



262B-Lecture 4

Date created: 2021.01.28
N. of Pages: 12

$\min f(x)$: Algorithm

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

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

special : gradient

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

find ? $\rightarrow g^{(k)}$: Approximate gradient

$$g^{(k)} = \nabla f(x^{(k)}) + \underbrace{e^{(k)}}_{\text{error}}$$

1 - $e^{(k)}$ is small relative to gradient

$$\|e^{(k)}\| \leq \|\nabla f(x^{(k)})\| \quad \forall k$$

downside: as we proceed, gradient gets

smaller & so the error must be small.

\Rightarrow Algorithm work!

$$\text{gradient direction} = \nabla f(x^{(k)})^T (-g^{(k)})$$

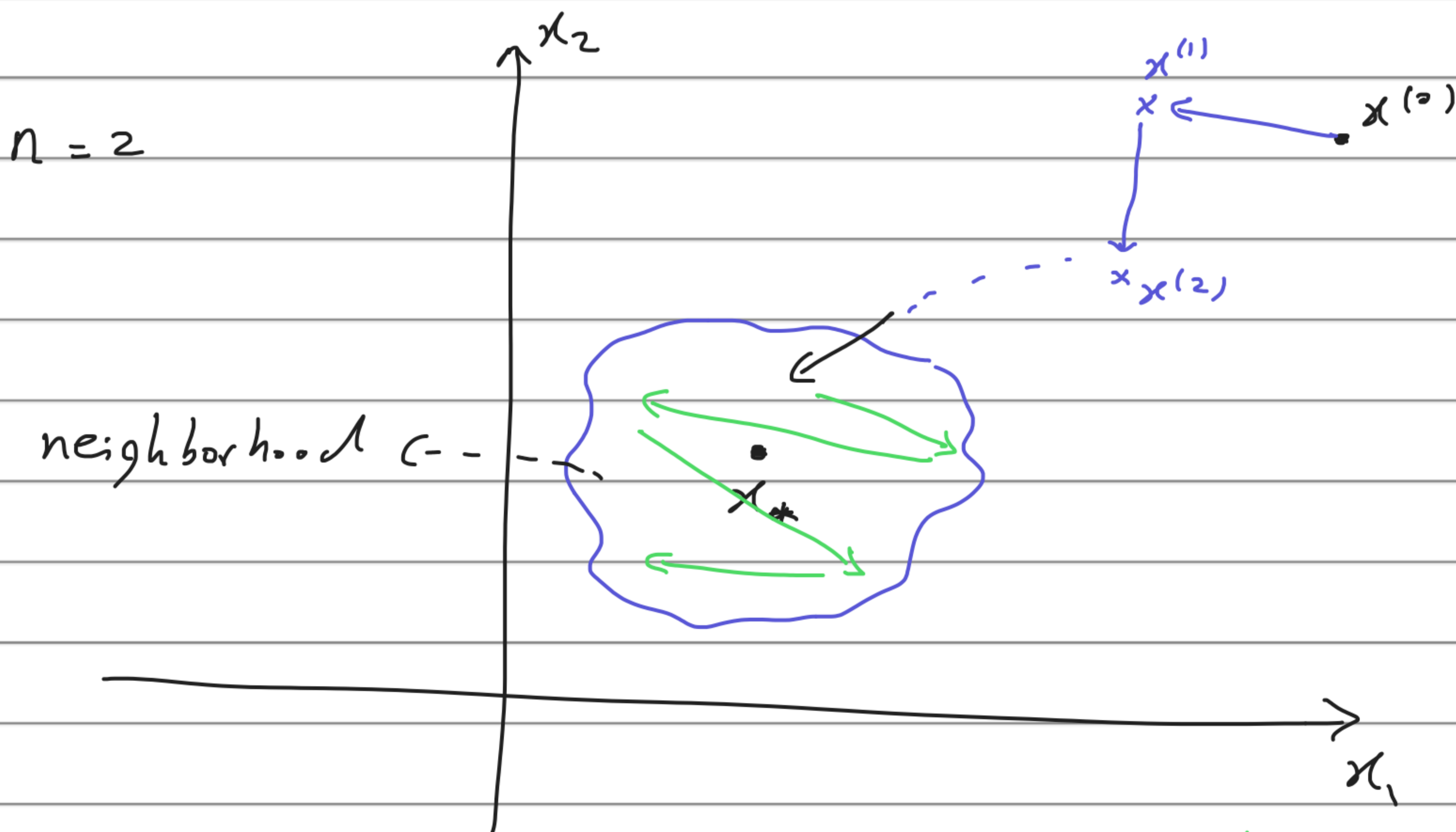
$$= -\|\nabla f(x^{(k)})\|^2 - \nabla f(x^{(k)})^T e^{(k)}$$

$$\leq -\|\nabla f(x^{(k)})\|^2 + \|\nabla f(x^{(k)})\| \cdot \|e^{(k)}\|$$

$$= \|\nabla f(x^{(k)})\| \left(\underbrace{+\|\nabla f(x^{(k)})\| - \|e^{(k)}\|}_{> 0 \text{ by assumption}} \right)$$

$\leq 0 \Rightarrow$ descent direction

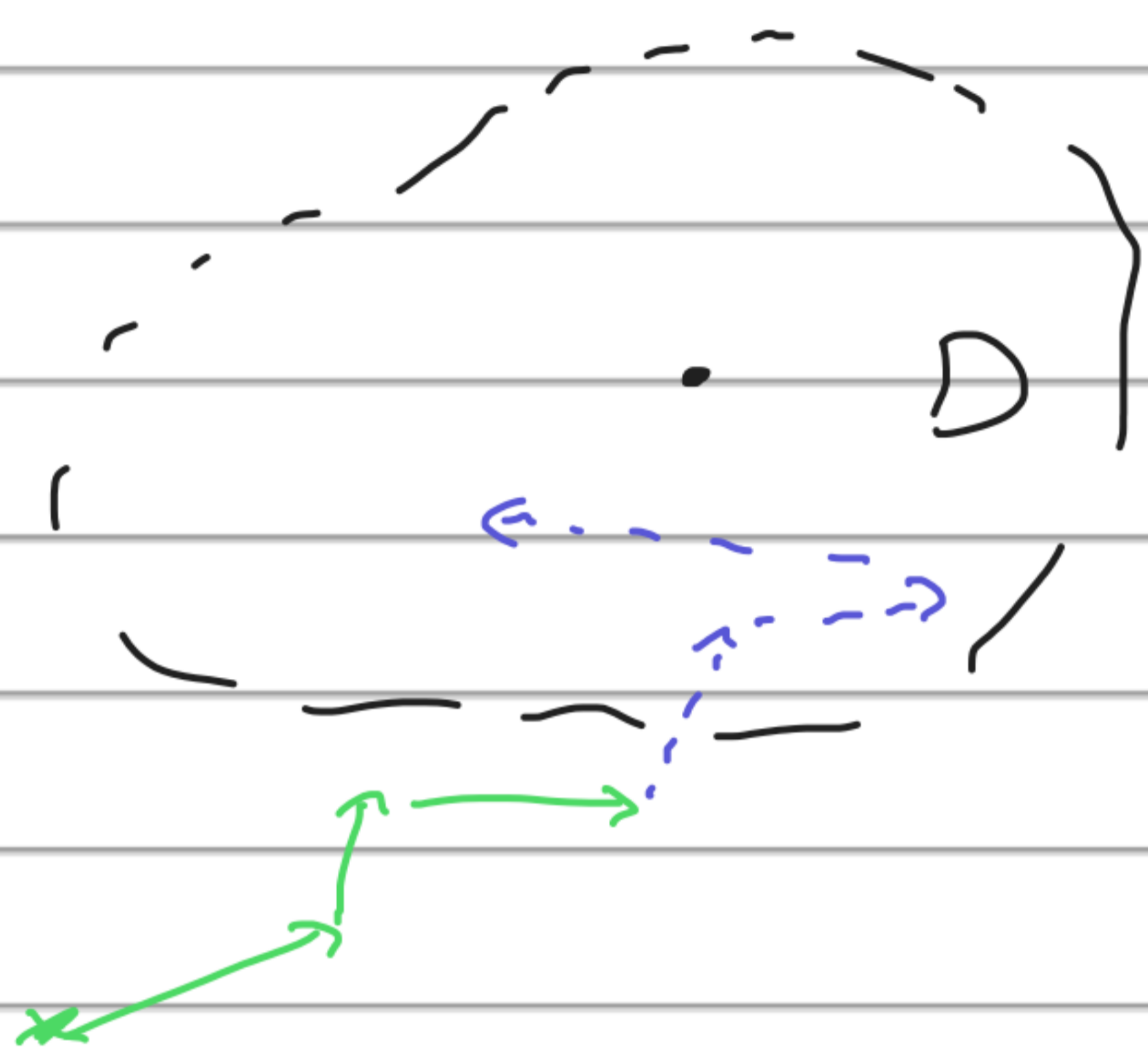
2. $e^{(k)}$ is bounded $\Rightarrow \|e^{(k)}\| \leq \delta \forall k$



erratic behavior

in the neighborhood

neighborhood: $\{x \mid \|\nabla f(x)\| \leq \delta\} = D$



outside of D:

descent algorithm

Assume: $x^{(k)} \notin D$

$$\text{gradient direction} = - \nabla f(x^{(k)})^T \underbrace{g^{(k)}}_{\nabla f(x^{(k)}) + e^{(k)}}$$

$$\leq - \underbrace{\|\nabla f(x^{(k)})\|^2}_{\text{blue}} + \underbrace{\|\nabla f(x^{(k)})\| \times \delta}_{\text{blue}}$$

$$= - \|\nabla f(x^{(k)})\| \left(\underbrace{+\|\nabla f(x^{(k)})\| - \delta}_{\text{green}} \right)$$

$\leq 0 \Rightarrow$ descent
direction

$$e^{(0)} \rightarrow e^{(1)} \rightarrow e^{(2)} \rightarrow \dots$$

special scenario: $e^{(k)} = \delta$

gradient alg.

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

perturbed gradient alg.

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

bias

min f(x)

min f(x) + $\underbrace{\delta^T x}_{\text{Linear}}$

3 - SGD : Stochastic gradient descent

$$g^{(k)} = \nabla f(x^{(k)}) + \underbrace{e^{(k)}}_{\text{random}}$$

min {f(x)}

$\mathbb{E}_w \{ F(x, w) \}$

uncertainty

(Stochastic or simulation optimization)

$$\nabla f(x) = \mathbb{E}_w \{ \nabla_x F(x, w) \}$$

hard to compute in real world

→ Simulate

Basic approach:

$$\text{gradient} : \nabla f(x^{(k)}) \xrightarrow{\text{approximate}} \nabla_x F(x^{(k)}, \underbrace{w^{(k)}}_{\text{one sample}})$$

$$\underbrace{e^{(k)}} = \nabla_x F(x^{(k)}, w^{(k)}) - \mathbb{E}_w (\nabla F(x^{(k)}, w))$$

Error: random variable

$$\mathbb{E}(e^{(k)}) = 0 \implies$$

$$\nabla f(x^{(k)}) (-g^{(k)}) < 0 \quad \underbrace{\text{on average}}$$

\implies SGD : $e^{(k)}$: i.i.d., zero mean, bounded

$$\text{Thm: } \alpha^{(k)} \rightarrow 0, \sum \alpha^{(k)} = \infty, \sum (\alpha^{(k)})^2 < \infty$$

\implies SGD: Converge to a stationary point.

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

Assume $\{\Delta x^{(k)}\}_{k=1}^{\infty}$ is gradient related.

And $\alpha^{(k)}$

- Exact line search
- Limited line search
- Armijo rule

③

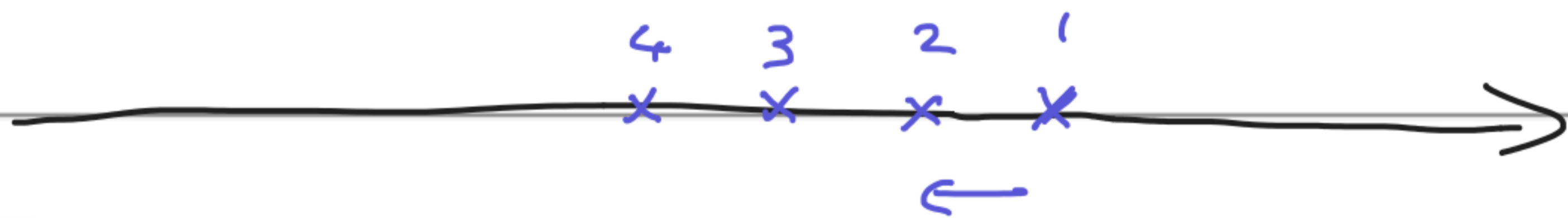
⇒ Every limit point of $\{x^{(k)}\}_{k=1}^{\infty}$ is a stationary point.

+1, -1, +1, -1, ...
 no convergence,
 two limit points: +1, -1

Proof of ③: (by contradiction)

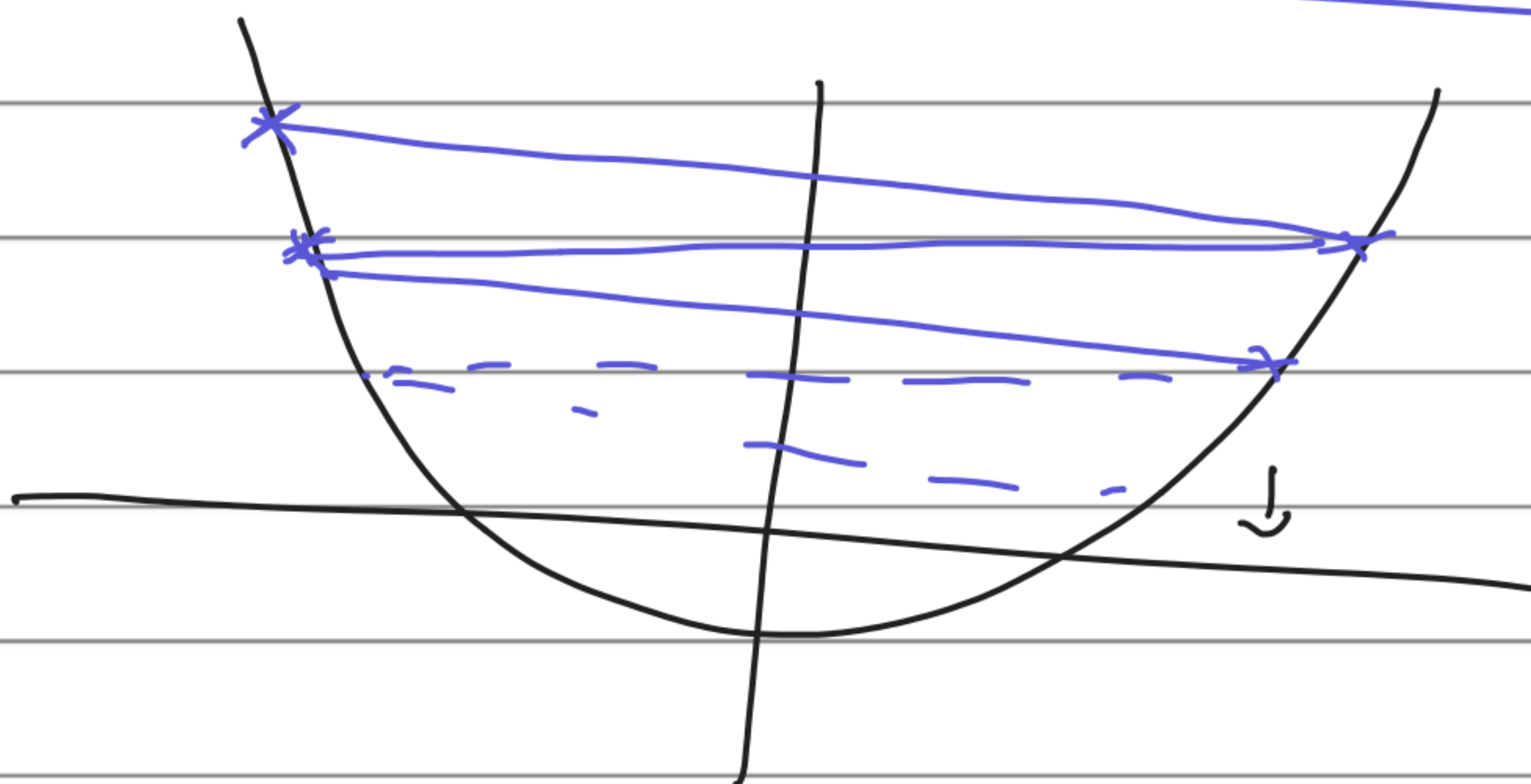
$$\begin{array}{ccccccc}
 x^{(0)} & \rightarrow & x^{(1)} & \rightarrow & x^{(2)} & \rightarrow & \dots \\
 \vdots & & \vdots & & \vdots & & \\
 f(x^{(0)}) & > & f(x^{(1)}) & > & f(x^{(2)}) & > & \dots
 \end{array}$$

⇒ a sequence in \mathbb{R} that is monotonic



$$1 - f(x^{(k)}) \rightarrow -\infty$$

$$2 - f(x^{(k)}) \rightarrow f_*$$



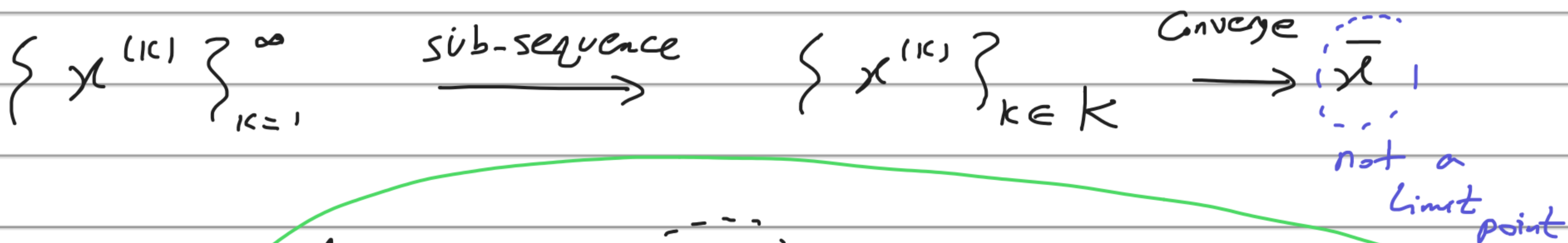
$$\begin{array}{l}
 x^{(k)} \rightarrow \{-1, +1\} \\
 f(x^{(k)}) \rightarrow 0
 \end{array}$$

$$f(x^{(k+1)}) - f(x^{(k)}) \leq \underbrace{\sigma}_{\sigma \in (0,1)} \underbrace{|\nabla f(x^{(k)})^T \Delta x^{(k)}|}_{\leq 0} \quad (*)$$

improvement

$k \rightarrow \infty : f(x^{(k)}) \rightarrow f_*$ \rightarrow left side $(*) \rightarrow 0$

Limit point: pick an arbitrary limit point \bar{x}



$(*) \Rightarrow \limsup_{\substack{k \rightarrow \infty \\ k \in K}} |\nabla f(x^{(k)})^T \Delta x^{(k)}| = 0$ ***

gradient related: $\limsup_{\substack{k \rightarrow \infty \\ k \in K}} \underbrace{|\nabla f(x^{(k)})^T \Delta x^{(k)}|}_{\leq 0} < 0$ ***

$(*)$, $(**)$ $\Rightarrow \limsup_{\substack{k \rightarrow \infty \\ k \in K}} \alpha^{(k)} = 0$

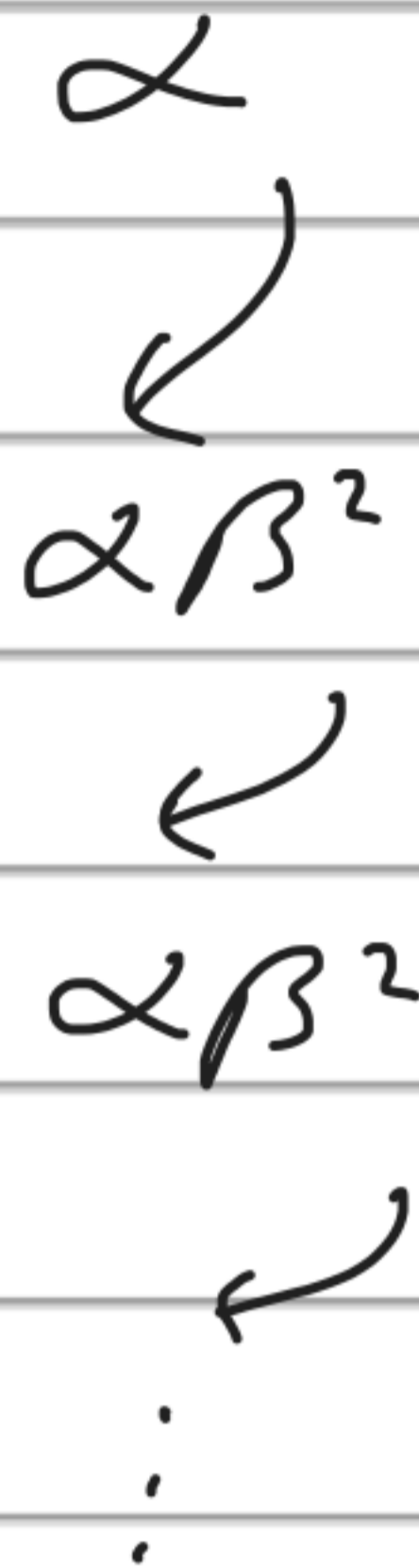
$\Rightarrow \alpha^{(k)} \rightarrow 0$ as $k \rightarrow \infty$
 $k \in K$

$$\lim_{\substack{k \rightarrow \infty \\ k \in K}} \Delta x^{(k)} = \overline{\Delta x}$$

$\{\Delta x^{(k)}\}_k$ is bounded,
so it has a convergent

subsequence with index set $\bar{K} \subseteq K$. assume $\bar{K} = k$ to simplify proof.

Backtracking:



$$\Rightarrow \alpha^{(k)} \rightarrow 0$$

$k: \text{large}$

$$\Rightarrow \exists \bar{k} \text{ s.t. } \forall k \geq \bar{k} : \alpha^{(k)} < \alpha$$

Armijo: $f(x^{(k)} + \alpha^{(k)} \Delta x^{(k)}) - f(x^{(k)}) \leq \sigma \nabla f(x^{(k)})^T (\alpha^{(k)} \Delta x^{(k)})$

①

but

$$f(x^{(k)} + \frac{\alpha^{(k)}}{\beta} \Delta x^{(k)}) - f(x^{(k)}) > \sigma \nabla f(x^{(k)})^T \left(\frac{\alpha^{(k)}}{\beta} \Delta x^{(k)} \right)$$

②

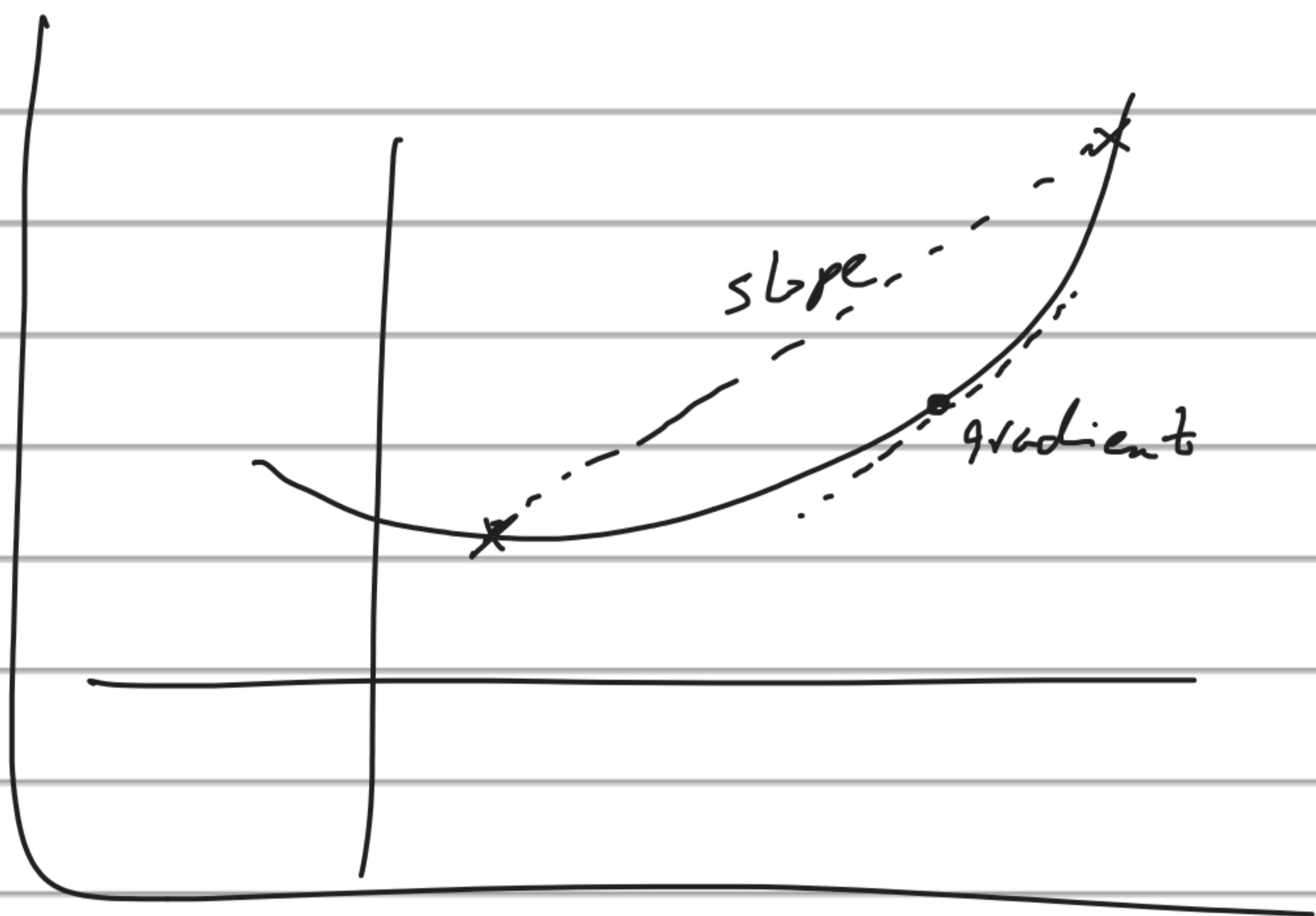


$$x \quad \alpha \beta^{t-1}$$

$$\checkmark \quad \alpha \beta^t = \alpha^{(k)}$$

$$f(x^{(k)} + \frac{\alpha^{(k)}}{\beta} \Delta x^{(k)}) - f(x^{(k)})$$

$$= \nabla f(x^{(k)} + \underbrace{\frac{\alpha^{(k)}}{\beta} \Delta x^{(k)}}_{\tilde{\alpha}^{(k)} \Delta x^{(k)}})^T \Delta x^{(k)}$$



$$\Rightarrow \nabla f(\bar{x}^{(k)} + \alpha^{(k)} \Delta \bar{x}^{(k)})^T \Delta \bar{x}^{(k)} > \sigma \nabla f(\bar{x}^{(k)})^T \Delta \bar{x}^{(k)}$$

$$k \rightarrow \infty, k \in \mathbb{K}$$

$$\Rightarrow \left[(1-\sigma) \nabla f(\bar{x})^T \Delta \bar{x} > 0 \right] \Rightarrow \nabla f(\bar{x})^T \Delta \bar{x} > 0$$

Limit point
Limit of direction

$$0 < \sigma < 1$$

$$\text{Also, } \nabla f(\bar{x})^T \Delta \bar{x} < 0 \text{ (descent)}$$

\Rightarrow Contradiction

Lipschitz continuity of gradient:

$$\exists L > 0 : \|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\|$$
$$\forall x, y \in \mathbb{R}^n$$

Thm: Lipschitz continuity of gradient
implies:

$$\underline{f(x+y)} \leq \underline{f(x)} + y^T \nabla f(x) + \frac{L}{2} \|y\|^2$$

$$\left(x=0 \Rightarrow f(y) \leq \underbrace{f(0) + y^T \nabla f(0)}_{\text{quadratic in } y} + \frac{L}{2} \|y\|^2 \right) \quad \forall x, y$$

Proof: $\underline{g(t)} = \underline{f(x+ty)}$, $t \in \mathbb{R}$

$$f(x+y) - f(x) = g(1) - g(0) = \int_0^1 \frac{dg(t)}{dt} dt$$

$$= \int_0^1 y^T \underline{\nabla f(x+ty)} dt = \int_0^1 y^T \nabla f(x) dt$$

$$+ \int_0^1 y^T (\nabla f(x+ty) - \nabla f(x)) dt$$

$$f(x+y) - f(x) = \int_0^1 y^T \nabla f(x) dt +$$

$$\int_0^1 y^T (\nabla f(x+ty) - \nabla f(x)) dt \leq$$

$$\int_0^1 \underbrace{y^T \nabla f(x)} dt + \int_0^1 \underbrace{\|y\| \cdot \|\nabla f(x+ty) - \nabla f(x)\|} dt$$

$$\leq L \underbrace{\|x+ty - x\|}_{t \|y\|}$$

$$= y^T \nabla f(x) + L \|y\|^2 \underbrace{\int_0^1 t dt}_{\frac{1}{2}}$$

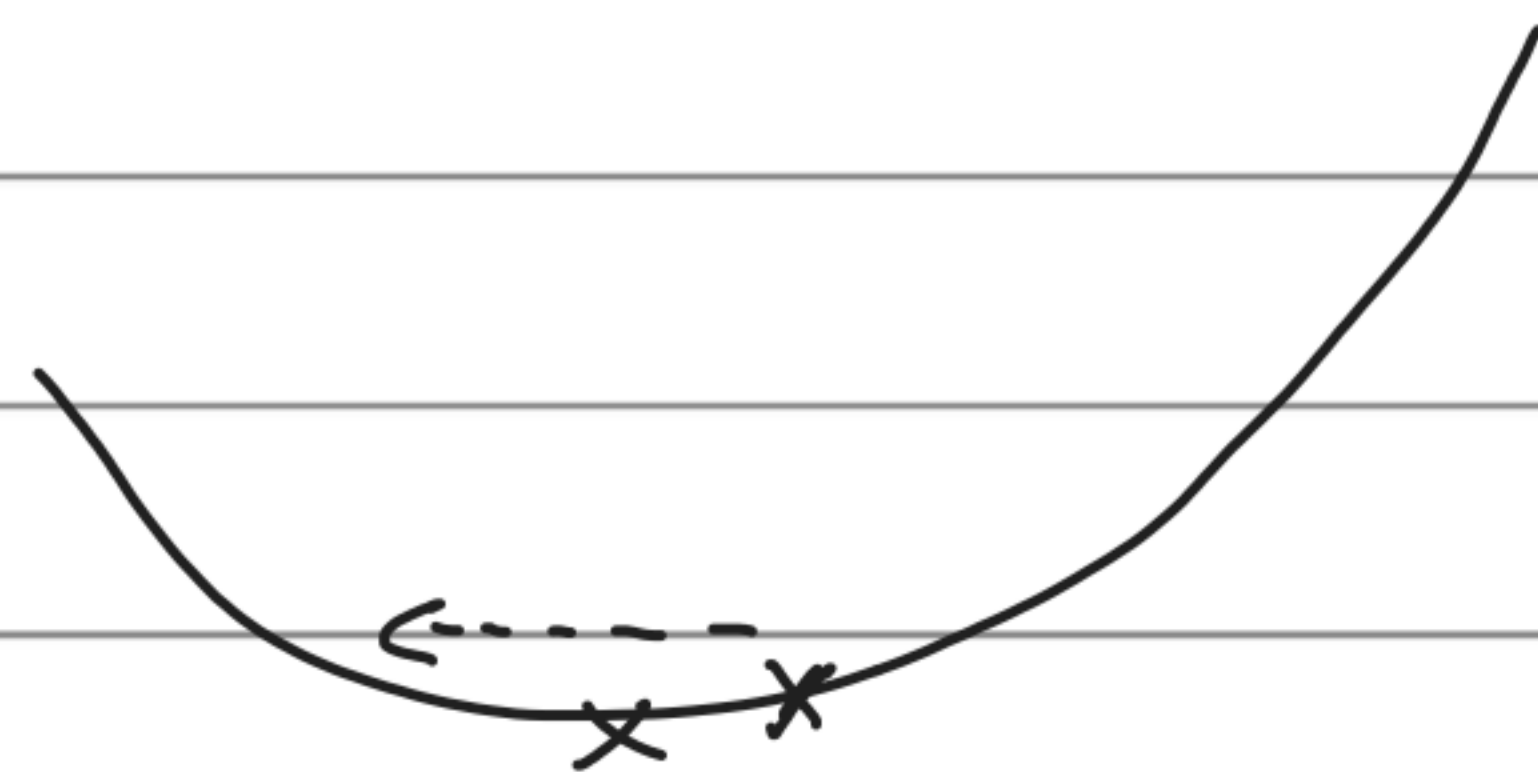
Rough idea: $x^{(k+1)} = x^{(k)} - \alpha^{(k)} \nabla f(x^{(k)})$

implement: $\alpha^{(k)} \quad ? \quad \rightarrow \quad \alpha^{(k)} = \alpha$

$\alpha < \frac{2}{L} \Rightarrow x^{(k)} \rightarrow \text{stationary point}$

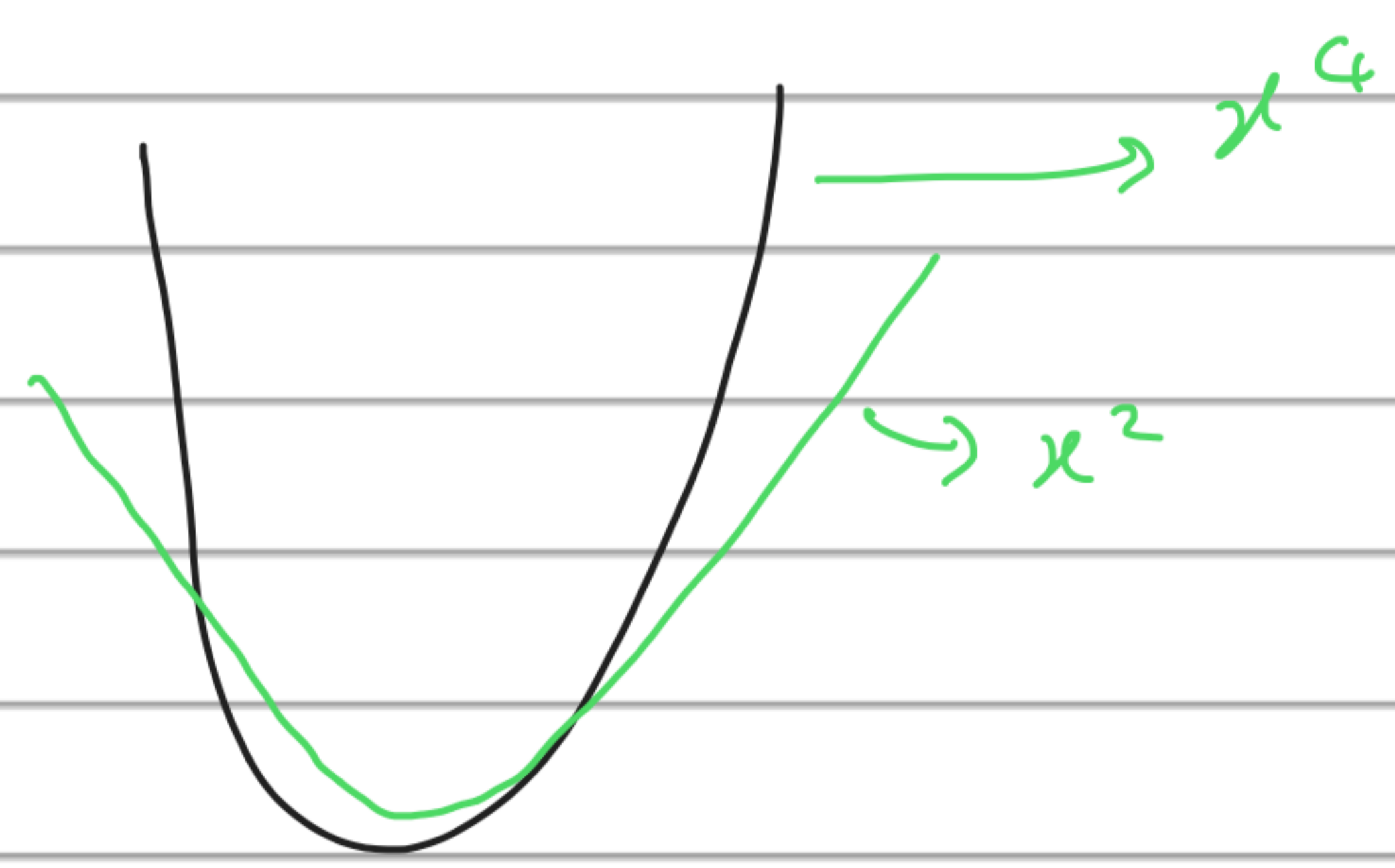
$x^{(k)} \rightarrow \text{close to } x_*, \quad \nabla f(x^{(k)}) \rightarrow 0$

$\underbrace{\alpha}_{\text{Constant}} \times \underbrace{\nabla f(x^{(k)})}_{\text{re-scaling}} \rightarrow 0$



$\mathcal{L} : f(x) \leq \text{quadratic function}$

$f(x) = x^4 :$



\Rightarrow local analysis: $f(x) \leq \text{quadratic function}$
in a neighborhood

