

Lec. 19 (11/20/2022)

## First Order Algorithms

only requires  
first derivative/gradient  
of the objective  
function, i.e.,  $\nabla f$

Gradient descent: (Most well-known 1<sup>st</sup> order algorithm)

Cauchy (1820)

$\sim 200$  yrs old!!

$$\min_{\underline{x} \in \mathcal{X}} f(\underline{x})$$

Suppose,  $\underline{x} \equiv \mathbb{R}^n$   
unconstrained

$$\underline{x}_{k+1} = \underline{x}_k - \eta_k \nabla f(\underline{x}_k),$$

where  $k = 0, 1, 2, \dots$

$\uparrow$   
iteration index

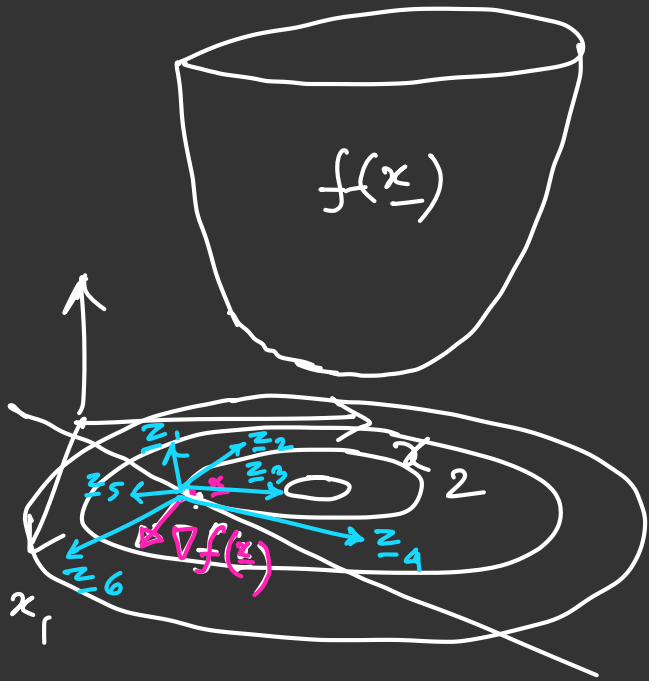
time step (may be constant  $\eta_k \equiv \eta$ )

Iterative algorithm

Main idea: Want a descent:

$$f(\underline{x}_{k+1}) \leq f(\underline{x}_k) \quad \forall k = 0, 1, 2, \dots$$

An algorithm that does this is called "Descent method"



1<sup>st</sup> order Taylor approximation  
of  $f(\underline{x} + \underline{z})$  around  $\underline{x}$  is

$$f(\underline{x} + \underline{z}) \approx \hat{f}(\underline{x} + \underline{z})$$

$$= f(\underline{x}) + \underbrace{\langle \nabla f(\underline{x}), \underline{z} \rangle}_{D_{\underline{z}} f(\underline{x})}$$

Directional derivative  
of  $f$  @  $\underline{x}$  in the  
direction  $\underline{z}$

$\therefore$  Any direction  $\underline{z}$  that makes  $D_{\underline{z}} f(\underline{x}) < 0$   
is a descent direction

$\therefore$  Any  $\underline{z}$  that makes an obtuse ( $> 90^\circ$ ) angle with  $\nabla f(\underline{x})$  is a descent direction.

So, out of all descent directions  $\underline{z}$ , what is the "steepest descent" direction?  
(makes  $D_{\underline{z}} f(\underline{x})$  most negative)

$$D_{\underline{z}} f(\underline{x}) = \langle \nabla f(\underline{x}), \underline{z} \rangle$$

$\therefore$  Most negative value of  $D_{\underline{z}} f(\underline{x})$  is achieved by  $\underline{z} = -\nabla f(\underline{x})$

$$\text{giving } \langle \nabla f(\underline{x}), -\nabla f(\underline{x}) \rangle = -\|\nabla f(\underline{x})\|_2^2$$

How fast can any gradient-based algorithm/  
first order algorithm converge?

---

Answer: Nemirovskii - Yudin (1983)  
optimal rate :  $\mathcal{O}(1/k^2)$

---

But gradient descent (particular 1<sup>st</sup> order  
algorithm) does NOT achieve this optimal  
1<sup>st</sup> order rate. Only achieves  $\mathcal{O}(1/k)$ .

Specifically, suppose  $\eta_k \equiv \eta > 0$

- $\boxed{If}$
- ①  $f$  is convex
  - ②  $f \in C^1$  ( $\Leftrightarrow$  continuously differentiable)
  - ③  $f$  is Lipschitz continuous with Lipschitz constant  $L$ :

$$\Leftrightarrow \|\nabla f(\underline{x}) - \nabla f(\underline{y})\|_2 \leq L \|\underline{x} - \underline{y}\|_2$$

$\boxed{Then}$  with  $\eta = \frac{1}{L}$ , generated sequence  $\{\underline{x}_k\}$   
using grad. descent next p. 8  $\forall \underline{x}, \underline{y} \in \text{dom}(f)$

satisfies :

$$f(\underline{x}_k) - f(\underline{x}_{opt}) \leq \underbrace{\frac{L}{2k} \|\underline{x}_0 - \underline{x}_{opt}\|_2^2}_{O(1/k)}$$

If  $f \in S_{L,m}^1$

$$\Leftrightarrow \textcircled{1} f \in C^1$$

$$\textcircled{2} f \text{ is convex}$$

$$\textcircled{3} \nabla f \text{ is Lipschitz with constant } L$$

$$\textcircled{4} f \text{ is } m\text{-strongly convex}$$

$$\Leftrightarrow f(\underline{y}) \geq f(\underline{x}) + \langle \nabla f(\underline{x}), \underline{x} - \underline{y} \rangle + \frac{m}{2} \|\underline{x} - \underline{y}\|_2^2$$

$$\text{Then } \eta_{opt} = \frac{2}{L+m}$$

optimal  $\uparrow$  constant stepsize convergence  
Still  $O(1/k)$  extra assumption

Gradient descent  
(1820, Cauchy)

$$\underline{x}_{k+1} = \underline{x}_k - \eta_k \nabla f(\underline{x}_k)$$

Descent  
method

$$\mathcal{O}(1/k)$$

Heavy Ball  
(1964, Polyak)

$$\underline{y}_{k+1} = \underline{x}_k - \eta_k \nabla f(\underline{x}_k)$$

$$\underline{x}_{k+1} = \underline{y}_{k+1} + \alpha_k (\underline{x}_k - \underline{x}_{k-1})$$



$$\underline{x}_{k+1} = \underline{x}_k - \eta_k \nabla f(\underline{x}_k) + \underbrace{\alpha_k (\underline{x}_k - \underline{x}_{k-1})}_{\text{momentum term}}$$

momentum  
term  
but may

Nesterov's  
accelerated (1983)  
Grad. descent

$$\underline{y}_{k+1} =$$

$$\underline{x}_k - \eta_k \nabla f(\underline{x}_k)$$

$$\underline{x}_{k+1} = \underline{y}_{k+1}$$

$$+ \alpha_k (\underline{y}_{k+1} - \underline{y}_k)$$

NOT descent  
method achieve  $\mathcal{O}(1/k^2)$



Heavy ball is NOT globally convergent  
even for strongly convex function.

In the domain in which it converges, it  
can beat gradient descent by achieving  $\mathcal{O}(1/k^2)$ .

---

## Handling Constraints

$$\boxed{\min_{\underline{x} \in \mathcal{X}} f(\underline{x})}$$

① Projected Gradient Descent:

$$\underline{x}_{k+1} = \text{proj}_{\mathcal{X}}^{\|\cdot\|_2} \left( \underline{x}_k - \eta_k \nabla f(\underline{x}_k) \right)$$



$$\underline{y}_{k+1} = \underline{x}_k - \eta_k \nabla f(\underline{x}_k)$$

$$\underline{x}_{k+1} = \text{proj}_{\mathcal{X}}^{\|\cdot\|_2}(\underline{y}_{k+1})$$

$$= \underset{\underline{x} \in \mathcal{X}}{\text{argmin}} \frac{1}{2} \|\underline{x} - \underline{y}_{k+1}\|_2^2$$

Next pg.

## ② Mirror Descent:

Constructs a "mirror function"

strictly convex function  $\Psi$   
related to the geometry  
of the constraint set  $\mathcal{X}$

$$\begin{cases} \nabla \Psi(\underline{y}_{k+1}) = \nabla \Psi(\underline{x}_k) - \eta_k \nabla f(\underline{x}_k) \\ \underline{x}_{k+1} = \underset{\mathcal{X}}{\text{proj}}^{D_\Psi}(\underline{y}_{k+1}) \end{cases}$$

where  $D_\Psi$  is the Bregman divergence induced  
by the mirror map/function  $\Psi$

Again,

$$\text{proj}_{\mathcal{X}}^{D_{\Psi}}(\underline{\eta}) := \underset{\underline{\xi} \in \mathcal{X}}{\text{argmin}} D_{\Psi}(\underline{\xi}, \underline{\eta})$$

projection of  $\underline{\eta}$  onto the set  $\mathcal{X}$

w.r.t.  $D_{\Psi}$

↑ not necessarily Euclidean distance