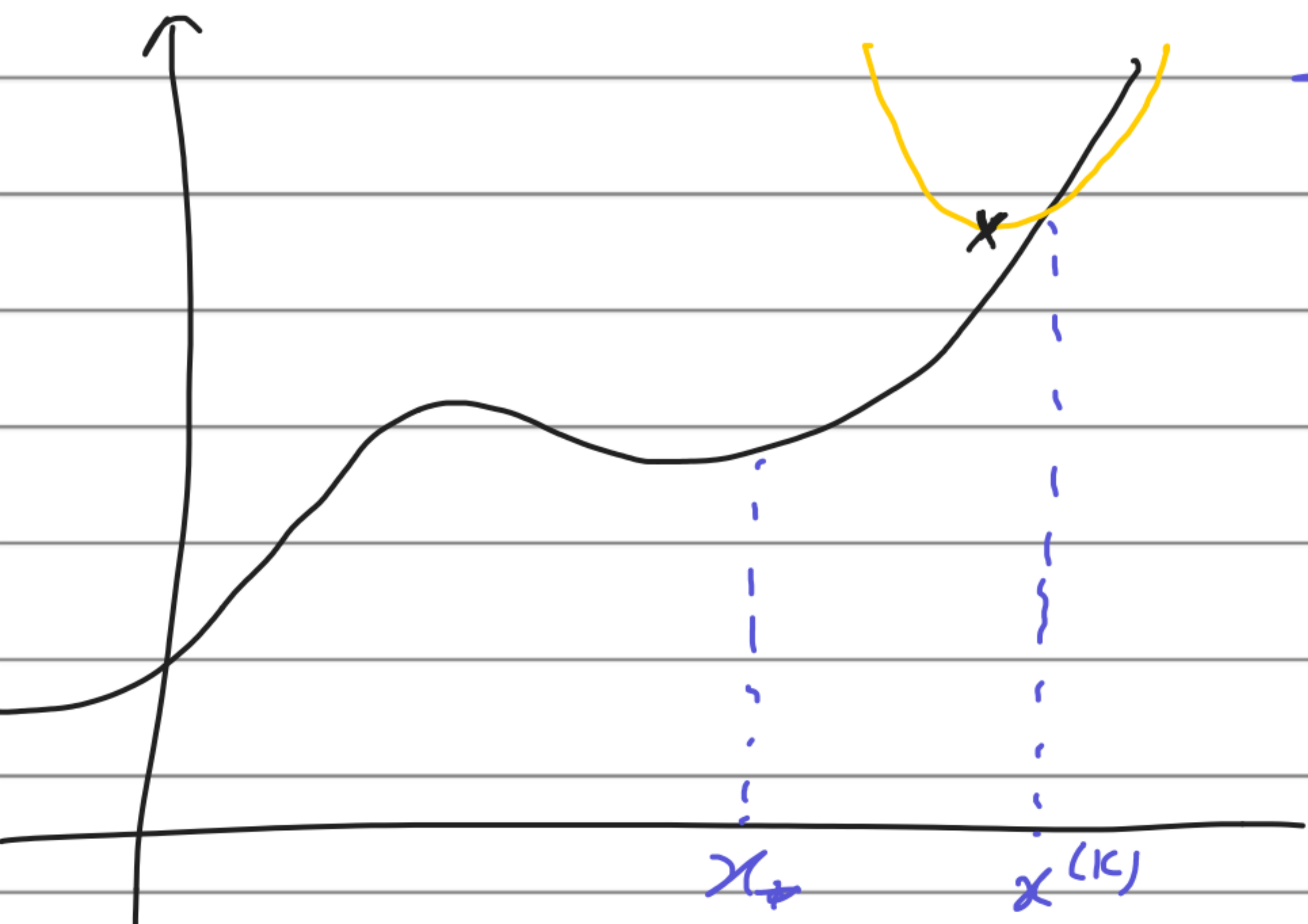
$$\min_{\Delta x^{(k)}} f(x^{(k)} + \Delta x^{(k)}) = \min_x f(x)$$

new point

approximate

$$\min_{\Delta x^{(k)}} f(x^{(k)}) + \nabla f(x^{(k)})^T \Delta x^{(k)}$$

$$+ \frac{1}{2} (\Delta x^{(k)})^T \underbrace{\nabla^2 f(x^{(k)})}_{P} \Delta x^{(k)}$$



$$\min_{\Delta x} r + q^T \Delta x + \frac{1}{2} \Delta x^T P \Delta x$$

$\downarrow$  Gradient = 0

$$\Delta x = -P^{-1} q$$

$$\Delta x^{(k)} = -P^{-1}q = -\underline{\underline{\nabla^2 f(x^{(k)})^{-1}}} \nabla f(x^{(k)})$$

$$= -\underline{\underline{D^{(k)}}} \nabla f(x^{(k)})$$

$$D^{(k)} = I \quad \text{or} \quad \nabla^2 f(x^{(k)})^{-1}$$

Gradient

Newton

↓  
slow  
(cheap)

between  
(--->)

↓  
fast  
(expensive)

$A^{-1} \rightarrow O(n^3)$

$$D^{(k)} = \text{diagonal} \quad \& \quad D_{ii}^{(k)} = \left( \frac{\partial^2 f(x^{(k)})}{\partial x_i \partial x_i} \right)^{-1}$$

cheap

$$\text{If } \nabla^2 f(x^{(k)}) \succ 0 \Rightarrow D^{(k)} \succ 0$$

$$\underbrace{D^{(0)}, D^{(1)}, \dots, D^{(p)}}_{\nabla^2 f(x^{(0)})^{-1}}, \quad \underbrace{D^{(p+1)}, D^{(p+2)}, \dots}_{\nabla^2 f(x^{(p)})^{-1}}$$

$$x^{(k+1)} = x^{(k)} + \alpha^{(k)} \Delta x^{(k)}$$

✓
?
✓

How to find stepsize?

1 - exact line search

2 - limited line search

3 - backtracking

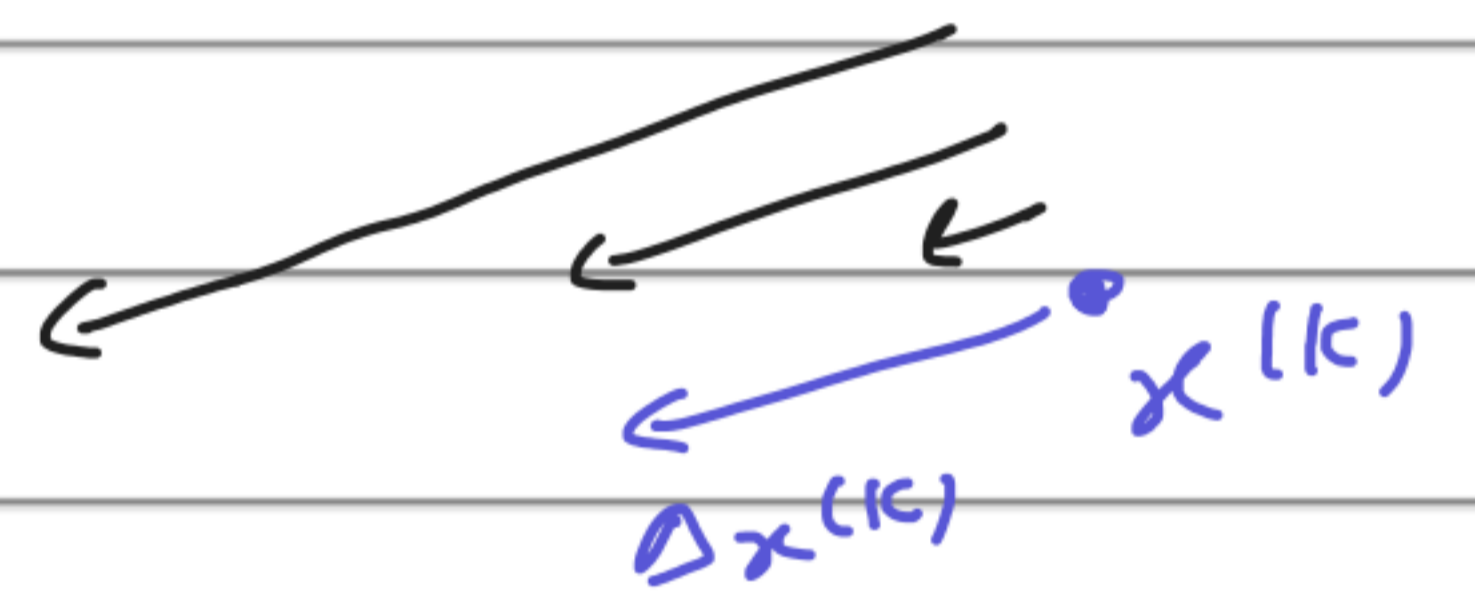
$$f(x^{(k+1)}) < f(x^{(k)}) \longrightarrow \text{best improvement}$$

reduce

$$\min_{\alpha} f(x^{(k)} + \alpha \Delta x^{(k)})$$

$\alpha$

$$\text{s.t. } \alpha \geq 0$$



$$\longrightarrow \alpha^{(k)} = \text{optimal } \alpha$$

exact line search

$$\min_x f(x)$$

$\longrightarrow$

$$x^{(0)}$$

$$\longrightarrow x^{(1)}$$

$$\longrightarrow x^{(2)}$$

$\longrightarrow \dots$

$\downarrow$

$$\min \alpha$$

$\downarrow$

$$\min \alpha$$

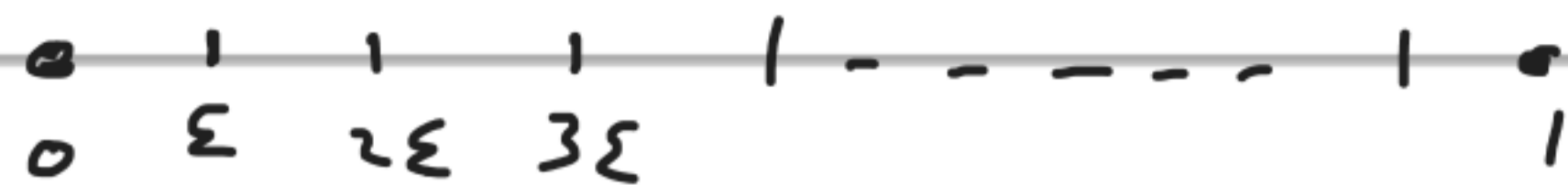
$\downarrow$

$$\min \alpha$$

$x \in \mathbb{R}^n \rightarrow$  multi-variate opt

$\alpha \in \mathbb{R} \rightarrow$  univariate opt

$$\min_{\alpha \in \mathbb{R}} g(\alpha) \quad \text{s.t.} \quad 0 \leq \alpha \leq 1$$



check  $g(\cdot)$  at all grid points to find

the best one.  $\rightarrow \frac{1}{\epsilon}$  evaluations

$$\min_{x \in \mathbb{R}^n} f(x) \quad \text{s.t.} \quad 0 \leq x_i \leq 1$$

$\Rightarrow$  gridding:  $(\frac{1}{\epsilon})^n$  evaluations

$$\epsilon = 0.5, \quad n = 400 \Rightarrow \# \text{ eval} \approx$$

# atoms in observable

universe

Limited line search:

$$\min f(x^{(k)} + \alpha \Delta x^{(k)}) \quad \text{s.t.} \quad 0 \leq \alpha \leq T$$

$\leftarrow$   
threshold

Backtracking : No optimization

Pick  $\alpha > 0$  ,  $0 < \beta < 1$

$\alpha$   
 $\alpha\beta$   
 $\alpha\beta^2$   
 $\alpha\beta^3$   
⋮  
↓  
0

→ pick the first one  
s.t.  $f(\text{new point}) < f(\text{old point})$

