

Lecture 12. Gradient descent method

Recall that x^* is optimal for $\min f(x) \iff \nabla f(x^*) = 0$.

Suppose $f(x) = ax^2 + bx + c$. $-f'(x)$ is the direction to optimal.

Intuitively, move from x to $x - t \cdot f'(x)$. when to stop?

Ideally, stop when $f'(x) = 0$. but impractical.

$|f'(x)| < \delta$. $|f(x_{\text{new}}) - f(x_{\text{old}})| < \delta$. # of iterations $< T$

Now consider $f: \mathbb{R}^n \rightarrow \mathbb{R}$. guess initial x_0 . then how to move?

Naïve idea: $x_{k+1} \leftarrow x_k - t \cdot d_k$. s.t. $f(x_{k+1}) < f(x_k)$.

For convex differential f , $-d_k$ is a descent direction $\iff d_k^T \nabla f(x_k) > 0$.

Simply select $d_k = \nabla f(x_k)$. Advantage: max rate descending direction.

Assume $\|d_k\| = 1$. $\nabla_{d_k} f(x_k) = \lim_{t \downarrow 0} (f(x_k + t \cdot d_k) - f(x_k)) / t = d_k^T \cdot \nabla f(x_k)$

By Cauchy-Schwarz, $d_k^T \cdot \nabla f(x_k) = \langle d_k, \nabla f(x_k) \rangle \leq \|d_k\| \cdot \|\nabla f(x_k)\|$

with equality iff $d_k = (\pm) \nabla f(x_k) / \|\nabla f(x_k)\|$. $\geq -\|d_k\| \cdot \|\nabla f(x_k)\|$.

Gradient descent: $x_{k+1} \leftarrow x_k - t \cdot \nabla f(x_k)$.

reasonable requirement: $f(x_{k+1}) < f(x_k)$. consider $f(x) = ax^2$. $a > 0$.

$f'(x) = 2ax \implies x_{k+1} = x_k - 2atx_k$. $\implies (1-2at)^2 x_k^2 < x_k^2 \implies t < \frac{1}{a}$.

consider $f(x) = x^T Q x$. ^{symmetric $Q \in \mathbb{R}^{n \times n}$} $\nabla f(x) = 2Qx$ $x_{k+1} = (I - 2tQ)x_k$.

$$f(x_{k+1}) = x_k^T Q x_k + 4t^2 (Qx_k)^T Q (Qx_k) - 4t (Qx_k)^T (Qx_k).$$

$$f(x_{k+1}) < f(x_k) \Leftrightarrow t \cdot (Qx_k)^T Q (Qx_k) < (Qx_k)^T (Qx_k)$$

Proposition. $\lambda_{\min} \|x\|_2^2 \leq x^T Q x \leq \lambda_{\max} \|x\|_2^2$. (so $t < 1/\lambda_{\max}$ suffice)

Proof: $Q \in \mathbb{R}^{n \times n}$ symmetric \Rightarrow orthogonally diagonalize $Q = U \Lambda U^T$

where $\Lambda = \text{diag}\{\lambda_1, \dots, \lambda_n\}$ and $U^T U = I$. let $x = Uy$.

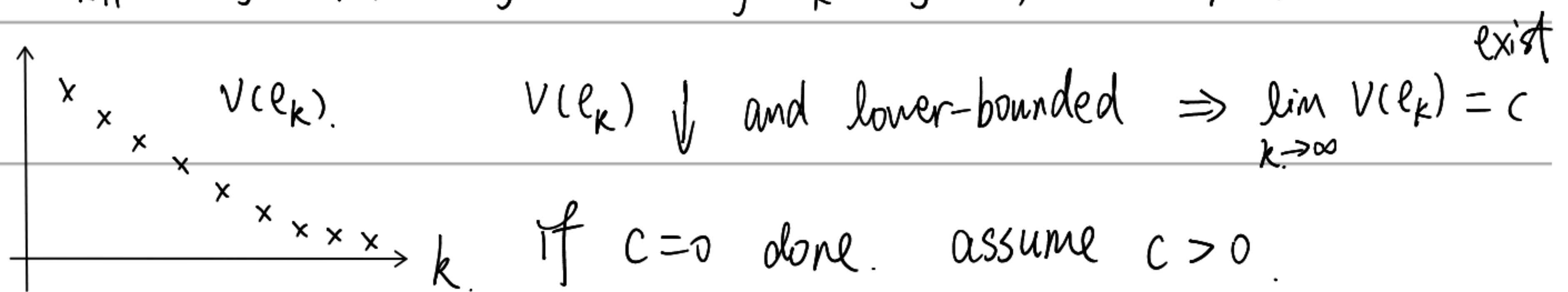
$$x^T Q x = y^T U^T Q U y = y^T \Lambda y = \sum_{i=1}^n \lambda_i y_i^2 \leq \max |\lambda_i| \cdot \|y\|_2^2.$$

$$\|y\|_2^2 = y^T y = y^T U^T U y = x^T x = \|x\|_2^2. \text{ Similarly for } \lambda_{\min}. \quad \square$$

Convergence: let $e_k = x_k - x^*$. and $v(e) = f(e + x^*) - f(x^*)$

$$v(e) = 2e^T Q x^* + e^T Q e. > 0 \text{ (} x^* \text{ optimal, } e \neq 0\text{). } v(0) = 0.$$

$$v(e_{k+1}) = f(x_{k+1}) - f(x^*) < f(x_k) - f(x^*) = v(e_k).$$



let $S = \{x: c \leq v(x) \leq v(e_0)\}$. S is compact. continuous function.

note that $e_{k+1} = e_k + x^* - 2tQ(e_k + x^*) - x^* \triangleq g(e_k)$

let $\delta = \min_{e \in S} |v(g(e)) - v(e)| > 0$ due to continuity and compactness

$$0 \notin S \Rightarrow V(g(e)) < V(e). \quad \forall e \in S \Rightarrow \delta > 0 \Rightarrow V(e_{k+1}) < V(e_k) - \delta. \quad \square$$

Lyapunov's global stability theorem in discrete time:

$$e_{k+1} = g(e_k) \quad \text{where } g: \mathbb{R}^n \rightarrow \mathbb{R}^n \text{ continuous and } g(0) = 0.$$

If there exists $V: \mathbb{R}^n \rightarrow \mathbb{R}$ continuous s.t. (Lyapunov function)

$$\textcircled{1} \quad V(0) = 0 \quad V(e) > 0 \quad \forall e \neq 0. \quad (\text{positivity})$$

$$\textcircled{2} \quad V(e) \rightarrow \infty \quad \text{as } \|e\| \rightarrow \infty \quad (\text{radical unboundedness})$$

$$\textcircled{3} \quad V(g(e)) < V(e). \quad \forall e \neq 0. \quad (\text{strict decrease})$$

Then $\forall e_0 \in \mathbb{R}^n$. we have $e_k \rightarrow 0$ as $k \rightarrow \infty$.

For gradient descent method. $e_k = x_k - x^*$. $g(e) = e - t \cdot \nabla f(e + x^*)$.

$$V(e) = f(e + x^*) - f(x^*). \quad \exists \text{ unique optimal } x^* \Rightarrow \text{convergence.}$$

If we are in math dept. then we are done! $x_k \rightarrow x^*$

But in CS dept. we'd like to ask the convergence rate. x_{k+1}

$$\|x_{k+1} - x^*\|^2 = \|x_k - t \nabla f(x_k) - x^*\|^2 = \|x_k - x^*\|^2 + t^2 \|\nabla f(x_k)\|^2$$

$$\nabla f(x_k)^T (x_k - x^*) \geq f(x_k) - f(x^*). \quad \text{by convexity.} \quad -2t \nabla f(x_k)^T (x_k - x^*)$$

$$\Rightarrow \|x_{k+1} - x^*\|^2 - \|x_k - x^*\|^2 \leq t^2 \|\nabla f(x_k)\|^2 + 2t (f(x^*) - f(x_k))$$

$$\Rightarrow \|x_T - x^*\|^2 - \|x_0 - x^*\|^2 \leq t^2 \sum_{k=0}^{T-1} \|\nabla f(x_k)\|^2 - 2t \sum_{k=0}^{T-1} (f(x_k) - f(x^*)).$$

$$\Rightarrow \sum_{k=0}^{T-1} (f(x_k) - f(x^*)) \leq \frac{1}{2t} (\|x_0 - x^*\|^2 - \|x_T - x^*\|^2) + \frac{t}{2} \sum_{k=0}^{T-1} \|\nabla f(x_k)\|^2.$$

Consider $f(x) = |x|$ with $|f'(x)| = 1$. oscillation between $-\frac{t}{2}$ and $\frac{t}{2}$.

Idea: bound $\|\nabla f(x_k)\|^2$ by $f(x_k) - f(x_{k+1})$

Note that $f(x_{k+1}) - f(x_k) \geq \nabla f(x_k)^T (x_{k+1} - x_k) = -t \|\nabla f(x_k)\|^2$.

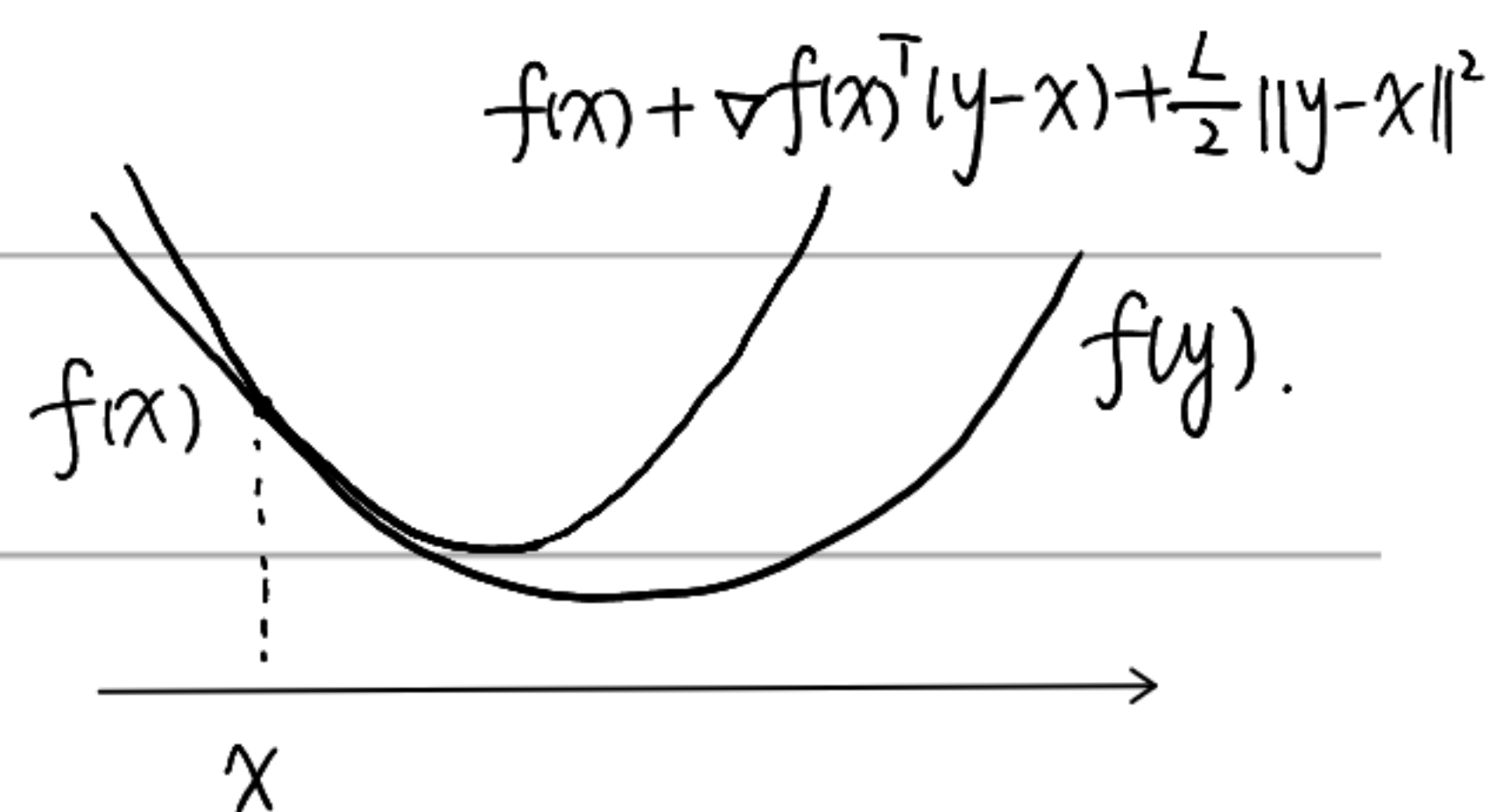
Hope LHS $\leq -\frac{t}{2} \|\nabla f(x_k)\|^2 \Leftrightarrow f(y) \leq f(x) + \nabla f(x)^T (y-x) + ?$

$$? = \frac{t}{2} \|\nabla f(x)\|^2 \quad \times$$

recall that $x_{k+1} - x_k = -t \nabla f(x)$

$$\text{so } ? = \frac{1}{2t} \|y - x\|^2.$$

let $? = \frac{L}{2} \|y - x\|^2$. choose $t < \frac{1}{L}$



$$f(y) \leq f(x) + \nabla f(x)^T (y-x) + \frac{L}{2} \|y-x\|^2 \Leftrightarrow f'(y) \leq f'(x) + L(y-x).$$

Lipschitz Continuity: $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ is Lipschitz continuous with $L > 0$.

if $\|f(x) - f(y)\| \leq L \|x - y\|$ $\forall x, y$. (or simply L -Lipschitz).

Example. $f(x) = w^T x$ is $\|w\|$ -Lipschitz since $|w^T(x-y)| \leq \|w\| \cdot \|x-y\|$.

Example. $f(x) = Qx$ with $Q \geq 0$ is $\lambda_{\max}(Q)$ -Lipschitz.

$$\text{recall } \|f(x) - f(y)\| = \|(x-y)^T Q^2 (x-y)\| \leq (\lambda_{\max}(Q^2))^{1/2} \|x-y\|.$$

remark: if $Q \neq 0$. then is $\sigma_{\max}(Q)$ -Lipschitz (singular value).

L -smooth: $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is L -smooth if ∇f is L -Lipschitz. i.e.,

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

Example. $f(x) = x^T Q x$ with $Q \succeq 0$ is $2\lambda_{\max}(Q)$ -smooth.

Lemma. a twice continuous differentiable $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is L -smooth

if $-L I \leq \nabla^2 f(x) \leq L I, \forall x$. or equivalently $|\lambda| \leq L, \forall$ eigenvalues λ .

Proof. " \Leftarrow ". assume $-L I \leq \nabla^2 f(x) \leq L I, \forall x$. $\exists z = x + \xi(y-x)$.

By the mean value theorem. $\nabla f(x) - \nabla f(y) = \nabla^2 f(z)(x-y)$

$$\|\nabla f(x) - \nabla f(y)\| = \|\nabla^2 f(z)(x-y)\| = (x-y)^T (\nabla^2 f(z))^2 (x-y) \leq L \|x-y\|.$$

" \Rightarrow ". assume f is L -smooth. fix $d \in \mathbb{R}^n$. let $g(\xi) = \nabla f(x + \xi d)$.

$$\Rightarrow \|g(\xi) - g(\omega)\| \leq \xi L \|d\| \Rightarrow \|(g(\xi) - g(\omega))/\xi\| \leq L \|d\|.$$

taking the limit and using the chain rule. $g'(\omega) = \nabla^2 f(x) \cdot d$

$$\Rightarrow \forall d \in \mathbb{R}^n \quad \|\nabla^2 f(x) d\| \leq L \|d\| \Rightarrow |\lambda| \leq L, \forall \text{ eigenvalue } \lambda. \quad \square$$

In particular, if f convex. f is L -smooth iff $\nabla^2 f(x) \leq L I$.

Lemma. f is L -smooth $\Rightarrow f(y) \leq f(x) + \nabla f(x)^T (y-x) + \frac{L}{2} \|x-y\|^2$.

Proof. Let $z(\xi) = x + \xi(y-x)$. $g(\xi) = f(z(\xi))$. so $g'(\xi) = \nabla f(z)^T (y-x)$

$$|g'(\xi) - g'(\omega)| = (\nabla f(z(\xi)) - \nabla f(z(\omega)))^T (y-x) \leq L \xi \|y-x\|^2.$$

$$f(y) - f(x) - \nabla f(x)^T (y-x) = g(1) - g(0) - g'(0) = \int_0^1 g'(\xi) - g'(0) d\xi \leq \frac{L}{2} \|y-x\|^2.$$

\square .

Recall our goal is to bound $\sum_{k=0}^{T-1} \|\nabla f(x_k)\|^2$.

$$f \text{ is } L\text{-smooth} \Rightarrow f(x_{k+1}) \leq f(x_k) + \nabla f(x_k)^T (x_{k+1} - x_k) + \frac{L}{2} \|x_{k+1} - x_k\|^2.$$

$$\Rightarrow f(x_{k+1}) \leq f(x_k) - t \|\nabla f(x_k)\|^2 + \frac{t^2 L}{2} \|\nabla f(x_k)\|^2. \text{ choose } 0 < t < 1/L.$$

$$\Rightarrow f(x_{k+1}) \leq f(x_k) - \frac{t}{2} \|\nabla f(x_k)\|^2. \Rightarrow \frac{t}{2} \sum_{k=0}^{T-1} \|\nabla f(x_k)\|^2 \leq f(x_0) - f(x_T).$$

$$\sum_{k=0}^{T-1} (f(x_k) - f(x^*)) \leq \frac{1}{2t} (\|x_0 - x^*\|^2 - \|x_T - x^*\|^2) + \frac{t}{2} \sum_{k=0}^{T-1} \|\nabla f(x_k)\|^2.$$

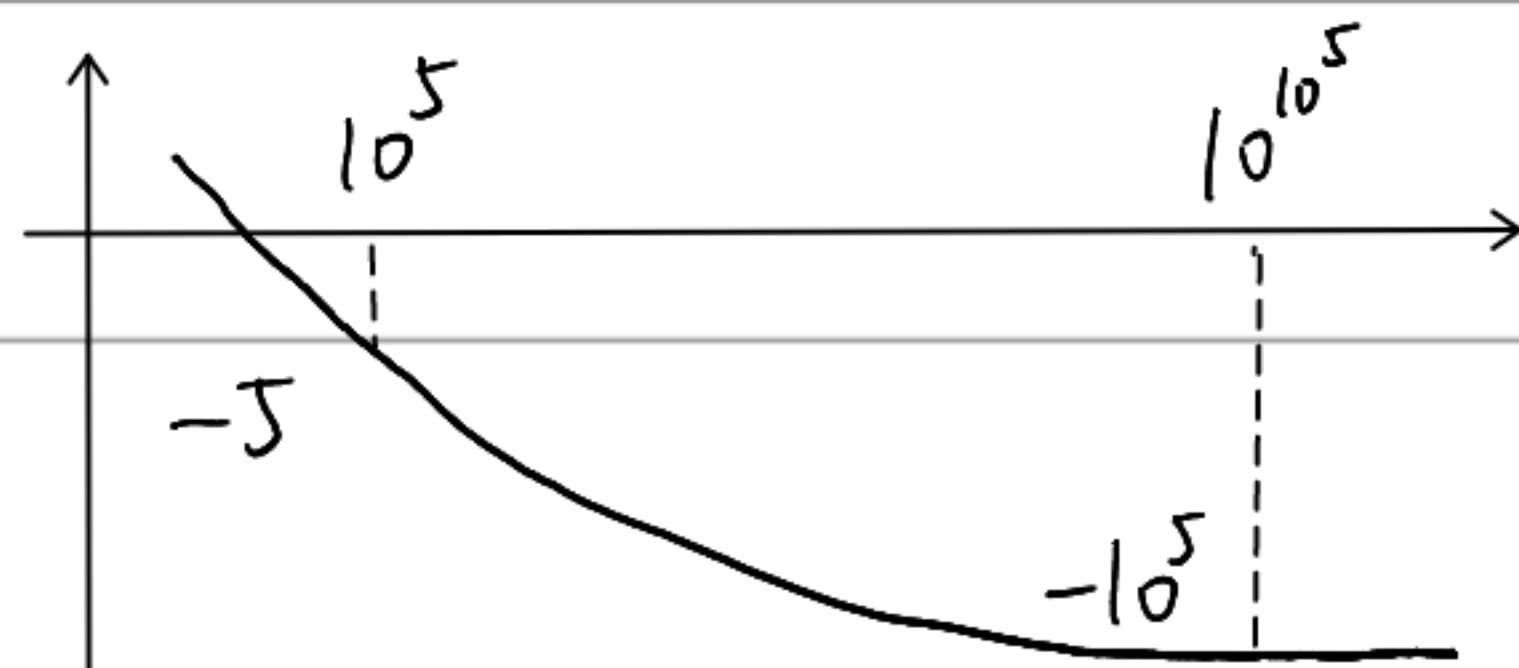
$$\Rightarrow \sum_{k=0}^{T-1} (f(x_k) - f(x^*)) \leq \frac{1}{2t} \|x_0 - x^*\|^2 + f(x_0) - f(x_T).$$

$$\Rightarrow \sum_{k=1}^T (f(x_k) - f(x^*)) \leq \frac{1}{2t} \|x_0 - x^*\|^2. \Rightarrow f(x_T) - f(x_0) \leq \frac{\|x_0 - x^*\|^2}{2tT}.$$

Remark: rate of convergence is $O(1/T)$; $T = O(1/\epsilon)$ to get ϵ -approximation.

Consider the following function:

$$f(x) = -\log x. \quad x < 10^{10^5} \text{ and } -10^5 \text{ on } \omega.$$



$$f'(x) = 1/x. \quad x_{k+1} \leftarrow x_k + t/x_k. \quad \text{stop at } |f'(x)| < 10^{-5}. \text{ Bad!}$$