

Lecture 10 - scribbles - convex opt

Friday, February 23, 2018

14:32

today: continue stochastic subgradient method (SGD)

[side note: $J(w) = \mathbb{E}_{(x,y) \sim p} J(x,y,w)$]

\downarrow
 $L(w)$ using calibration function

last time: we showed that $\min_{0 \leq t \leq T} \underbrace{(\mathbb{E} f(x_t) - f(x^*))}_{\triangleq \varepsilon_t} \leq \frac{B \sigma}{\sqrt{T+1}}$ when $\gamma_t = \frac{f_0}{B \sqrt{T+1}}$ (constant)

for stochastic subgradient method

⊛ can also show that with $\gamma_t = \frac{A}{\sqrt{t+1}}$, $\min_{0 \leq t \leq T} \varepsilon_t \leq O\left(\frac{\log(T+1)}{\sqrt{T+1}}\right)$

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

strongly convex case ($\mu > 0$)

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

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

multiply inequality by $\frac{(t+1)}{2}$
 (to get a telescoping sum)

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

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

$\frac{(t+1)}{2} \leq 1$

(4.2)

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

$$\Rightarrow \sum_{t=0}^T (t+1) \varepsilon_t \leq \frac{\mu}{4} \sum_{t=0}^T \underbrace{[(t+1)t r_t - (t+1)(t+2) r_{t+1}]}_{\substack{\geq U_t \\ U_0 - U_{T+1}}} + (T+1) \frac{B^2}{\mu}$$

$$\text{let } p_t \triangleq \frac{(t+1)}{S_T}$$

$$\sum_{t=0}^T S_T p_t \varepsilon_t \leq \frac{\mu}{4} \left[\overset{U_0 - U_{T+1}}{\downarrow} 0 - (T+1)(T+2) r_{T+1} \right] + \frac{(T+1)}{\mu} B^2$$

$$(\dagger) \quad \left[\sum_{t=0}^T p_t \varepsilon_t + \frac{\mu(T+1)(T+2)}{4 S_T} r_{T+1} \leq \frac{(T+1)}{\mu} \frac{B^2}{S_T} \right]$$

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

$$\text{by convexity, } f(\hat{x}_T) = f\left(\sum_t p_t x_t\right) \leq \sum_{t=0}^T p_t f(x_t)$$

$$\Rightarrow \mathbb{E} f(\hat{x}_T) - f(x^*) \leq \sum_{t=0}^T p_t \underbrace{[\mathbb{E} f(x_t) - f(x^*)]}_{\varepsilon_t} \stackrel{(\dagger)}{\leq} \frac{(T+1)}{\mu S_T} B^2$$

thus

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

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

$$\text{also: } \underbrace{\mathbb{E} \|x_{T+1} - x^*\|^2}_{r_{T+1}} \leq \frac{4 B^2}{\mu(T+2)} //$$

$$\text{note: } p_{t,T} = \frac{(t+1)}{S_T} = \frac{2(t+1)}{(T+1)(T+2)}$$

strategy:

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

$$\sum_t \delta_t \varepsilon_t \leq \sum_t \delta_t \left[\frac{r_t - r_{t+1}}{1 - \delta_t \mu} + \delta_t B^2 \right]$$

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

$$\Rightarrow \sum_{t=0}^T p_t = 1$$

$$\frac{(T+1)}{S_T} = \frac{2}{T+2}$$

Landscape of global convergence rates:

f is convex

$$f(x_t) - f(x^*) \leq \text{---}$$

$$r_0 \geq \text{dist}(x_0, x^*)$$

$$\mathbb{E} f(\hat{x}_t) - f(x^*) \leq \text{---} \quad (\text{stochastic setting})$$

assumptions	rate deterministic (batch)	stochastic setting	finite sum special case $\frac{1}{n} \sum_{i=1}^n f_i(x)$
1) non-smooth $\ \nabla f \ \leq B$	$O(\frac{Br_0}{\sqrt{t}})$ subgradient method	$O(\frac{Br_0}{\sqrt{t}})$	
2) smooth L -Lipschitz ∇f	$O(\frac{Lr_0^2}{t})$ gradient method	$O(\frac{\square}{\sqrt{t}})$ SGD	$O(\frac{\sqrt{\ln L}}{t})$ SAG/A
	$O(\frac{Lr_0^2}{t^2})$ Nesterov method (lower bound) "optimal method"		
f is μ -strongly convex			
3) non-smooth $\ \nabla f \ \leq B$	$O(\frac{B^2}{\mu t})$ subgradient method	$O(\frac{B^2}{\mu t})$	
smooth L -Lipschitz	$O(\exp(-\frac{\mu}{L}t))$ gradient method	$O(\frac{\square}{\mu t})$	$O(\exp(-\square t))$ SAG/A
	$O(\exp(-\sqrt{\frac{\mu}{L}}t))$ Nesterov method		

⊛ note: projecting gives the same rates
more generally, proximal gradient method as well

$$\min_x f(x) + h(x) \quad \text{"composite smooth optimization"}$$

smooth non-smooth

$$\text{constrained opt: } h(x) = S_C(x) \triangleq \begin{cases} +\infty & \text{if } x \notin C \\ 0 & \text{o.w.} \end{cases}$$

proximal gradient method: \uparrow prox step

related to step size...

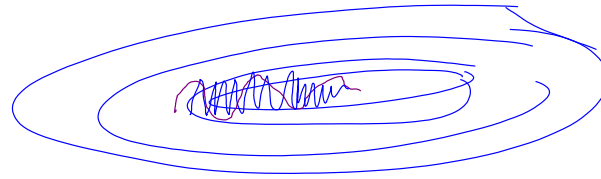
proximal gradient method:

related to step size...

$$x_{k+1} = \underset{x}{\operatorname{argmin}} f(x_k) + \langle \nabla f(x_k), x - x_k \rangle + \frac{L}{2} \|x - x_k\|^2 + h(x)$$

projected gradient method = proximal gradient method when $h(x) = \delta_C(x)$

example: $h(x) = \|x\|_1$



"linear coupling" paper