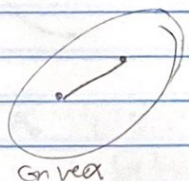
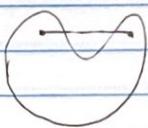


## Gradient Descent & SGD



convex



non convex

$K$  is convex

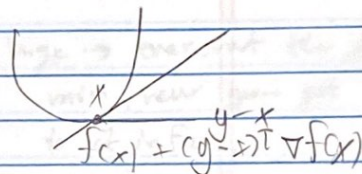
$$\Leftrightarrow \forall x, y. \lambda x + (1-\lambda)y \in K$$



This is convex

A function is convex  $\Leftrightarrow \forall x, y$ :

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



$$\nabla f(x) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}$$

• For diff  $f$ ,  $f(y) \geq f(x) + \nabla f(x)^T (y-x)$

• For twice diff function

$$(Hf)_{ij} = (\nabla^2 f)_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j}$$

$f$  is convex  $\Leftrightarrow Hf$  is positive semidefinite for all  $x$

$$\forall i, \lambda_i(A) \geq 0.$$

Def. A function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is  $\lambda$ -Lipchitz if  
 $\forall x, y \in \mathbb{R}^n, |f(x) - f(y)| \leq \lambda \|x - y\|_2$

Minimize a convex func:  $f: \mathbb{R}^n \rightarrow \mathbb{R}$

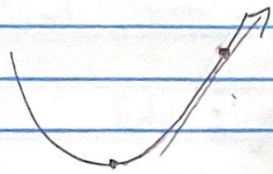
$$\nabla f(\vec{x}) = 0$$

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

$$\nabla f(x) = A x - b = 0$$

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





$$x_0 - \gamma \nabla f(x_0)$$

how far big?

step size  $\gamma$ :  $\uparrow$  small  $\rightarrow$  too long to reach local min

$\vec{x}_0$  = starting point  
for  $t=1 \dots T$ :

$$\vec{x}_t = \vec{x}_{t-1} - \gamma \nabla f(\vec{x}_{t-1})$$

• too large  $\rightarrow$  overshoot the min, never gonna get to it, in fact.

$$\text{return } \hat{x} = \frac{1}{T} \sum_t \vec{x}_t$$

Thm: Let  $x^*$  be the minimizer of  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ .  
If  $f$  is convex & diff and satisfies

$$\forall x \in \mathbb{R}^n, \|\nabla f(x)\|_2 \leq G, \text{ then}$$

\* setting  $T = \frac{1}{\epsilon} G^2 \|x_0 - x^*\|_2^2$  and  $\gamma_t = \gamma = \frac{\|x_0 - x^*\|_2}{G\sqrt{T}}$

$$\text{gives } f(\hat{x}) \leq f(x^*) + \epsilon$$

~~PF~~

Lemma A:

$$\sum_t f(x_t) \leq \left( \sum_t f(x^*) \right) + \frac{1}{2} G^2 t \gamma + \frac{1}{2\gamma} \|x_0 - x^*\|_2^2$$

Proof of Thm:

If we don't know about T  
know about T  
 $f(\hat{x}) = f(x^*) + O(\frac{1}{\sqrt{T}})$

$$f(\hat{x}) = f\left(\frac{1}{T} \sum_t f(x_t)\right) \stackrel{\text{by convexity of } f}{\leq} \frac{1}{T} \left( \sum_t f(x_t) \right)$$

$$\stackrel{\text{Lemma A}}{\leq} \frac{1}{T} (T f(x^*)) + \frac{1}{2} G^2 \gamma + \frac{1}{2\gamma} \|x_0 - x^*\|_2^2$$

$$\leq f(x^*) + \frac{G \|x_0 - x^*\|_2}{\sqrt{T}}$$

for  $\gamma$  plug in  $\gamma$   
for  $T$  plug in T

$$\leq f(x^*) + \frac{G \|x_0 - x^*\|_2}{\frac{1}{\epsilon} G \|x_0 - x^*\|_2} \leq f(x^*) + \epsilon$$



Proof lemma A

$$\text{Defn: } \phi_t = \frac{1}{2\eta} \|\vec{x}_t - x^*\|_2^2$$

you are f. from a only

$$\phi_0 = \frac{1}{2\eta} \|\vec{x}_0 - x^*\|_2^2$$

$$\text{Claim: } f(\vec{x}_{t+1}) + \phi_{t+1} - \phi_t \leq f(x^*) + \frac{1}{2} G^2 \eta$$

Proof

$$f(\vec{x}_t) + \phi_{t+1} - \phi_t = f(\vec{x}_t) + \frac{1}{2\eta} \left( \|\vec{x}_{t+1} - x^*\|_2^2 - \|\vec{x}_t - x^*\|_2^2 \right)$$

$$\begin{aligned} & \frac{1}{2\eta} \left( \|\underbrace{\vec{x}_{t+1} - \vec{x}_t}_{\Delta x} + \underbrace{\vec{x}_t - x^*}_z\|_2^2 - \|\vec{x}_t - x^*\|_2^2 \right) \\ &= (\Delta x + z)^T (\Delta x + z) - \|\vec{x}_t - x^*\|_2^2 \\ &= \|\Delta x\|_2^2 + 2\Delta x^T z + \|z\|_2^2 - \|\vec{x}_t - x^*\|_2^2 \end{aligned}$$

$$= f(\vec{x}_t) + \frac{1}{2\eta} \left( \|\Delta x\|_2^2 + 2\Delta x^T (\vec{x}_t - x^*) \right)$$

$$\Delta x = \vec{x}_{t+1} - \vec{x}_t = -\eta \nabla f(\vec{x}_t)$$

$$\begin{aligned} \text{plug } \Delta x \text{ in: } & \leq f(\vec{x}_t) + \frac{1}{2\eta} \left( \eta^2 \|\nabla f(\vec{x}_t)\|_2^2 - 2\eta \nabla f(\vec{x}_t)^T (\vec{x}_t - x^*) \right) \\ &= \frac{1}{2} \eta G^2 + f(\vec{x}_t) - \nabla f(\vec{x}_t)^T (\vec{x}_t - x^*) \end{aligned}$$

$$= \frac{1}{2} \eta G^2 + \underbrace{f(\vec{x}_t) + \nabla f(\vec{x}_t)^T (x^* - \vec{x}_t)}_{\text{linear approximation}}$$

$$\leq \frac{1}{2} \eta G^2 + f(x^*) \leq f(x^*)$$



Proof Lemma A

$$\begin{aligned}
 \sum_t f(\vec{x}_t) &= \left( \sum_t f(x^*) \right) + \frac{1}{2} G^T y + \frac{1}{2\eta} \|\vec{x}_0 - x^*\|_2^2 \\
 &= \left( \sum_t f(x^*) + \cancel{\phi_1} - \cancel{\phi_{T-1}} + \cancel{\phi_{T-1}} - \dots - \phi_0 \right) \\
 &\leq \left( \sum_t f(x^*) \right) + \frac{1}{2} G^T y \\
 \Rightarrow \sum_t f(\vec{x}_t) &\leq \left( \sum_t f(x^*) \right) + \frac{1}{2} G^T y + \cancel{\phi_0} - \cancel{\phi_0} \\
 &\leq \left( \sum_t f(x^*) \right) + \frac{1}{2} G^T y + \frac{1}{2\eta} \|\vec{x}_0 - x^*\|_2^2
 \end{aligned}$$

from lemma

SGD

minimize  $f(\vec{\theta}) = \frac{1}{n} \sum_{i=1}^n f_i(\vec{\theta})$

$(\vec{x}_i, y_i)$

$$f(\vec{\theta}) = \frac{1}{n} \sum_{i=1}^n (\vec{x}_i^T \vec{\theta} - y_i)^2$$

$\nabla f = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\vec{\theta})$  (n.d) work per iteration

$\mathbb{E}[f_i(\vec{\theta})] = \frac{1}{n} \sum_{i=1}^n f_i(\vec{\theta}) = f(\vec{\theta})$

$\mathbb{E}[\nabla f_i(\vec{\theta})] = \nabla f(\vec{\theta})$

$\mathbb{E}[\nabla f_i(\vec{x}_{t-1})]$

For  $t = 1 \dots T-1$

1. Pick  $i$  at random.

2.  $\vec{x}_t = \vec{x}_{t-1} - \eta \nabla f_i(\vec{x}_{t-1})$

per iter,  $n_{\text{row}}, O(n \cdot d) \text{ work} \leftarrow \text{SGD} \rightarrow \text{space reduction - (overall?)}$



Thm SGD. Sum

$\| \nabla f_i(x) \|_2 \leq G$ , then

$$\mathbb{E}[f(\bar{x})] \leq f(x^*) + \varepsilon$$

~~Clear~~

Expectation means nothing.

~~Min~~

• Minibatch SGD

pick  $B \subseteq [N]$ ,  $|B| = b$

$$\vec{x}_t = \vec{x}_{t-1} - \frac{\eta t}{b} \sum_{i \in B} \nabla f_i(\vec{x}_{t-1})$$

→ reducing a variance.

$$\mathbb{E} \left[ \frac{1}{b} \sum \nabla f_i(\vec{x}_{t-1}) \right] = \nabla f(\vec{x}_{t-1}), \text{Var: } \frac{1}{b}$$

$O(\sqrt{b\eta} + \frac{b\eta}{t})$  convergence.

if not batch  $O(\frac{1}{\sqrt{t}})$