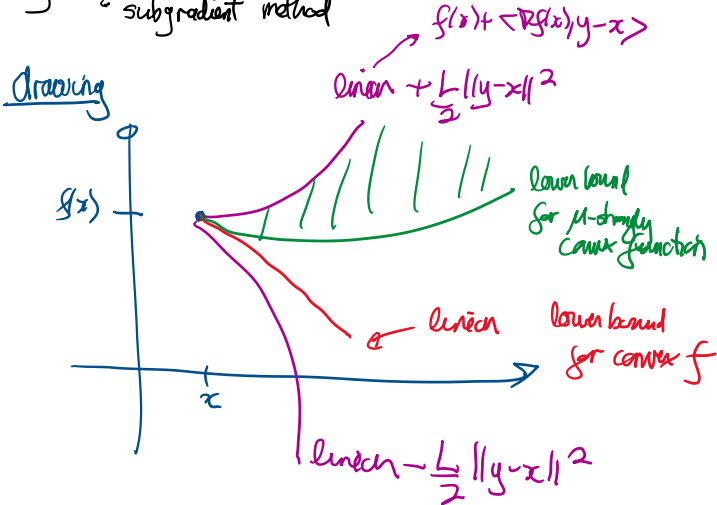


Lecture 10 - subgradient method

Thursday, February 15, 2024 9:35 AM

today:   
 • basic gradient methods   
 • subgradient method



when  $f$  is twice differentiable   
 $L = \sup_x (\lambda_{\max}(H(x)))$    
 $\mu = \inf_x (\lambda_{\min}(H(x)))$

$$f \text{ is } \mu\text{-strongly convex} \Leftrightarrow f(\cdot) - \frac{\mu \|\cdot\|^2}{2} \text{ is convex}$$

gradient descent:

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

a) when  $f$  is convex &  $L$ -smooth

$$f(x_t) - \min_x f(x) \leq O\left(\frac{L r_0^2}{t}\right)$$

$\stackrel{\Delta}{=} f^*$

↳ [see Nesterov book for proof]

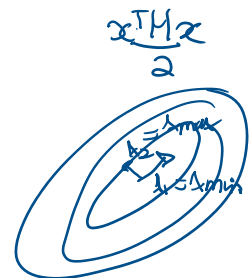
$\stackrel{\Delta}{=} \min_{x \in K} \|x_0 - x\|_2$    
 where  $r_0 \geq \text{dist}(x_0, X^*)$    
 $\uparrow$    
 $\text{argmin}_x f(x)$

"sublinear"

note: no guarantee on  $\text{dist}(x_t, X^*)$

(for general  $L$ -smooth convex  $f$ 's   
 for  $t \leq \text{dim}(x)$ )

↳ Nesterov lower bound



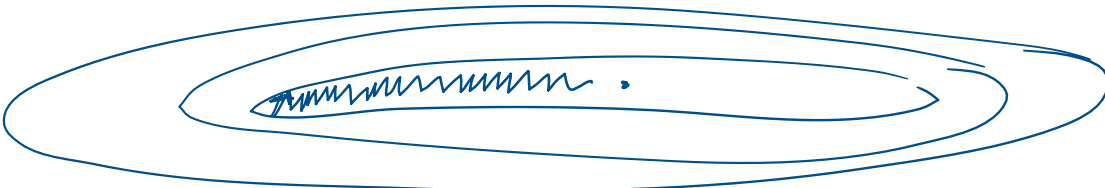
b) if  $f$  is  $\mu$ -strongly convex &  $L$ -smooth

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

"linear rate"

$$\exp\left(-\frac{t}{\kappa}\right)$$

$\frac{L}{\mu} \stackrel{\Delta}{=} \text{condition \# of } f \stackrel{\Delta}{=} \kappa$



Newton's method

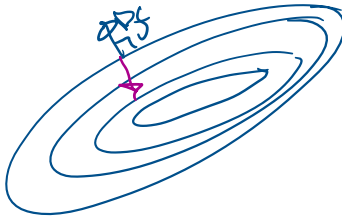
$$x_{k+1} = x_k - \gamma_k [H(x_k)^{-1}] \nabla f(x_k)$$

$$Hd = \nabla f(x_k)$$

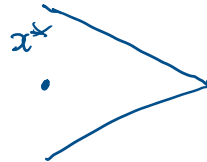
$$\min_d \|\nabla f(x_k) - Hd\|_2^2$$

Subgradient method

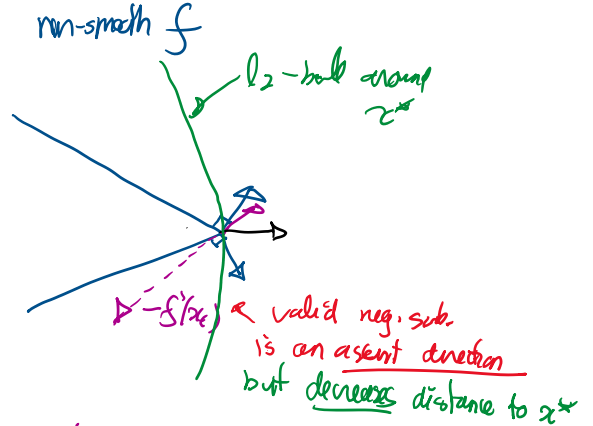
non-descent methods



$\rightarrow \nabla f$  is a descent direction



$-\nabla f(x_k)$  is not nec. a descent direction  
 a subgradient



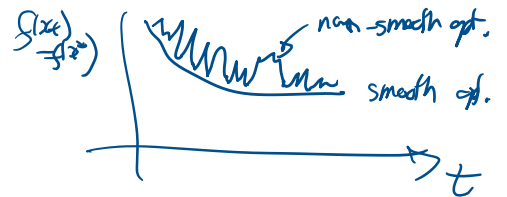
⊛ subgradient method is not nec. a descent method (on function values)

but  $-\nabla f(x_k)$  is a descent direction

on  $\|x(\gamma) - \tilde{x}\|_2^2$  for any  $\tilde{x}$  in sublevel set of  $x$   
 $L \triangleq x_k - \gamma \nabla f(x_k)$

$x(\gamma)$  gets closer to any  $\tilde{x}$  st.  $f(\tilde{x}) \leq f(x)$  for small enough  $\gamma$   
 thus get closer to any  $x^*$

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

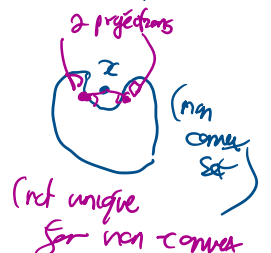


to stabilize:  $\rightarrow$  combine multiple pt.  $x_k$  to get  $\hat{x}_T$

argmin  $\sum_{k=1}^T f(x_k)$  [in batch setting]  
 weighted average  $\hat{x}_T = \sum_{k=1}^T p_k^{(T)} x_k$  [for stochastic setting or when  $f(x_k)$  is too expensive]

⊛ projection operator on a closed convex set  $C$   $P_C(x) \triangleq \text{argmin}_{y \in C} \|x - y\|_2^2$

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



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

(not unique for non-convex set)

$P_C(\cdot)$  is non-expansive i.e.  $\|P_C(x) - P_C(y)\|_2 \leq \|x - y\|_2 \quad \forall x, y$

• if  $y \in C$ , then  $P_C(y) = y$

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

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

$\Rightarrow$  this projection on  $C$  just makes iterates closer to  $x^*$

$$\left( \begin{array}{ccc} \min_x f(x) & \rightarrow & \min_{x \in C} f(x) \\ & & P_C(x) \end{array} \right)$$

10/28

stochastic subgradient method

some random objective function

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

computational assumptions { 2) projection on  $C$  is "cheap"  
3) we have a stochastic oracle which gives  $g(x, \xi)$  for random  $\xi$

$$\text{st. } \mathbb{E}_{\xi} [g(x, \xi) | x] = f'(x)$$

some subgradient of  $f$  at  $x$

[ examples:

a)  $f$  is differentiable in  $x$  & "well behaved"

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

"Leibniz rule"

$$\text{then } \mathbb{E}_{\xi} [\nabla_x h(x, \xi)] = \nabla_x [\mathbb{E}_{\xi} h(x, \xi)] = \nabla_x f(x)$$

b) ERM  
(finite sum)

$$f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x) \quad \text{e.g. } f_i(x) = \rho(x^{(i)}, y^{(i)}); \quad \text{"x" parameter}$$

$$= \mathbb{E}_{\xi} [f_{\xi}(x)] \quad \xi \in \xi_1, \dots, \xi_n \quad + \frac{\lambda}{2} \|x\|_2^2$$

$$h(x, \xi) \triangleq f_{\xi}(x) \quad \xi \in \xi_1, \dots, \xi_n$$

at step  $t$ , sample  $i_t \in \xi_1, \dots, \xi_n$

$$n(x, \xi) = \sum_{i=1}^n f_i(x) \quad \xi \in \{1, \dots, n\}$$

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

$$\text{use } g_t \triangleq g(x_t, i_t) \triangleq f_{i_t}'(x_t)$$

$$\text{here } \mathbb{E}_{\xi} [f_{i_t}'(x_t)] = \frac{1}{n} \sum_{i=1}^n f_i'(x_t) = f'(x_t)$$

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

↳ this replaces the Lipschitz gradient assump. for non-smooth & non-stochastic opt.

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

algorithm - (projected stochastic subgradient method)

$x_0 \in C$  initialization

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

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

let  $x_{t+1} = P_C [x_t - \delta_t g_t]$

↑  
step size

end

output  $\hat{x}_T \triangleq \sum_{t=0}^{T-1} \rho_{t,T} x_t$

where  $\rho_{t,T}$  are some ~~convex~~ combo. coeff.

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

"weighted average"

(w/ l.s.c.)  $\rightarrow O(\frac{1}{\sqrt{T}})$  rate

convergence proof:

important inequality:

$$f(y) \geq f(x) + \langle f'(x), y-x \rangle + \frac{\mu}{2} \|y-x\|^2 \quad (\mu \geq 0)$$

$$\Rightarrow \left( -\langle f'(x), x-y \rangle \leq -(f(x) - f(y)) + \frac{\mu}{2} \|y-x\|^2 \right) \quad (+)$$

$\forall x, y$

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

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

only feasible w.r.  $\tilde{x} \in C$

$$= \|x_t - \tilde{x}\|_2^2 + \delta_t^2 \|g_t\|^2 - 2\delta_t \langle g_t, x_t - \tilde{x} \rangle \quad [\text{valid } \forall \tilde{x} \in C]$$

any feasible  
pt.  $\tilde{x} \in C$

$$= \|x_t - \tilde{x}\|_2^2 + \delta_t^2 \|g_t\|^2 - 2\delta_t \langle g_t, x_t - \tilde{x} \rangle \quad [\text{valid } \forall \tilde{x} \in C]$$

$$\mathbb{E}[\|x_{t+1} - \tilde{x}\|^2 | x_t] \leq \|x_t - \tilde{x}\|_2^2 + \delta_t^2 \underbrace{\mathbb{E}[\|g_t\|^2 | x_t]}_{\leq B^2} - 2\delta_t \langle \underbrace{\mathbb{E}[g_t | x_t]}_{f'(x_t)}, x_t - \tilde{x} \rangle$$

$$\stackrel{\text{using (+)}}{\leq} \|x_t - \tilde{x}\|_2^2 + \delta_t^2 B^2 - 2\delta_t [f(x_t) - f(\tilde{x}) + \frac{\mu}{2} \|x_t - \tilde{x}\|^2]$$

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

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

↗ true even if  $\mu=0$ ; for  $\delta_t$  small enough

we have  $\mathbb{E}\|x_t - \tilde{x}\|^2 \stackrel{\triangleq}{=} r_t$  decreases for any  $\tilde{x} \in C$  st.  $f(\tilde{x}) \leq \mathbb{E}f(x_t)$

ie. we have  $r_{t+1} < r_t$