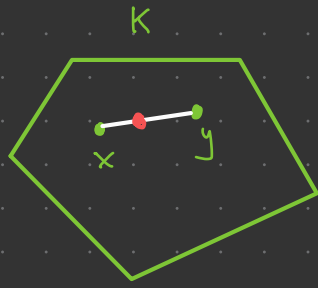


# ADFOCS 2021

Lecture 3: Gradient Descent  
Alejandro Cassis

Convexity:  $K \subseteq \mathbb{R}^n$  is convex if  $\forall x, y \in K \forall \alpha \in [0, 1]$   
 $(1-\alpha)x + \alpha y \in K$

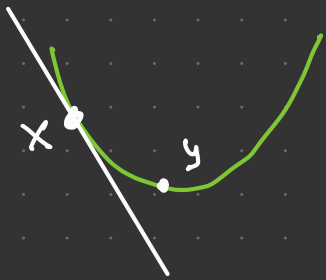


$f: K \rightarrow \mathbb{R}$  is convex if  $\forall x, y \in K \forall \alpha \in [0, 1]$   
 $f((1-\alpha)x + \alpha y) \leq (1-\alpha)f(x) + \alpha f(y)$



if  $f$  is differentiable:  $f(x + \alpha(y-x)) \leq (1-\alpha)f(x) + \alpha f(y)$

$$\Leftrightarrow \frac{f(x + \alpha(y-x)) - f(x)}{\alpha} \leq f(y) - f(x)$$



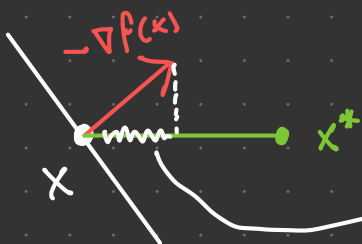
$$\text{as } \alpha \rightarrow 0: \nabla f(x)^T (y-x) \leq f(y) - f(x)$$

Usefulness: let  $x^* = \underset{x}{\operatorname{argmin}} f(x)$ . Plug  $y := x^*$ :

$$0 \leq \underbrace{f(x) - f(x^*)}_{\text{optimality gap}} \leq \nabla f(x)^T (x - x^*) = -\nabla f(x)^T (x^* - x)$$

optimality gap

$\rightarrow$  neg. gradient is positively correlated with direction to  $x^*$



progress is lower bounded by  $f(x) - f(x^*)$

GD Take 1: since moving in the negative gradient

direction brings us closer to  $x^*$ , this suggests

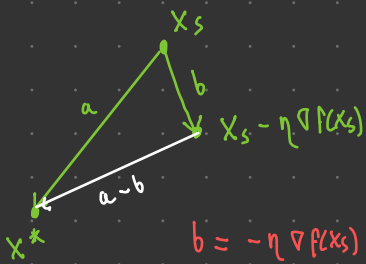
the following algorithm: (Assume  $K = \mathbb{R}^n$ )

let  $x_1 \in K$  be an initial point  
for  $s=1$  to  $T-1$ :

$$x_{s+1} := x_s - \eta \nabla f(x_s)$$

$\eta \in \mathbb{R}$   
learning rate  
or step size.

To analyze: use  $\|x^* - x_s\|^2$  as potential



change of potential:

$$\|a\|^2 - \|a-b\|^2 = \underbrace{2a^T b}_{\text{pos. curr. progress}} - \underbrace{\|b\|^2}_{\text{error term}}$$

$$\therefore \|x^* - x_s\|^2 - \|x^* - x_{s+1}\|^2 = 2(x^* - x_s)^T (-\eta \nabla f(x_s)) - \eta^2 \|\nabla f(x_s)\|^2$$

$$\Leftrightarrow \nabla f(x_s)^T (x_s - x^*) = \frac{\|x^* - x_s\|^2 - \|x^* - x_{s+1}\|^2}{2\eta} + \frac{\eta}{2} \|\nabla f(x_s)\|^2$$

**Theorem:** if  $\|\nabla f(x)\| \leq L \quad \forall x \in K$ , then

$$f\left(\frac{1}{T} \sum_{s=1}^T x_s\right) - f(x^*) \leq \frac{\|x_1 - x^*\|^2}{2\eta T} + \frac{\eta L^2}{2}$$

$$\begin{aligned} f\left(\frac{1}{T} \sum_{s=1}^T x_s\right) - f(x^*) &\leq \frac{1}{T} \sum_{s=1}^T f(x_s) - f(x^*) \leq \frac{1}{T} \sum_{s=1}^T \nabla f(x_s)^T (x_s - x^*) \\ &\leq \frac{\|x_1 - x^*\|^2}{2\eta T} + \frac{\eta L^2}{2} \quad \square \end{aligned}$$

**Theorem:** if  $\|\nabla f(x)\| \leq L \quad \forall x \in K$ , then

$$f\left(\frac{1}{T} \sum_{s=1}^T x_s\right) - f(x^*) \leq \frac{\|x_1 - x^*\|^2}{2\eta T} + \frac{\eta L^2}{2}$$

$R := \|x_1 - x^*\|^2$  "radius"

optimize  $\eta = \frac{R}{L\sqrt{T}} \Rightarrow f\left(\frac{1}{T} \sum_{s=1}^T x_s\right) - f(x^*) \leq \frac{RL}{\sqrt{T}}$

**Corollary:** Can find  $x$  st  $f(x) - f(x^*) \leq \xi$   
in  $T = O\left(\frac{R^2 L^2}{\xi^2}\right)$  iterations

**Remarks:**

-  $\|\nabla f(x)\| \leq L$  is equivalent to  $f$  being  $L$ -Lipschitz

e.g.  $\|\cdot\|_2$  is 1-Lipschitz  
i.e.  $\forall x, y \in K$   
 $|f(x) - f(y)| \leq L \|x - y\|$

- For some applications  $R, L$  are known.

Read convergence rate as  $1/\sqrt{T}$

- This rate is optimal in black box model

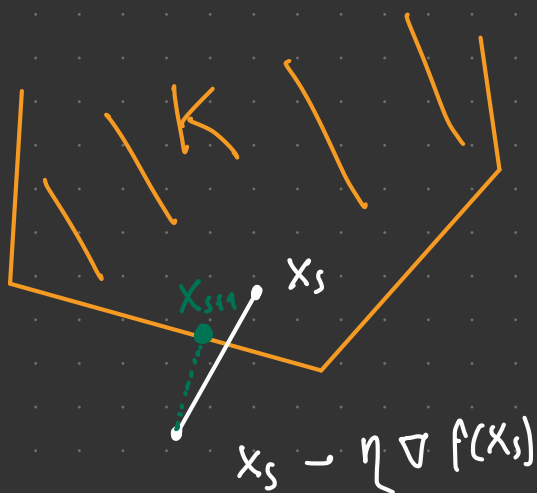
- Works for non-diff. fns. e.g.  $\|x\|_1$  or  $\|x\|_\infty$

What if  $K$  is not  $\mathbb{R}^n$ ?

Projected Gradient Descent:

$$x_{s+1} := \text{Proj}_K(x_s - \eta \nabla f(x_s))$$

where  $\text{Proj}_K(y) = \underset{x \in K}{\text{argmin}} \|x - y\|$



Same analysis goes through! Recall that potential

was  $\|x^* - x_s\|^2$ .

$$y_{s+1} := x_s - \eta \nabla f(x_s)$$

$$x_{s+1} := \text{Proj}_K(y_{s+1})$$

check:  $\|x^* - x_{s+1}\|^2 \leq \|x^* - y_{s+1}\|^2$

i.e: projection only brings us closer to  $x^*$ .

(only works for convex  $K$ !!!)

## Smoothness:

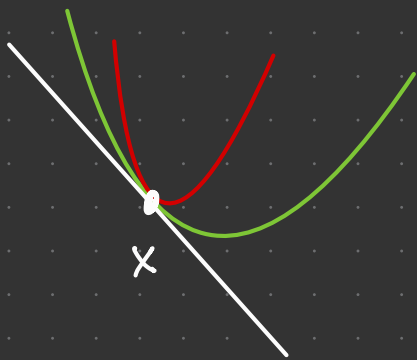
$f$  is  $\beta$ -smooth if its gradient is  $\beta$ -Lipschitz.

$$\text{i.e. } \forall x, y \in K \quad \|\nabla f(x) - \nabla f(y)\| \leq \beta \|x - y\|$$

if  $f$  is twice differentiable:  $\beta$ -smooth is equivalent to  $\nabla^2 f(x) \preceq \beta \cdot I \quad \forall x \in K$   
i.e.  $\lambda_{\max}(\nabla^2 f(x)) \leq \beta$ .

smoothness is useful to obtain upper bounds:

$$\left[ \text{lemma: } \forall x, y: f(y) \leq f(x) + \nabla f(x)^\top (y-x) + \frac{\beta}{2} \|x-y\|^2 \right]$$



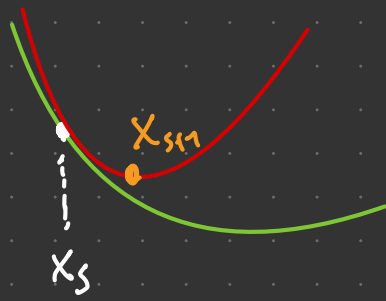
proof of lem: By Taylor's thm  $\exists z = (1-\alpha)x + \alpha y$  st

$$f(y) = f(x) + \nabla f(x)^\top (y-x) + \underbrace{\frac{1}{2} (y-x)^\top \nabla^2 f(z) (y-x)}_{\leq \frac{\beta}{2} \|y-x\|^2}$$

□

# How to exploit smoothness?

Minimize the upper bound!



$$x_{s+1} := \underset{z}{\operatorname{argmin}} \left\{ \underbrace{f(x_s) + \nabla f(x_s)^T (z - x_s) + \frac{\beta}{2} \|z - x_s\|^2}_{g(z) :=}$$

$$\nabla g(z) = \nabla f(x_s) + \beta (z - x_s) = 0$$

$$z = x_s - \frac{1}{\beta} \nabla f(x_s)$$

$$\Rightarrow \text{update rule: } x_{s+1} = x_s - \frac{1}{\beta} \nabla f(x_s)$$

what progress do we make!

$$\left[ \text{smoothness lemma: } f(x_{s+1}) - f(x_s) \leq -\frac{1}{2\beta} \|\nabla f(x_s)\|^2 \right]$$

proof: by smoothness:

$$f(z) \leq f(x_s) + \nabla f(x_s)^T (z - x_s) + \frac{\beta}{2} \|z - x_s\|^2$$

$$\text{plugging } z := x_{s+1} = x_s - \frac{1}{\beta} \nabla f(x_s)$$

$$\Rightarrow f(x_{s+1}) - f(x_s) \leq -\frac{1}{\beta} \|\nabla f(x_s)\|^2 + \frac{1}{2\beta} \|\nabla f(x_s)\|^2$$

□

# (Strong) Convexity:

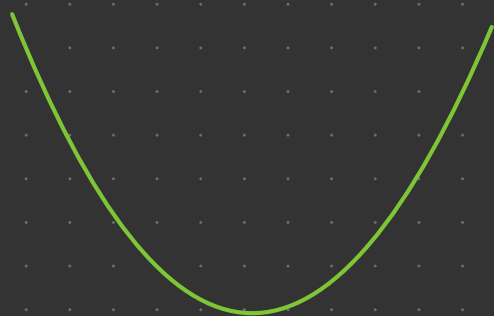
Recall that  $f$  convex gives us lower bounds:

$$\forall x, y \quad f(y) \geq f(x) + \nabla f(x)^T (y-x)$$

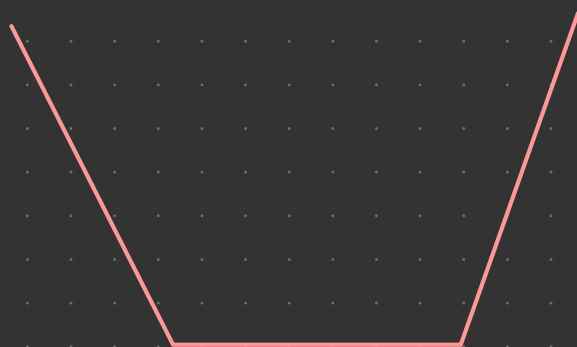
$f$  is  $\alpha$ -convex if  $f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{\alpha}{2} \|x-y\|^2$

[Similar as smoothness, for twice diff  $f$ 's we have that this is equivalent to  $\nabla^2 f(x) \succeq \alpha \cdot I \quad \forall x$ ]

- strongly convex functions have unique minimizers



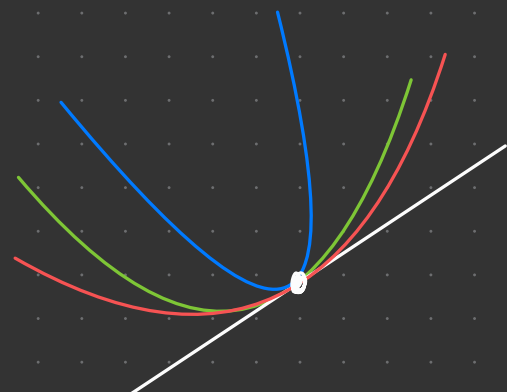
strongly convex



convex

$\beta$ -Smoothness +  $\alpha$ -convex  $\equiv$  quadratic upper and lower bounds

$$\forall x, y: \quad f(x) + \nabla f(x)^T (y-x) + \frac{\alpha}{2} \|x-y\|^2 \leq f(y) \leq f(x) + \nabla f(x)^T (y-x) + \frac{\beta}{2} \|x-y\|^2$$





$f$  is  $\alpha$ -convex if  $f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{\alpha}{2} \|x-y\|^2$

$$\left[ \alpha\text{-convex lemma: } \forall x \in \mathbb{R}^n \right. \\ \left. f(x) - f(x^*) \leq \frac{1}{2\alpha} \|\nabla f(x)\|^2 \right]$$

proof. by  $\alpha$ -convex:  $f(x^*) \geq f(x) + \nabla f(x)^T (x^* - x) + \frac{\alpha}{2} \|x^* - x\|^2$

$$\geq f(x) + \min_u \left\{ \nabla f(x)^T u + \frac{\alpha}{2} \|u\|^2 \right\}$$

minimizer is  $u = -\frac{1}{\alpha} \nabla f(x)$

$$= f(x) - \frac{1}{2\alpha} \|\nabla f(x)\|^2 \quad \square$$

$$\left[ \text{smoothness lemma: } f(x_{s+1}) - f(x_s) \leq -\frac{1}{2\beta} \|\nabla f(x_s)\|^2 \right]$$

$$\left[ \alpha\text{-convex lemma: } f(x) - f(x^*) \leq \frac{1}{2\alpha} \|\nabla f(x)\|^2 \right]$$

$$\text{Analyzing gradient descent: } x_{s+1} := x_s - \frac{1}{\beta} \nabla f(x_s)$$

$$\begin{aligned} f(x_{s+1}) - f(x^*) &\leq f(x_s) - f(x^*) - \frac{1}{2\beta} \|\nabla f(x_s)\|^2 \\ &\leq f(x_s) - f(x^*) - \frac{\alpha}{\beta} [f(x_s) - f(x^*)] \\ &= \left(1 - \frac{\alpha}{\beta}\right) [f(x_s) - f(x^*)] \\ &\leq \left(1 - \frac{\alpha}{\beta}\right)^s [f(x_1) - f(x^*)] \\ &\leq \exp\left(-\frac{\alpha}{\beta} \cdot s\right) [f(x_1) - f(x^*)] \end{aligned}$$

**Theorem:** Let  $f$  be  $\alpha$ -convex and  $\beta$ -smooth.  
Then,  $T$  steps of G.D. w/ step size  $\frac{1}{\beta}$  satisfy

$$f(x_{T+1}) - f(x^*) \leq \exp\left(-\frac{\alpha}{\beta} \cdot T\right) [f(x_1) - f(x^*)]$$

$$\Rightarrow \underbrace{\frac{\beta}{\alpha}}_K \log\left(\frac{f(x_1) - f(x^*)}{\epsilon}\right) \text{ to get } \epsilon\text{-close}$$

$K$ : condition number

**Theorem:** Let  $f$  be  $\alpha$ -convex and  $\beta$ -smooth.

Then,  $T$  steps of G.D. w/ step size  $\frac{1}{\beta}$  satisfy

$$f(x_{T+1}) - f(x^*) \leq \exp\left(-\frac{\alpha}{\beta} \cdot T\right) [f(x_1) - f(x^*)]$$

**Remarks:**

- Same result holds for projected G.D.  
analysis is a bit more involved see [Boeck 3.2]

- Progress is measured w/  $f(x_1) - f(x^*)$ .

Using  $\beta$ -smoothness:  $f(x_{T+1}) - f(x^*) \leq \beta \exp\left(-\frac{\alpha}{\beta} \cdot T\right) \|x_1 - x^*\|^2$

- Important application: solving  $Ax = b$

when  $\alpha I \preceq A \preceq \beta I$  e.g. Laplacian  $L$

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

$$\nabla f(x) = Ax - b$$

$\therefore$  unique minimizer satisfies  $\nabla f(x^*) = Ax^* - b = 0$

More in Rasmus' lectures!

# Excursion: Preconditioning

We saw that if  $\alpha I \preceq \nabla^2 f(x) \preceq \beta I$  and  $\frac{\beta}{\alpha}$  is "small" then G.D. converges exponentially.

→ Suppose that instead  $\alpha \cdot H \preceq \nabla^2 f(x) \preceq \beta \cdot H$   
for some p.d.  $H$ .

intuition:  $f$  is well condition in a different basis

Consider  $g(x) := f(Mx)$  where  $M$  is p.d.

$$\nabla g(x) = M^T \nabla f(Mx)$$

$$\nabla^2 g(x) = M^T \nabla^2 f(Mx) M$$

→ setting  $M := H^{-1/2}$  :  $\nabla^2 g(x) = H^{-1/2} \nabla^2 f(H^{-1/2}x) H^{-1/2}$

$$\therefore \alpha H \preceq \nabla^2 f(H^{-1/2}x) \preceq \beta \cdot H$$

recall  
 $A \succeq 0 \Leftrightarrow B^T A B \succeq 0$   
 $\forall B \in \mathbb{R}^{m \times n}$

$$\Leftrightarrow \alpha \cdot I \preceq \underbrace{H^{-1/2} \nabla^2 f(H^{-1/2}x) H^{-1/2}}_{= \nabla^2 g(x)} \preceq \beta \cdot I$$

i.e.  $g$  is  $\alpha$ -convex and  $\beta$ -smooth!

Gradient Descent on  $g$ : 
$$x_{s+1} = x_s - \frac{1}{\beta} H^{-\frac{1}{2}} \nabla f(H^{-\frac{1}{2}} x_s)$$

$$\Leftrightarrow H^{-\frac{1}{2}} x_{s+1} = \underbrace{H^{-\frac{1}{2}} x_s}_{y_s} - \frac{1}{\beta} H^{-1} \nabla f(\underbrace{H^{-\frac{1}{2}} x_s}_{y_s})$$

$$\therefore y_{s+1} = y_s - \frac{1}{\beta} H^{-1} \nabla f(y_s)$$

Newton's Method: 
$$x_{s+1} = x_s - (\nabla^2 f(x_s))^{-1} \nabla f(x_s)$$

i.e. preconditioning using the Hessian at  $x_s$ !

[ See Vishnoi's Chapter 9 for (much) more details on Newton's Method ]

## G.D. for smooth functions.

What if  $f$  is  $0$ -convex? i.e. only  $\beta$ -smooth.

Reduction to strongly convex case:

Fix  $x_1 \in \mathbb{R}^n$  and define

$$g(x) := f(x) + \underbrace{\frac{\alpha}{2} \|x_1 - x\|^2}_{\text{Regularization}}$$

check:  $g$  is  $\alpha$ -convex and  $(\alpha + \beta)$ -smooth

$\therefore$  can find  $x$  s.t.  $g(x) - g(x^*) \leq \xi$

in  $O\left(\frac{\alpha + \beta}{\alpha} \log\left(\frac{1}{\xi}\right)\right)$  iterations

let  $x^* = \underset{x}{\operatorname{argmin}} f(x)$ ,  $\tilde{x} = \underset{x}{\operatorname{argmin}} g(x)$

$$f(x) - f(x^*) = g(x) - g(x^*) + \frac{\alpha}{2} \left( \|x - x_1\|^2 - \|x^* - x_1\|^2 \right)$$

$$\leq \xi + \frac{\alpha}{2} R^2 \leq 2\xi$$

$$\alpha := \xi / R^2$$

$\therefore$  can find  $\xi$ -optimal pt in  $O\left(\frac{\beta \cdot R^2}{\xi} \cdot \log\left(\frac{f(x_1) - f(x^*)}{\xi}\right)\right)$

$$[R = \|x^* - x_1\|]$$

# Summary:

	Rate	# iters
L-Lipschitz	$LR/\sqrt{T}$	$L^2 R^2 / \epsilon^2$
$\beta$ -smooth	$\beta R^2 / T$	$\beta R^2 / \epsilon$
$\alpha$ -convex + $\beta$ -smooth	$R^2 \exp\left(-\frac{T}{K}\right)$ $(K = \frac{\beta}{\alpha})$	$K \log\left(\frac{R}{\epsilon}\right)$

- Various parameters:  $\beta, L, R$ . Adaptive G.D: obtain optimal rates w/o knowing these.

See Alina's Lectures!

- Rate for L-Lipschitz is optimal.

See [Bubeck 3.5] for lower bound.

$\beta$ -smooth can be improved to  $\frac{\beta R^2}{T^2}$   
 $\alpha$ -convex +  $\beta$ -smooth to  $R^2 \exp\left(-\frac{T}{\sqrt{K}}\right)$

via "acceleration". See Alina's L.2.

also Vishnoi's ch 8