

Lecture 11 - landscape of rates

Thursday, February 18, 2021 13:36

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

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

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

* 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 brute, let $\tilde{x}(x_0) = x^* = \operatorname{argmin}_{x \in X^*} \|x - x_0\|^2$]

let $r_t = \mathbb{E}[\|x_t - x^*\|^2]$

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

$$r_{t+1} \leq r_t - 2\alpha\mathcal{E}_t + \alpha^2 B^2 \quad \forall t$$

$$\Rightarrow 2\alpha\mathcal{E}_t \leq r_t - r_{t+1} + \alpha^2 B^2 \quad \forall t$$

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

$$2\left(\sum_{t=0}^T \alpha\right) \min_{0 \leq t \leq T} \mathcal{E}_t \leq r_0 + \left(\sum_{t=0}^T \alpha\right) B^2$$

a) \Rightarrow
$$\min_{0 \leq t \leq T} \mathcal{E}_t \leq \frac{r_0 + \left(\sum_{t=0}^T \alpha\right) B^2}{2\sum_{t=0}^T \alpha}$$

Note: $\min_{0 \leq t \leq T} \mathcal{E}_t \rightarrow 0$
 when $\frac{\sum_{t=0}^T \alpha^2}{\sum_{t=0}^T \alpha} \xrightarrow{T \rightarrow \infty} 0$

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

$$\Rightarrow \min_{t \leq T} \mathcal{E}_t \leq \frac{B r_0}{\sqrt{T+1}}$$

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

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

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

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

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

strongly convex case ($\mu > 0$)

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

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

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

multiply ineq. by $(t+1)$

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

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

telescoping sum? (trick)

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

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

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

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

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

ϵ_t (†) \dots

$\hat{x}_T = \sum_{t=1}^T p_t x_t$ (weighted avg.)
 $f(\hat{x}_T) \leq \sum p_t f(x_t)$

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

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

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

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

14h23

Landscape of global convergence rates

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

for stochastic setting $\mathbb{E}[f(x_t) - f(x^*)] \leq \dots$

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

assumptions	rate (deterministic/batch)	stochastic setting	finite sum special case $\frac{1}{n} \sum_{i=1}^n f(x_i)$
1) non-smooth $\ Df\ \leq B$	$O(\frac{Br_0}{\sqrt{t}})$ subgradient method	$O(\frac{Br_0}{\sqrt{t}})$	
2) smooth L -Lipschitz Df	$O(\frac{Lr_0^2}{t})$ gradient method $O(\frac{Lr_0^2}{t^2})$ Nesterov method matching lower bound "optimal method"	$O(\frac{Lr_0^2}{\sqrt{t}})$ SGD SAG/SAGA/SVRG \hookrightarrow "kaczmarz" version [Hofmann '14]	$O(\frac{\sqrt{L} r_0}{t})$
f is μ -strongly convex 3) non-smooth $\ Df\ \leq B$ 4) smooth L -Lipschitz	$O(\frac{B^2}{\mu t})$ subg. method $O(\exp(-\frac{\mu}{L} t))$ grad. $O(\exp(-\sqrt{\frac{\mu}{L}} t))$ Nesterov "optimal"	$O(\frac{B^2}{\mu t})$ $O(\frac{L}{\mu t})$	$O(\exp(-\min\{\frac{\mu}{L}, \frac{\mu}{2}\} t))$ SAG/SAGA/SVRG

interpolation regime $\mathbb{E}\|P_x h(x; \xi)\|^2 = 0$

\hookrightarrow overparameterized regime \Rightarrow get faster rates for SGD

* note: projecting gives the same rates

more generally, proximal gradient method is well

smooth non-smooth
 ↓ ↓

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

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:

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

* 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)

\rightarrow prox step becomes "soft-thresholding" operator

\rightarrow get same rate convergence as unconstrained

[accelerated prox gradient for l_1 = FISTA ; SOTA for deterministic Q -reg. problems (small n)]

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

$$\hat{J}(w) = R(w) + \frac{1}{n} \sum_{i=1}^n \ell(x^{(i)}, y^{(i)}; w)$$

say $R(w) = \frac{\lambda}{2} \|w\|_2^2$

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

CRF:

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

← neg. cond. log likelihood loss

here $\hat{J}_{\text{CRF}}(w)$ is L -smooth & λ -strongly convex

$$p_w(\tilde{y}|x) \propto \exp(s(\tilde{y}))$$

weighted avg. SGD \rightarrow get a rate of $O(\frac{\square}{\#t})$

what do we need to run SGD?

$$\nabla_w \mathcal{L}(x, y; w) = \frac{1}{\sum_{\tilde{y}} \exp(s(\tilde{y}))} \left[\frac{\partial}{\partial w} \exp(s(\tilde{y})) \right] [\phi(x, \tilde{y})] - \phi(x, y)$$

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

to run SGD:

$$\frac{\sum_y \exp(s(y))}{\sum_y \exp(s(y))} \rightarrow p(y|x)$$

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

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

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

maximal over \tilde{y}_c

Use sum-product alg on trees CG or junction tree alg for small tree width graph

Structured SVM

$$J_{\text{hinge}}(x^{(i)}, y^{(i)}; w) = \max_{\tilde{y} \in \mathcal{Y}} \langle w, \psi(\tilde{y}) \rangle + l(y^{(i)}; \tilde{y}) - \langle w, \psi(x^{(i)}; y^{(i)}) \rangle$$

let $l_i(\tilde{y}) \triangleq l(y^{(i)}; \tilde{y})$

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

$\psi_i(\tilde{y}) \triangleq \psi(x^{(i)}; y^{(i)}) - \psi(x^{(i)}; \tilde{y})$

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

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

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

"margin element"

note: if $\langle w, \psi_i(\tilde{y}^*) \rangle > 0$
 $\forall \tilde{y}^* \neq y^{(i)}$
 then $h_w(x^{(i)}) = y^{(i)}$

Structured SVM objective
 (non-smooth
 unconstrained form)

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

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

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

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

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

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

[batch subgradient]

argmax $l_i(\tilde{y}) - \langle w, \psi_i(\tilde{y}) \rangle$
 $\tilde{y} \in \mathcal{Y}$
 loss-augmented inference

convergence rate:

convergence rate:

here f is λ -strongly convex

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

then one can show that with $\gamma_t = \frac{2}{\lambda(t+2)}$
and $w_0 = 0$

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

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

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

, then $\|g_t\| \leq 4R^2 \leftarrow$ gives $\text{dim } \mathcal{B}^2$

[exercice: adapt App. A
of arxiv note [Lacoste-Julien](#)
Feb. 2012]