

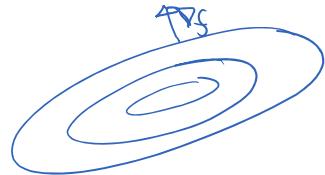
Lecture 10 - subgradient method

Thursday, February 7, 2019 13:37

today: subgradient Method

non descent methods

f smooth \Rightarrow smooth sublevel sets

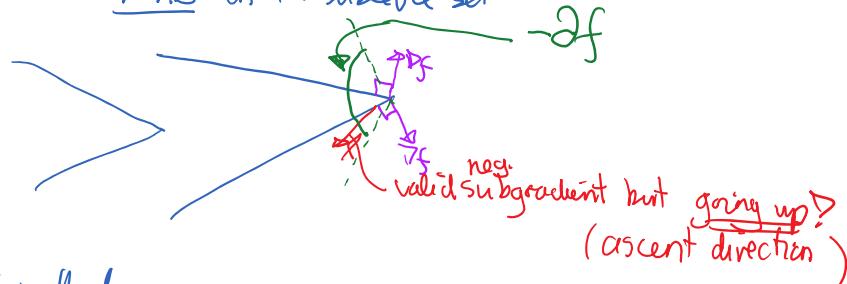


$-\nabla f$ is a descent direction

Non-smooth f

$-f'(x_t)$ is not nec. a descent direction
a subgradient

kinks in the sublevel set



④ subgradient method is not a descent method

but $-f'(x_t)$ is descent direction on $\|x(\gamma) - \tilde{x}\|^2$
for any \tilde{x} in sublevel set of x
(i.e. $\|x_t - \gamma f'(x_t) - \tilde{x}\|^2 \leq \|x_t - \tilde{x}\|^2$ for γ small enough)

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

and for any \tilde{x} s.t. $f(\tilde{x}) \leq f(x_t)$

thus get closer to any x^*

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

(to stabilize)

\Rightarrow combine multiple points x_t to get \hat{x}_T

$$\underset{\{x_t\}}{\operatorname{argmin}} f(x_t)$$

[in batch setting]

$$\text{weighted average } \hat{x}_T = \sum_t \beta_t x_t \quad \begin{array}{l} \text{weights} \\ \text{for stochastic setting} \\ \text{or when too} \\ \text{expensive to compute } f(x) \end{array}$$

* 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 a contraction ie $\|P_C(x) - P_C(y)\|_2 \leq \|x-y\|_2 \quad \forall x, y$

• if $y \in C$, then $P_C(y) = y$

and thus $\|P_C(x) - y\|_2 \leq \|x-y\|_2 \quad \forall y \in C$

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) $f \downarrow C$ are convex

- 2) projection on C is cheap
- 3) we have a stochastic oracle which gives a random $g(x, \xi)$
 s.t. $\mathbb{E}_\xi [g(x, \xi) | x] = f'(x)$
 $\qquad \qquad \qquad$ some subgradient of f at x
- [example : a) f is differentiable in x & "well behaved"
- $$g(x, \xi) \triangleq \nabla_x h(x, \xi)$$
- then $\mathbb{E}_\xi [\nabla_x h(x, \xi)] = f'(x)$ "Leibniz rule"

b) ERM example : $f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$ e.g. $f(x^{(i)}, y^{(i)}, \overset{\text{parameter}}{z}) + \frac{1}{2} \|x^{(i)}\|^2$

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

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

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

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

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

4) $\mathbb{E} \|g(x, \xi)\|_2^2 \leq B^2$ (finite variance condition) \leftarrow this replaces Lipschitz gradient assumption

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

← algorithm

Sufficient condition: $\|M(x, \xi)\| \leq B \quad \forall x, \xi$

algorithm

$x_0 \in C$ 'initialization'

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

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

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

\downarrow
stop size

end
output

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

"weighted average"

where p_t are some
convex comb. coeff.
 $\sum_{t=0}^T p_t = 1 \quad p_t \geq 0$

(4h3)

Convergence proof:

important inequality

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

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

use it on $-\langle f'(x_t), x_t - x^* \rangle$

$$y = z^*$$

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

$$\|x_{t+1} - \tilde{x}\|^2 \leq \|x_t - \gamma_t g_t - \tilde{x}\|^2$$

\uparrow
 any feasible
 i.e. $\tilde{x} \in C$

$$= \|x_t - \tilde{x}\|^2 + \gamma_t^2 \|g_t\|^2 - 2\gamma_t \langle g_t, x_t - \tilde{x} \rangle$$

\uparrow
 valid for all
 $\tilde{x} \in C$

$$\begin{aligned} \mathbb{E}[\|x_{t+1} - \tilde{x}\|^2 | x_t] &\leq \|x_t - \tilde{x}\|^2 + \gamma_t^2 \mathbb{E}[\|g_t\|^2 | x_t] - 2\gamma_t \underbrace{\langle \mathbb{E}[g_t | x_t], x_t - \tilde{x} \rangle}_{f'(x_t)} \\ &\stackrel{(\text{using } t)}{\leq} \|x_t - \tilde{x}\|^2 + \gamma_t^2 B^2 - 2\gamma_t [f(x_t) - f(\tilde{x}) + \frac{\mu}{2} \|x_t - \tilde{x}\|^2] \end{aligned}$$

$$\mathbb{E}[\mathbb{E}[\cdot | x_t]] = \mathbb{E}[\cdot] \quad (\text{also true with } \mu=0)$$

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

true even if $\mu=0$

$\left\{ \begin{array}{l} \gamma_t \text{ small enough, we have } \mathbb{E}\|x_t - \tilde{x}\|^2 \text{ decreasing} \\ \text{for any } \tilde{x} \in C \text{ s.t. } f(\tilde{x}) \leq \mathbb{E}f(x) \end{array} \right.$

\otimes non-strongly convex setting ($\mu=0$)

Let \tilde{x} be some min pt. \tilde{x} i.e. $f(x^*) = \min_{x \in C} f(x)$

$$\text{let } r_t \triangleq \mathbb{E}\|x_t - x^*\|^2$$

$$\text{Def } T_P \triangleq \mathbb{E}[f(x_t) - f(x^*)]$$

for better rate, let $\tilde{x}^* = \arg \min_{x \in X} \|x - x_0\|^2$

$$x^*$$

Let $\gamma_t = \mathbb{E}[\|x_t - \hat{x}_t\|^2]$

let $\mathbb{E}_t \triangleq \mathbb{E}[f(x_t) - f(\hat{x}_t)]$ expected suboptimality error

$$r_{t+1} \leq r_t - 2\gamma_t \varepsilon_t + \gamma_t^2 B^2$$

$$\Rightarrow 2\gamma_t [\mathbb{E}_t] \leq r_t - r_{t+1} + \gamma_t^2 B^2 \quad \forall t$$

* sum inequalities from $t=0$ to $t=T$

$$\Rightarrow 2 \sum_{t=0}^T \gamma_t \varepsilon_t \leq r_0 - r_{T+1} + \underbrace{\left(\sum_{t=0}^T (\gamma_t - \gamma_{t+1}) \right)}_{\text{telescoping sum}} \left(\sum_{t=0}^T \gamma_t^2 \right) B^2$$

$$2 \left(\sum_{t=0}^T \gamma_t \right) \min_{0 \leq t \leq T} \varepsilon_t \leq 2 \sum_t \gamma_t \varepsilon_t \leq \| \cdot \|$$

a)

$$\Rightarrow \boxed{\min_{0 \leq t \leq T} \varepsilon_t \leq r_0 + \frac{\left(\sum_{t=0}^T \gamma_t^2 \right) B^2}{2 \sum_{t=0}^T \gamma_t}}$$

note: $\min \varepsilon \rightarrow \infty$

when $\frac{\sum_{t=0}^T \gamma_t^2}{\sum_{t=0}^T \gamma_t} \rightarrow 0$

use $\gamma_t^* = \frac{r_0}{B\sqrt{T+1}}$ to minimize RHS

$$\Rightarrow \boxed{\min_{0 \leq t \leq T} \varepsilon_t \leq \frac{B r_0}{\sqrt{T+1}}}$$

b) for $\hat{x}_T = \sum_t p_t x_t$

by convexity

..... .. . / - .

Since f is convex, $f(\hat{x}_T) = f\left(\sum_t \gamma_t x_t\right) \leq \sum_t \gamma_t f(x_t)$

④ can also show with $\gamma_t = \frac{A}{\sqrt{T+1}}$, $\min_{C \in \mathcal{C}} \sum_{t=1}^T O\left(\frac{\log(T+1)}{\sqrt{T+1}}\right)$

and if set C is bounded, can show $O\left(\frac{\text{diam}(C)}{\sqrt{T+1}}\right)$ rate

strongly convex case ($\mu > 0$)

$$r_{t+1} \leq (1 - \mu \gamma_t) r_t - 2\gamma_t \varepsilon_t + \gamma_t^2 B^2$$

divide by this

$$\varepsilon_t \leq \frac{1}{2} (\gamma_t^2 - \mu) r_t - \frac{\gamma_t^{-1}}{2} r_{t+1} + \frac{\gamma_t B^2}{2}$$

use $\gamma_t = \frac{2}{\mu(t+2)}$

$$\gamma_t^{-1} = \frac{\mu(t+2)}{2}$$

Multiply ineq. by $(t+1)$

$$(t+1)\varepsilon_t \leq \frac{1}{2}(t+1) \left[\frac{t\mu + 2\mu - 2\mu}{2} \right] r_t - \frac{\mu(t+1)(t+2)}{4} r_{t+1} + \frac{(t+1)B^2}{2} \cancel{\frac{2}{\mu(t+2)}}$$

$$\leq \frac{\mu}{4} \left[\underbrace{(t+1) + r_t}_{\leq \bar{r}_t} - \underbrace{(t+1)(t+2)r_{t+1}}_{\bar{r}_{t+1}} \right] + \frac{B^2}{\mu} \leq \frac{B^2}{\mu}$$

$$\Rightarrow \sum_{t=0}^T (t+1)\varepsilon_t \leq \frac{\mu}{4} \left[\underbrace{u_0 - u_{T+1}}_0 \right] + (T+1) \frac{B^2}{\mu}$$

Let $\boxed{\rho_t \triangleq \frac{(t+1)}{S_T}}$

$$S_T \triangleq \sum_{t=0}^T (t+1) = \frac{(T+1)(T+2)}{2}$$

$$S_T \sum_{t=0}^T \rho_t \varepsilon_t \leq \frac{\mu}{4} [0 - (T+1)(T+2)r_{T+1}] + \frac{T+1}{\mu} B^2$$

$$\sum_{t=0}^T \sum_{t'=0}^{T+1} p_t c_t \leq \frac{\mu}{4} \left[(1 + \frac{1}{\mu})^2 - (1 + \frac{1}{\mu})^2 \right] \leq \frac{\mu}{4}$$

$$(\#) \quad \boxed{\sum_{t=0}^T p_t c_t + \frac{\mu}{4} \frac{(T+1)(T+2)}{S_T} R_{T+1} \leq \frac{(T+1)}{\mu} \frac{B^2}{S_T}}$$

let $\hat{x}_T \triangleq \sum_{t=0}^T p_t x_t$ (weighted avg.)

by convexity, $f(\hat{x}_T) = f\left(\sum_t p_t x_t\right) \leq \sum_t p_t f(x_t)$

$$\Rightarrow \mathbb{E}[f(\hat{x}_T) - f(x^*)] \leq \sum_t p_t \mathbb{E}[f(x_t) - f(x^*)] \stackrel{(\#)}{\leq} \frac{|T+1|}{\mu S_T} B^2$$

thus
$$\mathbb{E}[f(\hat{x}_T) - f(x^*)] \leq \frac{2B^2}{\mu(T+2)}$$

$$\frac{|T+1|}{S_T} = \frac{2}{T+2}$$

Vs. $O(\frac{1}{\sqrt{T}})$ rate when $\mu=0$

also
$$\mathbb{E}[\|\hat{x}_{T+1} - x^*\|^2] \leq \frac{4B^2}{\mu^2(T+2)}$$