

Optimization for Machine Learning

CSCI-599

Lecture 2: Gradient Descent

Sai Praneeth Karimireddy

USC – <https://spkreddy.org/optmlspring2025.html>

January 15, 2025



$$\min_{x \in \mathbb{R}^d} f(x) \rightarrow \text{CVX}$$

$$f^* = \inf_x f(x)$$

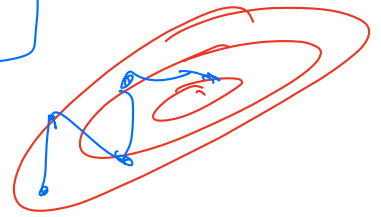
opt 1st $\nabla f(x^*) = 0$ $\| \nabla f(x^*) \|_2 \leq \epsilon$

2nd $f(x^*) - f^* \leq 0 ?$

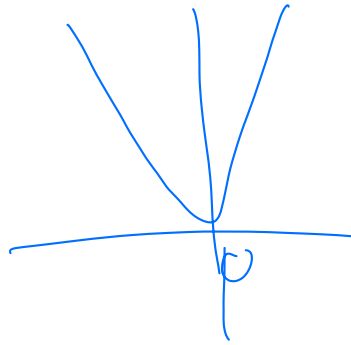
\hat{x} Suppose f^* exists
 $(f(\hat{x}) - f^* \leq \epsilon) \checkmark$

Suppose $\exists x^*$ s.t. $f(x^*) = f^*$

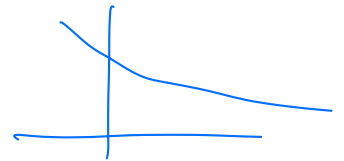
$$\| \hat{x} - x^* \|_2 \leq \epsilon$$



(a)

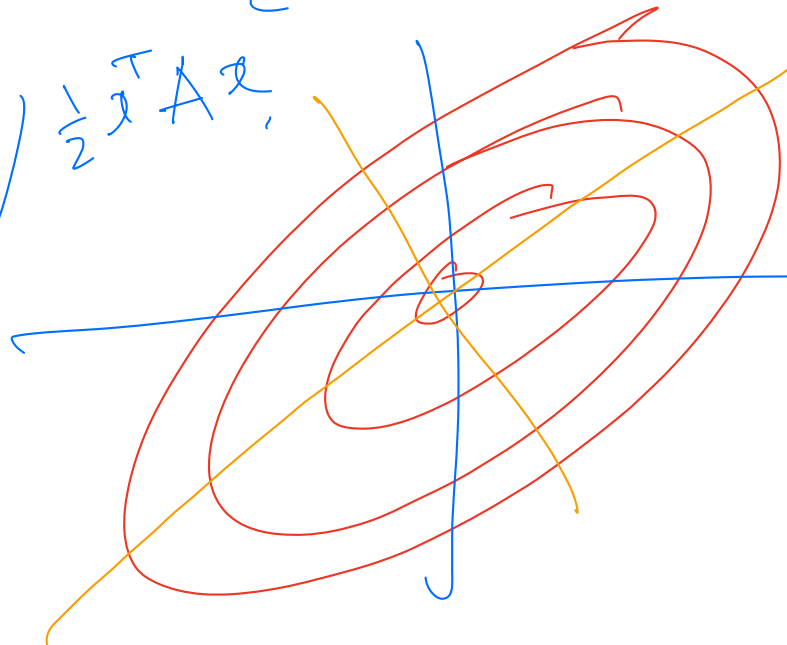
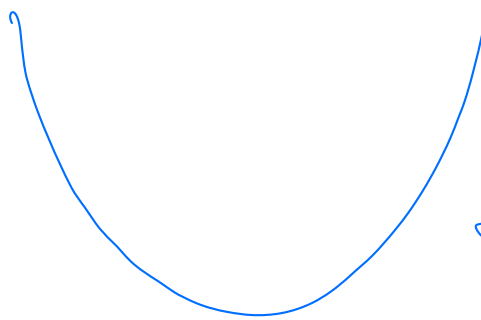


$$\inf_x f(x) = x$$



$$e^{-x}$$

$$\frac{1}{2} x^T A x$$



By Taylor's Expansion

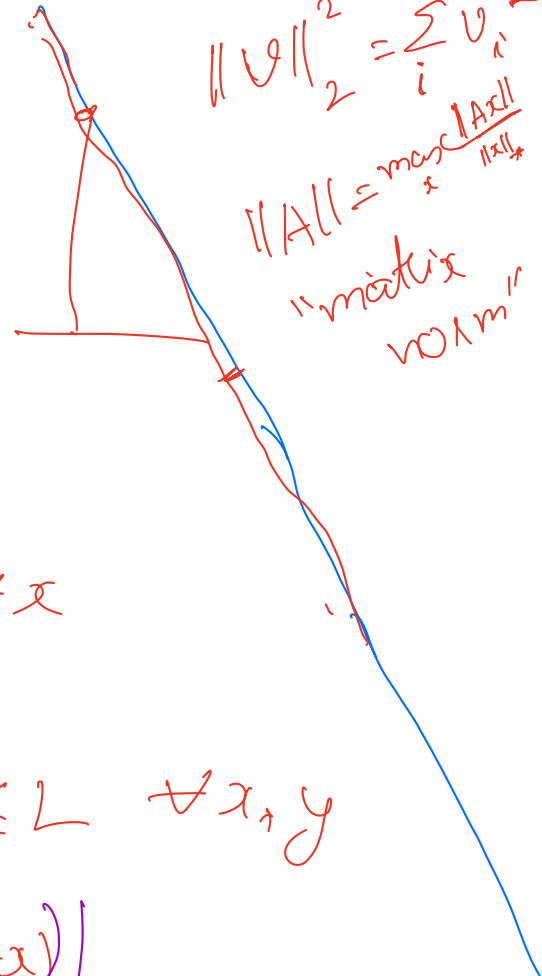
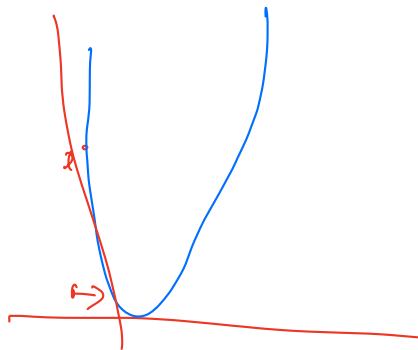
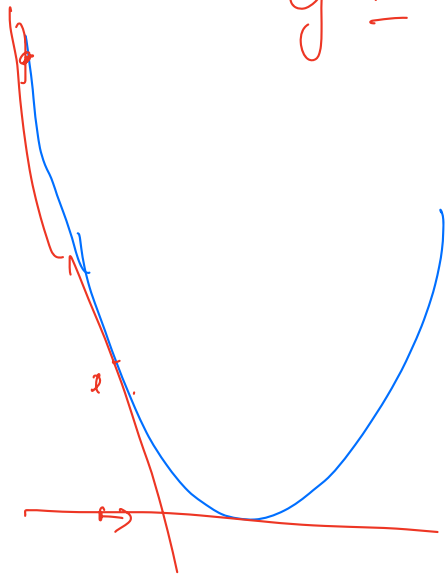
for a "close enough" y ,

$$f(y) \approx f(x) + \nabla f(x)^T (y-x)$$

$\nearrow -\nabla f(x)$

$$y = x - \gamma \nabla f(x)$$

\hookrightarrow step-size, learning rate



$$\|v\|_2^2 = \sum_i v_i^2$$

$$\|A\| = \max_x \frac{\|Ax\|}{\|x\|}$$

"matrix norm"

Def: Smoothness

2nd order $\| \nabla^2 f(x) \|_2 \leq L \quad \forall x$ \rightarrow norm

1st order $\frac{\| \nabla f(y) - \nabla f(x) \|_2}{\|y-x\|_2} \leq L \quad \forall x, y$

0th order $\frac{|f(y) - (f(x) + \nabla f(x)^T (y-x))|}{\|y-x\|_2^2} \leq \frac{L}{2} \quad \forall x, y$

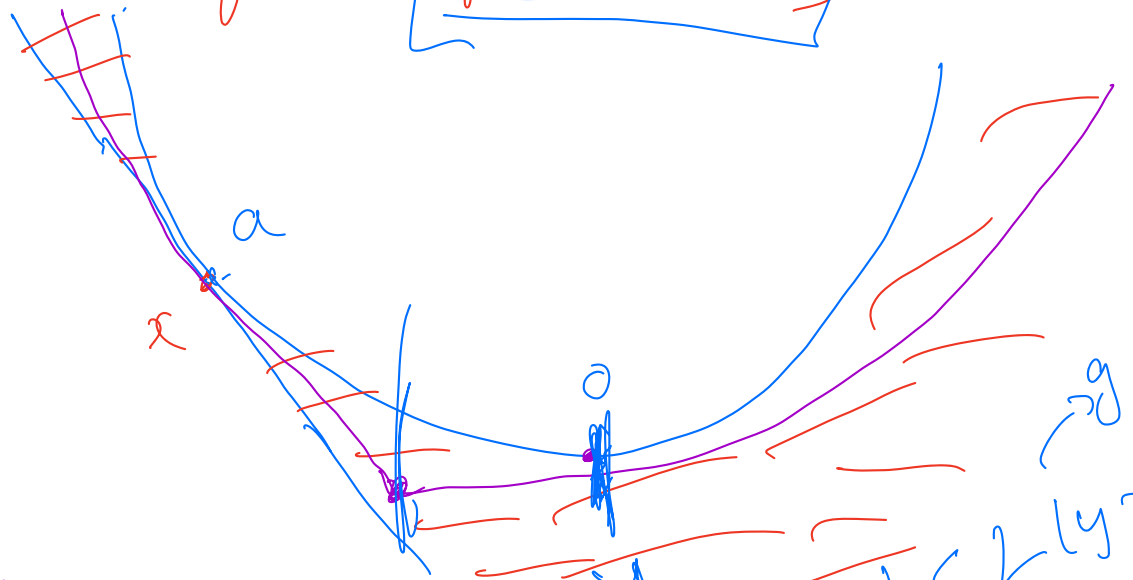
$$\frac{f(y) - f(x)}{\|y-x\|_2} - \frac{\nabla f(x)^T (y-x)}{\|y-x\|_2} \approx \frac{\| \nabla f(y) - \nabla f(x) \|}{\|y-x\|}$$

$$f(y) \leq \underbrace{f(x) + \nabla f(x)^T (y-x)}_{g(y)} + \frac{L}{2} \|y-x\|^2 \quad \forall x, y \in \mathcal{D}$$

$$\nabla g(y) = \nabla f(x) + L(y-x), \nabla^2 g(y) = L$$

f is convex & L -smooth $\Rightarrow \nabla g(x) = \nabla f(x)$

$$0 \leq f(y) - \underbrace{(f(x) + \nabla f(x)^T (y-x))}_{g(y)} \leq \frac{L}{2} \|y-x\|^2$$



$$\frac{\|\nabla f(y) - \nabla f(x)\|}{\|y-x\|} \leq L$$

$$|a| \leq L \|y-x\|$$

$$|a-\delta| = \frac{|a+\delta|}{f(x)}$$

$$\frac{|a-\delta|}{\|y-x\|} \geq L$$

Claim

Minimizer of the upper quadratic never "overshoots"

$\Rightarrow \delta = \frac{1}{L}$ never "overshoots"

$$y = x - \frac{1}{L} \nabla f(x)$$

$$f(x) = \frac{1}{2} x^T A x - b^T x$$

np-inv
np-solve

$$\nabla f(x) = \frac{1}{2} (A + A^T) x - b$$

A is symmetric
& inv
 $x^* = A^{-1} b$

$$\nabla f(x) = Ax - b$$

$$x_0 = 0$$

$$x_{t+1} = x_t - \gamma (Ax_t - b) \quad - \textcircled{1}$$

setting γ

$$\nabla^2 f(x) = A$$

$$L = \|A\|_2$$

$$\Rightarrow \gamma = \frac{1}{\|A\|_2}$$

Unrolling $\textcircled{1} \Rightarrow$

$$x_{t+1} = (I - \gamma A) x_t + \gamma b$$

$$x_{t+1} = \underline{\hspace{10em}}$$

$$\|x_{t+1} - A^{-1} b\|_2 \leq ???$$

Claim if f is convex and L -smooth

$$x_{t+1} = x_t - \frac{1}{L} \nabla f(x_t)$$

$$\textcircled{1} \quad f(x_t) - f^* \leq O\left(\frac{1}{t}\right)$$

$\textcircled{2}$ if f is μ -strongly-convex \Rightarrow

$$f(x_t) - f^* \leq O\left(\exp\left(-\frac{\mu}{L}t\right)\right)$$

Crash Course on Spectral Norm

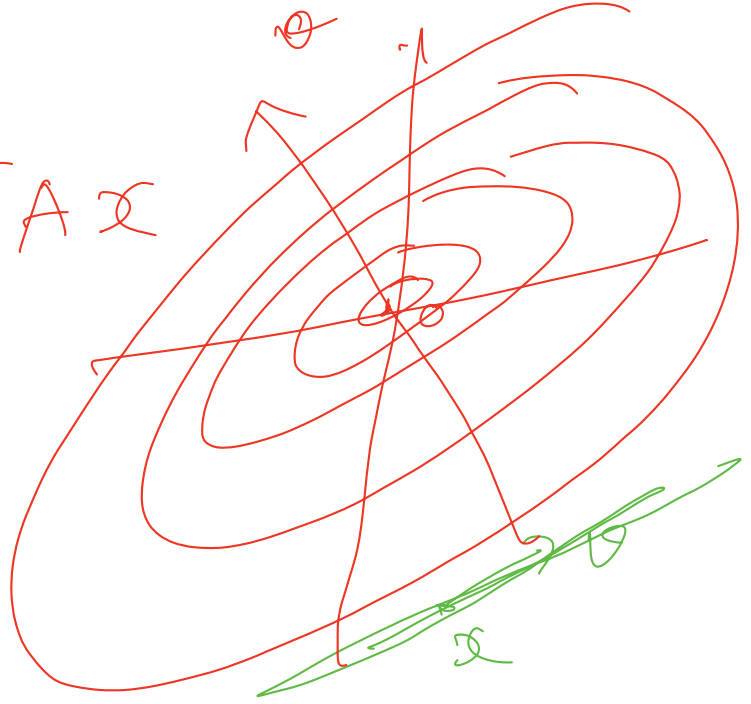
→ symmetric p.s.d. finite dim.

$$A \quad d \times d$$

$$f(x) = \frac{1}{2} x^T A x$$

$$\nabla f(x) = Ax$$

$$\nabla^2 f(x) = A$$



$$\rightarrow g_\theta(\lambda) = f(x + \lambda \theta) =$$

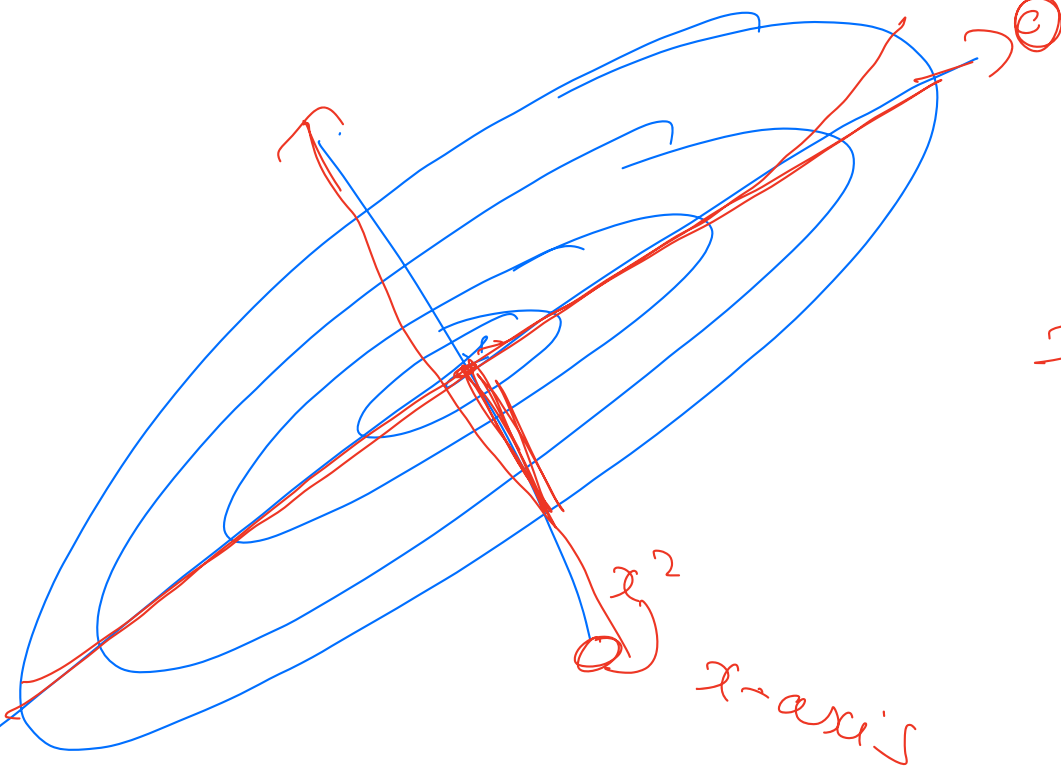
$$\frac{1}{2} x^T A x + Ax^T \theta \lambda + \frac{1}{2} \lambda^2 \underbrace{\theta^T A \theta}_{=}$$

$$\rightarrow \frac{d g_\theta(\lambda)}{d \lambda} = Ax^T \theta$$

$$\frac{d^2 g_\theta(\lambda)}{d \lambda^2} =$$

$$\theta^T A \theta$$

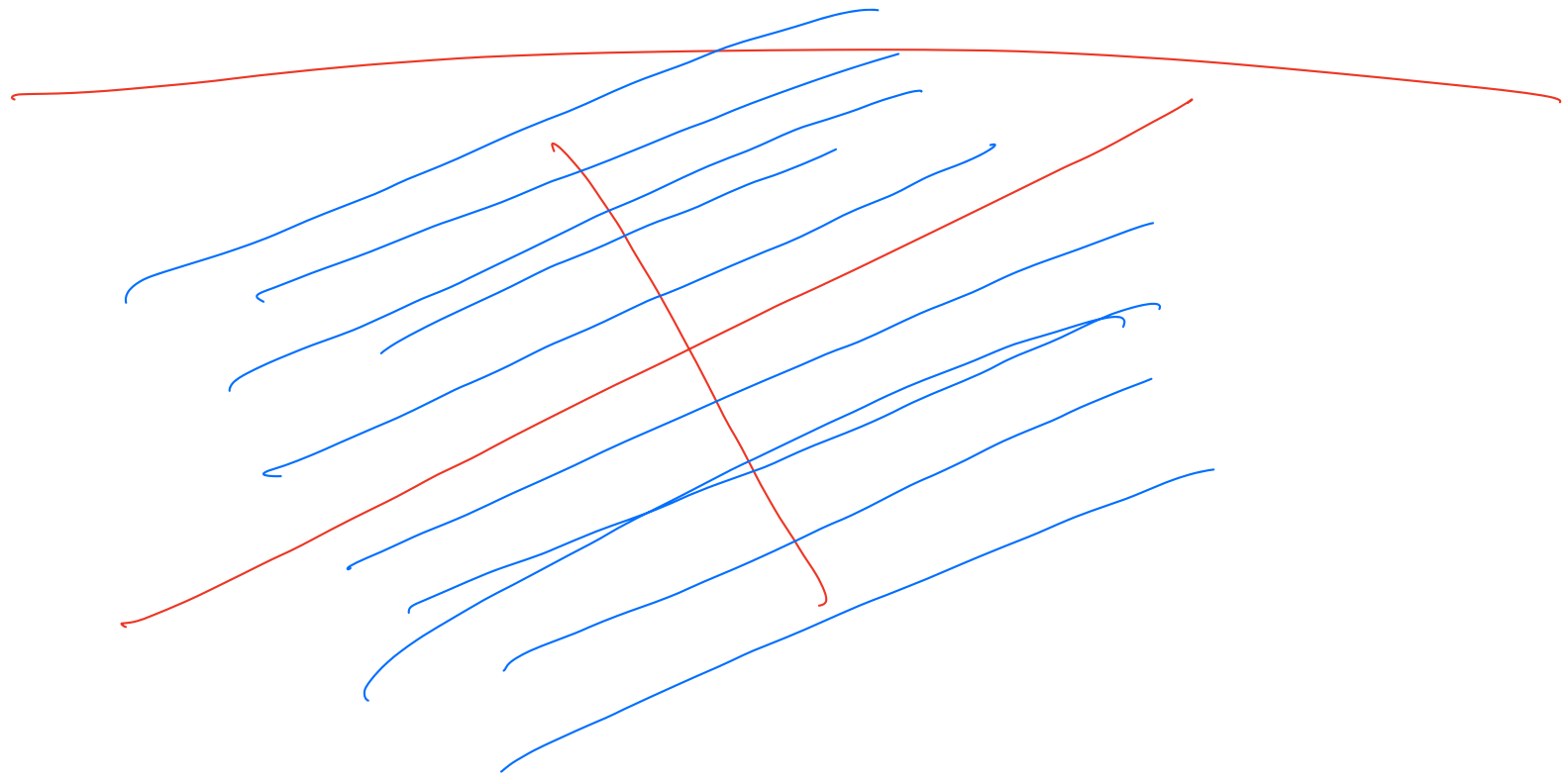
min eigenvalue



$$\|A\|_2 = 2$$

= max eigenvalue

min eigenvalue



Jan 27

A convex function f is L -smooth if

(i) $\|\nabla^2 f(x)\|_2 \leq L \quad \forall x$

(ii) $\|\nabla f(x) - \nabla f(y)\|_2 \leq L\|x-y\|_2 \quad \forall x, y$

(iii) $f(y) \leq f(x) + \nabla f(x)^T(y-x) + \frac{L}{2}\|y-x\|_2^2$

\Rightarrow $x_{t+1} = x_t - \frac{1}{L} \nabla f(x_t)$

Thm Claim

$f(x_{t+1}) - f^* \leq$ some decreasing function of t

$f(x) = \frac{1}{2} x^T (A^T A) x - b^T x$

$x_{t+1} = x_t - \gamma (A^T A x_t - b)$

$A^T A$ is symm, pd

$\Rightarrow x^* = (A^T A)^{-1} b$

$\therefore \nabla f(x) = (A^T A)x - b$

at $x^* \quad \nabla f(x^*) = 0$

Observation 1

$$f(x) - f^* = \frac{1}{2} \|A(x - x^*)\|^2 \rightarrow A^T A (x - x^*)$$

$$\nabla f(x) = (A^T A)x - b = A^T A x - \cancel{A^T A(A^T A)^{-1} b} = A^T A x - b$$

$$\|x_{t+1} - x^*\|_2^2 = \|x_t - \gamma(A^T A x_t - b) - x^*\|_2^2$$

$$= \|x_t - x^*\|_2^2 - 2\gamma (A^T A x_t - b)^T (x_t - x^*) \leq 0$$

$$+ \gamma^2 \|A^T A x_t - b\|_2^2$$

Attempt 1

$$-2\gamma (A^T A x_t - b)^T (x_t - (A^T A)^{-1} b)$$

$$-2\gamma (A^T A x_t - b)^T (A^T A)^{-1} (A^T A) (x_t - (A^T A)^{-1} b)$$

$$-2\gamma \left[(A^T A x_t - b)^T \right] \left[(A^T A)^{-1} \right] \left[(A^T A x_t - b) \right] \geq 0$$

$$\leq -2\gamma \lambda_{\min} \|A^T A x_t - b\|^2$$

λ_{\min}
P.d M
 $M \succeq \lambda_{\min} I$

$$\|x_{\text{reg}} - x^*\|^2 \leq \|x_{\text{reg}} - x^*\|^2 \quad (M - \lambda_{\min} I) \succeq 0$$

$$- 2\sigma \lambda_{\min} \|A^T A x_{\text{reg}} - b\|^2$$

$$+ \sigma^2 \|A^T A x_{\text{reg}} - b\|^2$$

Attempt 2

$$f(x_{\text{reg}}) - f^* = \|A(x_{\text{reg}} - x^*)\|^2$$

$$\|A(x_{\text{reg}} - (A^T A)^{-1} b)\|^2$$

$$-2\sigma \underbrace{(A^T A x_{\text{reg}} - b)^T}_{\downarrow} \underbrace{(x_{\text{reg}} - (A^T A)^{-1} b)}_{\downarrow} + \sigma^2 \|A^T A x_{\text{reg}} - b\|^2$$

$$(A^T A x_{\text{reg}} - b)^T (A^T A)^{-1} (A^T A x_{\text{reg}} - b)$$

$$= A(x_{\text{reg}} - (A^T A)^{-1} b)^T (A^T A)^{-1} A(x_{\text{reg}} - (A^T A)^{-1} b)$$

Strong Convexity

- $\nabla^2 f(x) \geq 0$ - convex

- $\nabla^2 f(x) > 0$ - strict convex

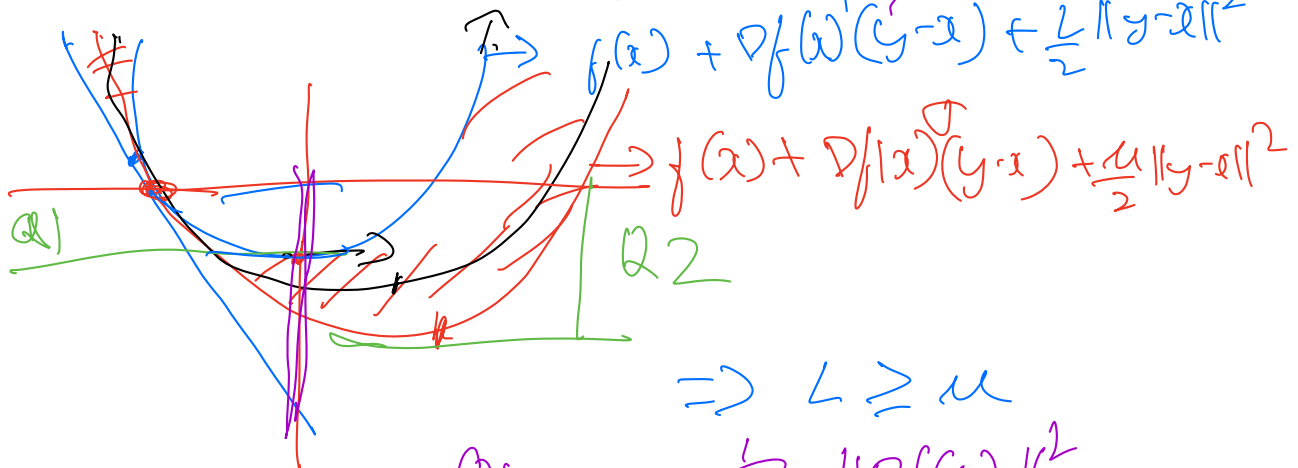
- $\nabla^2 f(x) \geq \mu I$ - μ -strongly convex
 $\mu > 0$

- $\langle \nabla f(y) - \nabla f(x), y - x \rangle \geq \mu \|y - x\|^2$

$\Rightarrow \|\nabla f(y) - \nabla f(x)\| \geq \mu \|y - x\|$

$\frac{L}{2} \|y - x\|^2 \Rightarrow \left| f(y) - \left[f(x) + \nabla f(x)^T (y - x) \right] \right| \geq \frac{\mu}{2} \|y - x\|^2$

f is L -smooth $\Leftarrow \mu$ -strongly convex $\rightarrow y = x - \frac{1}{L} \nabla f(x)$



$\Rightarrow L \geq \mu$

$Q1 = \frac{1}{2L} \|\nabla f(x)\|^2$

$Q2 = \frac{1}{2\mu} \|\nabla f(x)\|^2$

Claim: a) $f(x_{t+1}) \leq f(x_t) - \textcircled{Q1}$

b) $f(x^*) - f(x_t) \leq -\textcircled{Q2}$ \square

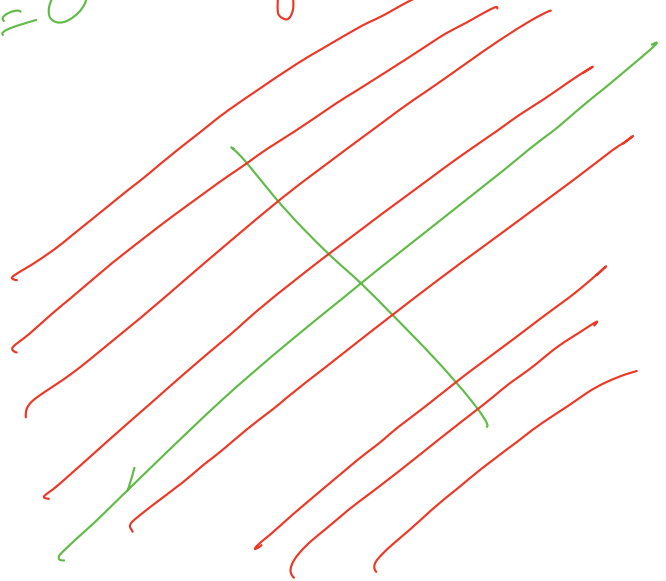
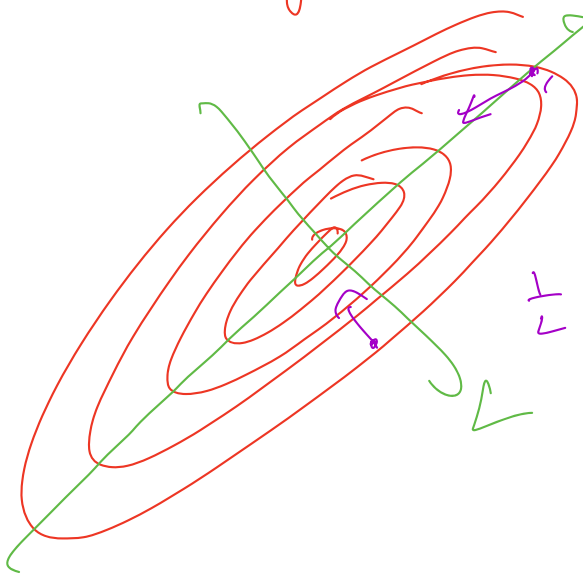
$$\begin{aligned} f(x_{t+1}) - f^* &\leq f(x_t) - f^* - \frac{1}{2L} \|\nabla f(x_t)\|^2 \\ &= f(x_t) - f^* - \frac{\mu}{L} \left(\frac{1}{2\mu} \|\nabla f(x_t)\|^2 \right) \\ &\leq f(x_t) - f^* - \frac{\mu}{L} (f(x_t) - f^*) \end{aligned}$$

$$\begin{aligned} (f(x_{t+1}) - f^*) &\leq \left(1 - \frac{\mu}{L}\right) (f(x_t) - f^*) \\ &\leq \left(1 - \frac{\mu}{L}\right)^t (f(x_0) - f^*) \end{aligned}$$

Vis

$f(x) = \frac{1}{2} x^T A x \rightarrow \text{p.s.d}$
 $\mu = 0$

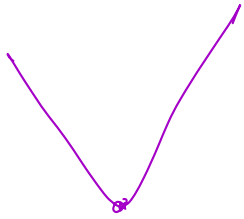
$D^2 f(x) = A$



f is non-smooth, but convex

$$\|x_{t+1} - x^*\|^2 = \|x_t - \sigma g_t - x^*\|^2$$

$$x_{t+1} = x_t - \sigma \underbrace{\frac{\nabla f(x_t)}{\|g_t\|}}_{g_t}$$



$$= \|x_t - x^*\|^2$$

$$- 2\sigma \underbrace{g_t^T (x_t - x^*)}_{\geq 0}$$

$$+ \sigma^2 \|g_t\|^2$$

Claim

$$\nabla f(x_t)^T (x_t - x^*) \geq f(x_t) - f^*$$

Proof

$$f(y) \geq f(x_t) + \nabla f(x_t)^T (y - x_t)$$

$$y = x^*$$

$$f^* \geq f(x_t) + \nabla f(x_t)^T (x^* - x_t)$$

$$\Rightarrow \|x_{t+1} - x^*\|^2 \leq \|x_t - x^*\|^2 - 2\sigma (f(x_t) - f^*) + \sigma^2 \underbrace{\|\nabla f(x_t)\|_2^2}_{\leq \sigma^2}$$

$$\frac{1}{T} \sum_k 2\sigma (f(x_k) - f^*) \leq \frac{1}{2} \sum_k (\|x_k - x^*\|^2 - \|x_{k+1} - x^*\|^2) + T\sigma^2 G^2$$

$$= \frac{1}{2} (\|x_0 - x^*\|^2 - \|x_{T+1} - x^*\|^2) + T\sigma^2 G^2$$

$$\gamma = \sqrt{\frac{D^2}{T G^2}} \leq \frac{1}{2} (\|x_0 - x^*\|^2 + T\sigma^2 G^2) \leq \frac{D^2}{\gamma} + T\sigma^2 G^2 = 2DG\sqrt{T}$$

$$\frac{2}{T} \sum_k f(x_k) - f^* \leq \frac{2DG\sqrt{T}}{T} = \frac{2DG}{\sqrt{T}}$$

$$\bar{x}_T = \frac{1}{T} \sum_k x_k$$

$$f(\bar{x}_T) - f^* \leq \frac{1}{T} \sum_k f(x_k) - f^* \leq \frac{2DG}{\sqrt{T}}$$

Assumption $\|x_0 - x^*\|^2 \leq D$
 $\|\nabla f(x_k)\|^2 \leq G^2$

$$\|\nabla f(x)\|_2 \leq C \quad \forall x$$

$$|f(y) - f(x)| \leq C \|y - x\|_2$$

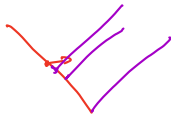
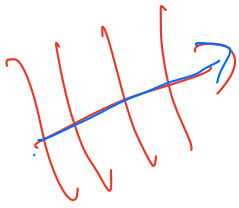
C-Lip

Recap

① f is L -smooth + μ -S.C
 $\Rightarrow f(x_t) - f^* \leq \mathcal{O}\left(\left(1 - \frac{\mu}{L}\right)^t\right)$

② f is G -Lip + CVX +
 $\|x_0 - x^*\|_2 \leq 1 \Rightarrow$
 $f(x_T) - f^* \leq \mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$

③ f is L -smooth + CVX \Rightarrow
 $f(x_T) - f^* \leq \mathcal{O}\left(\frac{1}{T}\right)$

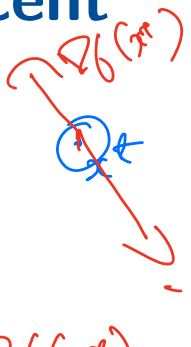


Chapter 2

Gradient Descent

$$x^{t+1} = x^t - \nabla^2 f(x^t)^{-1} \nabla f(x^t)$$

$$\nabla f(x^t) \neq 0$$



$$f(y) \approx f(x^t) + \nabla f(x^t)^T (y - x^t)$$

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

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

The Algorithm

$$f(y) = f(x^*) + f'(x^*)(y-x^*) + \frac{1}{2} f''(x^*) \cdot (y-x^*)^2$$

Get near to a minimum \mathbf{x}^* / close to the optimal value $f(\mathbf{x}^*)$?

(Assumptions: $f : \mathbb{R}^d \rightarrow \mathbb{R}$ convex, differentiable, has a global minimum \mathbf{x}^*)

Goal: Find $\mathbf{x} \in \mathbb{R}^d$ such that

$$f(\mathbf{x}) - f(\mathbf{x}^*) \leq \varepsilon.$$

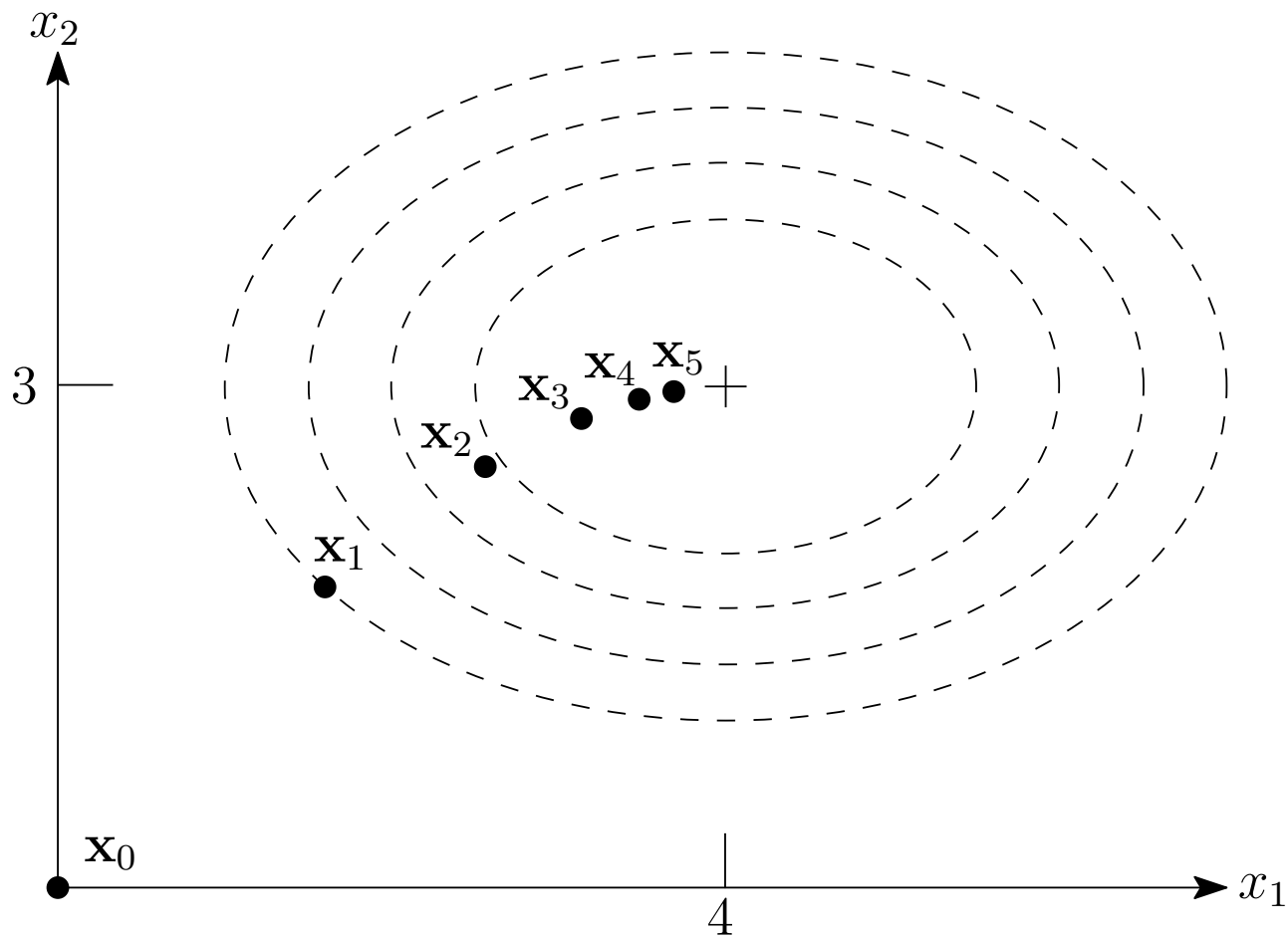
Note that there can be several global minima $\mathbf{x}_1^* \neq \mathbf{x}_2^*$ with $f(\mathbf{x}_1^*) = f(\mathbf{x}_2^*)$.

Iterative Algorithm: choose $\mathbf{x}_0 \in \mathbb{R}^d$.

$$\mathbf{x}_{t+1} := \mathbf{x}_t - \gamma \nabla f(\mathbf{x}_t),$$

for **timesteps** $t = 0, 1, \dots$, and **stepsize** $\gamma \geq 0$.

Example



$$f(x_1, x_2) := 2(x_1 - 4)^2 + 3(x_2 - 3)^2, \mathbf{x}_0 := (0, 0), \gamma := 0.1$$

Vanilla analysis

How to bound $f(\mathbf{x}_t) - f(\mathbf{x}^*)$?

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \gamma \mathbf{g}_t$$

- ▶ Abbreviate $\mathbf{g}_t := \nabla f(\mathbf{x}_t)$ (gradient descent: $\mathbf{g}_t = (\mathbf{x}_t - \mathbf{x}_{t+1})/\gamma$).

$$\mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*) =$$

$$\frac{1}{2\gamma} \cdot \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 = \frac{1}{2\gamma} \left(\|\mathbf{x}_t - \gamma \mathbf{g}_t - \mathbf{x}^*\|^2 \right)$$
$$\frac{1}{2\gamma} \left(\|\mathbf{x}_t - \mathbf{x}^*\|^2 - 2\gamma \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*) + \gamma^2 \|\mathbf{g}_t\|^2 \right)$$

Vanilla analysis

How to bound $f(\mathbf{x}_t) - f(\mathbf{x}^*)$?

- ▶ Abbreviate $\mathbf{g}_t := \nabla f(\mathbf{x}_t)$ (gradient descent: $\mathbf{g}_t = (\mathbf{x}_t - \mathbf{x}_{t+1})/\gamma$).

$$\mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*) = \frac{1}{\gamma} (\mathbf{x}_t - \mathbf{x}_{t+1})^\top (\mathbf{x}_t - \mathbf{x}^*).$$

- ▶ Apply $2\mathbf{v}^\top \mathbf{w} = \|\mathbf{v}\|^2 + \|\mathbf{w}\|^2 - \|\mathbf{v} - \mathbf{w}\|^2$ to rewrite

$$\mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*) =$$

Vanilla analysis

How to bound $f(\mathbf{x}_t) - f(\mathbf{x}^*)$?

$$g_t = \nabla f(\mathbf{x}_t)$$

- ▶ Abbreviate $\mathbf{g}_t := \nabla f(\mathbf{x}_t)$ (gradient descent: $\mathbf{g}_t = (\mathbf{x}_t - \mathbf{x}_{t+1})/\gamma$).

$$\mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*) = \frac{1}{\gamma} (\mathbf{x}_t - \mathbf{x}_{t+1})^\top (\mathbf{x}_t - \mathbf{x}^*).$$

- ▶ Apply $2\mathbf{v}^\top \mathbf{w} = \|\mathbf{v}\|^2 + \|\mathbf{w}\|^2 - \|\mathbf{v} - \mathbf{w}\|^2$ to rewrite

$$\begin{aligned} \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*) &= \frac{1}{2\gamma} (\|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 + \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2) \\ &= \frac{\gamma}{2} \|\mathbf{g}_t\|^2 + \frac{1}{2\gamma} (\|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2) \end{aligned}$$

- ▶ Sum this up over the first T iterations:

$$\sum_{t=0}^{T-1} \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*) =$$

$$\left(\|\mathbf{x}_0 - \mathbf{x}^*\|^2 - \|\mathbf{x}_T - \mathbf{x}^*\|^2 \right) / \gamma$$

Vanilla analysis

How to bound $f(\mathbf{x}_t) - f(\mathbf{x}^*)$?

- ▶ Abbreviate $\mathbf{g}_t := \nabla f(\mathbf{x}_t)$ (gradient descent: $\mathbf{g}_t = (\mathbf{x}_t - \mathbf{x}_{t+1})/\gamma$).

$$\mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*) = \frac{1}{\gamma} (\mathbf{x}_t - \mathbf{x}_{t+1})^\top (\mathbf{x}_t - \mathbf{x}^*).$$

- ▶ Apply $2\mathbf{v}^\top \mathbf{w} = \|\mathbf{v}\|^2 + \|\mathbf{w}\|^2 - \|\mathbf{v} - \mathbf{w}\|^2$ to rewrite

$$\begin{aligned} \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*) &= \frac{1}{2\gamma} (\|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 + \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2) \\ &= \frac{\gamma}{2} \|\mathbf{g}_t\|^2 + \frac{1}{2\gamma} (\|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2) \end{aligned}$$

- ▶ Sum this up over the first T iterations:

$$\sum_{t=0}^{T-1} \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*) = \frac{\gamma}{2} \sum_{t=0}^{T-1} \|\mathbf{g}_t\|^2 + \frac{1}{2\gamma} (\|\mathbf{x}_0 - \mathbf{x}^*\|^2 - \|\mathbf{x}_T - \mathbf{x}^*\|^2)$$

Vanilla analysis II

Use first-order characterization of convexity: $f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}), \forall \mathbf{x}, \mathbf{y}$

▶ with $\mathbf{x} = \mathbf{x}_t, \mathbf{y} = \mathbf{x}^*$:

Vanilla analysis II

Use first-order characterization of convexity: $f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}), \forall \mathbf{x}, \mathbf{y}$

► with $\mathbf{x} = \mathbf{x}_t, \mathbf{y} = \mathbf{x}^*$:

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*)$$

giving

$$\sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq$$

Vanilla analysis II

Use first-order characterization of convexity: $f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}), \forall \mathbf{x}, \mathbf{y}$

- ▶ with $\mathbf{x} = \mathbf{x}_t, \mathbf{y} = \mathbf{x}^*$:

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*)$$

giving

$$\sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq \frac{\gamma}{2} \sum_{t=0}^{T-1} \|\mathbf{g}_t\|^2 + \frac{1}{2\gamma} \|\mathbf{x}_0 - \mathbf{x}^*\|^2,$$

(Handwritten green annotations: a large circle around the sum of squared gradients, and a smaller circle around the initial distance term. Below the first term is $\leq TB^2$ and below the second is $\leq R^2$)

an upper bound for the **average error** $f(\mathbf{x}_t) - f(\mathbf{x}^*)$ over the steps

- ▶ last iterate is not necessarily the best one
- ▶ stepsize is crucial

$$\gamma = \frac{R}{B\sqrt{T}} \Rightarrow \leq \frac{\gamma TB^2}{2} + \frac{R^2}{2\gamma}$$

(Handwritten green annotations: the final inequality above)

Lipschitz convex functions: $\mathcal{O}(1/\varepsilon^2)$ steps

Assume that all gradients of f are bounded in norm.

- ▶ Equivalent to f being Lipschitz (Theorem ??; Exercise ??).
- ▶ Rules out many interesting functions (for example, the “supermodel” $f(x) = x^2$)

Theorem

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex and differentiable with a global minimum \mathbf{x}^* ; furthermore, suppose that $\|\mathbf{x}_0 - \mathbf{x}^*\| \leq R$ and $\|\nabla f(\mathbf{x})\| \leq B$ for all \mathbf{x} . Choosing the stepsize

$$\gamma := \frac{R}{B\sqrt{T}},$$

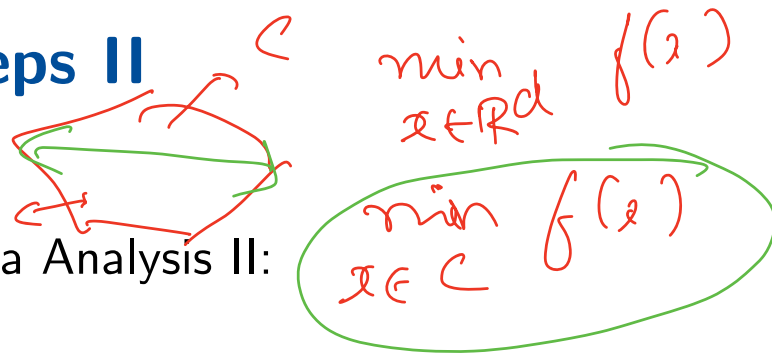
gradient descent yields

$$\frac{1}{T} \sum_{t=0}^{T-1} f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{RB}{\sqrt{T}}.$$

f is B-Lip

$$g(x) = \|x\|, \quad \frac{1}{2}x^T A x - b^T x$$

Lipschitz convex functions: $\mathcal{O}(1/\varepsilon^2)$ steps II



Proof.

- ▶ Plug $\|\mathbf{x}_0 - \mathbf{x}^*\| \leq R$ and $\|\mathbf{g}_t\| \leq B$ into Vanilla Analysis II:

$$\sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq \frac{\gamma}{2} \sum_{t=0}^{T-1} \|\mathbf{g}_t\|^2 + \frac{1}{2\gamma} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \leq \frac{\gamma}{2} B^2 T + \frac{1}{2\gamma} R^2.$$

- ▶ choose γ such that

$$q(\gamma) = \frac{\gamma}{2} B^2 T + \frac{R^2}{2\gamma}$$

is minimized.

- ▶ Solving $q'(\gamma) = 0$ yields the minimum $\gamma = \frac{R}{B\sqrt{T}}$, and $q(R/(B\sqrt{T})) = RB\sqrt{T}$.
- ▶ Dividing by T , the result follows.

$f(x) = \|x\|_C$ is convex \square

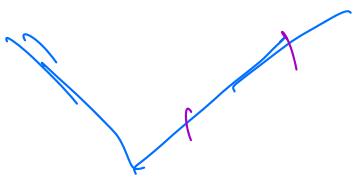
Limitations

$$\|x\|_2 \leq 1$$

- ① How to be sure $\|x_0 - x^*\| \leq R$ w/o knowing x^* ? Need to constrain x .
- ② Gradient ^{may be} undefined for a non-smooth f

- ③ How to know $\|Df(x)\|_2 \leq B \forall x$?

$f(x) = \|x\|_1$, Eg 1 $|f(y) - f(x)| \leq B \|y - x\|_2$



$$(\|y\|_1 - \|x\|_1) \leq B \|y - x\|_2$$

$$\|y\|_1 - \|x\|_1 \leq \|y - x\|_1 \leq \sqrt{2} \|y - x\|_2$$

$$\|z\|_1 \leq \sqrt{2} \|z\|_2$$

Eg 2

$$f(x) = \frac{1}{2} x^T A x$$
$$\|x\|_2 \leq 1$$

$$f(x) = x^2 \quad |x| \leq 1$$
$$Df(x) = 2x \leq 2$$

$$\|Df(x)\|_2 = \|Ax\|_2 \leq \|A\|_2 \|x\|_2 \leq \|A\|_2$$

ToDo

- ① Run GD on constrained domains.
- ② Make GD work when func is not diff.

Lipschitz convex functions: $\mathcal{O}(1/\varepsilon^2)$ steps III

$$T \geq \frac{R^2 B^2}{\varepsilon^2} \quad \Rightarrow \quad \text{average error} \leq \frac{RB}{\sqrt{T}} \leq \varepsilon.$$

Advantages:

- ▶ dimension-independent (no d in the bound)!
- ▶ holds for both average, or best iterate

In Practice:

What if we don't know R and B ? → **Exercise ??** (having to know R can't be avoided)

Smooth functions

“Not too curved”

Definition

Let $f : \text{dom}(f) \rightarrow \mathbb{R}$ be differentiable, $X \subseteq \text{dom}(f)$, $L \in \mathbb{R}_+$. f is called **smooth** (with parameter L) over X if

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2, \quad \forall \mathbf{x}, \mathbf{y} \in X.$$

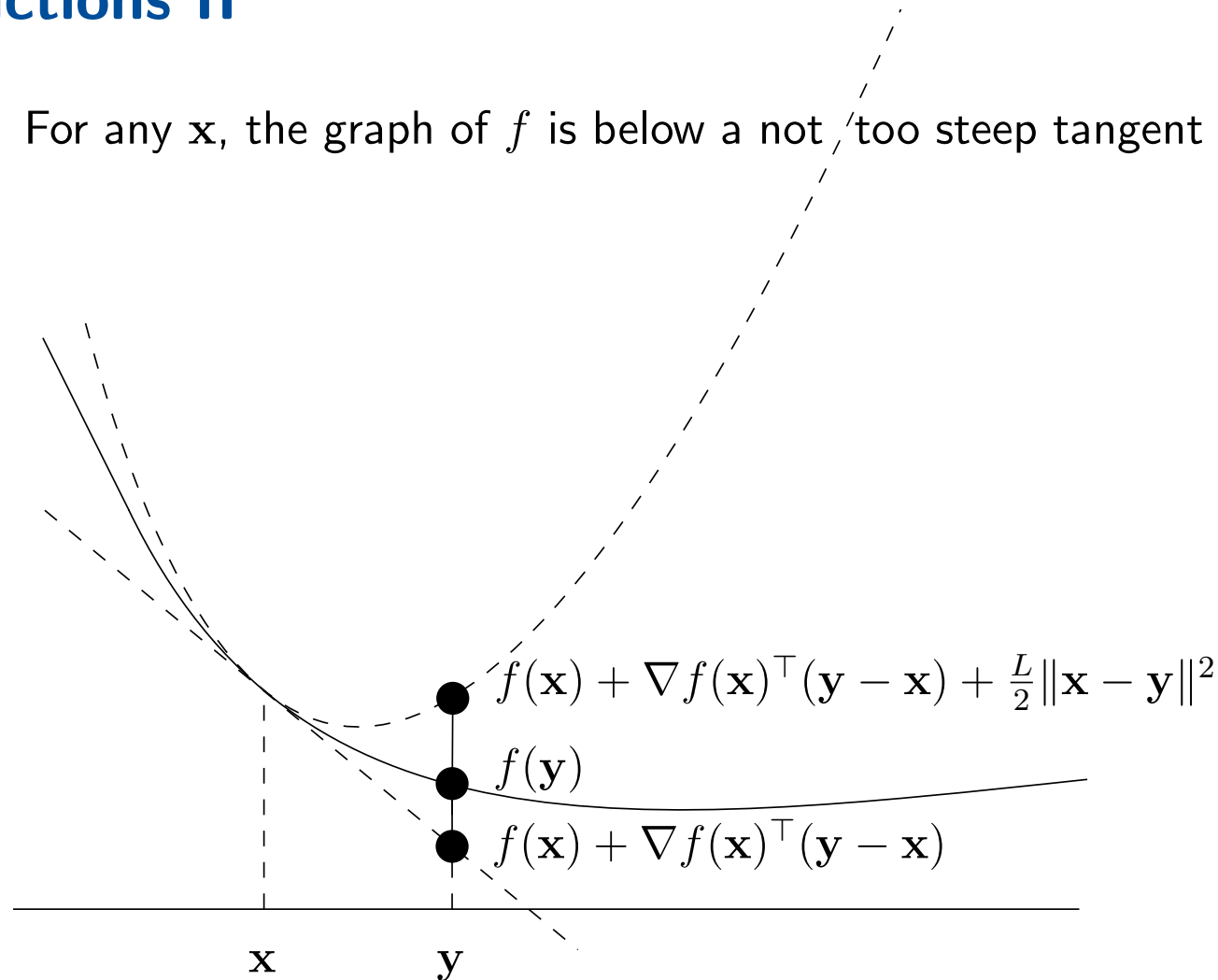
f smooth $:\Leftrightarrow f$ smooth over \mathbb{R}^d .

Definition does not require convexity (useful later)

$$\forall \mathbf{y}, \mathbf{x}, \quad (\nabla f(\mathbf{y}) - \nabla f(\mathbf{x}))^\top (\mathbf{y} - \mathbf{x}) \leq L \|\mathbf{y} - \mathbf{x}\|_2^2$$
$$\|\nabla^2 f(\mathbf{x})\|_2 \leq L \quad \forall \mathbf{x}$$

Smooth functions II

Smoothness: For any \mathbf{x} , the graph of f is below a not too steep tangent paraboloid at $(\mathbf{x}, f(\mathbf{x}))$:



Smooth functions III

- ▶ In general: quadratic functions are smooth (**Exercise ??**).
- ▶ Operations that preserve smoothness (the same that preserve convexity):

Lemma (Exercise ??)

- (i) *Let f_1, f_2, \dots, f_m be functions that are smooth with parameters L_1, L_2, \dots, L_m , and let $\lambda_1, \lambda_2, \dots, \lambda_m \in \mathbb{R}_+$. Then the function $f := \sum_{i=1}^m \lambda_i f_i$ is smooth with parameter $\sum_{i=1}^m \lambda_i L_i$.*
- (ii) *Let f be smooth with parameter L , and let $g(\mathbf{x}) = A\mathbf{x} + \mathbf{b}$, for $A \in \mathbb{R}^{d \times m}$ and $\mathbf{b} \in \mathbb{R}^d$. Then the function $f \circ g$ is smooth with parameter $L\|A\|^2$, where $\|A\|$ is the **spectral norm** of A (Definition ??).*

Smooth vs Lipschitz

- ▶ Bounded gradients \Leftrightarrow Lipschitz continuity of f
- ▶ Smoothness \Leftrightarrow Lipschitz continuity of ∇f (in the convex case).

Lemma

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex and differentiable. The following two statements are equivalent.

- (i) f is smooth with parameter L .
- (ii) $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|$ for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$.

Proof in lecture slides of L. Vandenberghe, <http://www.seas.ucla.edu/~vandenbe/236C/lectures/gradient.pdf>.

Sufficient decrease

Lemma

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be differentiable and smooth with parameter L . With stepsize

$$\gamma := \frac{1}{L},$$

gradient descent satisfies

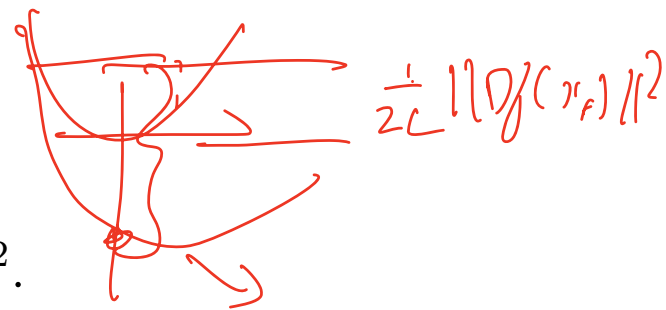
$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2, \quad t \geq 0.$$

Remark

More specifically, this already holds if f is smooth with parameter L over the line segment connecting \mathbf{x}_t and \mathbf{x}_{t+1} .

Sufficient decrease II

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \frac{1}{2L} \nabla f(\mathbf{x}_t) \quad f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2.$$



Proof.

Use smoothness and definition of gradient descent ($\mathbf{x}_{t+1} - \mathbf{x}_t = -\nabla f(\mathbf{x}_t)/L$):

$$\begin{aligned} f(\mathbf{x}_{t+1}) &\leq f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_{t+1} - \mathbf{x}_t) + \frac{L}{2} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \\ &= f(\mathbf{x}_t) - \frac{1}{L} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2 \\ &= f(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2 \end{aligned}$$

Sufficient decrease II

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2.$$

\Rightarrow

$$\|\nabla f(\mathbf{x}_t)\|_2^2 \leq 2L (f(\mathbf{x}_t) - f(\mathbf{x}_{t+1}))$$

Proof.

Use smoothness and definition of gradient descent ($\mathbf{x}_{t+1} - \mathbf{x}_t = -\nabla f(\mathbf{x}_t)/L$):

$$\begin{aligned} f(\mathbf{x}_{t+1}) &\leq f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_{t+1} - \mathbf{x}_t) + \frac{L}{2} \|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 \\ &= f(\mathbf{x}_t) - \frac{1}{L} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2 \\ &= f(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2. \end{aligned}$$

□

Smooth convex functions: $\mathcal{O}(1/\varepsilon)$ steps

Theorem

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex and differentiable with a global minimum \mathbf{x}^* ; furthermore, suppose that f is smooth with parameter L . Choosing stepsize

$$\gamma := \frac{1}{L},$$

gradient descent yields

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{L}{2T} \|\mathbf{x}_0 - \mathbf{x}^*\|^2, \quad T > 0.$$

Smooth convex functions: $\mathcal{O}(1/\varepsilon)$ steps II

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{L}{2T} \|\mathbf{x}_0 - \mathbf{x}^*\|^2, \quad T > 0.$$

Proof.

Vanilla Analysis II:

$$\sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq \frac{\gamma}{2} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2\gamma} \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

This time, we can bound the squared gradients by sufficient decrease:

$$\frac{1}{2L} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 \leq$$

Handwritten derivation showing the relationship between the squared gradient and the function value decrease:

$$\frac{1}{2L} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 \leq \frac{1}{2L} \sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}_{t+1}))$$

The above equation is written in red ink. Below it, the expression $f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t)$ is written in blue ink, with a green line connecting the $f(\mathbf{x}_t)$ term in the sum above to the $f(\mathbf{x}_t)$ term in the expression below, and another green line connecting the $f(\mathbf{x}_{t+1})$ term in the sum above to the $f(\mathbf{x}_{t+1})$ term in the expression below.

Smooth convex functions: $\mathcal{O}(1/\varepsilon)$ steps II

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{L}{2T} \|\mathbf{x}_0 - \mathbf{x}^*\|^2, \quad T > 0.$$

Proof.

Vanilla Analysis II:

$$\sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq \frac{\gamma}{2} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2\gamma} \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

This time, we can bound the squared gradients by sufficient decrease:

$$\frac{1}{2L} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 \leq \sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}_{t+1})) = \underbrace{f(\mathbf{x}_0)} - \underbrace{f(\mathbf{x}_T)}.$$

Smooth convex functions: $\mathcal{O}(1/\varepsilon)$ steps III

Putting it together with $\gamma = 1/L$:

$$\begin{aligned} \sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) &\leq \frac{1}{2L} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \\ &\leq \underbrace{f(\mathbf{x}_0) - f(\mathbf{x}_T)}_{\geq f(\mathbf{x}^*)} + \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2. \end{aligned}$$

Rewriting:

$$\begin{aligned} &\leq \underline{f(\mathbf{x}_0) - f(\mathbf{x}^*)} \\ \Rightarrow \frac{1}{T} \sum_{t=1}^T f(\mathbf{x}_t) - f(\mathbf{x}^*) &\leq \frac{L}{T} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \end{aligned}$$

Smooth convex functions: $\mathcal{O}(1/\varepsilon)$ steps III

Putting it together with $\gamma = 1/L$:

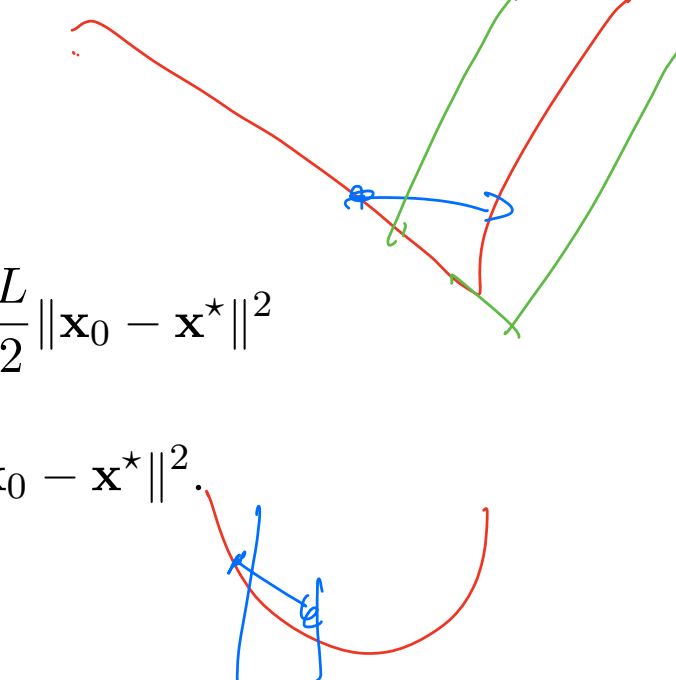
$$\begin{aligned} \sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) &\leq \frac{1}{2L} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \\ &\leq f(\mathbf{x}_0) - f(\mathbf{x}_T) + \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2. \end{aligned}$$

Rewriting:

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{1}{T} \sum_{t=1}^T (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq \frac{1}{T} \left(\frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \right)$$

As last iterate is the best (sufficient decrease!):

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{L}{2T} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \leq \frac{LR^2}{2T} = \frac{1}{2L} \left(\frac{L}{R} \|\nabla f(\mathbf{x}_T)\| \right)^2$$



$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2$$

Smooth convex functions: $\mathcal{O}(1/\varepsilon)$ steps III

Putting it together with $\gamma = 1/L$:

$$\begin{aligned}\sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) &\leq \frac{1}{2L} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \\ &\leq f(\mathbf{x}_0) - f(\mathbf{x}_T) + \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2.\end{aligned}$$

Rewriting:

$$\sum_{t=1}^T (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

As last iterate is the best (sufficient decrease!):

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{1}{T} \left(\sum_{t=1}^T (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \right) \leq \frac{L}{2T} \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

□

Smooth convex functions: $\mathcal{O}(1/\varepsilon)$ steps IV

$$R^2 := \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

$$T \geq \frac{R^2 L}{2\varepsilon} \quad \Rightarrow \quad \text{error} \leq \frac{L}{2T} R^2 \leq \varepsilon.$$

- ▶ $50 \cdot R^2 L$ iterations for error 0.01 ...
- ▶ ... as opposed to $10,000 \cdot R^2 B^2$ in the Lipschitz case

In Practice:

What if we don't know the smoothness parameter L ?

→ **Exercise ??**

Lemma 2.4

$$f(x) = x^T Q x + b^T x + c$$

$\rightarrow \text{Sym, psd}$

is $2\|Q\|_2$ smooth

Proof \rightarrow

$$\nabla^2 f(x) = 2Q$$

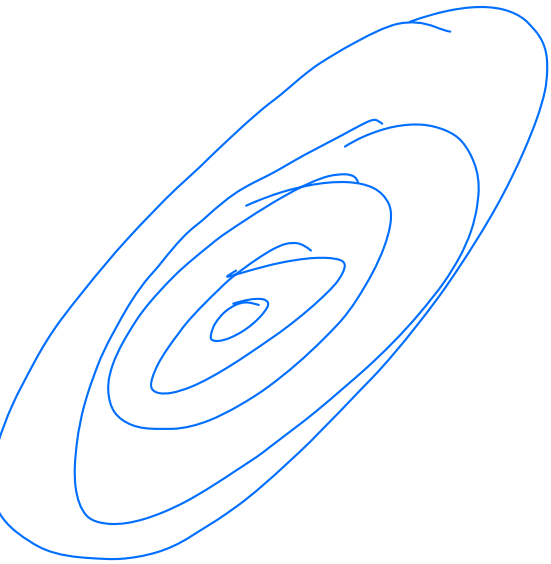
$$\|\nabla^2 f(x)\|_2 \leq \|2Q\|_2 = 2\|Q\|_2$$

\rightarrow $\|\nabla f(y) - \nabla f(x)\|_2 =$

$$\|2Qy - b - (2Qx - b)\|_2$$

$$= \|2Q(y-x)\|_2$$

$$\leq 2\|Q\|_2 \|y-x\|_2$$



Prove

i) $\{f_i \text{ is } L_i\text{-Smooth}\}$

$$g(x) = \sum_{i=1}^m \lambda_i f_i \text{ is } \left(\sum_{i=1}^m \lambda_i L_i \right)\text{-Smooth}$$

$$\nabla^2 g(x) = \left\| \sum_{i=1}^m \lambda_i \nabla^2 f_i(x) \right\|_2$$

$$\leq \sum_{i=1}^m \lambda_i \|\nabla^2 f_i(x)\|_2$$

$$\leq \sum_{i=1}^m \lambda_i L_i$$

ii) $g(x) = f(Ax+b)$ is $2\|A\|_2^2$ smooth

$$\nabla g(x) = A^T \nabla f(Ax+b)$$

$$\|\nabla^2 g(x)\| = \|A^T \nabla^2 f(Ax+b) A\|_2$$

$$= \max_v \frac{\|A^T (\nabla^2 f(Ax+b) A) v\|_2}{\|v\|_2}$$

$$\leq \max_{\psi} \underbrace{\|A\|_2}_{\|0\|_2} \|\nabla^2 f(Ax+b) \underbrace{A\psi}_{\psi^2}\|_2$$

$$\leq \max_{\psi} \underbrace{\|A\|_2 \|\nabla^2 f(Ax+b)\|_2}_{\|0\|_2} \|A\psi\|_2$$

$$\leq \max_{\psi} \underbrace{\|A\|_2^2 \|\nabla^2 f(Ax+b)\|_2}_{\|0\|_2} \|\psi\|_2$$

$$= \|A\|_2^2 \|\nabla^2 f(Ax+b)\|_2$$

Ex

$$f(x) = x^4$$

$$x \in (-a, a)$$

$$\nabla^2 f(x) = 12x^2 \leq 12a^2 \quad \forall x \in (-a, a)$$

$$\|\nabla f(x)\| = |4x^3| \leq 4a^3$$