

Lecture 9 - convex optimization

Thursday, February 11, 2021 13:35

- today:
- furnish calibration
- start convex optimization

optimization normalization for calibration function

Setup let $s(x, \tilde{y})$ be of the form $F\Theta(x) \quad \Theta(x) \in \mathbb{R}^r$

$$s(x, \cdot) = F\Theta(x) \quad \underbrace{\Theta(\cdot)}_{\in \mathcal{H} \subset \text{RKHS}} \quad \underbrace{\Theta}_{\text{optimization variables}}$$

$$\mathcal{J}(\Theta) = \mathbb{E}_{(x,y) \sim p} \mathcal{J}(x, y; \Theta)$$

$$\text{run projected kernelized SGD on } \mathcal{J}(\Theta) \text{ ie. } \Theta^{(t+1)} = \underset{\mathcal{H}}{\text{P}}[\Theta^{(t)} - \gamma \nabla_{\Theta} \mathcal{J}(x^{(t)}, y^{(t)}; \Theta)]$$

projection
↓
 $(x^{(t)}, y^{(t)}) \stackrel{\text{iid}}{\sim} P$

ball of radius D around 0

$$\nabla_{\Theta} \mathcal{J}(x^{(t)}, y^{(t)}; \Theta) = F^T \nabla_{\Theta} \mathcal{J}(x^{(t)}, y^{(t)}; s) \Phi(x^{(t)})^T$$

$\Gamma \times \mathbb{R}$

feature map of \mathcal{H}

$$\Theta^{(t)} = F^T \sum_t \alpha_t \Phi(x^{(t)})^T$$

$\Gamma \times \mathbb{R}$

$$\Theta^{(t)}(x) \rightsquigarrow \langle \Phi(x^{(t)}), \Phi(x) \rangle$$

$K(x^{(t)}, \cdot)$

Convergence result:

(Thm. 5) if $\|\Theta^*\|_{\text{HS}} \leq D$ \downarrow \mathcal{J} is convex and differentiable

and if $\mathbb{E}_{(x,y) \sim p} \|\nabla_{\Theta} \mathcal{J}(x, y; \Theta)\|_{\text{HS}}^2 \leq M^2$

$\Theta^{[n]} = \frac{1}{n} \sum_{t=1}^n \Theta^{(t)}$ then averaged projected SGD with step-size $\gamma = \frac{2D}{M\sqrt{n}}$

$$\boxed{\mathbb{E}[\mathcal{J}(\Theta^{[n]})] - \mathcal{J}(\Theta^*) \leq \frac{2DM}{\sqrt{n}}}$$

gives

projected SGD with step-size $\gamma = \frac{2D}{M\sqrt{n}}$

$$\boxed{\mathbb{E}[\mathcal{J}(\Theta^{[n]})] - \mathcal{J}(\Theta^*) \leq \frac{2DM}{\sqrt{n}}}$$

Thm. 6 Learning complexity

Let G minimize $L(\theta)$ with $\|\theta\|_{HS} \leq D$

$$\text{choosing } n \geq \frac{4D^2M^2}{H(\epsilon)^2} \text{ implies } E[L(\theta^{(n)})] \leq L(\theta^*) + \epsilon$$

define a meaningful scale

in the paper: we compute $D f M \leq H(\epsilon)$ for specific losses l and the quadratic \mathcal{L} to get sample complexity



moral here:

- * some losses are harder than others (worst case sample complexity)
 - [ℓ_0 -loss is difficult in general
"too harsh of a loss"]

* have linked computation to statistical performance in consistency framework

↳ convex surrogate losses

* (non-parametric analysis) could handle dependence on x using RKFs

Caveats:

- * distribution-free result (i.e. worst case over all distributions)

→ still need more theory?

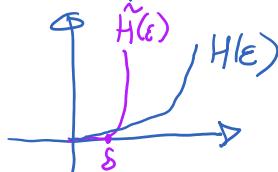
(e.g. role of $p(y|x)$)

or other surrogates? ▷

models $\|\theta\|_{HS} \leq D$ constraint

could be big for "bad p"
[→ no free lunch]

* follow-up: inconsistent surrogate loss with computational/statistical advantages [Neuens 2018]



14/20

part II: convex optimization

$$\text{motivation: } \min_w \frac{1}{2} \|w\|^2 + \sum_{i=1}^n l(x^{(i)}, y^{(i)}; w)$$

convex surrogate loss

convex analysis recap:

$$f: \mathbb{R}^d \rightarrow \mathbb{R}$$

$$f: \mathbb{R}^d \rightarrow \mathbb{R}$$

f is convex \Leftrightarrow

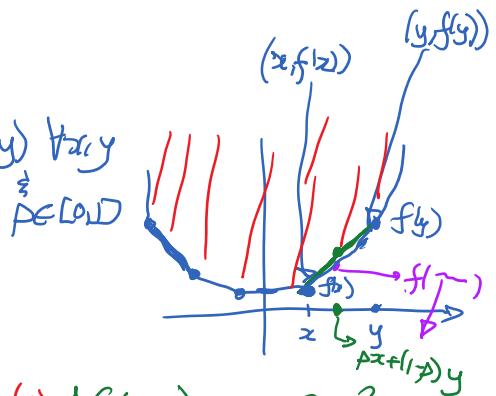
$$f(\rho x + (1-\rho)y) \leq \rho f(x) + (1-\rho)f(y) \quad \forall x, y$$

Convex combination
between $x \& y$

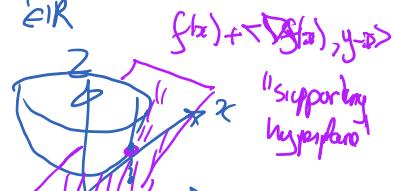
$$y + f(x-y)$$

$$y + f(x-y)$$

$$y + f(x-y)$$



$$\text{epigraph}(f) \triangleq \{(x, y) : y \geq f(x)\}$$



$$|x| = \max\{x, -x\}$$

* If f is differentiable at x and convex

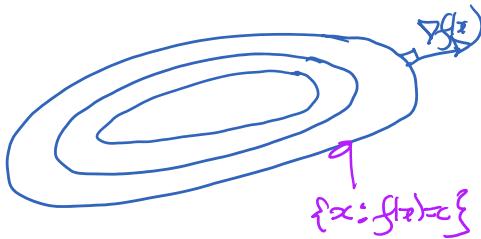
$$\Rightarrow f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle$$

(suppose f is convex)

subdifferentiable

Subgradient v of f at x : $v \in \partial f(x)$

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



$$\partial f(x) = \text{conv}\{\nabla f_i(x) : i \in \arg\max_i f_i(x)\}$$

when $f(x) = \max_i f_i(x)$ where f_i is differentiable

$$\partial f(x) = \text{conv}\{\nabla f_i(x) : i \in \arg\max_i f_i(x)\}$$

(Danzkin's theorem)

Danzkin's theorem: https://en.wikipedia.org/wiki/Danzkin%27s_theorem

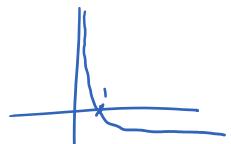
Clarke's subdifferential \rightarrow nice gen. to non-convex

$$f: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{-\infty\}$$

$$\text{dom}(f) \triangleq \{x \in \mathbb{R}^d : f(x) < \infty\}$$

e.g. $-\log x$ is convex

$$\text{dom}(f) = \{x \in \mathbb{R} : x > 0\}$$



$$\Rightarrow \min_x f(x) = \min_{x \in \text{dom}(f)} f(x)$$

Some standard assumptions:

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

strong convexity constant

$\langle v, y-x \rangle$ for any $v \in \partial f(x)$

$$f \text{ is } \mu\text{-strongly convex} \Leftrightarrow f - \frac{\mu}{2} \| \cdot \|_2^2 \text{ is convex}$$

$$f \text{ is } L\text{-smooth} \text{ i.e. } f \text{ is } L\text{-Lipschitz obs. gradient } \forall x \text{ (with respect to norm } \|\cdot\| \text{)}$$

$$\Leftrightarrow \|\nabla f(x) - \nabla f(y)\|_2 \leq L \|x-y\| \quad \forall x, y$$

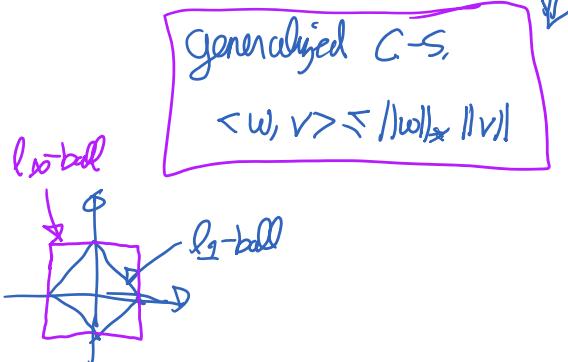
$$(\|\cdot\|_p)_* = \|\cdot\|_q$$

$$\text{where } \frac{1}{p} + \frac{1}{q} = 1$$

$$p=2 \Rightarrow q=2$$

$$p=1 \Rightarrow q=\infty$$

$$\|w\|_* \stackrel{\text{"dual norm"}}{\triangleq} \sup_{\|v\| \leq 1} \langle w, v \rangle$$



Fundamental descent lemma:

when ∇f is L -Lipschitz (Lemma holds even if f is not convex)

$$*\boxed{f(y) \leq f(x) + \langle \nabla f(x), y-x \rangle + \frac{L}{2} \|y-x\|^2 \quad \forall x, y}$$

$$*\boxed{f(x-\gamma \nabla f(x)) \leq f(x) - \gamma \underbrace{\langle \nabla f(x), \nabla f(x) \rangle}_{\|\nabla f(x)\|_2^2} + \frac{L\gamma^2}{2} \|\nabla f(x)\|_2^2}$$

$$= f(x) - \underbrace{\left[\gamma \left(1 - \frac{L}{2} \right) \right]}_{>0} \|\nabla f(x)\|_2^2$$

$$\Leftrightarrow \boxed{0 < \gamma < \frac{2}{L}}$$

$$\gamma_0 \Leftrightarrow \boxed{0 < \gamma < \frac{\alpha}{L}}$$

→ minimize RHS with respect to γ

gives $\boxed{\gamma^* = \frac{1}{L}}$

$$f(y_{\gamma^*}) \leq f(x) - \frac{1}{2} \frac{\|\nabla f(x)\|^2}{L}$$

Proof intuition of descent lemma:

Integral form of remainder

think of 2nd order Taylor expansion

$$f(y) = f(x) + \langle \nabla f(x), y - x \rangle + \frac{1}{2} \int_{y=x}^1 \langle y-x, H(x+\delta(y-x))y-x \rangle d\delta$$

Hessian of f

\int_L L-smooth
twice diff. \Rightarrow top eigenvalue of H $\leq L$
in absolute value

$$\sqrt{H} \leq \lambda_{\max}(H) \|v\|_2^2 \leq L \|y-x\|^2$$