

today:

- finish calibration
- start convex optimization

optimization normalization for calibration function $\xrightarrow{\text{ker}}$ $\Theta \rightarrow \mathbb{R} \times \mathbb{R}$

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

$$S(x, \cdot) = F\Theta(x) \quad \Theta(\cdot) \in \mathcal{H} \leftarrow \text{RKHS}$$

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

$$\text{run projected (kernelized) SGD on } \mathcal{J}(\Theta) \quad \text{i.e. } \Theta^{(t+1)} = P_{\mathcal{B}}(\Theta^{(t)} - \gamma \nabla_{\Theta} \mathcal{J}(x^{(t)}, y^{(t)}; \Theta^{(t)}))$$

ball of radius D around 0

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

$F \times \mathbb{R}$

feature map of \mathcal{H}

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

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

$$= F^T \sum_t \alpha_t(\cdot) \underbrace{\langle \Phi(x^{(t)}), \Phi(x) \rangle}_{K(x^{(t)}, x)} \rightarrow \| \Theta \|_{HS}^2 = \sum_{j=1}^r \| \Theta_j \|_{RKHS}^2$$

convergence result:
(Thm. 5)

if $\|\Theta^*\|_{HS} \leq D$ & \mathcal{J} is convex & differentiable

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

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

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

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

Thm. 6 Learning complexity

Let Θ^* minimize $L(\Theta)$ with $\|\Theta^*\|_{HS} \leq D$

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

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

offers a meaningful scale

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

moral here:

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

[$l=1$ loss is difficult in general]
"too harsh of a loss"

* have linked computation to statistical performance in consistency framework

↳ convex surrogate loss

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

consists:

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

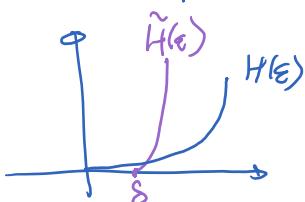
→ still need more theory?

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

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

could be big for "bad"
↳ no free lunch

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



15h18

part II: convex optimization

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

↳ convex surrogate loss

convex analysis recap:

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

f is convex \Leftrightarrow

$\dots \dots \dots \rightarrow \dots \dots \dots \rightarrow \dots \dots \dots$

\Leftrightarrow epigraph(f) is convex



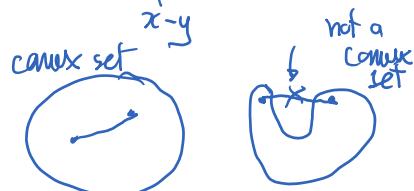
f is convex \Leftrightarrow

$$f(px + (1-p)y) \leq pf(x) + (1-p)f(y) \quad \forall x, y$$

convex combination
between $x \notin y$

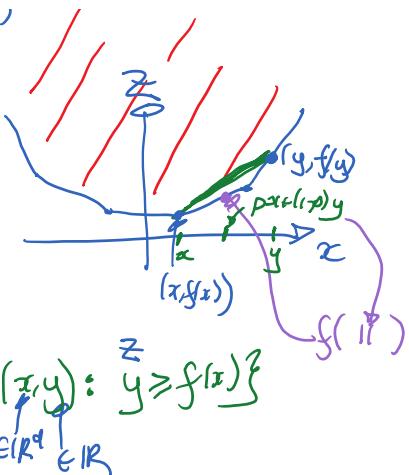
$$y + p(x-y)$$

$p \in [0, 1]$



"is convex"

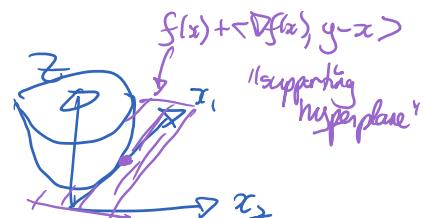
$$\nexists p \in [0, 1]$$



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

* f is differentiable at x and convex

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

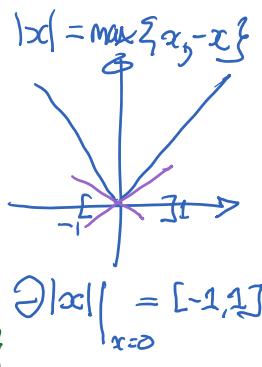
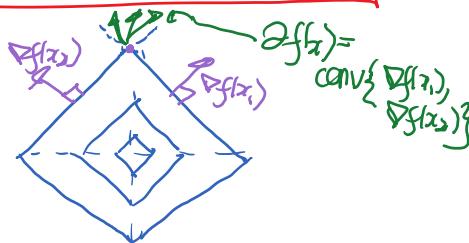
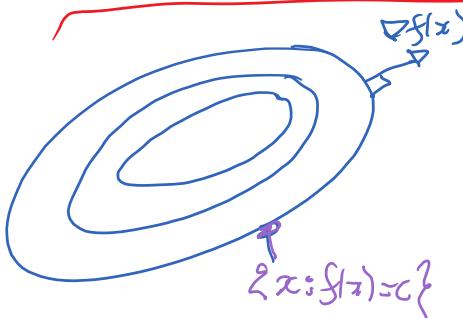


(suppose f is convex)

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

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

subdifferential



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

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

(Danzkin's thm.)

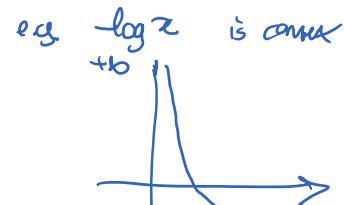
Danzkin's theorem: https://en.wikipedia.org/wiki/Danzkin%27s_theorem

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

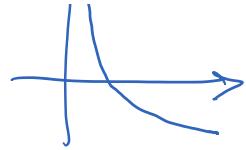
$$f: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{-\infty\}$$

"extended reals"

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



$f(x) = \text{convex function}$



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

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

Some standard assumptions

f is μ -strongly convex $\Leftrightarrow f(y) \geq f(x) + \underbrace{\langle \nabla f(x), y-x \rangle}_{\langle v, y-x \rangle} + \frac{\mu}{2} \|y-x\|^2$
 strong convexity constant

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

f is L -smooth i.e. f has L -Lipschitz continuous gradient ∇f (with respect to norm $\|\cdot\|$)
 $\Leftrightarrow \|\nabla f(x) - \nabla f(y)\|_* \leq L \|x-y\| \quad \forall x, y$

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

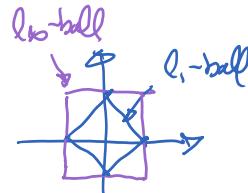
$$\text{where } \frac{1}{p} + \frac{1}{q} = 1$$

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

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

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

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

apply lemma with $y = x - \gamma \nabla f(x)$

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

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

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

→ minimize RHS with respect to γ

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

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

proof intuition of descent Lemma:

think of 2nd order Taylor expansion

$$f(y) = f(x) + (\nabla f(x)) \cdot (y-x) + \frac{1}{2} \int_{\gamma=0}^1 \underbrace{\langle y-x, H(x+\gamma(y-x)) \cdot (y-x) \rangle}_{\text{Hessian of } f} d\gamma$$

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

f is L -smooth \Leftrightarrow top eigenvalue of $H \leq L$
 $\&$ twice diff. in absolute value

$$v^T H v \leq \lambda_{\max}(H) \|v\|_2^2$$

* in absolute value