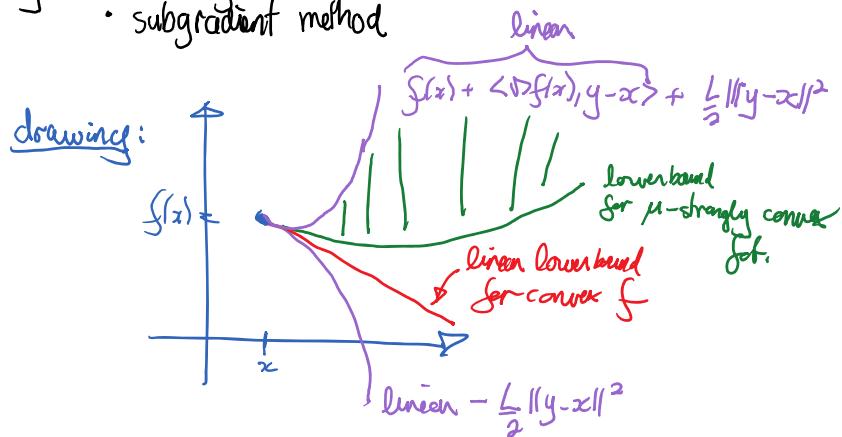


Lecture 10 - subgradient method

Thursday, February 16, 2023 1:35 PM

- Today:
 - basic gradient methods
 - subgradient method



Hessian

If f is twice differentiable

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

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

gradient descent:

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

a) when f is convex & L-smooth

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

$\underbrace{x}_{\in f^*}$

[See Nesterov book for proof]

"sublinear"

Note: no guarantee on $\text{dist}(x_t, x^*)$
 (for general L-smooth convex fns
 for $t \leq \text{dim}(x_t)$)

\hookrightarrow Nesterov lower bound

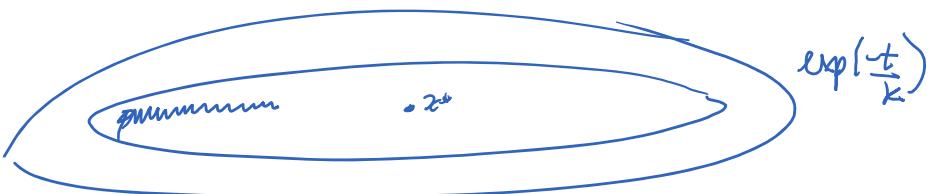
$$\begin{cases} \min_{x \in X^*} \|x_0 - x\|_2 \\ \text{where } r_0 \geq \text{dist}(x_0, X^*) \\ \text{augmin}_x f(x) \end{cases}$$

b) if f is μ -strongly convex & L-smooth

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

"linear rate"

$$\frac{L}{\mu} \triangleq \underline{\text{condition \# of } f} = k \geq 1$$



Newton's method

$$x_{t+1} = x_t - \gamma_t [H(x_t)]^{-1} Df(x_t)$$

{ linear rate sketch: $f(x_\gamma) \leq f(x) - \frac{1}{2L} \|Df(x)\|^2$

$\gamma = \frac{1}{L}$

μ -strongly convex f.

$f(x) - f^* \leq \frac{1}{2\mu} \|Df(x)\|^2$

$f(x_{t+1}) - f^* \leq f(x_t) - f^* - \frac{1}{2L} \|Df(x_t)\|^2$

$\triangleq h_{t+1}$

$- \|Df\|^2 \leq 2\mu (f(x) - f^*)$

$\leq f(x_t) - f^* - \frac{\mu}{2} (f(x_t) - f^*)$

$\triangleq h_t$

$h_{t+1} \leq (1 - \frac{\mu}{L}) h_t$

$\leq (1 - \frac{\mu}{L})^t h_0$

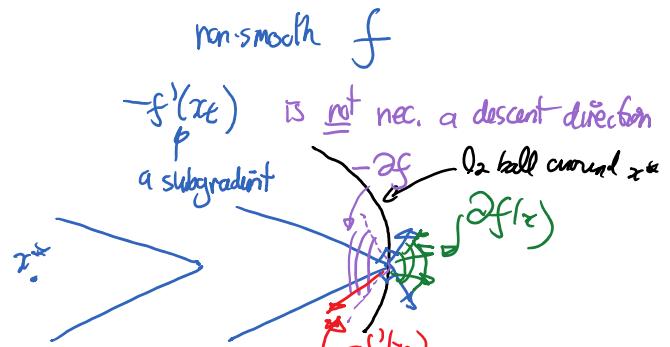
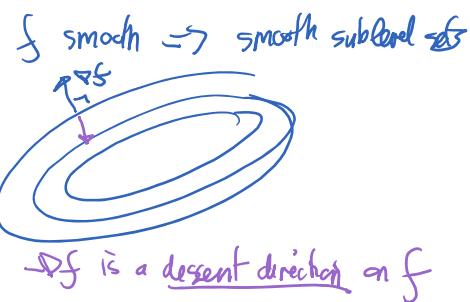
$(1 - \frac{\mu}{L})^t \leq \exp(-\frac{\mu t}{L}) \quad \forall t$

$\leq \exp(-\frac{\mu t}{L}) h_0 //$

14h19

Subgradient method

non-descent methods



② subgradient method is not nec. a descent method

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

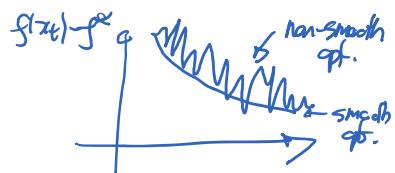
on $\|x(\gamma) - \tilde{x}\|_2^2$ for any \tilde{x}
in sublevel set at x

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

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

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



* in full gradient optimization, just goes up & down!

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



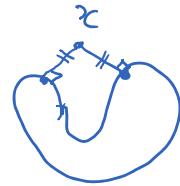
$$\underset{\sum x_t = \text{const}}{\operatorname{argmin}} f(x_t) \quad [\text{in batch setting}]$$

$$\text{weighted average } \hat{x}_T = \sum_t p_t^{(T)} x_t \quad \text{convex weights}$$

④ projection operator on
a closed convex set C

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

"Euclidean projection" of x on C



$P_C(\cdot)$ 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 $P_C(y) = y$

$$\underset{x \in C}{\operatorname{min}} f(x) \quad P_C(x^*) = x^*$$

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

\Rightarrow thus projection on C just makes iterates closes to x^*

stochastic subgradient method

setup: want to solve $\underset{x \in C}{\operatorname{min}} f(x)$ 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 ξ

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

stochastic subgradient
of f at x

[examples:

a) if h is diff. in $x \notin$ well behaved

$$g(x, \xi) \triangleq \nabla_x h(x, \xi) \quad \text{"Subgrad"}$$

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

b) ERM example $f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$ e.g. $f_i(w) = \ell(x^{(i)}, y^{(i)}; w) + \frac{\lambda \|w\|_2^2}{2}$

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

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

at step t , sample $\bar{w}_t \sim \xi \in \{1, \dots, n\}$

$$\text{use } g_t \triangleq g(x_t, \bar{w}_t) \triangleq f'_{\bar{w}_t}(x_t)$$

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

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

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

end
output

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

"weighted avg."

where $p_{t,T}$ are some
convex comb. coefs.

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

$\rightarrow O(\frac{1}{\sqrt{T}})$ rate

convergence proof:

important inequality (by μ -strong convexity)

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

$$\langle -f'(x), y-x \rangle \leq -(f(x) - f(y) + \frac{\mu}{2} \|y-x\|^2) \quad (+) \quad \forall x \in C$$

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

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

$$\begin{aligned} \underset{\text{any } \tilde{x} \in C}{=} & \|x_t - \tilde{x}\|^2 + \gamma_t^2 \|g_t\|^2 - 2\gamma_t \langle g_t, x_t - \tilde{x} \rangle \quad (\text{valid if } \tilde{x} \in C) \\ & = \|x_t - \tilde{x}\|^2 + \gamma_t^2 \|g_t\|^2 - 2\gamma_t \langle g_t, x_t - \tilde{x} \rangle \end{aligned}$$

$$\mathbb{E}[\|x_{t+1} - \tilde{x}\|^2 | x_t] \leq \|x_t - \tilde{x}\|^2 + \gamma_t^2 \mathbb{E}[\|g_t\|^2 | x_t] - 2\gamma_t \langle \mathbb{E}[g_t | x_t], x_t - \tilde{x} \rangle$$

Using (+)

$$\leq \|x_t - \tilde{x}\|^2 + \gamma_t^2 B^2 - 2\gamma_t [f(x_t) - f(\tilde{x}) + \frac{\mu}{2} \|x_t - \tilde{x}\|^2]$$

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

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

is true even if $\mu=0$

let $r_t \triangleq \mathbb{E}\|x_t - \tilde{x}\|^2$; for γ_t small enough

r_t decreases with t for any $\tilde{x} \in C$ s.t.

$$f(\tilde{x}) \leq \mathbb{E}[f(x_t)]$$

i.e. we have $r_{t+1} \leq r_t$