

AA Lecture 21

3/29.

Gradient Descent

$$\min_{x \in \mathbb{R}^d} f(x).$$

Taylor expansion:

$$f(x+\delta) = f(x) + \nabla f(x)^T \delta + \delta^T \nabla^2 f(y) \cdot \delta \cdot \frac{1}{2} \quad \text{(TE)}$$

Def: f is β -smooth: $\lambda_{\max}(\nabla^2 f(y)) \leq \beta$.

$$\text{(TE)}: f(x+\delta) \leq \underbrace{f(x) + \nabla f(x)^T \delta + \frac{1}{2}\beta \|\delta\|^2}_{\delta \text{ that is min}}$$

$$\delta \text{ that is min} \uparrow \text{ is } \delta = -\frac{1}{\beta} \cdot \nabla f(x).$$

Fix $x^t \in \mathbb{R}^d$

$$x^{t+1} = x^t - \frac{1}{\beta} \nabla f(x^t).$$

$$f(x^{t+1}) \leq f(x^t) - \frac{1}{\beta} \cdot \nabla f(x^t)^T \cdot \nabla f(x^t) + \frac{\beta}{2} \cdot \frac{1}{\beta^2} \|\nabla f(x^t)\|^2$$

$$= f(x^t) - \frac{1}{2\beta} \|\nabla f(x^t)\|^2$$

$$< f(x^t) \text{ if } \nabla f(x^t) \neq 0.$$

When is $\nabla f(x) = 0$?

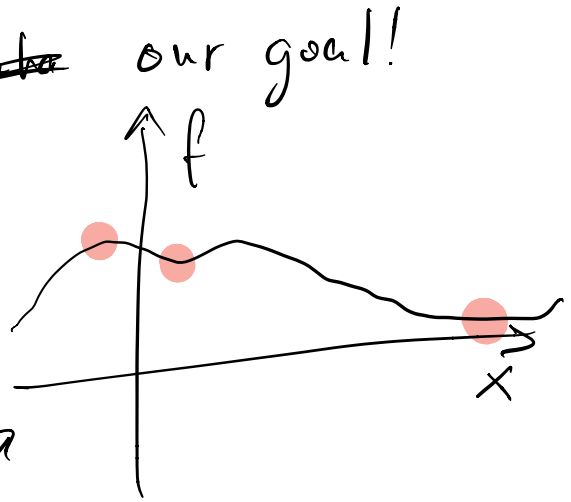
1. Global min \leftarrow ~~is~~ our goal!

2. Local min

3. A gl/loc. max

4. Saddle point

(in some directions \nearrow
in others \searrow)



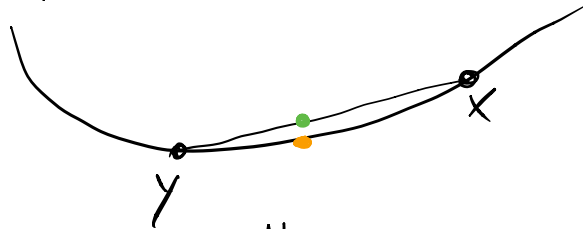
Notes: usually easy to get out of
#3, #4.

#2 \rightarrow non-convex opt.

Def: f is convex iff $\lambda \nabla^2 f(y) \succeq 0$.

Fact: f is convex iff $\forall x, y \in \mathbb{R}^d, \alpha \in [0, 1]$

$$f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y).$$



Claim: if $f(x)$ is convex, then

$\nabla f(x) = 0 \Rightarrow x$ is global min.

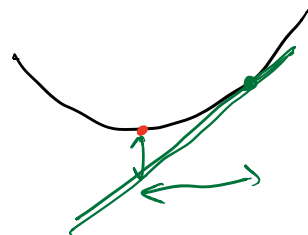
pf: fix x s.t. $\nabla f(x) = 0$.

$$\text{TE: } f(x+\delta) = f(x) + \underbrace{\nabla f(x)^T \cdot \delta}_0 + \underbrace{\frac{1}{2} \delta^T \nabla^2 f(y) \cdot \delta}_{\geq 0}$$

$$\geq f(x).$$

Assumption #2: f is convex. \square

$\Rightarrow \nabla f(x^t) = 0 \Rightarrow$ done!
 $x^t = \text{gl. min.}$



$\nabla f(x^t) \neq 0 \Rightarrow \text{GD is making progress.}$

$$f(x^{t+1}) \leq f(x^t) - \frac{1}{2\beta} \cdot \|\nabla f(x)\|^2.$$

Goal: find some x' s.t.

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

\nwarrow global min

$x' \rightarrow$ ϵ -approx. sol.

Relate $f(x^t) - f(x^*)$ to $\|\nabla f(x)\|$:

$$\begin{aligned}
 f(x^*) &= f(x + (x^* - x)) \\
 &\stackrel{TE}{=} f(x) + \nabla f(x)^T \cdot (x^* - x) + \\
 &\quad \frac{1}{2} (x^* - x)^T \cdot \nabla^2 f(y) \cdot (x^* - x)
 \end{aligned}$$

$$\begin{aligned}
 &\geq f(x) + \nabla f(x)^T \cdot (x^* - x) \\
 \nabla f(x)^T \cdot (x - x^*) &\geq f(x) - f(x^*)
 \end{aligned}$$

$\frac{1}{p} \leq \|\nabla f\|_q$
 $\|g\|_p$
 $\frac{1}{p} + \frac{1}{q} = 1$
 Hölder inequality.

$$f(x) - f(x^*) \leq \|\nabla f(x)\| \cdot \|x - x^*\|.$$

$$\|\nabla f(x)\| \geq \frac{f(x) - f(x^*)}{\|x - x^*\|}$$

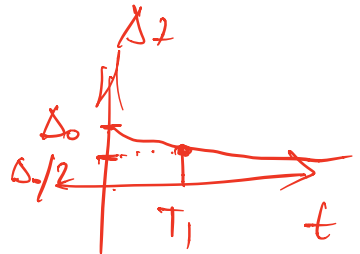
For example: if x^t s.t. $f(x^t) - f(x^*) \geq \epsilon$
 $\Rightarrow \|\nabla f(x)\| \geq \frac{\epsilon}{\|x - x^*\|}$

Thm: $f(x^t) - f(x^*) \leq \epsilon$ after

$$T = O\left(\beta \cdot \frac{D^2}{\epsilon}\right) \text{ iterations}$$

$$\begin{aligned}
 D &\triangleq \max_x \|x - x^*\| \\
 x &: f(x) \leq f(x^*)
 \end{aligned}$$

Pf:



$$\text{Let } \Delta_t = f(x^t) - f(x^*) \geq \epsilon.$$

$$\|\nabla f(x^t)\| \geq \frac{\Delta_t}{\|x^t - x^*\|} \geq \frac{\Delta_t}{D}.$$

$$\text{Let } T_1 = \# \text{ iterations until } \Delta_{T_1} \leq \frac{\Delta_0}{2}.$$

Before reaching $t = T_1$:

$$\|\nabla f(x^t)\| \geq \frac{\Delta_0/2}{D}.$$

In each iteration $f(x^t) - f(x^*)$

drops by $\geq \frac{1}{2\beta} \cdot \|\nabla f(x^t)\|^2$

$$\geq \frac{1}{2\beta} \cdot \frac{\Delta_0^2}{4D^2} = \frac{\Delta_0^2}{8D^2\beta}.$$

$$\Rightarrow T_1 \leq \frac{\Delta_0/2}{\Delta_0^2/8D^2\beta} = 4D^2\beta/\Delta_0.$$

$$\text{Let } T_2 = \# \text{ iterations until } \Delta_{T_1+T_2} \leq \frac{\Delta_{T_1}}{2}.$$

$$\text{Similarly, } T_2 \leq 4D^2\beta/\Delta_{T_1} = 4D^2\beta \cdot 2/\Delta_0.$$

$$T_3 : T_3 \leq 4D^2\beta \cdot \frac{4}{\Delta_0}$$

...

$$T_k : T_k \leq 4D^2\beta \cdot 2^k / \Delta_0 \Rightarrow f(x^{T_1, \dots, T_k}) - f(x^*) \leq \frac{\Delta_0}{2^k} \leq \epsilon.$$

At time $T = T_1 + T_2 + \dots + T_k$, value:

$$f(x^T) - f(x^*) \leq \epsilon$$

Total time: $T \equiv T_1 + T_2 + \dots + T_k$

$$= 4D^2\beta \cdot \left[\frac{1}{\Delta_0} + \frac{2}{\Delta_0} + \dots + \frac{2^k}{\Delta_0} \right]$$

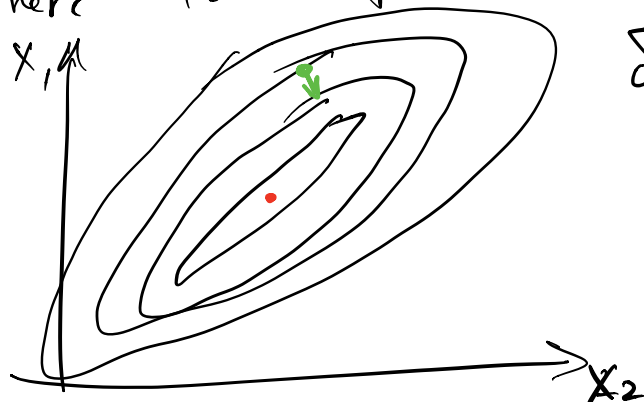
$$= 4D^2\beta \left[\frac{1}{\Delta_0} + \frac{2}{\Delta_0} + \dots + \frac{1}{\epsilon} \right]$$

$$\leq 4D^2\beta \cdot 2/\epsilon. \quad \square$$

$$T \leq O(\beta D^2 / \epsilon).$$

$$f(x^T) - f(x^*) \leq \epsilon.$$

Example where lack of smoothness is bad:



$$\sigma = -\frac{1}{\beta} \cdot \nabla$$

Def: d -strongly convex for $d > 0$ iff

$$\lambda_{\min}(\nabla^2 f(y)) \geq d.$$