



262B-Lecture 5

Date created: 2021.02.03
N. of Pages: 15

Lipschitz continuity of gradient:

$$\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\| \quad \forall x, y$$

Thm $\implies f(x+y) \leq f(x) + \nabla f(x)^T y + \frac{L}{2} \|y\|^2 \quad \forall x, y$

Constant stepsize:

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

$$x^{(k+1)} = x^{(k)} + \alpha^{(k)} \Delta x^{(k)} \quad \text{where } \{\Delta x^{(k)}\} \text{ is}$$

gradient related. Assume Lipschitz continuity

of gradient holds. Pick $\varepsilon \in (0, 1]$ and

define:
$$\bar{\alpha}^{(k)} = \frac{|\nabla f(x^{(k)})^T \Delta x^{(k)}|}{L \|\Delta x^{(k)}\|^2}$$

\implies every limit point of $\{x^{(k)}\}$ is a

stationary point if $\alpha^{(k)} \in [\varepsilon, (2-\varepsilon)\bar{\alpha}^{(k)}]$.

$(0, 2\bar{\alpha}^{(1)}) \leftarrow$ pick arbitrary stepsize

$(0, 2\bar{\alpha}^{(2)}) \leftarrow \dots$

\vdots

$(0, 2\bar{\alpha}^{(k)}) \leftarrow \dots$

\vdots

proof:

$$f(x^{(k+1)}) = f(x^{(k)} + \underbrace{\alpha^{(k)} \Delta x^{(k)}}_{\substack{\text{previous thm} \\ \leq}}) \leq f(x^{(k)}) + \underbrace{\nabla f(x^{(k)})^T (\alpha^{(k)} \Delta x^{(k)})}_{\substack{\text{negative,} \\ \text{descent direction}}} + \frac{L}{2} \|\alpha^{(k)} \Delta x^{(k)}\|^2$$

$$\leq f(x^{(k)}) - \alpha^{(k)} |\nabla f(x^{(k)})^T \Delta x^{(k)}|$$

$$+ \frac{L}{2} \|\alpha^{(k)} \Delta x^{(k)}\|^2$$

$$= f(x^{(k)}) - \underbrace{\alpha^{(k)}}_{\geq \varepsilon} \left(|\nabla f(x^{(k)})^T \Delta x^{(k)}| - \frac{L}{2} \alpha^{(k)} \|\Delta x^{(k)}\|^2 \right)$$

$$\Rightarrow \left(f(x^{(k+1)}) - f(x^{(k)}) \right) \leq \frac{-\varepsilon^2}{2} |\nabla f(x^{(k)})^T \Delta x^{(k)}| \geq \frac{1}{2} \varepsilon |\nabla f(x^{(k)})^T \Delta x^{(k)}|$$

Armijo rule,
enough improvement

Gradient algorithm:

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

$$\Rightarrow \bar{\alpha}^{(k)} = \frac{|\nabla f(x^{(k)})^\top \Delta x^{(k)}|}{L \times \|\Delta x^{(k)}\|^2} = \frac{1}{L}$$

$$\Rightarrow \alpha^{(k)} \in \left[\varepsilon, \frac{2-\varepsilon}{L} \right], \quad \varepsilon > 0$$
$$\rightarrow \left(0, \frac{2}{L} \right)$$

Thm: If stepsize $\alpha^{(k)} = \alpha = \text{constant}$,

then gradient algorithm works if

$$\alpha \in \left(0, \frac{2}{L} \right)$$

Another version: diminishing step size,

If 1) Lipschitz continuity

2) $\alpha^{(k)} \rightarrow 0$ as $k \rightarrow \infty$

3) $\sum_{k=1}^{\infty} \alpha^{(k)} \rightarrow \infty \Rightarrow$ every

limit point is a station
point. any

What if L doesn't exist?

$f(x) = x^4 \rightarrow$ grows faster than a quadratic function

So far, we defined L globally.

What if we define it locally?

Previously: $\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\|$

$\forall x, y \in \mathbb{R}^n$

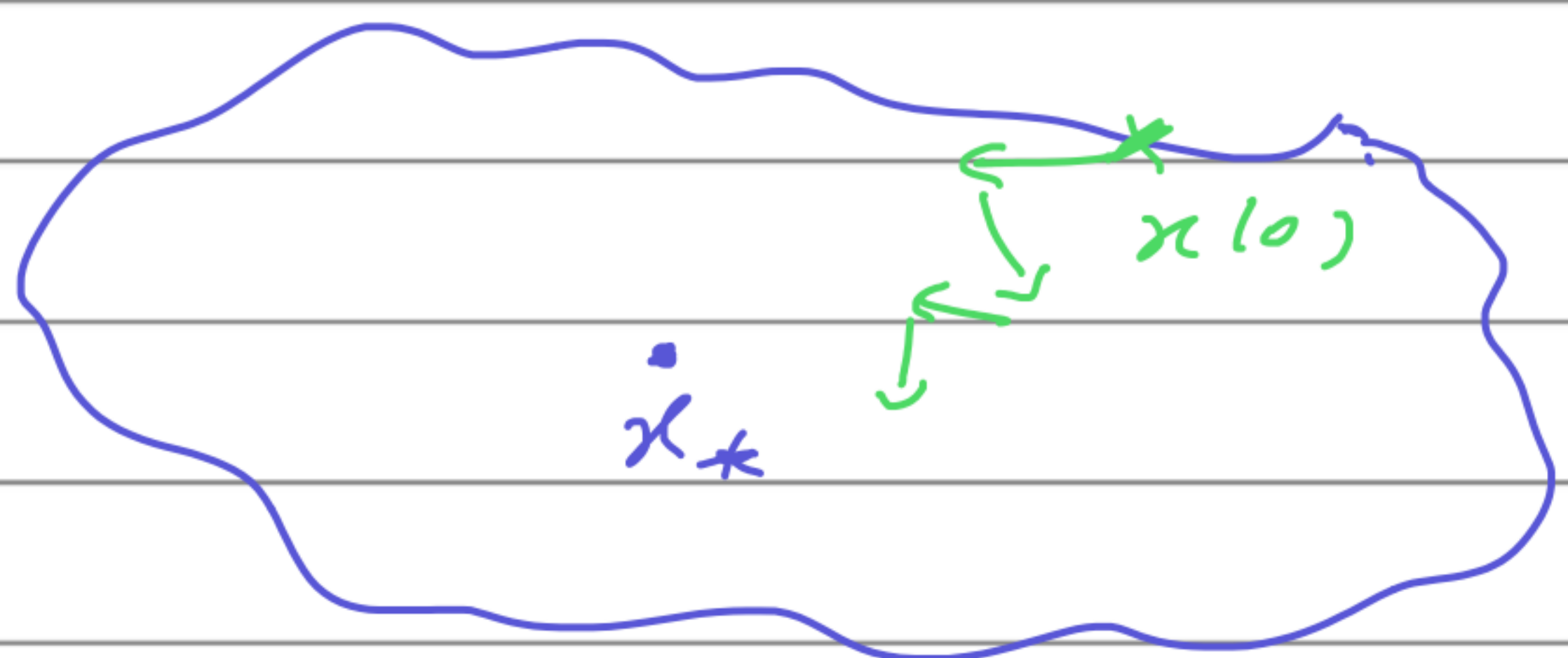
Now: $\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\| \quad \forall x, y \in D$

$D =$ a neighborhood of a stationary point x_* .

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

Property $\textcircled{*}$

sublevel set: $\Sigma = \{x \mid f(x) \leq f(x^{(0)})\}$



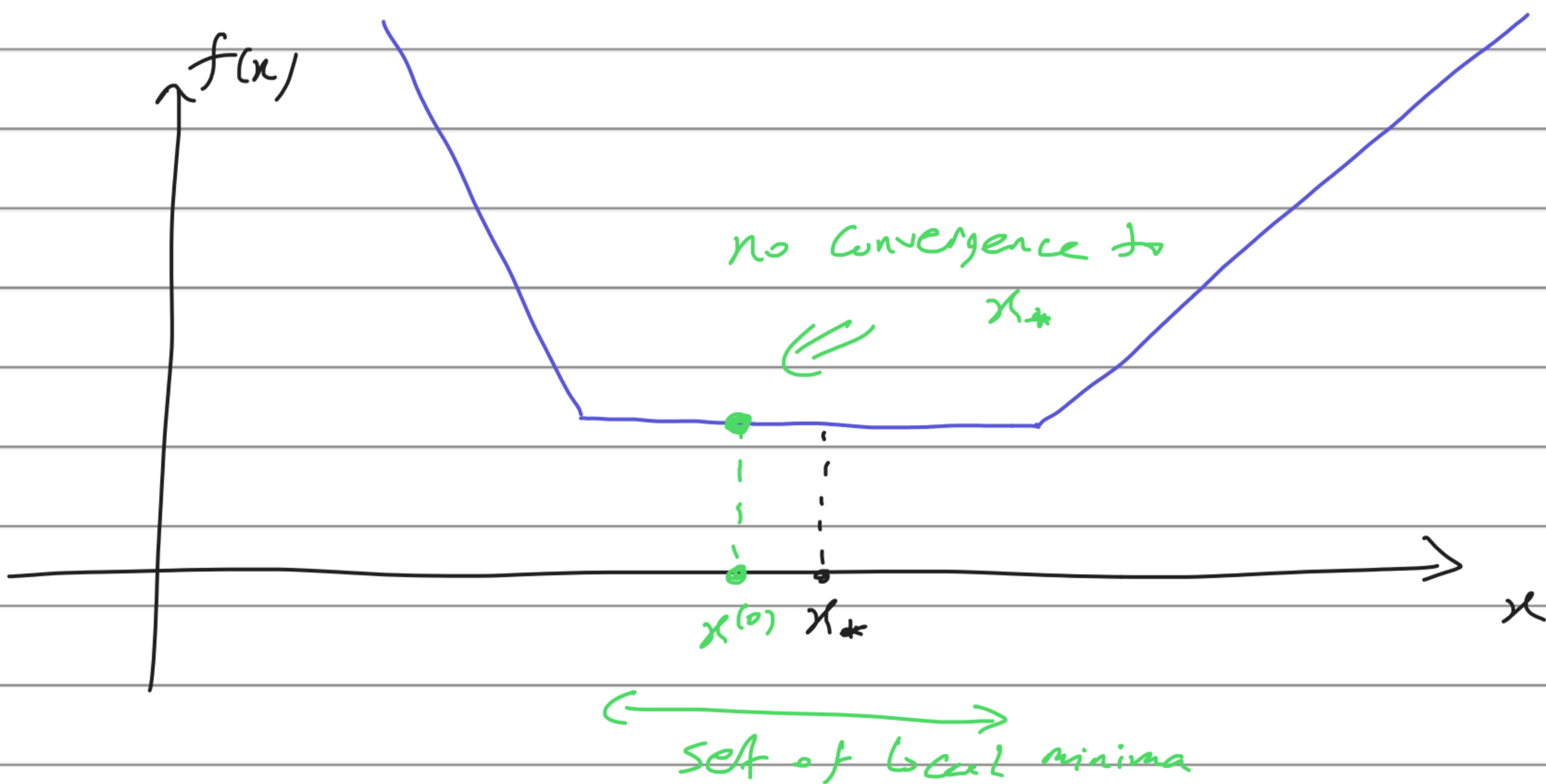
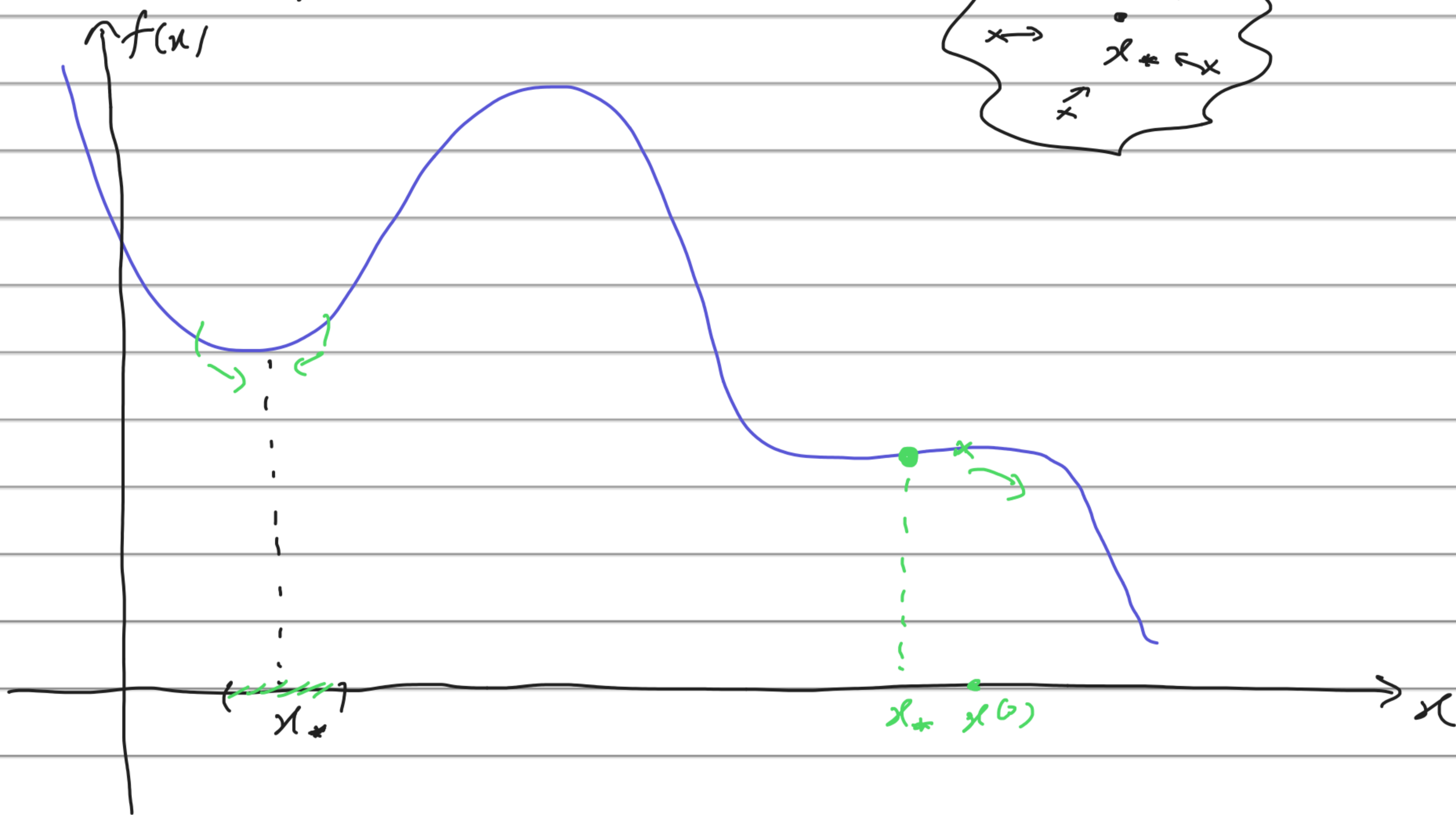
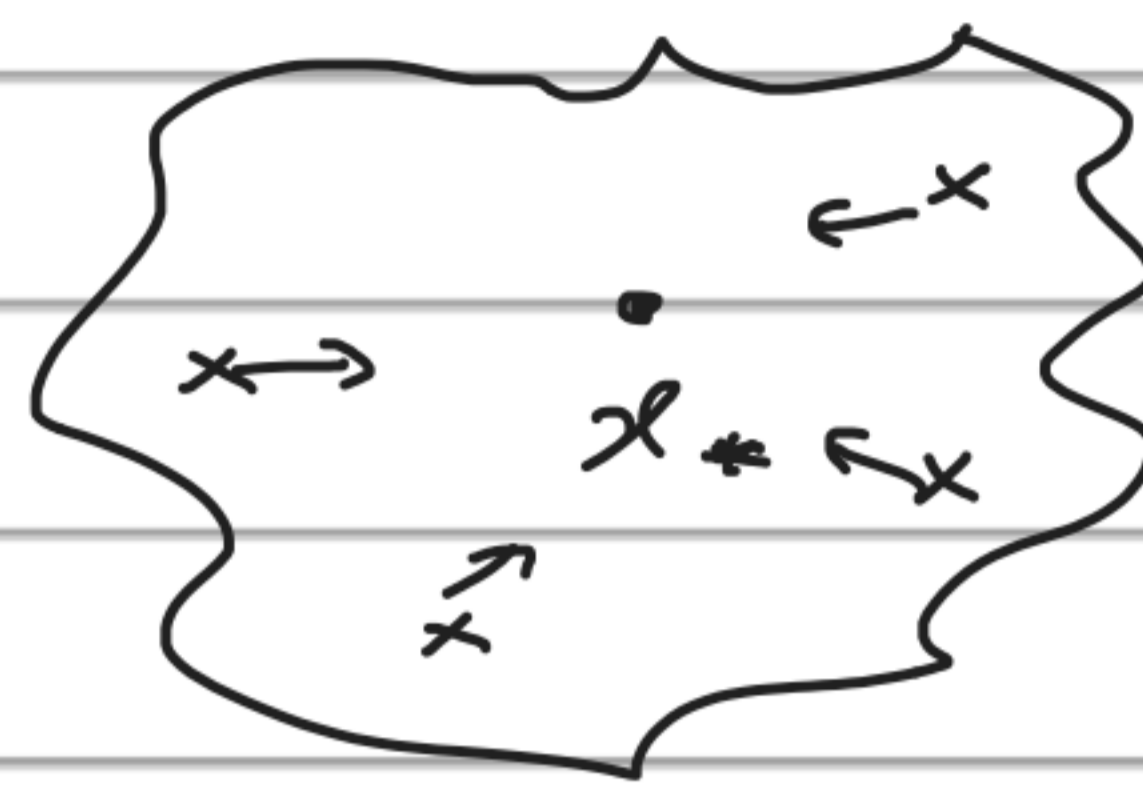
$\Rightarrow D: \underline{\Sigma}$

(a little bit bigger than Σ)

Capture theorem: If $x^{(0)}$ is close enough

to an isolated local min x_* , then

$$\{x^{(k)}\} \rightarrow x_*$$



Convergence rate

Alg : $x^{(0)} \rightarrow x^{(1)} \rightarrow x^{(2)} \rightarrow \dots \rightarrow x_*$

a bunch of points in \mathbb{R}^n

Assume: $\{x^{(k)}\} \rightarrow$ unique limit point = x_*

Measure of speed: \rightarrow scalar

error: 1 - $e(x) = \|x - x_*\|$

2 - $e(x) = |f(x) - f(x_*)|$

Asymptotic convergence rate

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

beginning unimportant

tail matters.

$\{e^{(k)}\}$: Alg 1 versus Alg 2

1 : 1 0.9 0.8 0.2 10^{-2} ...

2 : 1 0.7 0.6 0.3 10^{-3} ...

Need a baseline for comparison.

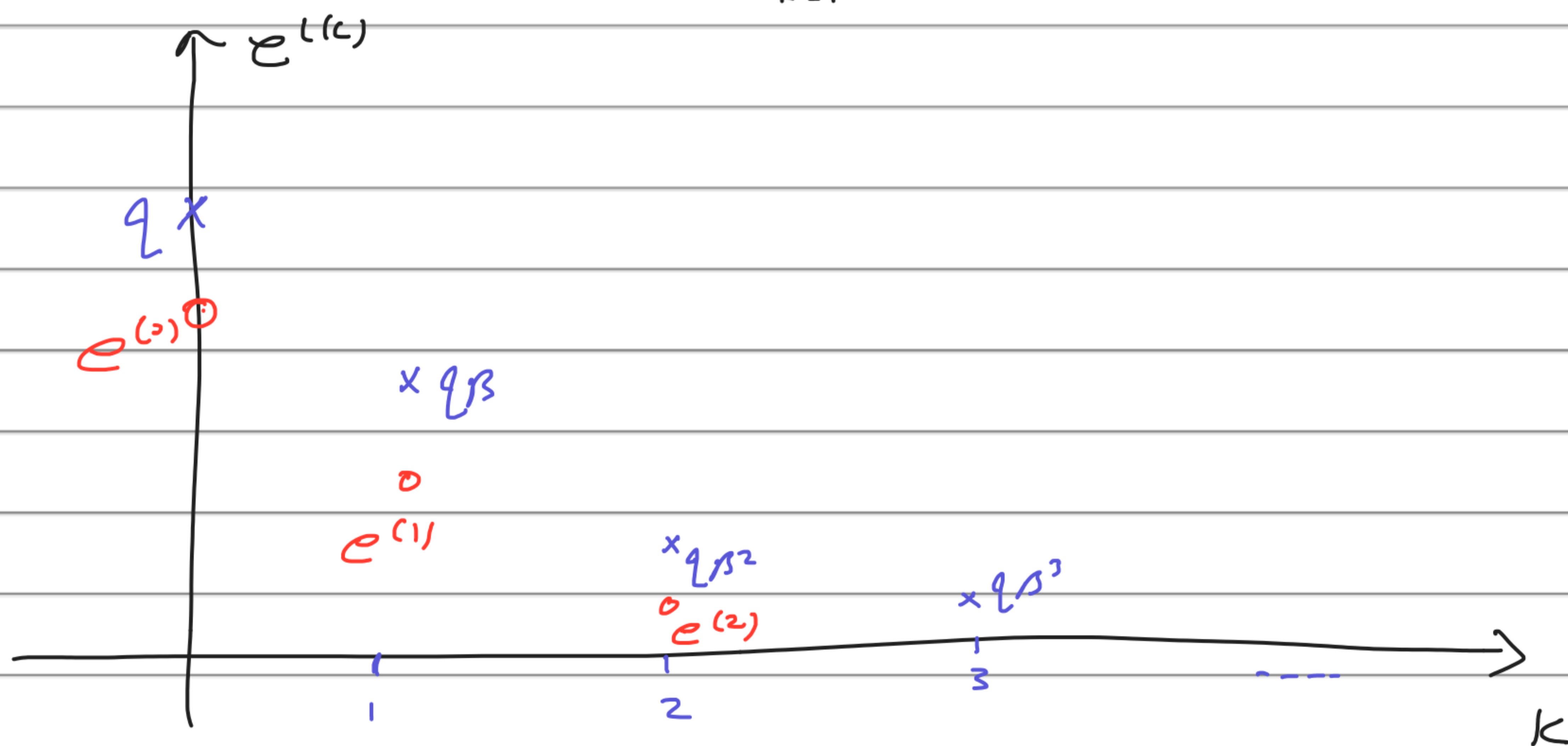
Baseline: pick $\beta \in (0, 1)$

generate a sequence: $1 \rightarrow \beta \rightarrow \beta^2 \rightarrow \beta^3 \rightarrow \dots$

geometric
sequence

Rescale $\{e^{(k)}\} \rightarrow$ pick $q > 0$

Compare $\{e^{(k)}\}_{k=1}^{\infty}$ against $\{q\beta^k\}_{k=1}^{\infty}$



$$e^{(k)} = e(x^{(k)}) \leq q\beta^k \quad \forall k \quad \text{--- (1)}$$

error sequence baseline

$$e(x^{(k)}) \leq \tilde{q}\beta^k \quad \forall k \geq \bar{k} \quad \text{--- (2)}$$

large enough

If that happens $\rightarrow \{e(x^{(k)})\}$

converges linearly or geometrically

with factor β .

Special case: $e(x^{(k)}) = \alpha \beta^k$

$$\Rightarrow \frac{e(x^{(k+1)})}{e(x^{(k)})} = \beta \Rightarrow$$

error reduces by factor β at
each iteration.

Thm: $\{e^{(k)}\}$ converges linearly if

$$\lim_{k \rightarrow \infty} \frac{e^{(k+1)}}{e^{(k)}} < 1$$

with what factor? $\beta = \lim_{k \rightarrow \infty} \frac{e(x^{(k+1)})}{e(x^{(k)})}$

Alg 1: $\longrightarrow \beta = 0.9$ slow

Alg 2: $\longrightarrow \beta = 0.1$ fast

Thm: As β decreases, algorithm gets faster.

What if $\beta = \lim_{k \rightarrow \infty} \frac{e(x^{(k+1)})}{e(x^{(k)})} = 0$?

That means the sequence $\{e^{(k)}\}$ is

faster than $\{\beta^k\}$ for every $0 < \beta < 1$.

\Rightarrow Converges super linearly

Def: $\{e(x^{(k)})\}$ converges super linearly

with order p if

$p > 1$

$$e(x^{(k+1)}) \leq q \times \beta^{(p^k)}$$

for some $q > 0$ &

$(\forall k \text{ or large enough } k) \quad 0 < \beta < 1$

$$\beta = 0.5, \quad q = 1, \quad k = 10$$

<p>Linear</p> <p>↓</p> $e(x^{(10)}) \leq (0.5)^{10}$ <p style="text-align: center;">small</p>	}	<p>Super Linear $p=2$</p> <p>↓</p> $e(x^{(10)}) \leq (0.5)^{\underbrace{(2^{10})}_{1024}}$ <p style="text-align: center;">Extremely small</p>
---	---	--

special case:

$$e(x^{(k)}) = q \times \beta^{(p^k)} \quad \forall k$$

$$\Rightarrow \underbrace{e(x^{(k+1)})}_{\text{new error}} = q \times \beta^{(p^{k+1})} = q \times (\beta^{(p^k)})^p$$

$$= q \times \left(\frac{e(x^{(k)})}{q} \right)^p$$

$$= \underbrace{q^{(1-p)}}_{\text{any number}} \times \underbrace{e(x^{(k)})^p}_{\text{old error}}$$

Thm: $\{e^{(k)}\}$ converges superlinearly

with order p if

$$\lim_{k \rightarrow \infty} \frac{e(x^{(k+1)})}{e(x^{(k)})^p} < \infty$$

note:

since $x^{(k)} \rightarrow x_*$,

we have

$$e(x^{(k)}) \rightarrow 0$$

$p=2$: Quadratic convergence

gradient alg \rightarrow Linear convergence

Newton's alg \rightarrow quadratic convergence

Role of Condition number:

$Q > 0 \rightarrow$ Condition number (c.d.)

$$= \frac{\max \text{eig } Q}{\min \text{eig } Q} = \frac{\lambda_{\max}(Q)}{\lambda_{\min}(Q)}$$

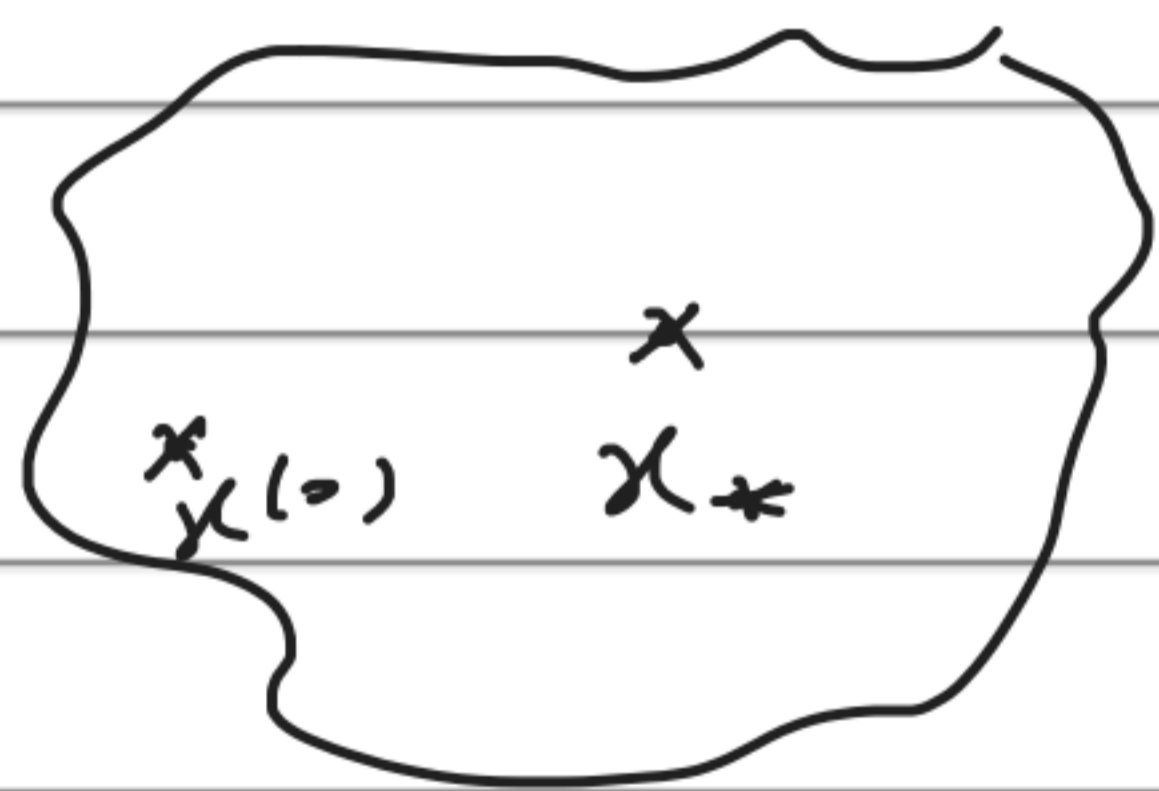
$$\geq 1$$

$f: \mathbb{R}^n \rightarrow \mathbb{R}$, x_* satisfies SOC (sufficient)

\Rightarrow strict local min
isolated

$$\begin{aligned} \Rightarrow f(x) &= f(x_*) + \cancel{\nabla f(x_*)^T (x - x_*)} \\ &+ \frac{1}{2} (x - x_*)^T \nabla^2 f(x_*) (x - x_*) + O(\|x - x_*\|^3) \end{aligned}$$

Now, if x is close to x_* ,



Capture them

$$\nabla^2 f(x_*) \rightarrow \mathbb{Q}$$

$$x - x_* \rightarrow x$$

\Rightarrow right side $\mathbb{Q} \succ 0 \sim$

$$\left(\frac{1}{2} x^T \mathbb{Q} x \right)$$

Min $f(x)$ where $f(x) = \frac{1}{2} x^T \mathbb{Q} x$

$$\Rightarrow \nabla f(x) = \mathbb{Q} x, \quad \nabla^2 f(x) = \mathbb{Q} \succ 0$$

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

$$\begin{aligned} \Rightarrow x^{(k+1)} &= x^{(k)} - \alpha^{(k)} Q x^{(k)} \\ &= \underbrace{(I - \alpha^{(k)} Q)}_{\text{matrix}} x^{(k)} \end{aligned}$$

$*$

error: $e(x^{(k)}) = \|x^{(k)} - x_*\| = \|x^{(k)}\|$

$$* \Rightarrow \|x^{(k+1)}\|^2 = (x^{(k)})^T \underbrace{(I - \alpha^{(k)} Q)^2}_{\geq 0} x^{(k)}$$

$$\Rightarrow \underbrace{\|x^{(k+1)}\|^2}_{e(x^{(k+1)})^2} \leq \lambda_{\max}(I - \alpha^{(k)} Q)^2 \underbrace{\|x^{(k)}\|^2}_{e(x^{(k)})^2}$$

$$x^T A x \geq \lambda_{\min}(A) \|x\|^2$$

$$\lambda_{\min}(A) \|x\|^2$$

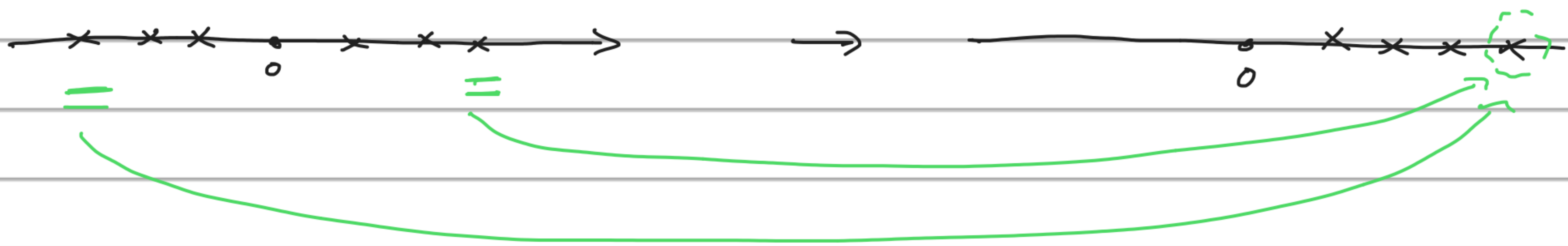
$$\Rightarrow \frac{e(x^{(k+1)})}{e(x^{(k)})} \leq \sqrt{\lambda_{\min}(I - \alpha^{(k)} Q)^2}$$

eigs of Q : $\lambda_1, \lambda_2, \dots, \lambda_n$

eigs of $I - \alpha^{(k)} Q$: $1 - \alpha^{(k)} \lambda_1, 1 - \alpha^{(k)} \lambda_2, \dots, 1 - \alpha^{(k)} \lambda_n$

eigs of $(I - \alpha^{(k)} Q)^2$: $(1 - \alpha^{(k)} \lambda_1)^2, (1 - \alpha^{(k)} \lambda_2)^2, \dots$

Eigs of $I - \alpha^{(k)} Q$ \rightarrow Eigs of $(I - \alpha^{(k)} Q)^2$



min eig of Q : m (little)

max eig of Q : M (capital)

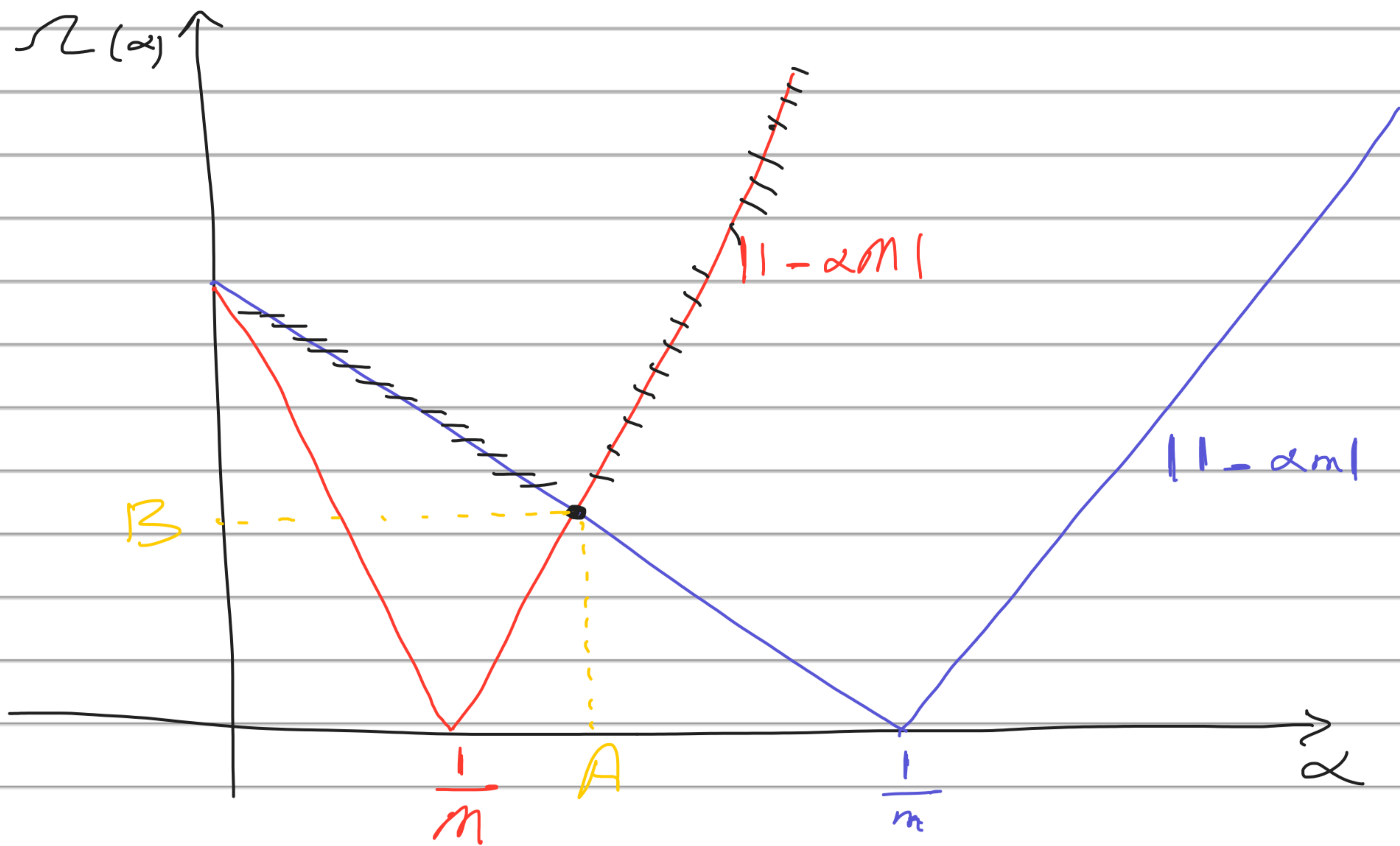
$$\Rightarrow \lambda_{\max} (I - \alpha^{(k)} Q)^2 = \max \left\{ (1 - \alpha^{(k)} m)^2, (1 - \alpha^{(k)} M)^2 \right\}$$

$$\Rightarrow \frac{\rho(x^{(k+1)})}{\rho(x^{(k)})} \leq \max \left\{ |1 - \alpha^{(k)} m|, |1 - \alpha^{(k)} M| \right\}$$

How to pick $\alpha^{(k)}$: two methods

method 1: Find $\alpha^{(k)}$ to minimize

upper bound. $\rightarrow \min_{\alpha^{(k)}} \Omega$



$$A = \frac{2}{m + m_c}, \quad B = \frac{M - m}{M + m} = \frac{\text{c.d.}(\mathcal{Q}) - 1}{\text{c.d.}(\mathcal{Q}) + 1}$$

Optimize upper bound to find stepsize:

$$\left\{ e^{(k)} \right\}_{k=1}^{\infty} \rightarrow \frac{e(x^{(k+1)})}{e(x^{(k)})} \leq \underbrace{\frac{\text{c.d.}(\mathcal{Q}) - 1}{\text{c.d.}(\mathcal{Q}) + 1}}_{0 < \beta < 1}$$

⇒ Gradient alg: Linear convergence
with rate depending on c.d.