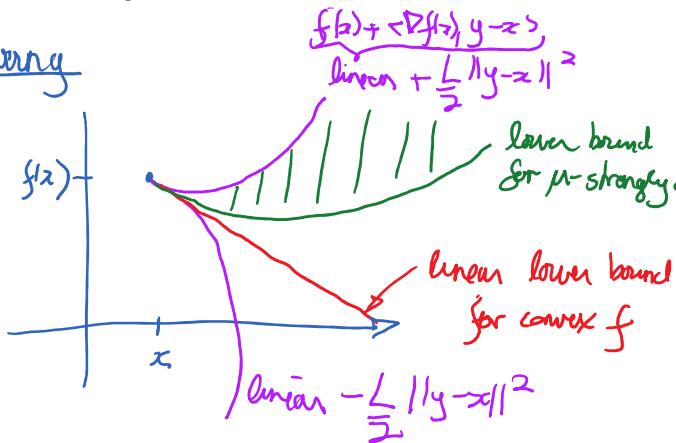


## Lecture 10 - subgradient method

Tuesday, February 16, 2021 14:33

- today:
  - basic gradient methods
  - subgradient method

drawing



when  $f$  is twice differentiable

$$L = \sup_x (\lambda_{\max}(H(x)))$$

$$\mu = \inf_x (\lambda_{\min}(H(x)))$$

$$f \text{ is } \mu\text{-strongly convex} \iff f - \frac{\mu}{2} \|x\|^2 \text{ is convex}$$

gradient descent:

$$x_{t+1} = x_t - \gamma Df(x_t) \quad \gamma = \frac{1}{L}$$

a) when is  $f$  convex  $\nless L$ -smooth

$$f(x_t) - \min_x f(x) \leq O\left(\frac{Lr_0^2}{t}\right)$$

$\underbrace{x}_{\triangleq x^*}$

"sublinear"

[see Nesterov book for proof]

$$\triangleq \underset{\bar{x} \in X^*}{\operatorname{argmin}} \|x_0 - \bar{x}\|_2$$

$$\text{where } r_0 \geq \operatorname{dist}(x_0, X^*)$$

$$\operatorname{argmin}_x f(x)$$

$$\frac{x^T H x}{2}$$

note: no guarantee on  $\operatorname{dist}(x_t, X^*)$   
(for general  $L$ -smooth for  $t \leq \dim(X)$ )  
convex sets

$\hookrightarrow$  Nesterov lower bound



$$\lambda_2 \rightarrow \lambda_{\max}$$

$$\lambda_1 \rightarrow \lambda_{\min}$$

b) if  $f$  is  $\mu$ -strongly convex  $\nless L$ -smooth

$$f(x_t) - f(x^*) \leq O\left(\exp\left(-\frac{\mu}{L} t\right)\right) \quad \text{"linear rate"}$$

$$\frac{L}{\mu} \triangleq \underline{\text{condition # of } f}$$





Newton's method

$$x_{t+1} = x_t - \Delta_t [H(x_t)]^{-1} \nabla f(x_t)$$

Subgradient method

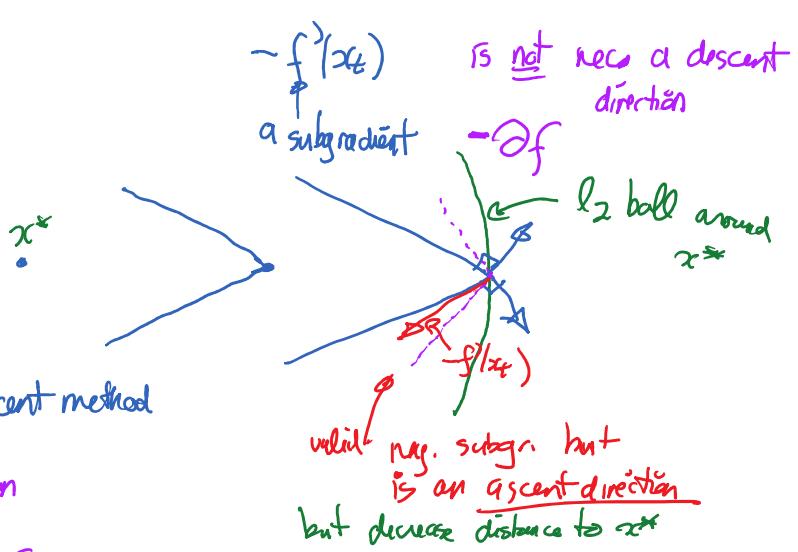
non-descent methods

$f$  smooth  $\Rightarrow$  smooth sublevel sets



$-\nabla f$  is a descent direction

Non-smooth  $f$



(\*) Subgradient method is not nec a descent method

but  $-f'(x_t)$  is a descent direction

on  $\|x(\bar{x}) - \bar{x}\|_2^2$  for any  $\bar{x}$   
in sublevel set at  $x$

$$x_t - \gamma f'(x_t)$$

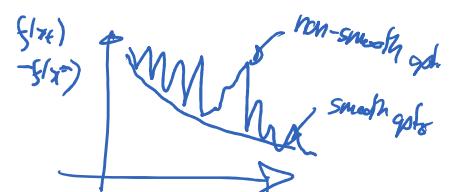
$x(\gamma)$  gets closer to any  $\bar{x}$  s.t.  $f(\bar{x}) \leq f(x_t)$  for small enough  $\gamma$

thus get closer to any  $x^*$

\* in non-smooth optimization,  $f(x_t)$  can go up & down

to stabilize

$\Rightarrow$  combine multiple pt.  $x_t$  to get  $\hat{x}_T$



augment  $f(x_t)$   $\{x_t\}_{t=1}^T$  [in batch setting]

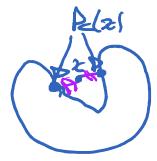
weighted average  $\hat{x}_T = \sum_t p_t^{(T)} x_t$  [for stochastic setting or]

(\*) projection operator on a  
closed convex set  $C$

$$P_C(x) \triangleq \underset{y \in C}{\operatorname{argmin}} \|x - y\|_2^2$$

convex weights

when too expensive  
to compute  $f(x)$



"Euclidean projection" of  $x$  on  $C$

$P_C(\cdot)$  is non-expansive

$$\text{i.e. } \|P_C(x) - P_C(y)\|_2 \leq \|x - y\|_2 \quad \forall x, y$$

• if  $y \in C$ , then  $P_C(y) = y$

$$\min_{x \in C} f(x) \rightarrow x^* \quad P_C(x^*) = x^*$$

$$\|P_C(y) - x^*\|_2 \leq \|y - x^*\|_2$$

$\rightarrow$  thus projection on  $C$  just make iterates  
closer to  $x^*$

15h33

### Stochastic Subgradient Method

Setup: want to solve  $\min_{x \in C} f(x)$  where

$$f(x) \triangleq \mathbb{E}_{\xi} [h(x, \xi)]$$

Assumptions: 1)  $f$  &  $C$  are convex

2) projection on  $C$  is "cheap"

3) we have a stochastic oracle which gives  $g(x, \xi)$  for random  $\xi$

$$\text{s.t. } \boxed{\mathbb{E}_{\xi} [g(x, \xi) | x] = f'(x)}$$

some subgradient  
of  $f$  at  $x$

[examples:

a)  $f$  is diff. in  $x$  { "well behaved"

$$g(x, \xi) \triangleq P_x h(x, \xi)$$

then  $\mathbb{E}_{\xi} [P_x h(x, \xi)] = \nabla_x [\mathbb{E}_{\xi} [h(x, \xi)]] = \nabla f(x)$  "Leibniz rule" parameter

b) ERM example

$$f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x) \quad \text{o.g. } f_i(x) = \mathcal{J}(x^{(i)}, y^{(i)}, \|x\|)$$

(finite sum)

$$h(x, \xi) \triangleq f_{\xi}(x)$$

$\xi \in \{1, \dots, n\}$

at step  $t$ , sample  $i_t \stackrel{\text{uniform}}{\sim} \{1, \dots, n\}$

$$\text{use } g_t \triangleq g(x_t, i_t) \triangleq \hat{f}_{i_t}(x_t)$$

$$\text{here } \mathbb{E}_g[\hat{f}_{i_t}|x] = \frac{1}{n} \sum_{i=1}^n \hat{f}_i(x) = f'(x_t)$$

$$+ \frac{\Delta}{2} \|x\|^2$$

4)  $\mathbb{E} \|g(x, \xi)\|_2^2 \leq B^2$  (finite variance condition)

↳ This replaces the Lipschitz gradient ass.

for non-smooth/stochastic opt.

[sufficient condition  $\|h'(x, \xi)\| \leq B \quad \forall x, \xi$ ]

Algorithm - (projected stochastic sub. method)

$x_0 \in C$  initialization

for  $t=0, \dots, T-1$

let  $g_t$  be  $g(x_t, \xi_t)$  [from oracle]

let  $x_{t+1} = P_C[x_t - \gamma_t g_t]$  ↳ step-size

end  
output

$$\hat{x}_T \triangleq \sum_{t=0}^T p_{t,T} x_t$$

"weighted average"

where  $p_{t,T}$  are some convex (onto) coef.

$$\sum_t p_{t,T} = 1 \quad p_{t,T} \geq 0$$

→  $O(\frac{1}{\sqrt{T}})$  rate

Convergence proof:

Important inequality

$$s(y) \geq f(x) + \langle f'(x), y-x \rangle + \frac{\mu}{2} \|y-x\|^2 \quad (\text{H2O})$$

$$f(y) \geq f(x) + \langle f'(x), y-x \rangle + \frac{\mu}{2} \|y-x\|^2 \quad (\text{MZO})$$

$$\rightarrow -\langle f'(x), x-y \rangle \leq -(f(x) - f(y) + \frac{\mu}{2} \|y-x\|^2) + \underbrace{\forall x \in g}_{y}$$

$x_{t+1} = P_C(x_t - \gamma_t g_t)$  by def.

$$\|x_{t+1} - \tilde{x}\|_2^2 \stackrel{\text{by } P_C}{\leq} \|x_t - \gamma_t g_t - \tilde{x}\|_2^2$$

$$\begin{aligned} \text{any feasible pt. } \tilde{x} \in C \\ &= \|x_t - \tilde{x}\|^2 + \gamma_t^2 \|g_t\|^2 - 2\gamma_t \langle g_t, x_t - \tilde{x} \rangle \quad [\text{wth } \tilde{x} \in C] \end{aligned}$$

$$\begin{aligned} \mathbb{E}[\|x_{t+1} - \tilde{x}\|^2 | x_t] &\leq \|x_t - \tilde{x}\|^2 + \gamma_t^2 \mathbb{E}[\|g_t\|^2 | x_t] - 2\gamma_t \langle \mathbb{E}[g_t | x_t], x_t - \tilde{x} \rangle \\ &\stackrel{\text{using (t)}}{\leq} \|x_t - \tilde{x}\|^2 + \gamma_t^2 B^2 - 2\gamma_t [f(x_t) - f(\tilde{x}) + \frac{\mu}{2} \|x_t - \tilde{x}\|^2] \end{aligned}$$

$$\mathbb{E}[\mathbb{E}[ \cdot | x_t]] = \mathbb{E}[\cdot]$$

$$\mathbb{E}[\|x_{t+1} - \tilde{x}\|^2] \leq (1-\gamma_t) \mathbb{E}[\|x_t - \tilde{x}\|^2] - 2\gamma_t [\mathbb{E}f(x_t) - f(\tilde{x})] + \gamma_t^2 B^2$$

↑ true even if  $\mu=0$ ; for  $\gamma_t$  small enough

we have  $\mathbb{E}\|x_t - \tilde{x}\|^2 \underset{\equiv r(t)}{\text{decreases}}$  for any  $\tilde{x} \in C$  s.t.  $f(\tilde{x}) \leq \mathbb{E}f(x_t)$

i.e. we have  $r(t_{m+1}) \leq r(t_m)$