

CS 331, Fall 2024  
Lecture 16 (10/23)

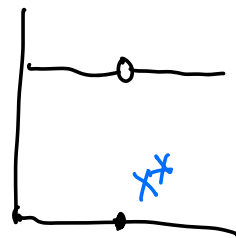
Today: - Convexity  
& Smoothness  
- Gradient descent

## Convexity & Smoothness (Part VI, Section 3.1)

---

Last time: Unstructured  $\min_{x \in \mathcal{X}} f(x)$

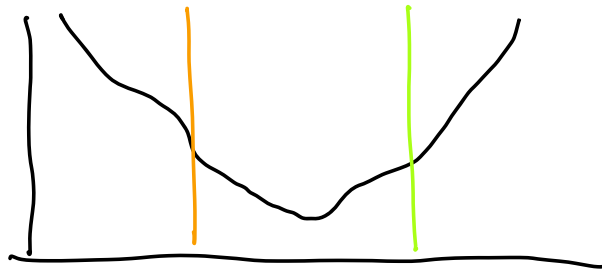
... is impossible!



With structure, Unimodal  $\min_{x \in \mathcal{X}} f(x)$

$\exists$  algorithms: is possible when  $\mathcal{X} \subseteq \mathbb{R}$

(ternary search)



What about high dimensions?  $x \in \mathbb{R}^d$

Can't quite ternary search...

Under stronger condition convexity

( $\Rightarrow$  unimodal)

We can use

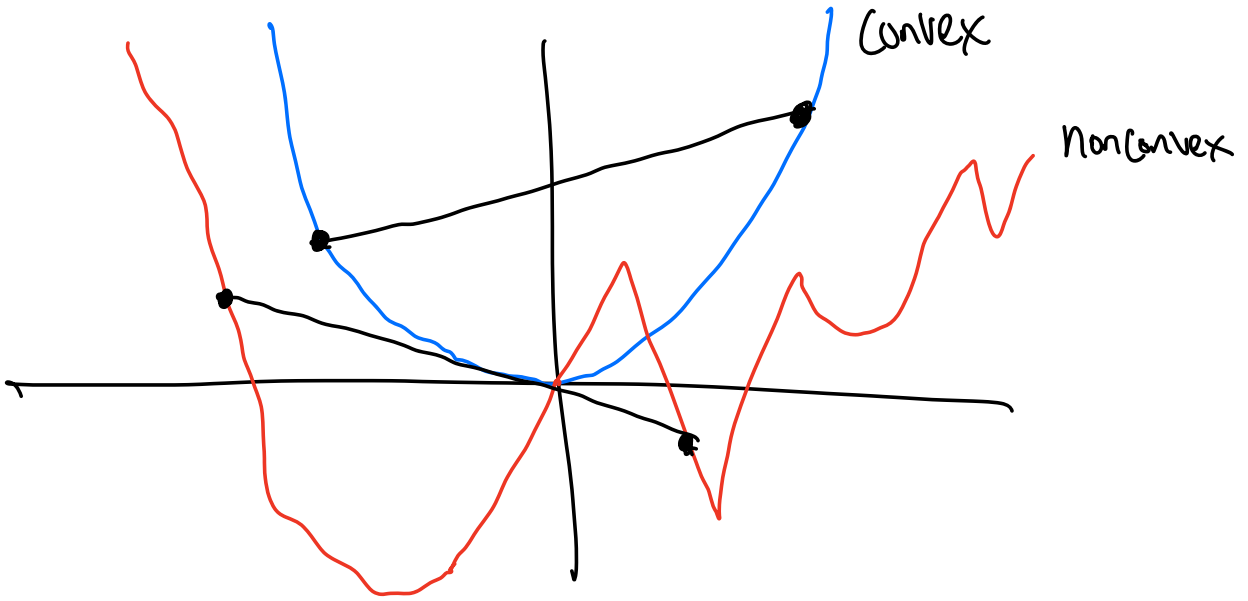


This is the most important algo in modern ML. It's how ChatGPT learned.

Today: GD in 1-d,  $x = \mathbb{R}$

# Convexity

$f$  is convex if  $\geq$ /below the line.



Other convex / non convex:



✓

✗

✓

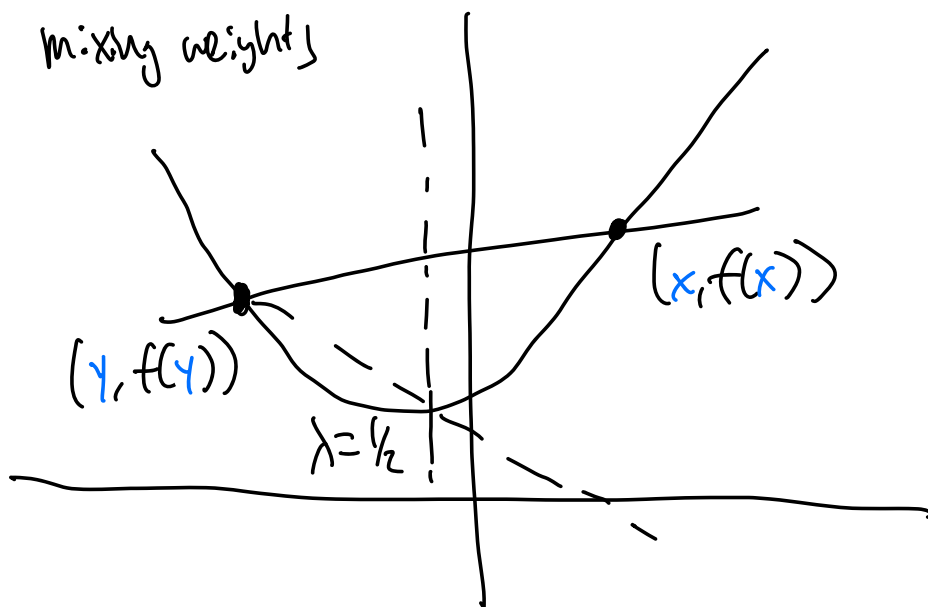
✗

("at the line")

Formally,  $\forall \lambda \in [0, 1]$

$$f((1-\lambda)x + \lambda y) \leq (1-\lambda)f(x) + \lambda f(y)$$

↑      ↑  
mixing weights



Another fact:

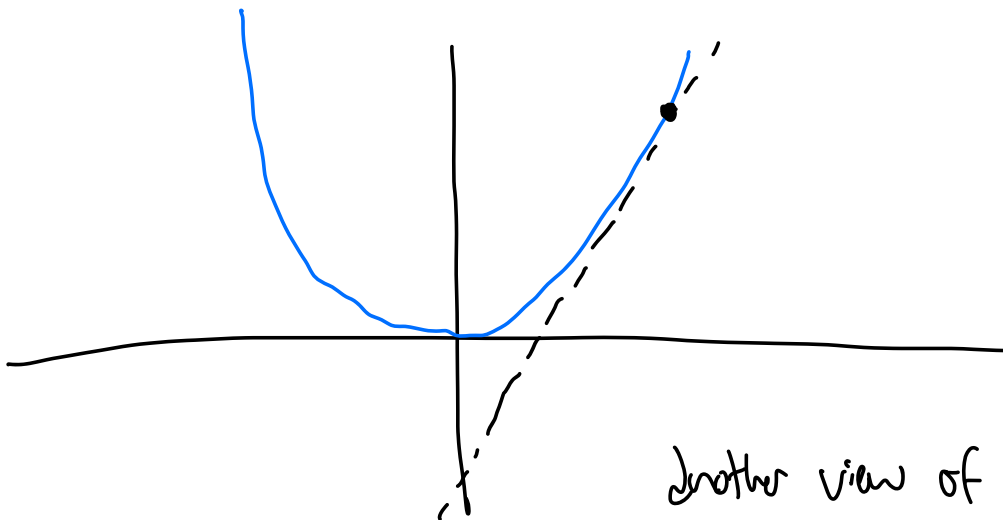
$$f(y) \geq f(x) + \underbrace{\frac{f(x + \lambda(y-x)) - f(x)}{\lambda}}_{\lim_{\lambda \rightarrow 0} = f'(x)(y-x)}$$

In summary  $\forall x, y$

$$f(y) \geq f(x) + f'(x)(y-x)$$



"Euler approximation"



Another view of convexity...  
above the tangent line

Why good? Binary search!

$$f(x^*) \geq f(x) + f'(x)(x^* - x)$$

$f'(x) > 0$ : go left      else: go right

Another equivalent condition:  $f'' \geq 0$ .

e.g.

"The derivative (gradient) increases!"

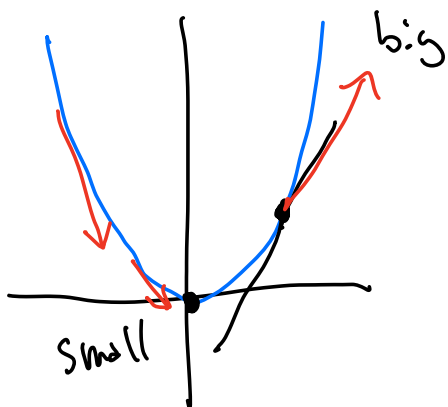
Issue 1: What if  $f'$  doesn't exist?

Issue 2: What about  $\mathbb{R}^d$ ???

Idea: gradient descent! in high-d,  
(ignore 1)  $\nabla f(x)$

$$x \leftarrow x - \eta f'(x)$$

step size



Why scale w  $f'(x)$ ?

further from  $x^* \Rightarrow |f'|$  bigger

We would love to just simulate a rolling ball ( $n \rightarrow 0$ ). Unfortunately, we must discretize ☹

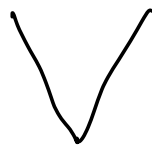
Important to pick  $N$  carefully.

To make it work, we need...

Smoothness



Smooth



non-smooth

$f$  is  $L$ -Lipschitz if

$$|f(x) - f(y)| \leq L|x - y|$$

nearby points have close values.

$f$  is  $L$ -smooth if "no corners"

$$|f'(x) - f'(y)| \leq L |x - y|$$

Gradient queries Why ok?

$$\left| f'(x) - \underbrace{\frac{f(x+\delta) - f(x)}{\delta}}_{\text{sim. with function queries}} \right| = O(\delta) \text{ if smooth.}$$

$\delta \rightarrow 0$

sim. with function queries

GD Take 1: Quadratics (Part VI, Section 3.2)

In this section,  $f = q$  a quadratic

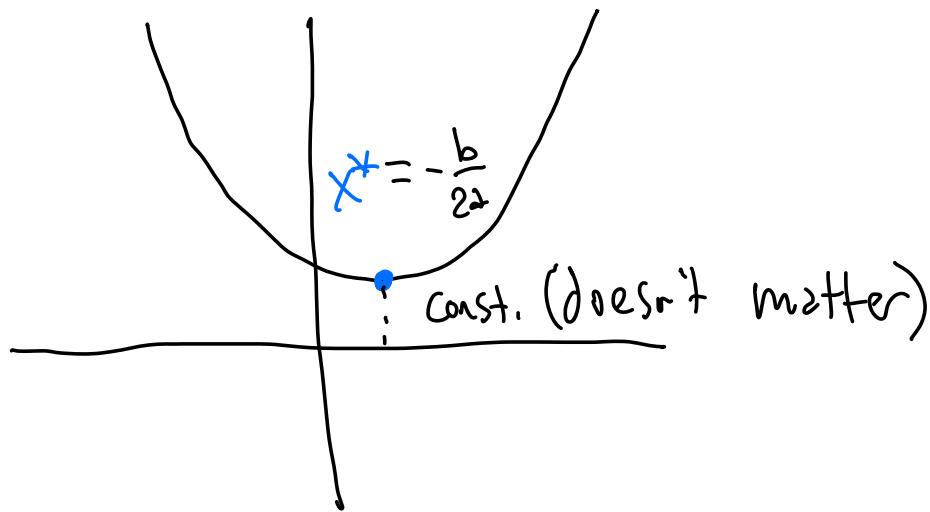
$$q(x) = ax^2 + bx + c, \quad a > 0$$



What is  $x^*$ ?

Complete the square:

$$ax^2 + bx + c = a \left( x + \frac{b}{2a} \right)^2 + \text{const.}$$



Suboptimality @  $x$ :

$$a \left( x + \frac{b}{2a} \right)^2 - a \left( x^* + \frac{b}{2a} \right)^2$$

$$= a \left( x - x^* \right)^2 \quad \text{★}$$

What is  $q'$ ?

$$q'(x) = \frac{d}{dx} (2x^2 + bx + c)$$
$$= 2ax + b = 2a(x - x^*)$$

Sanity check:  $q'(x^*) = 0$

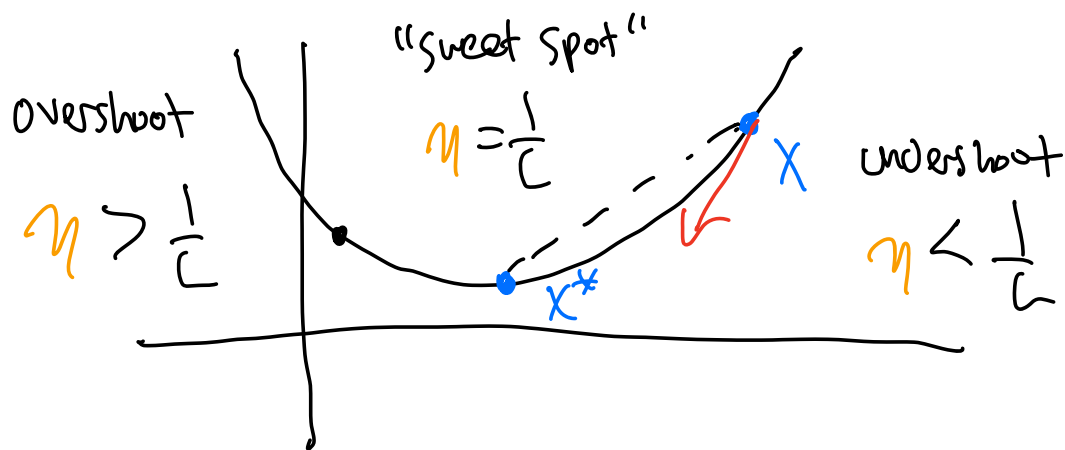
Hence,  $L = 2a$ -smooth:

$$|q'(x) - q'(y)| = 2a|x - y|$$

What step size?

$$x - \eta f'(x) = x - \eta \cdot 2a(x - x^*)$$

$$\eta = \frac{1}{2a} = \frac{1}{L} \Rightarrow \text{step to } x^*$$



full characterization of quadratics

(G) Take 2: Smoothness (Part VI, Section 3.3)

---

$$|f'(x) - f'(y)| \leq L |x - y|$$

Why is it good? Upper bounded by  $q$

Intuition:  $q'' = L$ ,  $f'' \leq L$

$$\lim_{y \rightarrow x} \frac{f'(x) - f'(y)}{x - y}$$

Just as best linear fit

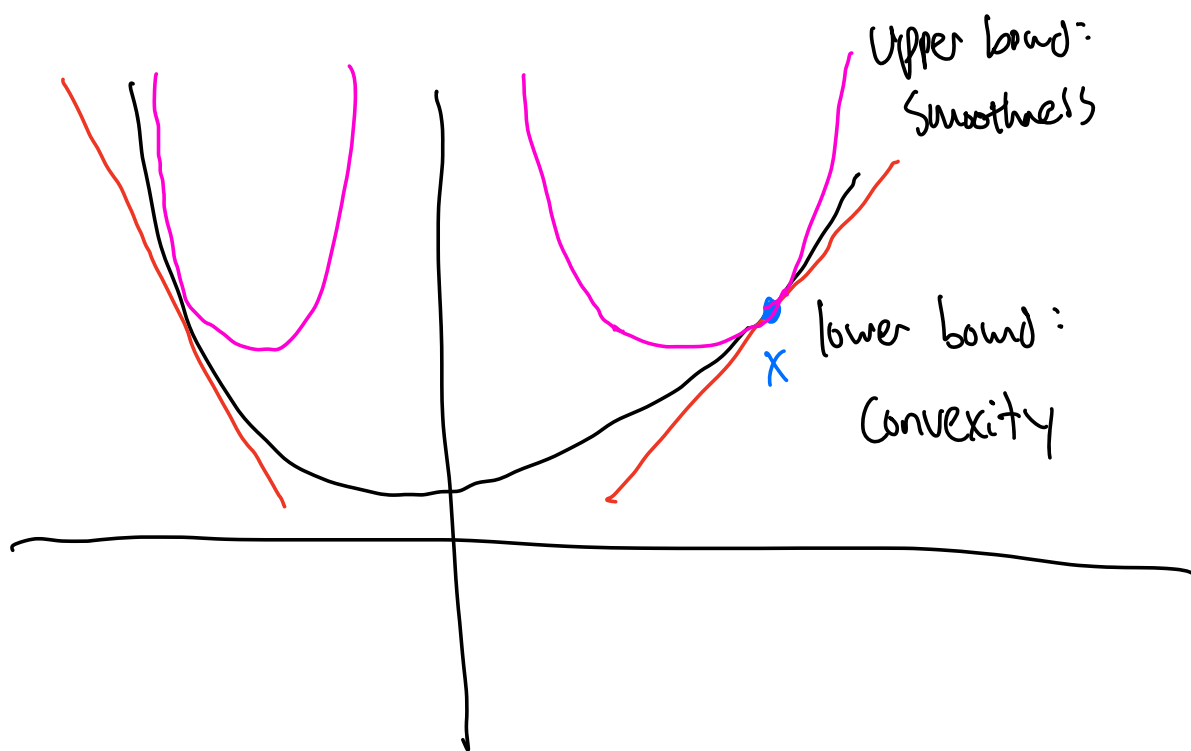
$$f(y) \approx f(x) + f'(x)(y-x)$$

Best quadratic fit

$$f(y) \approx f(x) + f'(x)(y-x) + \frac{1}{2} \underbrace{f''(x)}_{\leq L} (y-x)^2$$

Can be made rigorous:  $\forall x, y$

$$f(y) \leq f(x) + f'(x)(y-x) + \frac{L}{2} (y-x)^2 \quad \left. \vphantom{f(y)} \right\} q(y)$$



Note: @  $x$ ,  $q' = f' = l'$

Also,  $q$  is  $L$ -Smooth.

$$\begin{aligned} \text{Hence, } x &\leftarrow x - \frac{1}{L} f'(x) \\ &= x - \frac{1}{L} q'(x) \end{aligned}$$

minimizes  $q$  @ each step.

# Critical points

L-smooth, not convex

Algo: assume  $f(x_0) - f(x^*) \leq \Delta$

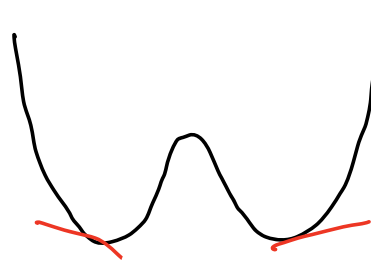
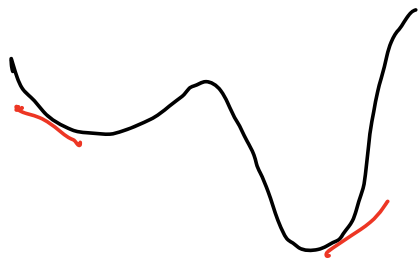
iterate  $\forall t \in [T]$

$$x_t \leftarrow x_{t-1} - \frac{1}{L} f'(x_{t-1})$$

Claim:  $T \geq \frac{2L\Delta}{\epsilon^2} \Rightarrow \min_t |f'(x_t)| \leq \epsilon$   
 $\approx$  local minimum

Proof: Suppose otherwise. Progress / iter  
 $\geq \frac{1}{2L} |f'(x_t)|^2 \geq \frac{\epsilon^2}{2L}$  ~~(A)~~

After  $T$  iters, more than  $\Delta$  progress ~~(B)~~



modern neural nets?

# Global optimality

Algo: assume  $|x_0 - x^*| \leq 1$ ,  
 $f$  is  $L$ -smooth & convex  
iterate  $\forall t \in [T]$

$$x_t \leftarrow x_{t-1} - \frac{1}{L} f'(x_{t-1})$$

Helper claim:  $|x_t - x^*| \leq |x_0 - x^*|$   
(see notes) never overshoots.

Claim: After  $T \geq \frac{L^2}{\epsilon^2}$  iters,  
 $f(x_t) \leq \underbrace{f(x^*)}_{\text{global opt}} + \epsilon$

Proof:  $\Delta \leq \frac{L}{2}$ ,  $\overbrace{\leq 1}^{\leq 1}$

$$f(x_0) \leq f(x^*) + \frac{L}{2} (x_0 - x^*)^2$$

Use critical points,  $T \geq \frac{L^2}{\epsilon^2} = \frac{2L\Delta}{\epsilon^2}$

$$f(x) - f(x^*) \leq f'(x)(x - x^*)$$

(Convexity)  $\leq |f'(x)| \leq \epsilon$

... and beyond

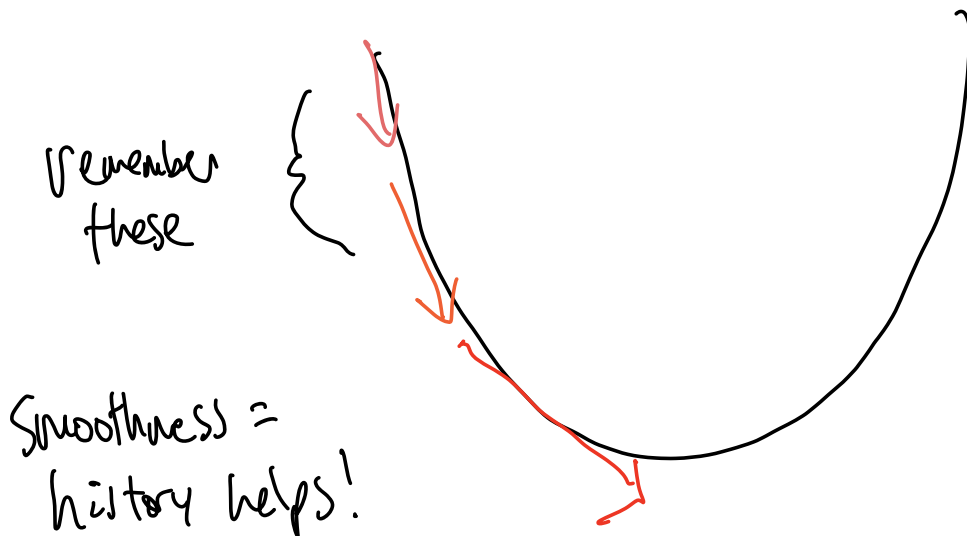
We proved  $T \geq \frac{L^2}{\epsilon^2}$  suffices today.

1) Same algo but smarter analysis  $\Rightarrow \frac{L}{\epsilon}$

2) Optimal algo:  $\sqrt{\frac{L}{\epsilon}}$  (Nesterov)



Key idea: momentum



Why is this important?

for that GPT,  
 $\delta = 10^{12}$

Generalizes verbatim to  $\mathbb{R}^d$

(critical points & global optima)

$$\nabla f(x) = \left( \frac{\partial f}{\partial x_1}(x), \dots, \frac{\partial f}{\partial x_d}(x) \right)$$

"gradient"

"hint" of where to find  $x^*$