

Lecture 11 - landscape of rates

Tuesday, February 11, 2020 14:31

today: finish proof for SGD

- rate landscape
- CRF / structured SVM optimization

$$\mathbb{E} L_{t+1}(x_t) = \mathbb{E} L_t$$

(more like a mean)

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

✓ true even if $\mu=0$

; for γ_t small enough,
we have $\mathbb{E} \|x_t - \tilde{x}\|^2$ decreases for any $\tilde{x} \in C$ st.
 $f(\tilde{x}) \leq \mathbb{E} f(x_t)$
 $\frac{\partial}{\partial t} r(t)$ with t
we have $r(t+1) \leq r(t)$

② non-strongly convex setting ($\mu=0$)

set \tilde{x}^* to be some minimizer of f i.e. $f(x^*) = \min_{x \in C} f(x)$
in ineq. above

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

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

$$\text{let } \mathcal{E}_t \triangleq \mathbb{E} [f(x_t) - f(\tilde{x})], \quad \text{expected suboptimality error}$$

$$\Gamma_{t+1} \leq \Gamma_t - 2\gamma_t \mathcal{E}_t + \gamma_t^2 B^2 \quad \forall t$$

$$\Rightarrow 2\gamma_t [\mathcal{E}_t] \leq \Gamma_t - \Gamma_{t+1} + \gamma_t^2 B^2 \quad \forall t$$

$$\Rightarrow 2 \sum_{t=0}^T \gamma_t \mathcal{E}_t \leq \underbrace{\Gamma_0 - \Gamma_{T+1}}_{\substack{\text{telescoping sum} \\ \sum_{t=0}^T (\Gamma_t - \Gamma_{t+1})}} + \left(\sum_{t=0}^T \gamma_t^2 \right) B^2$$

$$2 \left(\sum_t \gamma_t \right) \min_t \mathcal{E}_t \leq 2 \sum_t \gamma_t \mathcal{E}_t \leq \Gamma_0 - \Gamma_{T+1}$$

a) \Rightarrow

$$\min_{0 \leq t \leq T} \mathcal{E}_t \leq \Gamma_0 + \frac{\left(\sum_{t=0}^T \gamma_t^2 \right) B^2}{2 \sum_{t=0}^T \gamma_t}$$

use $\gamma_t = \frac{\Gamma_0}{B\sqrt{T+1}}$

to minimize RHS

note: $\min_t \mathcal{E}_t \rightarrow 0$

$$\text{when } \begin{aligned} \sum_{t=0}^T \gamma_t^2 &\xrightarrow{T \rightarrow \infty} 0 \\ \sum_{t=0}^T \gamma_t &\xrightarrow{T \rightarrow \infty} 0 \end{aligned}$$

$$\min_{0 \leq t \leq T} \mathcal{E}_t \leq \frac{\Gamma_0}{\sqrt{T+1}}$$

b) for $\hat{x}_T = \frac{1}{T} \sum_t x_t$

Jensen

since f is convex, $f(\bar{x}_T) = f\left(\sum_t p_t x_t\right) \leq \sum_t p_t f(x_t)$

④ can also show that with $\gamma_t = \frac{1}{\sqrt{t+1}}$, $\min_{t \in T} \leq O\left(\frac{\log|T|}{\sqrt{T}}\right)$

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

strongly convex case ($\mu > 0$)

$$\begin{aligned} r_{t+1} &\leq (1 - \mu \gamma_t) r_t - \underbrace{\gamma_t \epsilon_t + \gamma_t^2 B^2}_{\text{divide}} \\ &\Rightarrow \epsilon_t \leq \frac{1}{2} (\gamma_t^{-1} - \mu) r_t - \frac{\gamma_t^{-1}}{2} r_{t+1} + \frac{\gamma_t^2 B^2}{2} \end{aligned}$$

$$\text{use } \left| \begin{array}{l} \gamma_t = \frac{2}{\mu(t+2)} \\ \gamma_t^{-1} = \frac{\mu(t+2)}{2} \end{array} \right.$$

Multiply ineq. by $(t+1)$

$$\begin{aligned} (t+1)\epsilon_t &\leq \frac{1}{2} (t+1) \left(\frac{\mu(t+2) - 2\mu}{2} \right) r_t - \frac{\mu(t+1)(t+2)}{4} r_{t+1} + \frac{(t+1)\gamma_t^2 B^2}{2} \\ &\leq \frac{\mu}{4} \left[\underbrace{t(t+1)r_t}_{\equiv u_t} - \underbrace{(t+1)(t+2)r_{t+1}}_{u_{t+1}} \right] + \frac{B^2}{\mu} \end{aligned}$$

(sum ineq.)

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

telescoping sum

$$\text{Let } \boxed{P_T \triangleq \frac{(T+1)}{S_T}}$$

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

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

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

$$\frac{\mu}{2} r_{T+1} \quad \frac{2}{T+2} \frac{B^2}{\mu}$$

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

$$f(\hat{x}_T) \stackrel{\text{convex}}{\leq} \sum_t p_t f(x_t)$$

$$\epsilon_t \quad \dots \quad \downarrow$$

$$f(\hat{x}_T) \stackrel{\text{convex}}{\leq} \sum_t p_t f(x_t)$$

$$\Rightarrow \mathbb{E}[f(\hat{x}_T) - f(x^*)] \leq \sum_t p_t \underbrace{\mathbb{E}[f(x_t) - f(x^*)]}_{\epsilon_t} \stackrel{(\dagger)}{\leq} \frac{2B^2}{T+2}$$

thus $\boxed{\mathbb{E} f(\hat{x}_T) - f(x^*) \leq \frac{2B^2}{T+2}}$ vs. $O(\frac{1}{\sqrt{T}})$ rate when $\mu=0$

also $\boxed{\mathbb{E} \|x_{T+1} - x^*\|^2 \leq \frac{4B^2}{T+2}}$

Landscape of global convergence rates

f is convex rate on suboptimality $f(x_t) - f(x^*) \leq \dots$

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

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

assumptions	rate (deterministic/batch)	stochastic setting	finite sum special case $\frac{1}{n} \sum_i^n f(x_i)$
1) non-smooth $\ \nabla f\ \leq B$	$O\left(\frac{Br_0}{\sqrt{t}}\right)$ subgradient method	$O\left(\frac{Br_0}{\sqrt{t}}\right)$	
2) smooth L -Lipschitz ∇f	$O\left(\frac{Lr_0^2}{t}\right)$ gradient method $O\left(\frac{Lr_0^2}{\sqrt{t}}\right)$ Nesterov method <small>(matching lower bound)</small>	$O\left(\frac{\square}{\sqrt{t}}\right)$ SGD <small>"optimal method"</small>	$O\left(\frac{\sqrt{n} L}{t}\right)$ SAG/SAGA SVRG <small>↳ Replicates version</small> <small>[Hoffman et al.]</small>
3) non-smooth $\ \nabla f\ \leq B$ μ -strongly convex	$O\left(\frac{B^2}{\mu t}\right)$ Subgradient meth.	$O\left(\frac{B^2}{\mu t}\right)$	
4) smooth L -Lipschitz	$O\left(\exp\left(-\frac{r_0^2}{L^2 t}\right)\right)$ gradient method $O\left(\exp\left(-\frac{r_0^2}{L^2 t}\right)\right)$ Nesterov "optimal"	$O\left(\frac{\square}{\mu t}\right)$	$O\left(\exp\left(-\min\left\{\frac{1}{L}, \frac{1}{L^2 t}\right\}\right)\right)$ SAG/SAGA/SVRG

15h35

⊕ note: projecting gives the same rates

more generally, proximal gradient method as well

smooth

non-smooth

Setup: "composite smooth opt." $\min_x f(x) + h(x)$

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

proximal gradient method:

$$x_{t+1} = \underset{x}{\operatorname{argmin}} \underbrace{f(x_t) + \langle \nabla f(x_t), x - x_t \rangle + \gamma \frac{1}{2} \|x - x_t\|^2}_{\frac{\gamma L}{2} \|x - (x_t - \frac{1}{\gamma L} \nabla f(x_t))\|^2 + \text{const.}} + h(x)$$

* if $h = S_C$ \Rightarrow projected gradient method

- but can also run on other "simple" h , e.g., $h(x) = \|\omega\|_2$ (Lasso type problem)

[accelerated prox gradient]
 \Rightarrow FISTA ; SOTA for deterministic L1-reg. problems
(small n)

optimization of $\hat{S}(w)$

$$\hat{S}(w) = R(w) + \frac{1}{n} \sum_{i=1}^n S(x^{(i)}, y^{(i)}; w) \quad \text{say } R(w) = \frac{\lambda}{2} \|w\|^2$$

$$\text{recall } h_w(x) = \underset{y \in \mathcal{Y}}{\operatorname{argmax}} \langle w, \varphi(x, \tilde{y}) \rangle$$

CRF: $\mathcal{L}_{\text{CRF}}(x, y; w) \stackrel{?}{=} \log \left(\sum_{\tilde{y}} \exp(\langle w, \varphi(x, \tilde{y}) \rangle) \right) - \langle w, \varphi(x^{(i)}, y^{(i)}) \rangle$ conditioned neg. log-likelihood loss
 here $\hat{S}_{\text{CRF}}(w)$ is L-smooth & λ -strongly convex $p(y|x) \propto \exp(S(y))$

weighted avg-SGD \rightarrow get a rate of $O(\frac{1}{\epsilon n t})$

What do we need?
 to run SGD

$$\begin{aligned} \text{compute } D_w \hat{S}(x, y; w) &= \frac{1}{n} \sum_{i=1}^n \frac{\sum_{\tilde{y}} \exp(s(\tilde{y}))}{\sum_{\tilde{y}} \exp(s(\tilde{y}))} [\varphi(x, \tilde{y})] - \varphi(x, y) \\ &\rightarrow p_{\theta}(\tilde{y}|x) \\ &= \mathbb{E}_{\tilde{y}|x; \theta} [\varphi(x, \tilde{y})] - \varphi(x, y) \end{aligned}$$

$$\text{CRF: } \varphi(x, \tilde{y}) = \sum_{c \in \mathcal{C}} \varphi_c(x, \tilde{y}_c)$$

, marginal over \tilde{y}_c .

$$\text{CRF: } \varphi(x, \tilde{y}) = \sum_{c \in C} \varphi_c(x, \tilde{y}_c)$$

$$\text{then } \mathbb{E}_{\tilde{y}|x} [\varphi(x, \tilde{y})] = \sum_{c \in C} \mathbb{E}_{\tilde{y}_c|x} [\varphi_c(x, \tilde{y}_c)]$$

use sum-product alg
on trees e.g.
or junction tree alg
for small tree with graph

structured SVM:

$$\text{SVMage}(x^{(i)}, y^{(i)}; w) = \max_{\tilde{y} \neq y} \langle w, \varphi(x^{(i)}, \tilde{y}) \rangle - \langle w, \varphi(x^{(i)}, y^{(i)}) \rangle$$

$$\text{let } l_i(\tilde{y}) \triangleq \ell(y^{(i)}, \tilde{y})$$

$$H_i(w) \triangleq \text{SVMage}(x^{(i)}, y^{(i)}; w) \rightarrow \max_{\tilde{y} \in \mathcal{Y}} l_i(\tilde{y}) - \underbrace{\langle w, \mathbf{1}_{\tilde{y}} \rangle}_{m(\tilde{y})}$$

$$N_i(\tilde{y}) \triangleq \ell(x^{(i)}, y^{(i)}) - \varphi(x^{(i)}, \tilde{y})$$

$$H_i(w; \tilde{y}) \triangleq l_i(\tilde{y}) - \langle w, N_i(\tilde{y}) \rangle \quad \text{note: if } \langle w, N_i(\tilde{y}) \rangle > 0 \quad \forall \tilde{y} \neq y^{(i)}$$

$$\text{then } h_w(x^{(i)}) = y^{(i)}$$

Structured SVM objective
(non-smooth unconstrained)
form

$$\min_w \frac{\|h(w)\|^2}{2} + \sum_{i=1}^n H_i(w)$$

★ this fits stochastic sub. method framework : $f(w) = \mathbb{E}_i h(w_i)$ where $h(w_i) \triangleq \frac{\|h(w)\|^2}{2} + H_i(w)$

now a subgradient of $h(w_i)$

$$h'(w_i) = \lambda w - N_i(\hat{y}_i|w)$$

$$\mathbb{E}_i h'(w_i) = \lambda w - \sum_{i=1}^n N_i(\hat{y}_i|w) = f'(w)$$

[batch subgradient]

$$\stackrel{\triangle}{=} \max_{\tilde{y} \in \mathcal{Y}} l_i(\tilde{y}) - \langle w, N_i(\tilde{y}) \rangle$$

loss-augmented influence

convergence rate:

have f is Δ -strongly convex

$$\Rightarrow \text{example: } \varphi(x, \tilde{y}) = \sum_{c \in C} \varphi_c(x, \tilde{y}_c)$$

$$\text{Suppose that } \|N_i(\tilde{y})\| \leq R \quad \forall i, \forall \tilde{y}$$

$$\|\varphi(x, \tilde{y})\|_2 \leq \sum_{c \in C} \|\varphi_c(x, \tilde{y}_c)\|_2$$

then one can show that with $\gamma_t = \frac{2}{\lambda(t+2)}$, then $\mathbb{E}\|g_t\|^2 \leq 4R^2 \leftarrow \text{gives our } B^2$
and $w_0 = 0$

[exercise: adapt App.A of Armin note Larochelle et al. 2012]

$$\rightarrow O\left(\frac{B^2}{\Delta t}\right) \text{ rate}$$