

today: • finish calibration
• start convex optimization

optimization normalization for calibration set.

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

$$s(x, \cdot) = F\theta(x) \in \mathbb{R}^k$$

$\theta_j(\cdot) \in \mathcal{H} \leftarrow$ RKHS
optimization variables

$$\theta_j(x) = \langle \theta_j, \Phi(x) \rangle_{\mathcal{H}}$$

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

$(x^{(t)}, y^{(t)}) \sim p$

run projected kernelized SGD on $J(\theta)$

$$\text{i.e. } \theta^{(t+1)} = \mathcal{P}_D(\theta^{(t)} - \gamma \nabla_{\theta} J(x^{(t)}, y^{(t)}; \theta^{(t)}))$$

projection ball of radius D around 0

$$\nabla_{\theta} J(x^{(t)}, y^{(t)}; \theta) \approx$$

$r \times k$

$$F^T \sum_{k=1}^K \nabla_{\theta} J(x^{(k)}, y^{(k)}; s^{(k)}) \Phi(x^{(k)})^T$$

$F\theta(x^{(k)})$ feature map of \mathcal{H}
 $k(x^{(k)}, \cdot)$

$$\theta^{(n)} = F^T \sum_{t=0}^{n-1} \alpha_t \nabla_{\theta} J(x^{(t)}, y^{(t)}; s^{(t)})$$

$$\alpha_t = \gamma \nabla_{\theta} J(x^{(t)}, y^{(t)}; s^{(t)})$$

$$\theta^{(n)}(x) \rightarrow \langle \Phi(x^{(k)}, \Phi(x) \rangle = k(x^{(k)}, x)$$

convergence result
(thm. 5)

if $\|\theta^*\|_{\mathcal{H}} \leq D$ & J is convex and differentiable
and if $\mathbb{E}_{(x,y) \sim p} \|\nabla_{\theta} J(x,y;\theta)\|_{\mathcal{H}}^2 \leq M^2$

then averaged projected SGD with step size $\gamma = \frac{\partial D}{M\sqrt{n}}$

$$\theta^{[n]} = \frac{1}{n} \sum_{t=1}^n \theta^{(t)}$$

gives

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

$$\|\theta\|_{\mathcal{H}}^2 = \sum_{j=1}^r \|\theta_j\|_{\mathcal{H}}^2$$

thm. 6 Learning complexity

let θ^* minimize $L(\theta)$ with $\|\theta^*\|_{\mathcal{H}} \leq D$

"complexity measure" of θ^*

$$\text{choosing } n \geq \frac{4DM^2}{\epsilon^2}$$

implies that $\mathbb{E}[L(\theta^{[n]})] \leq L(\theta^*) + \epsilon$

defines a meaningful scale

in the paper: we compute D & M & $H(\epsilon)$ for

defines a meaningful scale

in the paper: we compute $D \approx M \approx H(\epsilon)$ for specific losses Q and the quadratic Q to get sample complexity

⊛ moral here:

* some losses Q are harder than others (worst case sample complexity)

[0-1 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 RKHS

concepts:

• distribution-free result (ie. worst case overall distribution.)

→ still need more theory

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

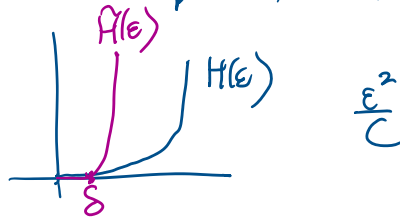
models $\| \cdot \|_{H_S} \leq D$ constraint

call to big for "bad p "
↳ no free lunch

(related to kernel choice?)

⊛ follow-up: inconsistent surrogate loss with computational/statistical advantages

[NeuIPS 2018]



10h17

part II: Convex optimization

motivation: $\min_w \lambda \frac{\|w\|^2}{2} + \frac{1}{n} \sum_{i=1}^n \mathcal{L}(x^{(i)}; y^{(i)}; w)$

convex surrogate loss

convex analysis recap:

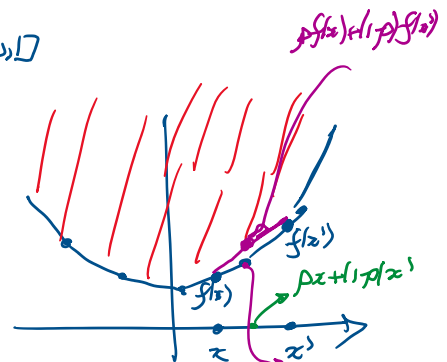
$f: \mathbb{R}^d \rightarrow \mathbb{R}$

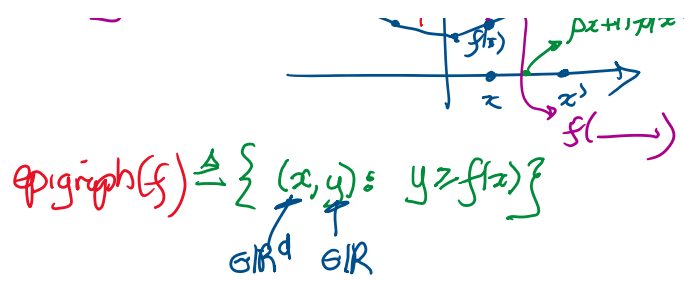
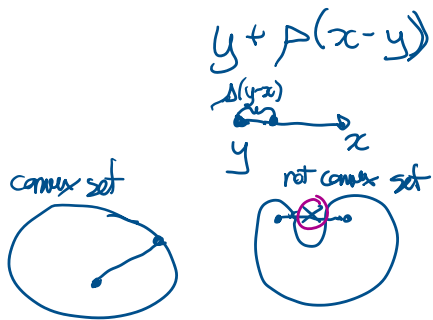
$\forall x, y \in \mathcal{D}, p \in [0, 1]$

f is convex $\Leftrightarrow f(pz + (1-p)y) \leq pf(x) + (1-p)f(y)$

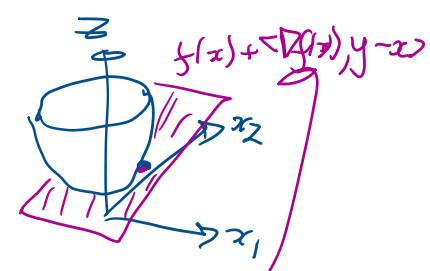
convex combination between x & y

$y + p(x - y)$





* If f is differentiable at x and convex
 $\Rightarrow f(y) \geq f(x) + \langle \nabla f(x), y-x \rangle \quad \forall y$

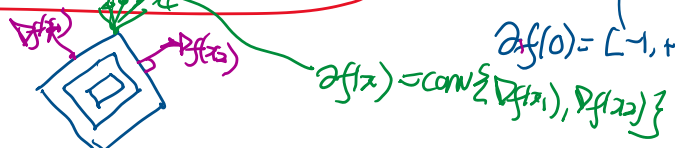
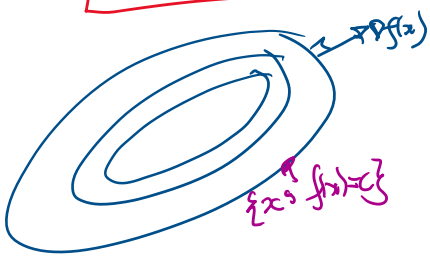
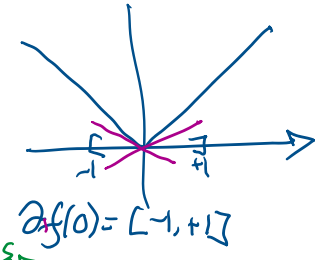


(suppose f is convex)

subdifferential

subgradient v of f at x : $v \in \partial f(x)$
 $\Leftrightarrow \forall y \in \text{dom}(f), f(y) \geq f(x) + \langle v, y-x \rangle$

"supporting hyperplane"
 $|x| = \max\{x, -x\}$



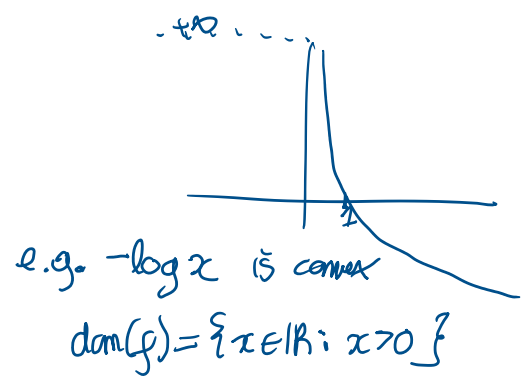
when $f(x) = \max_i f_i(x)$ where f_i is differentiable
 $\partial f(x) = \text{conv} \{ \nabla f_i(x) : i \in \text{argmax}_j f_j(x) \}$
 (Donskin's thm)

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

Clarke's subdifferential \rightarrow nice gen. to non-convex
 (limits of gradients for locally Lipschitz \rightarrow a.e. differentiable)

* extended reals convention for convex fct.:

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



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

Some standard assumptions:

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

$\forall x, y \in \text{dom}(f)$
 $\langle v, y-x \rangle$
for any $v \in \partial f(x)$

strong convexity constant

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

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

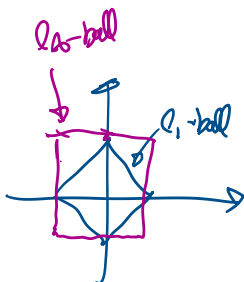
$$\iff \|\nabla f(x) - \nabla f(y)\|_* \leq L \|x-y\| \quad \forall x, y$$

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

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

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

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



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

Generalized CS.

$$\langle w, v \rangle \leq \|w\|_* \|v\|$$

Fundamental descent Lemma:

When ∇f is L -Lipschitz (lemma holds even if f is not convex)

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

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

$\delta \delta$

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

→ minimize RHS with respect to δ

guess $\delta^* = \frac{1}{L}$

$$f(y^*) \leq f(x) - \frac{1}{2L} \|\nabla f(x)\|_2^2$$

$$\boxed{f(y_*) \leq f(x) - \frac{1}{2L} \|\nabla f(x)\|_2^2}$$

$$f(x) - f(x^*) \leq \frac{\|\nabla f(x)\|_2^2}{2\mu}$$

$$f(x_{t+1}) - f(x^*) \leq f(x_t) - f(x^*) - \frac{\mu}{L} (f(x_t) - f(x^*))$$

$$\epsilon_{t+1} \leq \left(1 - \frac{\mu}{L}\right) \epsilon_t$$

$$(1-x)^t \leq \exp(-xt) \quad \forall x$$

$$\epsilon_t \leq \left(1 - \frac{\mu}{L}\right)^t \epsilon_0 \leq \exp\left(-\frac{\mu}{L}t\right) \epsilon_0 \quad \text{"linear rate"}$$

proof intuition of descent lemma:

think of 2nd order Taylor expansion

$$f(y) = f(x) + \langle \nabla f(x), y-x \rangle + \frac{1}{2} \int_{\gamma=0}^1 \langle y-x, \nabla^2 f(x+\gamma(y-x)) \rangle y-x \, d\gamma$$

integral form of remainder

Hessian of f

f is L -smooth
↳ twice diff.

⇒ for e-value of $H \leq L$
in absolute value

$$\leq L \|y-x\|^2$$

$$v^T H v \leq \lambda_{\max}(H) \|v\|^2$$