# Overview of week

HW end of week / early next grades

==A4 due Monday 11/2==

today : Optimization

rest of week : nonlinear models

Friday : project

---

Today's goals : SGD w/ mini-batches

convex sets & functions

why we like convexity

# SGD w/ minibatch

data pts ↙

$$I_t = \text{uniform random index } \{1, \ldots, n\}$$

$$\vec{d}_t = -\nabla \ell_{I_t}(\vec{w}_t)$$

$$\mathbb{E}[\vec{d}_t] = -\nabla C(\vec{w}_t)$$

$$\vec{w}_{t+1} = \vec{w}_t + h_t \vec{d}_t$$

$$= -\nabla \left( \underbrace{\frac{1}{n} \sum_{i=1}^{n} \ell_i(\vec{w}_t)}_{\color{red}\text{average of } n} \right)$$

mini-batching: variance decreases $\sim 1/B$

$\Rightarrow$ closer to average

Pick $B$ random indices (w/ replacement)

$$I_t = \{I_{t1}, I_{t2}, \ldots, I_{tB}\} \qquad |I_t| = B$$

$$\vec{d}_t = -\frac{1}{\textcircled{B}} \sum_{J \in I_t} \nabla \ell_I(\vec{w}_t) \qquad (\text{SGD}: B=1)$$

↳ # of data pts in batch

# Practical considerations, advantages of SGD

- less computation than GD $\qquad \mathcal{O}(B)$ vs. $\mathcal{O}(n)$
  memory ⌉ <span style="color:red">important for GPU and NNs</span>

- parallelizes easily
  - different processes / computers working on different batches (Hogwild!)

- special sauce of SGD for NNs
  - "implicit bias" (bias-var) of SGD small $\|\vec{w}\|$
    similar to ridge
  - noise helps avoid local min
  - debatable how important

<span style="color:magenta">depends (convex)</span>
<span style="color:magenta">GD $\|\vec{w}_t - \vec{w}^*\| \leq \frac{1}{t}$</span>
<span style="color:magenta">SGD $\|\vec{w}_t - \vec{w}^*\| \leq \frac{1}{\sqrt{t}}$</span>

Disadvantages:
- more iterates than GD
- noisier trajectories, not always descent

# Convex sets :

**Defn** A set $K \subseteq \mathbb{R}^d$ is <u>convex</u> if

$$\underbrace{(1-t)\vec{x} + t\vec{y}}_{\text{line between } \vec{x} \text{ and } \vec{y}} \in K \text{ for any } \vec{x}, \vec{y} \in K, \ 0 \leq t \leq 1$$

K
convex

$\vec{x}$
$t=0$
$t=\frac{1}{2}$
$\vec{y}$
$t=1$

K not
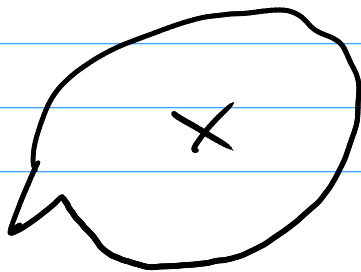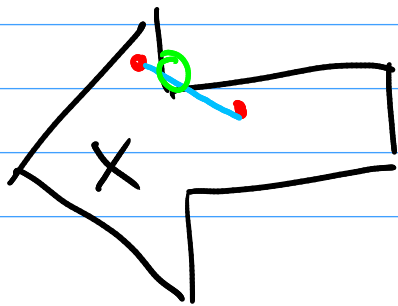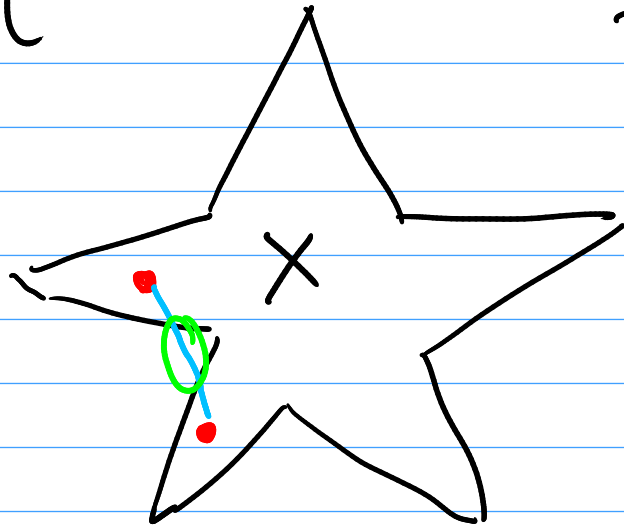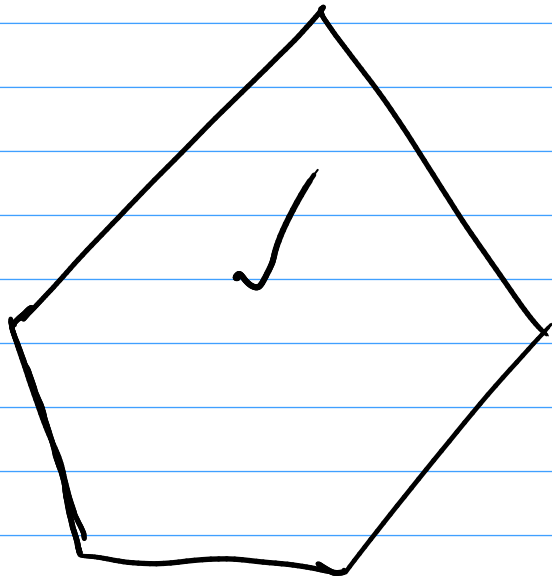convex

$\notin K$

Is it convex?

$$\mathbb{R}_{\geq 0}^d = \{\vec{x} : x_i \geq 0\}$$

$$(1-t)\underbrace{\vec{x}}_{\geq 0} + \underbrace{t\vec{y}}_{\geq 0} = \vec{z} \geq 0$$

$\checkmark$

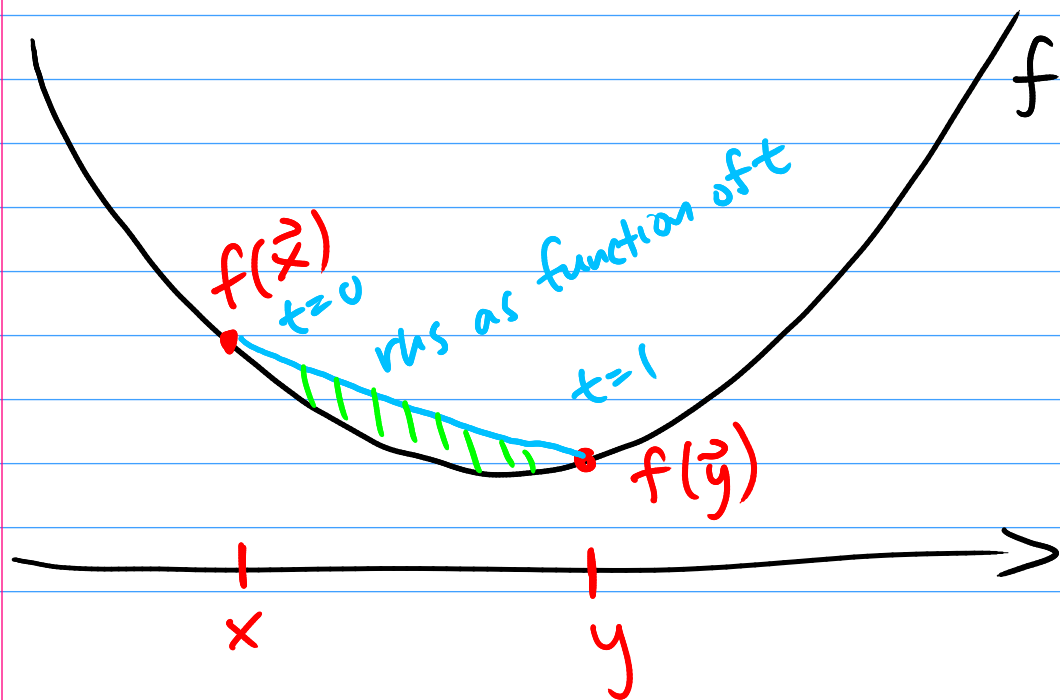✓ (checkmark inside square)

X (in star)

X (in arrow)

X (in blob)

half-space

$\checkmark$

K

# Convex functions

Defn $f: \mathbb{R}^d \longrightarrow \mathbb{R}$ is convex iff

$$f\left((1-t)\vec{x} + t\vec{y}\right) \leq (1-t)f(\vec{x}) + t f(\vec{y})$$

$\underbrace{\qquad\qquad}$
any pt. on line

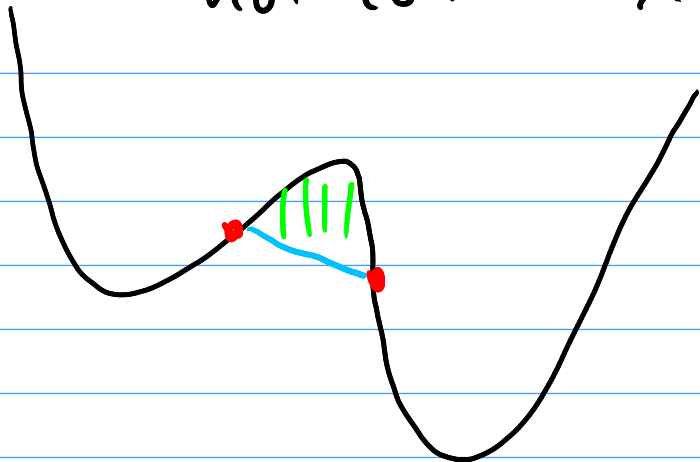for any $\vec{x}, \vec{y} \in \text{dom}(f)$, $0 \leq t \leq 1$

If $f$ is convex:

"$f$ lies underneath line segment connecting any two points"
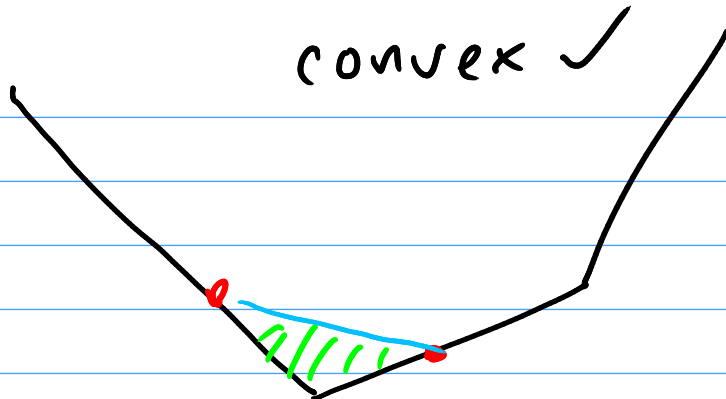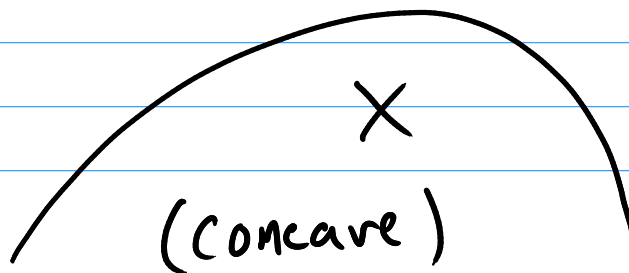


$f(\vec{x})$
$t=0$
rhs as function of $t$
$t=1$
$f(\vec{y})$

$x$
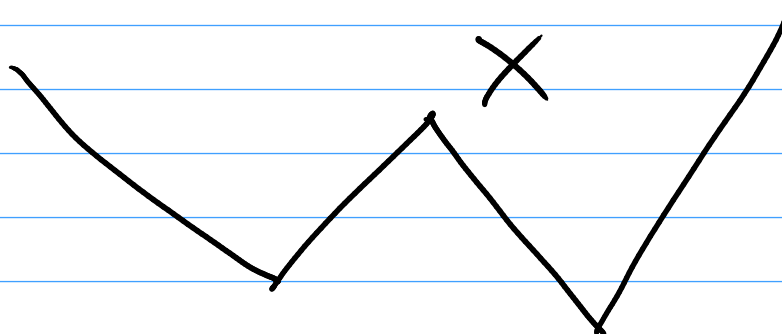$y$

$f$

not convex X

convex ✓

$e^{-x}$ ✓

X

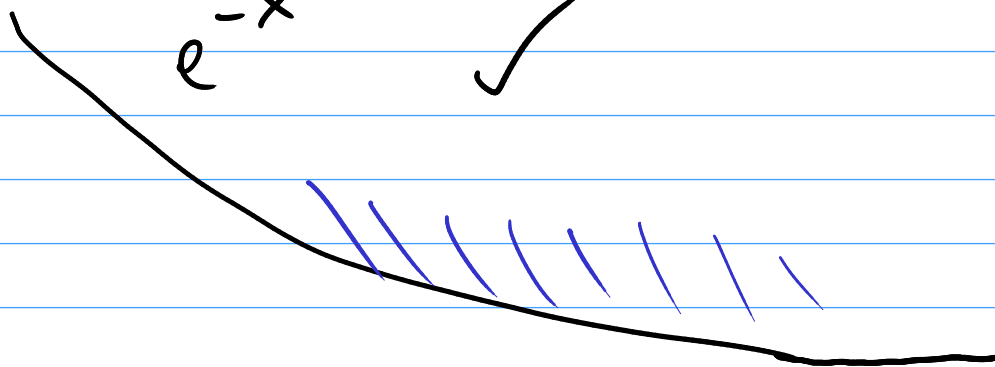(concave) X

# Why do we like them?

- all local minima are global
- algorithms are efficient and find these minima
  - GD, SGD
- A4 $\Big\{$
  - coordinate descent
  - tricks for nonsmooth $\quad\vee\quad \ell_1$-norm
    "proximal"
    "sub-gradient"
  - accelerated versions $\quad 1/t \longrightarrow 1/t^2$
    averaging

ex/ $\quad \|X\vec{w} - \vec{y}\|^2 \qquad$ convex $\qquad\qquad L(\vec{w}) + \lambda R(\vec{w})$

$\|\vec{w}\| \qquad$ (real) norms

$f(\vec{w}) + g(\vec{w}) \qquad f, g$ convex