

today : finish calibration  
start convex opt.

Optimization normalization for calibration (cf.)

Setup let  $S(x, y)$  be of the form  $F(\theta(z))$   $G(z) \in \mathbb{R}^r$

$S(x, \cdot) = F(\theta(x))$   $\theta(\cdot) \in \mathcal{H} \leftarrow \text{RHS}$   
 $\in \mathbb{R}^k$  optimization variables

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

run projected SGD on  $L(\theta)$  (kernelized) i.e.  $\theta^{(t+1)} = \underset{\theta \in \mathcal{H}}{\text{Proj}} \left( \theta^{(t)} - \gamma \nabla_{\theta} S(x^{(t)}, y^{(t)}; \theta^{(t)}) \right)$   
 (ball of radius  $D$  around  $\theta$ )

$\nabla_{\theta} S(x^{(t)}, y^{(t)}; \theta) = F^T \nabla_S S(x^{(t)}, y^{(t)}; s) \Phi(x^{(t)})^T$   
 $r \times b$   $\uparrow$  feature map of  $\mathcal{H}$   
 $K(x^{(t)}, \cdot)$

convergence result: (thm. 5) if  $\|\theta^*\|_{\mathcal{H}} \leq D$   $S$  is convex and differentiable  
 and if  $\mathbb{E}_{(x,y) \sim p} \|\nabla_{\theta} S(x,y;\theta)\|_{\mathcal{H}}^2 \leq M^2$

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

gives  $\mathbb{E}[L(\bar{\theta}^n)] - L(\theta^*) \leq \frac{2DM}{\sqrt{n}}$  (convergence result)  
 $\uparrow$   
 $\frac{1}{n} \sum_{t=1}^n \theta^{(t)}$

thm. 6 Learning complexity

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

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

define a meaningful scale

in the paper we compute  $D \lesssim M \lesssim \sqrt{H(\epsilon)}$  for specific losses  $S$  and the quadratic  $L$  to get sample complexity

⊗ Moral here:  
 \* some losses are harder than others (worst case sample complexity)

\* some losses are harder than others (worst case sample complexity)  
 [0-1 loss is difficult in general]

\* have linked computation to statistical performance in consistency framework  
 ↳ convex surrogate loss

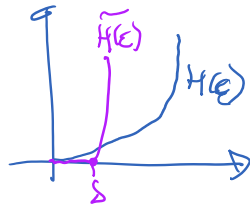
\* could handle dependence on  $x$  using RKHS

cautions:

• distribution free result (i.e. worst case over all distribution)  
 modulo  $\mathcal{H}^{\infty}(\mathcal{X}, \mathcal{Y})$  constraint

→ still need more theory! (e.g. rate of  $p(y|x)$ ? or other surrogates?)

\* follow-up paper: inconsistent surrogate loss with computational/statistical advantage NIPS 2015



15h38

part II: convex optimization

Motivation:  $\min_w \frac{\|w\|_2^2}{2} + \frac{1}{n} \sum_{i=1}^n \ell(x^{(i)}, y^{(i)}; w)$

convex surrogate loss

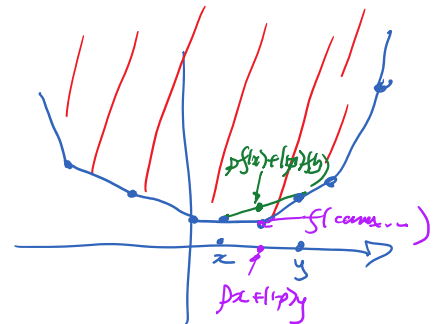
Convex analysis recap:

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

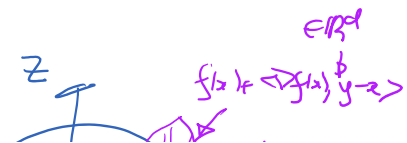
$f$  is convex  $\Leftrightarrow$

$f(px + (1-p)y) \leq pf(x) + (1-p)f(y)$   $\forall x, y \in \text{dom}(f)$

convex combination between  $x$  &  $y$   
 $y + p(x-y)$

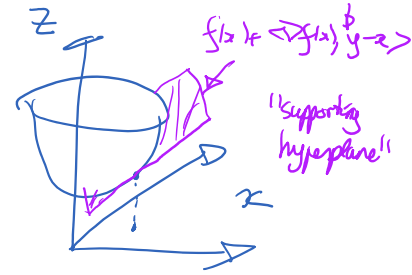


epigraph  $(f) \triangleq \{ (x, y) : y \geq f(x) \}$   
 $x \in \mathbb{R}^d, y \in \mathbb{R}$



\* if  $f$  is differentiable at  $x$  and convex

$$\Rightarrow f(y) \geq f(x) + \langle \nabla f(x), y-x \rangle \quad \underline{\underline{by}}$$

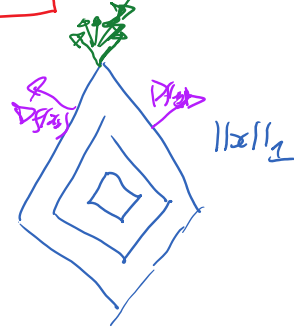
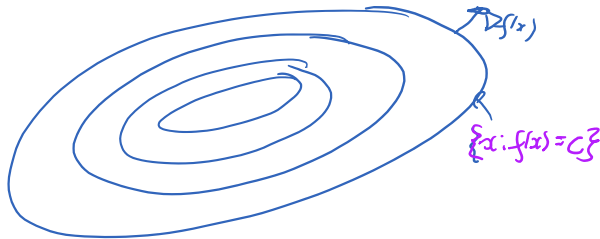
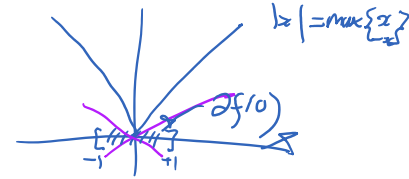


(suppose  $f$  is convex)

subdifferentiable

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

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



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

$$\partial f(x) = \text{conv}\{ \nabla f_{i^*}(x) : i^* \in \text{argmax}_i f_i(x) \}$$

(Danskin's thm.)

Danskin's theorem : [https://en.wikipedia.org/wiki/Danskin%27s\\_theorem](https://en.wikipedia.org/wiki/Danskin%27s_theorem)

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

some standard assumptions :

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

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

$\forall x, y \in \text{dom}(f)$

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

strong convexity constant  $\uparrow$

$f$  is  $L$ -smooth i.e.  $f$  has  $L$ -Lipschitz gradient  $\forall x$  (with ref. norm  $\|\cdot\|_p$ )

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

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

generalized C-S.

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

$$(\|\cdot\|_p)^* = \|\cdot\|_q \quad \text{where } \frac{1}{p} + \frac{1}{q} = 1$$

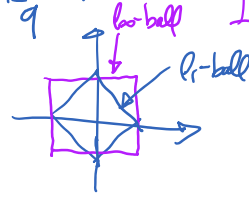
$p=2 \Rightarrow q=2$

$\uparrow$  ball  $\downarrow$  -ball

$$(\|\cdot\|_p) = (\|\cdot\|_q) \text{ where } \frac{1}{p} + \frac{1}{q} = 1$$

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

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



$$|\langle w, v \rangle| \leq \|w\|_p \|v\|_q$$