

Lecture 4 - scribbles

Wednesday, January 18, 2017
10:27

today: optimization (of $J(w)$)

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

$$\text{structured SVM: } R(w) = \frac{\lambda \|w\|^2}{2} \quad \mathcal{S}(x, y, w) = \langle w, \phi(x, y) \rangle$$

$$\mathcal{J}(x^{(i)}, y^{(i)}, w) = \max_{\tilde{y} \in \mathcal{Y}(x^{(i)})} [\langle w, \phi(x^{(i)}, \tilde{y}) \rangle + \ell(y^{(i)}, \tilde{y})] - \langle w, \phi(x^{(i)}, y^{(i)}) \rangle$$

(structured hinge loss)

$$\text{let } \ell_i(\tilde{y}) \triangleq \ell(y^{(i)}, \tilde{y}) \quad \mathcal{Y}_i \triangleq \mathcal{Y}(x^{(i)})$$

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

$$H_i(w) \triangleq \mathcal{J}(x^{(i)}, y^{(i)}, w) = \max_{\tilde{y} \in \mathcal{Y}_i} [\ell_i(\tilde{y}) - \langle w, \psi_i(\tilde{y}) \rangle]$$

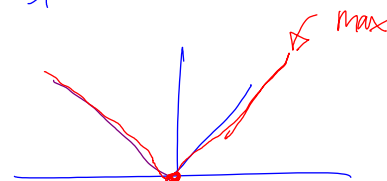
structured SVM objective (non-smooth unconstrained form)

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

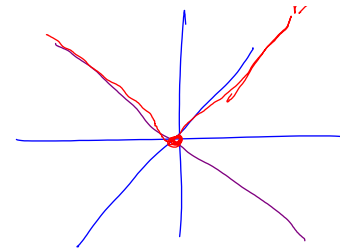
regularization hyperparameter (somewhat akin to $(1 - y_i w^T x_i)_+$)

$\max_{\tilde{y} \in \mathcal{Y}_i} [\ell_i(\tilde{y}) - \langle w, \psi_i(\tilde{y}) \rangle]$

$$|x| = \max \begin{cases} x \\ -x \end{cases}$$



$$|x| = \max \begin{cases} x \\ -x \end{cases}$$



convex analysis recap

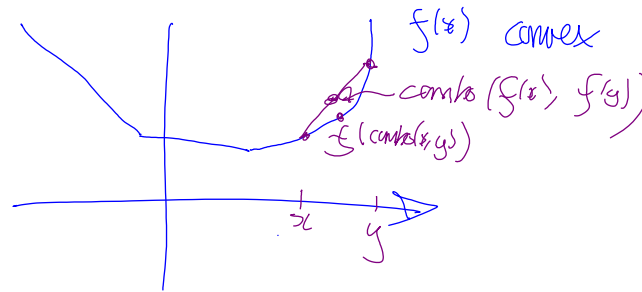
f is convex

$$\Leftrightarrow f(\underbrace{px + (1-p)y}_{\text{conv-combo}(x,y)}) \leq \underbrace{pf(x) + (1-p)f(y)}_{\text{conv-combo}(f(x), f(y))}$$

$p \in [0,1]$

subgradient: $v \in \partial f(x)$, " v is a subgradient of f at x "

$$\Leftrightarrow \forall y \in \text{dom } f \quad f(y) \geq f(x) + \langle v, y - x \rangle$$



standard assumptions on f :

smooth f :

∇f is L -Lipschitz continuous i.e. $\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\|$ Hessian

[if f is twice differentiable; you have $L = \max_{x \in \text{dom } f} \|H(x)\|$]

non-smooth f :

∂f is bounded i.e. $\|v\| \leq B$ for all subgradients v of f

strong convexity of f :

$$f \text{ is } \mu\text{-strongly convex} \Leftrightarrow f(x) \geq f(y) + \langle \nabla f(y), x - y \rangle + \frac{\mu}{2} \|x - y\|^2 \quad \forall x, y \in \text{dom } f$$

$$f \text{ is } \mu\text{-strongly convex} \Leftrightarrow f(y) \geq f(x) + \underbrace{\langle \nabla f(x), y-x \rangle}_{\leq V, y-x \text{ for any } v \in \partial f(x)} + \frac{\mu}{2} \|x-y\|^2 \quad \forall y \in \text{dom } f$$

strong convexity constant

$$[f \text{ is } C^2, \quad \mu = \lambda_{\min}(\text{Hessian}(f))]$$

minimum e-value

properties:

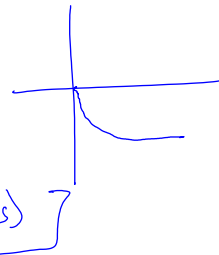
$$\nabla f \text{ L-Lipschitz} \Rightarrow f(y) \leq f(x) + \langle \nabla f(x), y-x \rangle + \frac{L}{2} \|y-x\|^2$$

$$f \text{ } \mu\text{-strongly convex} \Rightarrow f(y) \geq f(x) + \langle \nabla f(x), y-x \rangle + \frac{\mu}{2} \|y-x\|^2$$

[question: if $f(x)$ is bounded and convex,
 $\Rightarrow f$ is L-lipschitz]

answer: no; counterexample $f(x) = -\sqrt{x}$ on $[0,1]$

(f is not Lipschitz continuous)



Landscape of convergence rates: \rightarrow for first order methods

example:

convergence result: $\underbrace{f(x_k) - f(x^*)}_{\text{function suboptimality}} \leq \frac{\text{const.}}{\sqrt{k}} \quad (\text{non-smooth rate})$

\nearrow
set of solutions

suppose $\text{dist}(x_0, X^*) \leq \underline{r_0}$

assumptions	rate deterministic	rate stochastic $\rightarrow \frac{1}{n} \sum_{i=1}^n f_i(x)$
1) non-smooth $\ \nabla f\ \leq B$	$O\left(\frac{Br_0}{\sqrt{k}}\right)$ subgradient method	$O\left(\frac{1}{\sqrt{k}}\right)$ stochastic subgradient method

f is μ strongly convex	2) smooth L -Lipschitz	$O\left(\frac{Lr^2}{k}\right)$ gradient method $O\left(\frac{Lr^2}{k^2}\right)$ Nesterov method	$O\left(\frac{1}{\sqrt{k}}\right)$ SGD	
	3) non-smooth $\ x\ \leq B$	$O\left(\frac{B^2}{\mu k}\right)$ subgradient method	$O\left(\frac{B}{\mu k}\right)$	$\mathbb{E}_{\xi} f(x, \xi)$
	4) smooth L -Lipschitz	$O(\exp(-\frac{\mu}{L} k))$ gradient method $O(\exp(-\sqrt{\frac{\mu}{L}} k))$ Nesterov method	true expectation $O\left(\frac{1}{k}\right)$ SGD but finite sum faster rate $\exp(-\text{something} \cdot k)$ using SAG, etc. (stochastic averaged gradient method)	$\sum_{i=1}^n S_i(x)$ $\rightarrow \frac{1}{n} \sum_{i=1}^n S_i(x)$

⊗ note: projecting gives same rate \rightarrow reason is contraction property of projection
 e.g. gradient method (unconstrained optimization) $x_{t+1} = x_t - \frac{1}{L} \nabla f(x_t)$
 i.e. $\|P_C(y) - P_C(x)\|_2 \leq \|y - x\|_2$

projected gradient method for $\min_{x \in C} f(x)$: $x_{t+1} = P_C \left[x_t - \frac{1}{L} \nabla f(x_t) \right]$

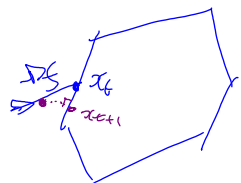
Euclidean projection on C i.e.

$$P_C(y) \triangleq \arg \min_{x \in C} \|x - y\|_2$$

minimizes

$$\|P_C(y) - P_C(x^*)\| = \|P_C(y) - x^*\| \leq \|y - x^*\|$$

\Rightarrow so rates transfer b "projected" versions



complexity in optimization

→ "work" to get to ϵ suboptimality $O(\frac{1}{\sqrt{\epsilon}})$ rate $\leadsto O(\frac{1}{\epsilon^2})$ # of iterations

if want to compare methods, include also "cost" of iterations

stochastic subgradient method

say want to solve $\min_{x \in C} f(x)$ where $f(x) \triangleq \mathbb{E}_{z \sim P} [h(x, z)]$

assumption: 1) projection on C is cheap

2) f is convex in x

3) we have stochastic oracle which gives g_t at time t , a random direction

such that
$$\mathbb{E}[g_t | x_t] = \underbrace{f'(x_t)}_{\text{subgradient of } f}$$

 $\mathbb{E}[g_t | \text{"past"}]$

[e.g. if f is differentiable in x and "well behaved"]

$$\nabla_x h(x, z_t) \text{ for } z_t \sim P \text{ then } \mathbb{E}_z [\nabla h(x, z)] = \nabla f(x)$$

more specifically if $f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$

then pick i_t unif. at random $1, \dots, n$

and letting $g_t = \tilde{g}_t(x_t)$

$$\text{then } \mathbb{E}[g_t | x_t] = \frac{1}{n} \sum_{i=1}^n f'_i(x_t) = f'(x_t)$$

$$4) \mathbb{E} \|g_t\|^2 \leq B^2 \quad [\text{finite variance condition}]$$

($\|f'_i(x)\| \leq B$ is sufficient for this) there

algorithm:

- $x_0 \in C$ initialization
- for $t=0, \dots, T-1$

get g_t from oracle \swarrow step-size

$$\text{let } x_{t+1} = \mathcal{P}_C[x_t - \gamma_t g_t]$$

• output $\hat{x}_T = \sum_{t=0}^T \lambda_t x_t$ where $\lambda_t \propto t$

"weighted average"

ie. $\lambda_t = \frac{2t}{T(T+1)}$

pointers:

- book for rates of convergence & lower bound:
 - Nesterov, "Introductory Lectures on Convex Optimization", 2004
- proof for convergence of weighted average stochastic subgradient method:
 - Lacoste-Julien, Schmidt, Bach, "A simpler approach to obtaining an $O(1/t)$ convergence rate for the projected stochastic subgradient method", arXiv:1212.2002
<https://arxiv.org/abs/1212.2002>