

\* HW 5 due Nov 4

\* Today: GD for Smooth fns, Online gradient descent

## Smoothness and Strong-Convexity

Roughly, smoothness is an upper bound on 2nd derivative  
strong-convexity is a lower " " "

$\beta$ -Smoothness:

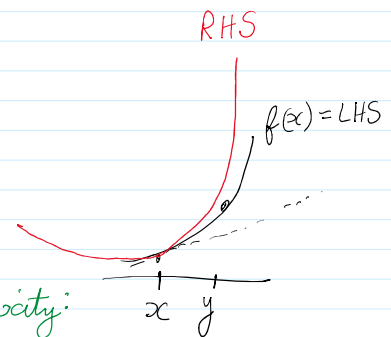
$$(1) \quad \nabla^2 f(x) \preceq \beta \quad \leftarrow \text{all e.v.s} \leq \beta$$

$$(2) \quad \|\nabla f(x) - \nabla f(y)\| \leq \beta \cdot \|x - y\|$$

$$(3) \quad f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\beta}{2} \|y - x\|^2$$

Recall 1st-order defn of convexity:

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle$$



$\alpha$ -Strong Convexity:

$$(1) \quad \nabla^2 f(x) \succeq \alpha \quad \leftarrow \text{all e.v.s} \geq \alpha$$

$$(2) \quad \|\nabla f(x) - \nabla f(y)\| \geq \alpha \cdot \|x - y\|$$

$$(3) \quad f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\alpha}{2} \|y - x\|^2$$

E.g.:  $f(x) = \|Ax - b\|^2 / 2$   
 $\nabla f(x) = Ax - b$ ,  $\nabla^2 f(x) = A^T A$   
 $\beta = \text{Largest e.v. of } A^T A$   
 $\alpha = \text{Min e.v. of } A^T A$

What can be shown:

	Function class	# iterations (ignoring non- $\epsilon$ terms)	Step Size
→ 1)	Convex	$1/\epsilon^2$	$1/\sqrt{T}$
→ 2)	Strongly-convex	$1/\epsilon$	$1/(\alpha t)$
→ 3)	Smooth	$1/\epsilon$	$1/\beta$
→ 4)	Str-convex + Smooth	$\text{polylog}(1/\epsilon)$	$1/\beta$

Last lecture  
HW5  
Today

## Gradient Descent for $\beta$ -Smooth functions

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

**Thm:** Grad descent on  $\beta$ -smooth convex fn  
and given  $x_0$  s.t.  $\|x_0 - x^*\| \leq D$  gives

$$f(x_T) - f(x^*) \leq \beta \frac{D^2 (1 + \ln T)}{2T}$$

**Remark:** (1) The  $\ln T$  can be removed with a better analysis.  
(2) This result is about  $x_T$  (and not  $\sum_{t=1}^T x_t$ )

**Proof:** Define potential  $\phi_t := t \cdot [f(x_t) - f(x^*)]$

Want to show  $\phi_T$  is small

$$\begin{aligned} \phi_{t+1} - \phi_t &= (t+1) [f(x_{t+1}) - f(x^*)] - t [f(x_t) - f(x^*)] \\ &= (t+1) \left[ \underbrace{f(x_{t+1}) - f(x_t)}_{\text{Claim 1}} + \underbrace{[f(x_t) - f(x^*)]}_{\text{Claim 2}} \right] \\ &\stackrel{\text{Claim 1}}{\leq} \underbrace{-\frac{\|\nabla f(x_t)\|^2}{2\beta}}_{\text{Claim 2}} \leq \underbrace{\frac{(t+1)\|\nabla f(x_t)\|^2}{2\beta}}_{\text{Claim 2}} + \underbrace{\frac{D^2\beta}{2(t+1)}}_{\text{Claim 2}} \end{aligned}$$

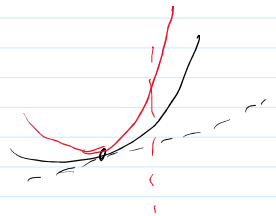
Assuming claims,

$$\begin{aligned} \underbrace{\phi_T - \phi_0}_{= T[f(x_T) - f(x^*)]} &= \sum_{t=0}^{T-1} (\phi_{t+1} - \phi_t) \leq \frac{D^2\beta}{2} \sum_{t=0}^{T-1} \frac{1}{t+1} \leq \frac{D^2\beta}{2} (1 + \ln T) \end{aligned}$$

**Claim 1:**  $f(x_{t+1}) - f(x_t) \leq -\frac{\|\nabla f(x_t)\|^2}{2\beta}$

**Pf:**  $LHS \stackrel{\text{smoothness}}{\leq} \langle \nabla f(x_t), x_{t+1} - x_t \rangle + \frac{\beta}{2} \|x_{t+1} - x_t\|^2$

$$= -\eta \|\nabla f(x_t)\|^2 + \frac{\beta}{2} \eta^2 \|\nabla f(x_t)\|^2 = RHS \text{ for } \eta = 1/\beta$$



**Claim 2:**  $f(x_t) - f(x^*) \leq \frac{t+1}{\beta} \frac{\|\nabla f(x_t)\|^2}{2} + \frac{D^2\beta}{2(t+1)}$

Claim 2: 
$$f(x_t) - f(x^*) \leq \frac{t+1}{\beta} \frac{\|\nabla f(x_t)\|^2}{2\beta} + \frac{D^2\beta}{2(t+1)}$$

Pf: LHS  $\stackrel{\text{Convexity}}{\leq} \langle \nabla f(x_t), x_t - x^* \rangle \leq \|\nabla f(x_t)\| \cdot \|x_t - x^*\|$

$$\leq \underbrace{\|\nabla f(x_t)\| \sqrt{\frac{t+1}{\beta}}}_a \cdot \underbrace{\frac{D}{\sqrt{(t+1)/\beta}}}_b$$

Recall  $(a-b)^2 \geq 0$   
 $\Rightarrow \frac{a^2 + b^2}{2} \geq ab$

$$\leq \frac{1}{2} \left[ \|\nabla f(x_t)\|^2 \frac{t+1}{\beta} + \frac{D^2\beta}{t+1} \right] = \text{RHS}$$

## Online Convex Optimization / Online Learning

Example (Regression): Given  $T$  labeled samples  $(a_1, b_1) \dots (a_T, b_T)$   
 where  $a_t \in \mathbb{R}^n$  and  $b_t \in \mathbb{R}$

Find  $x$  s.t.  $a_t x \approx b_t \leftarrow \text{offline optim / Learning}$

Formally, 
$$\min_x \sum_{t=1}^T (a_t x - b_t)^2 = \min_x \underbrace{\|Ax - b\|^2}_{f(x)}$$

Gradient descent will find  $x$  s.t.  $f(x) - f(x^*) \leq \epsilon$ .

What if the samples not given upfront and arrive over time?

## Problem Model

- 1) We are given a convex body  $K \subseteq \mathbb{R}^n$
  - 2)  $T$  rounds:
    - (a) Alg plays  $x_t \in K$
    - (b) Convex cost function  $f_t: \mathbb{R}^n \rightarrow \mathbb{R}$  is revealed
- E.g.  $f_t(x) = (a_t x - b_t)^2$

(b) Convex cost function  $f_t : \mathbb{R}^n \rightarrow \mathbb{R}$  is revealed

$$\text{Goal: } \min \sum_{t=0}^{T-1} f_t(x_t)$$

$$\text{Benchmark: } \min_{x^* \in K} \sum_{t=0}^{T-1} f_t(x^*)$$

Thm:

Online gradient descent guarantees

$$\frac{1}{T} \sum_{t=0}^{T-1} [f_t(x_t) - f_t(x^*)] \leq \frac{DG}{\sqrt{T}} \quad \text{where } D = \text{Diameter}(K)$$

$$\& G = \max_t \|\nabla f_t(x)\| \text{ for } x \in K$$

$$x_{t+1} = x_t - \eta \nabla f_t(x_t)$$