

Lecture 11 - landscape of rates

Tuesday, February 12, 2019 14:29

today:

- rate landscape
- CRF / structured SVM optimization

Landscape of global convergence rates

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

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

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

assumptions	rate deterministic (batch)	stochastic setting	finite sum special case $\sum_i f_i(x)$
1) non-smooth bounded $\ \nabla f\ \leq B$ subgradient	$O\left(\frac{B r_0}{\sqrt{t}}\right)$ 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 $O\left(\frac{L r_0^2}{t^2}\right)$ (lower bound) Nesterov method "optimal method"	$O\left(\frac{\square}{\sqrt{t}}\right)$ SGD	$O\left(\frac{\sqrt{n} L}{t}\right)$ SAG/SAGA SVRG \rightarrow Nesterov's style
f is μ -strongly convex { 3) non-smooth $\ \nabla f\ \leq B$ { 4) smooth L -Lipschitz	$O\left(\frac{B^2}{\mu t}\right)$ subgradient method $O\left(\exp(-\frac{\mu}{2}t)\right)$ gradient method $O\left(\exp(-\frac{\mu}{2}t)\right)$ Nesterov	$O\left(\frac{B^2}{\mu t}\right)$ $O\left(\frac{\square}{\mu t}\right)$	$O\left(\exp(-\min\left\{\frac{1}{n}, \frac{\mu}{2}\right\} t)\right)$ SAG/SAGA

(L-Lipschitz)

- "PCT" method
 $O(\exp(-\frac{1}{2}t))$ Nesterov
(lower bound)

ut)

- "I.C. on the side"

⊗ note: projecting gives the same rates

more generally, proximal gradient method
as well

setup: 'composite smooth optimization'

$$\min_x \underbrace{f(x)}_{\text{smooth}} + \underbrace{h(x)}_{\text{non-smooth}}$$

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

proximal gradient method: prox step

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

if $h = S_C$

$$\frac{1}{2} \|x - (x_t - \frac{1}{2} \nabla f(x_t))\|^2 + \text{cst.}$$

⇒ projected gradient method

= prox grad. method when $h(x) = S_C(x)$
• but can also run on other h , e.g. $h(x) = \|x\|_1$

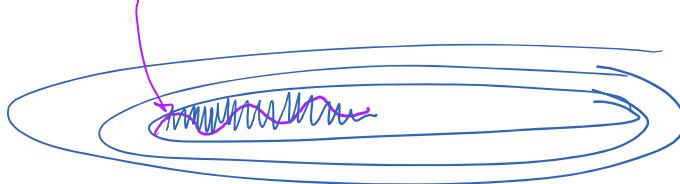
(Lasso type problem)

[accelerated prox gradient

= FISTA ; state-of-the-art

for deterministic Q_1 -reg. problems]

Aside: Nesterov method behavior



15h16

Optimization of $\hat{f}(w)$

$$\hat{f}(w) = R(w) + \frac{1}{n} \sum_{i=1}^n \ell(x^{(i)}, y^{(i)}; w) \quad \text{say } R(w) = \frac{\lambda \|w\|_2^2}{2}$$

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

CRF: $\mathcal{L}_{\text{CRF}}(x, y; w) \triangleq \log \left(\sum_{\tilde{y}} \exp(\langle w, \ell(x, \tilde{y}) \rangle) \right) - \langle w, \ell(x^{(i)}, y^{(i)}) \rangle$ & negative conditional log-likelihood loss

here $\hat{f}(w)$ is L-smooth & λ -strongly convex

weighted average SGD \rightarrow get a rate of $O(\frac{1}{\sqrt{t}})$

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

what do we need? compute $\nabla_w \mathcal{L}_{\text{CRF}}(x, y; w) = \frac{1}{\sum_{\tilde{y}} \exp(s(\tilde{y}))} \sum_{\tilde{y}} \exp(s(\tilde{y})) (\ell(x, \tilde{y}) - \ell(x, y))$

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

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

maximal over \tilde{y}_c

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

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

marginal over \tilde{y}_c

use sum-product alg
on trees e.g.
or function tree alg.
on small treewidth graphs

Structured SUM:

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

$$\text{let } \ell_i(\tilde{y}) \triangleq \ell(y^{(i)}; \tilde{y})$$

$$H_i(w) \triangleq S_{\text{hinge}}(x^{(i)}, y^{(i)}; w) = \max_{\tilde{y} \in \mathcal{Y}} \ell_i(\tilde{y}) - \langle w, \psi_i(\tilde{y}) \rangle$$

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

$$\begin{aligned} &\hookrightarrow \langle w, \psi_i(\tilde{y}) \rangle = m(\tilde{y}) \quad (\text{"margin"}) \\ H_i(w; \tilde{y}) &\triangleq \ell_i(\tilde{y}) - \langle w, \psi_i(\tilde{y}) \rangle \end{aligned}$$

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

Structured SUM
objective
(non-smooth unconstrained form)

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

⊗ This fits stochastic subgradient method framework: $f(w) = \mathbb{E}_i h(w, i)$ where $h(w, i) \triangleq \frac{1}{2} \|w\|^2 + H_i(w)$

now a subgradient of $h(w, i)$: $h'(w, i) = \Delta w - \Psi_i(\hat{y}_i(w))$

$$\hat{y}_i(w) \triangleq \arg \max_{\tilde{y}} \ell_i(\tilde{y}) - \langle w, \psi_i(\tilde{y}) \rangle$$

$$\hat{y} \triangleq \operatorname{argmax}_{y \in \mathcal{Y}} \varphi_i(\tilde{y}) - \langle w, \varphi(y) \rangle$$

loss-augmented inference

[sidenote : soon, we will see that

$$w^* = \frac{1}{\lambda n} \sum_{i=1}^n \sum_{y \in \mathcal{Y}_i} \alpha_i^*(\tilde{y}) \varphi_i(\tilde{y})$$

← true both for CRF & structured SVM

for SVM

$$\left\{ \begin{array}{l} \alpha_i^*(\tilde{y}) = 0 \text{ if } H_i(w; \tilde{y}) < \max_{y \in \mathcal{Y}_i} H_i(w; y) = H_i(w) \\ \alpha_i^*(\tilde{y}) > 0 \text{ "support vectors"} \end{array} \right.$$

Convergence rate :

here f is λ -strongly convex

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

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

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

then one can show that with $\gamma_t = \frac{2}{\lambda(t+2)}$, then $\|\varphi_t\|^2 \leq 4R^2$ or this gives B^2
and $w_0 = 0$

[exercise: adapt appendix A of Axiv note Lacoste-Julien et al. 2012]

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

other approaches to optimize SUM struct

(UP)
unconstrained primal

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

[unconstrained]
non-smooth

(PQP)
primal QP
→ quadratic program

$$\min_{w \in \mathbb{R}^n} \frac{1}{2} \|w\|^2 + \sum_{i=1}^n \xi_i$$

[constrained formulation]
 $\xi_i \geq H_i(w; \tilde{y}) \quad \forall \tilde{y} \in \mathcal{Y}_i \nsubseteq \mathcal{X}_i$ (smooth) convex QP
with exp. # of linear constraints

1) generic approach to use convexity of loss-augmented decoding: [Taskar et al. ICML 2005]

Idea: here, we suppose that loss-augmented decoding can be expressed as a "compact" maximization problem of a concave fct,

$$\text{i.e. } H_i(w) = \max_{\tilde{y} \in \mathcal{Y}_i} l_i(\tilde{y}) - \langle w, \eta_i(\tilde{y}) \rangle = \max_{z \in Z} g_i(w; z)$$

cts. \nearrow convex
discrete

where g_i is concave in z and convex in w

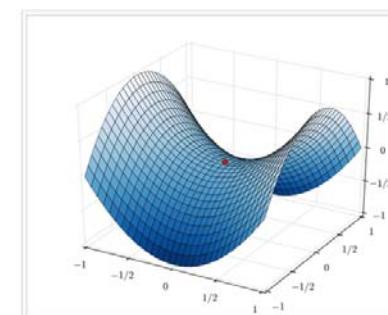
- Z: • should not depend on w
- should have a tractable size

a) saddle point formulation:

$$\min_w \max_{z_i \in Z_i} \frac{1}{2} \|w\|^2 + \sum_{i=1}^n g_i(w; z_i)$$

$$\min_w \max_z \mathcal{L}(w, z)$$

convex-concave
saddle point problem

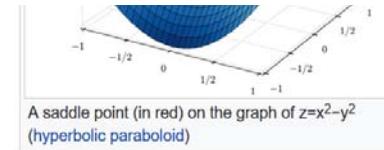


$$\min_w \max_z S(w, z)$$

Convex-concave
saddle point problem

(under reg. conditions $\min_w \max_z = \max_z \min_w \rightarrow$ "saddle point")

$$\forall z \quad S(w^*, z) \leq S(w^*, z^*) \leq S(w, z^*) \quad \forall w$$



in general

$$\min_w \max_z \geq \max_z \min_w$$

circular dependence
dependence
it might not exist?

Standard algorithms

Extragradient algorithm:

"lookahead step"

$$\begin{pmatrix} \tilde{w}_{t+1} \\ \tilde{z}_{t+1} \end{pmatrix} = \begin{pmatrix} w_t \\ z_t \end{pmatrix} + \delta t \begin{pmatrix} -\nabla_w S(w_t, z_t) \\ \nabla_z S(w_t, z_t) \end{pmatrix}$$

$$\begin{pmatrix} w_{t+1} \\ z_{t+1} \end{pmatrix} = \begin{pmatrix} w_t \\ z_t \end{pmatrix} + \delta t \begin{pmatrix} -\nabla_w S(\tilde{w}_{t+1}, \tilde{z}_{t+1}) \\ \nabla_z S(\tilde{w}_{t+1}, \tilde{z}_{t+1}) \end{pmatrix}$$

Converges $O(\frac{1}{t})$ for convex-concave game

applied to structured SUM
[Taskar et al., JMLR 2006]

