

Lecture 11 - landscape of rates

Tuesday, February 20, 2024 9:28 AM

- today:
- finish for SGD
 - rate landscape
 - CRF / structured SVM optimization

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

$\forall \tilde{x} \in C$

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

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

[for better rate, let $\tilde{x}(x_0) \triangleq x^* = \operatorname{argmin}_{x \in X^*} \|x - x_0\|^2$]

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

let $\epsilon_t \triangleq \mathbb{E} [f(x_t) - f(x^*)]$ expected suboptimality error

$$r_{t+1} \leq r_t - 2\delta_t \epsilon_t + \delta_t^2 B^2 \quad \forall t$$

$$\Rightarrow 2\delta_t [\epsilon_t] \leq r_t - r_{t+1} + \delta_t^2 B^2 \quad \forall t$$

$$\Rightarrow 2 \sum_{t=0}^T \delta_t \epsilon_t \leq \underbrace{r_0 - r_{T+1}}_{\text{telescoping sum: } \sum_{t=0}^T (r_t - r_{t+1})} + \underbrace{\left(\sum_{t=0}^T \delta_t^2 \right)}_{\downarrow} B^2$$

$$2 \left(\sum_{t=0}^T \delta_t \right) \min_{0 \leq t \leq T} \epsilon_t \leq r_0 + \left(\sum_{t=0}^T \delta_t^2 \right) B^2$$

$$\min_{0 \leq t \leq T} \epsilon_t \leq \frac{r_0 + \left(\sum_{t=0}^T \delta_t^2 \right) B^2}{\sum_{t=0}^T \delta_t}$$

note: $\sum_{t=0}^T \delta_t \xrightarrow{T \rightarrow \infty} \infty$
 if $\sum_{t=0}^T \delta_t^2 \xrightarrow{T \rightarrow \infty} 0$
 $\sum_{t=0}^T \delta_t$

then $\min_{t \leq T} \epsilon_t \rightarrow 0$

a) use $\delta_t = \frac{r_0}{B\sqrt{t+1}}$ to minimize RHS

$$\|x_0 - x^*\|^2$$

$$\Rightarrow \min_{t \leq T} \epsilon_t \leq \frac{B r_0}{\sqrt{T+1}}$$

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

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

$$\mathbb{E}f(\hat{x}_T) - f^* \leq \sum_t p_t \underbrace{(\mathbb{E}f(x_t) - f^*)}_{\varepsilon_t}$$

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

and if set C is bounded

can show $O\left(\frac{\text{diam}(C)}{\sqrt{T+1}}\right)$ rate

Strongly convex case ($\mu > 0$)

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

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

\downarrow
 $\frac{\mu(t+2)}{2}$

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

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

multiply ineq. by $(t+1)$

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

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

$\leq \frac{B^2}{\mu}$

telescoping sum? (trick)

(sum ineq.)

$$\Rightarrow \sum_{t=0}^T \frac{(t+1)}{4} \varepsilon_t \leq \frac{\mu}{4} [U_0 - U_{T+1}] + (T+1) \frac{B^2}{\mu}$$

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

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

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

$$\sum_{t=0}^T p_t E_t + \underbrace{\frac{\mu (T+2)}{4} \Gamma_{T+1}}_{\frac{\mu}{2} \Gamma_{T+1}} \leq \underbrace{\frac{(T+1) B^2}{S_T \mu}}_{\frac{2}{T+2} \frac{B^2}{\mu}} \quad (\neq)$$

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

$$\mathbb{E} f(\hat{x}_T) - f^* \leq \sum_{t=0}^T p_t E_t \leq \frac{2B^2}{(T+2)\mu} \quad (\neq)$$

by convexity

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

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

also

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

$$\mathbb{E} E_{T+1} \leq \frac{L}{2} \mathbb{E} \|x_{T+1} - x^*\|^2 \rightarrow \frac{2(L)}{\mu^2} \frac{B^2}{(T+2)}$$

10h24

Landscape of global convergence rates

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

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

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

assumptions	rate (deterministic / self kahn)	stochastic setting	finite sum special case $\frac{1}{n} \sum_{i=1}^n f_i(x)$
1) non-smooth $\ g\ \leq B$	$O\left(\frac{B r_0}{\sqrt{t}}\right)$ subgradient method	$O\left(\frac{B r_0}{\sqrt{t}}\right)$	
2) smooth L -Lipschitz ∇f	$O\left(\frac{L r_0^2}{t}\right)$ gradient method $O\left(\frac{L r_0^2}{t^2}\right)$ Nesterov method "optimal method" "matching lower bound"	$O\left(\frac{L}{\sqrt{t}}\right)$ SGD	$O\left(\frac{\sqrt{\ln L}}{t}\right)$ SAG/SAGA/SVRG "loopless" [Hoffman's al.]
f is μ -strongly convex	3) non-smooth $\ g\ \leq B$ $O\left(\frac{B^2}{\mu t}\right)$ subgradient method	$O\left(\frac{B^2}{\mu t}\right)$	
4) smooth L -Lipschitz	$O\left(\exp\left(-\frac{\mu}{L} t\right)\right)$ grad. $O\left(\exp\left(-\frac{\mu}{L} t\right)\right)$ Nesterov	$O\left(\frac{L}{\mu t}\right)$	$O\left(\exp\left(-\min\left\{\frac{\mu}{L}, \frac{\mu}{n}\right\} t\right)\right)$ SA/SAGA/SVRG

$\left. \begin{array}{l} \text{L-Lipschitz} \\ \cup (\exp(-\frac{\mu t}{L})) \text{ grad.} \\ \cup (\exp(-\sqrt{\frac{\mu}{L}} t)) \text{ Nesterov} \\ \text{"optimal"} \end{array} \right\}$

$\cup (\frac{1}{\mu t})$

$\cup (\exp(-\frac{1}{n} \sum_{i=1}^n \mu_i t))$
 SA(B)/SABA/SVRG

interpolation regime $\mathbb{E}_{\xi} \|\nabla_x h(x^*, \xi)\|^2 = 0$

↳ overparameterized regime

⇒ get faster rates for SGD

⊛ rate: projecting gives the same rates

more generally, proximal gradient method as well

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

smooth ↓ non-smooth ↓

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

proximal gradient method:

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

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

* but can also run on other "simple" h e.g. $h(x) = \|x\|_1$ (Lasso type problem)
 ↳ prox step becomes "soft-thresholding" operator

→ get same rate convergence as no h (unconstrained or no $\| \cdot \|_2$ etc.)

[accelerated prox. gradient for \mathcal{D}_1 = FISTA ; sorta for deterministic \mathcal{D}_1 -reg. problems (small n)]

optimization of $\hat{\mathcal{J}}(w)$

$\hat{\mathcal{J}}(w) = R(w) + \frac{1}{n} \sum_{i=1}^n \mathcal{L}(x^{(i)}, y^{(i)}; w)$ say $R(w) = \frac{\lambda}{2} \|w\|_2^2$

recall $h(w) = \arg \max \langle w, \varphi(x, y) \rangle$

recall $h_w(x) = \arg \max_{\tilde{y} \in \mathcal{Y}} \langle w, \phi(x, \tilde{y}) \rangle$

CRF: $\mathcal{L}_{\text{CRF}}(x, y; w) \triangleq \log \left(\sum_{\tilde{y}} \exp(\langle w, \phi(x, \tilde{y}) \rangle) \right) - \langle w, \phi(x, y) \rangle$

here $\mathcal{L}_{\text{CRF}}(w)$ is L -smooth & λ -strongly convex
 weighted avg. SGD \rightarrow get a rate of $O\left(\frac{1}{\lambda t}\right)$

what do we need to run SGD?

$$\nabla_w \mathcal{L}(x, y; w) = \frac{1}{\sum_{\tilde{y}} \exp(\langle w, \phi(x, \tilde{y}) \rangle)} \sum_{\tilde{y}} \exp(\langle w, \phi(x, \tilde{y}) \rangle) (\phi(x, \tilde{y}) - \phi(x, y))$$

$\rightarrow p_w(\tilde{y}|x)$

$$= \mathbb{E}_{\tilde{y}|x; w} [\phi(x, \tilde{y})] - \phi(x, y)$$

CRF: $\phi(x, \tilde{y}) = \sum_{c \in \mathcal{C}} \phi_c(x, \tilde{y}_c)$

then $\mathbb{E}_{\tilde{y}|x} [\phi(x, \tilde{y})] = \sum_{c \in \mathcal{C}} \mathbb{E}_{\tilde{y}_c|x} [\phi_c(x, \tilde{y}_c)]$

\uparrow marginal over \tilde{y}_c

"marginalization oracle"
 use sum-product of ϕ_c on trees e.g.
 or junction tree alg. for small treewidth graphs

Structured SVM

$$\mathcal{L}_{\text{hinge}}(x, y; w) = \max_{\tilde{y} \in \mathcal{Y}} \langle w, \phi(x, \tilde{y}) \rangle + \ell(y, \tilde{y}) - \langle w, \phi(x, y) \rangle$$

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

$H_i(w) \triangleq \mathcal{L}_{\text{hinge}}(x^{(i)}, y^{(i)}; w)$

$\psi_i(\tilde{y}) \triangleq \phi(x^{(i)}, y^{(i)}) - \phi(x^{(i)}, \tilde{y})$

$H_i^0(w; \tilde{y}) \triangleq l_i(\tilde{y}) - \langle w, \psi_i(\tilde{y}) \rangle$

"hinge loss elements"

$H_i(w) = \max_{\tilde{y}} H_i^0(w; \tilde{y})$

$\max_{\tilde{y} \in \mathcal{Y}} l_i(\tilde{y}) - \langle w, \psi_i(\tilde{y}) \rangle$

$m_i(\tilde{y})$ "margin element"

note: if $\langle w, \psi_i(\tilde{y}) \rangle \gg 0$ $\forall \tilde{y} \neq y^{(i)}$

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

structured SVM objective

(non-smooth unconstrained form)

$$\min_w \frac{\lambda \|w\|_2^2}{2} + \frac{1}{n} \sum_{i=1}^n H_i(w)$$

use this obj. to solve with Frank-Wolfe: $(w^{(t)}) = \mathbb{E}_i h(w^{(t-1)})$

(unconstrained form)

⊗ this fits stochastic subg. method framework $= f(w) = \mathbb{E}_i h(w, i)$

$$\text{where } h(w, i) \triangleq \frac{\lambda \|w\|^2 + M_i(w)}{2}$$

now a subgradient of $h(w, i)$

$$h'(w, i) = \lambda w - \Psi_i(\hat{y}_i(w))$$

$$\hat{y}_i(w) \triangleq \underset{\tilde{y} \in \mathcal{Y}}{\text{argmax}} \ell_i(\tilde{y}) - \langle w, \Psi_i(\tilde{y}) \rangle$$

loss-augmented inference

$$\mathbb{E}_i h'(w, i) = \lambda w - \frac{1}{n} \sum_{i=1}^n \Psi_i(\hat{y}_i(w)) = f'(w)$$

[batch subgradient]

convergence rate

here f is λ -strongly convex

suppose that $\|\Psi_i(\tilde{y})\| \leq R \quad \forall i, \tilde{y} \in \mathcal{Y}$

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

$$\leadsto O\left(\frac{R^2}{\lambda t}\right)$$

[exercise: adapt App. A of arxiv note Lacoste-Julien + al. 2012]