

today: finish calibration for
convex optimization

(continuation)

thm. 6: learning complexity

let G^* minimize $L(\epsilon)$ with $\|G^*\|_{H_S} \leq D$

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

define
a meaningful
scale

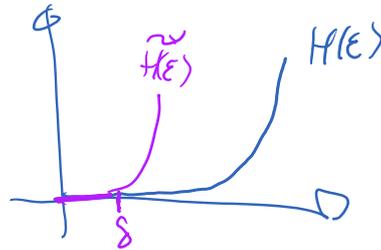
in the paper: we compute D & M & $H(\epsilon)$ for specific losses l
and the quadratic L
to get sample complexity

(*) Moral here:

- * 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 distributions)
→ still need more theory! (e.g. role of poly(x)?
or other surrogates)

* follow-up paper: inconsistent surrogate losses
with computational/statistical advantage NIPS 2018



15h15

PART II: convex optimization

convex surrogate loss

PART II: convex optimization

convex surrogate loss

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

convex analysis recap:

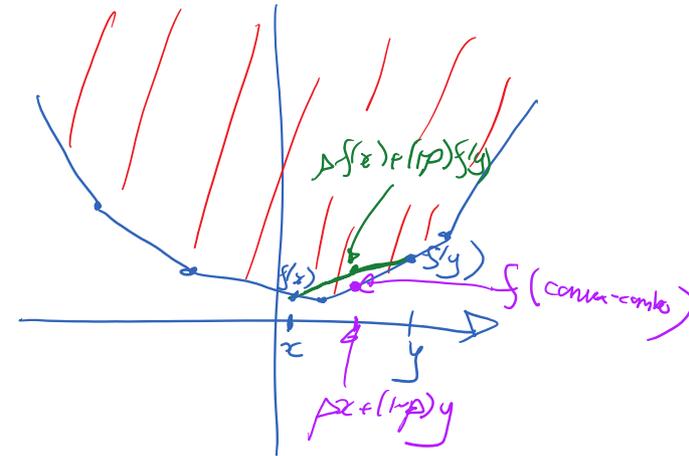
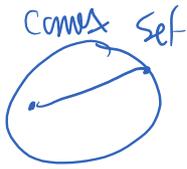
$f: \mathbb{R}^d \rightarrow \mathbb{R}$

f is convex \Leftrightarrow

$f(\underbrace{px + (1-p)y}_{\text{convex combination between } x \text{ \& } y}) \leq pf(x) + (1-p)f(y) \quad \forall x, y, p \in [0, 1]$

convex combination between x & y

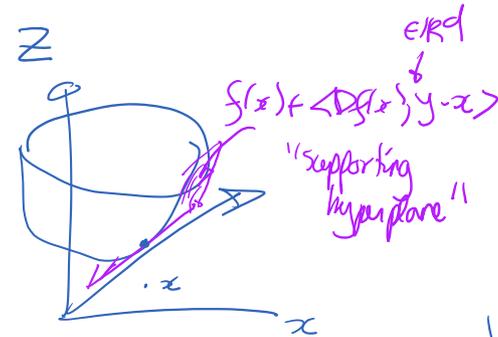
$\rightarrow y + p(x - y)$



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

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



(suppose f is convex)

subdifferential at $x = \{ \text{subgradients} \}$

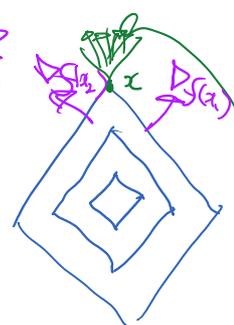
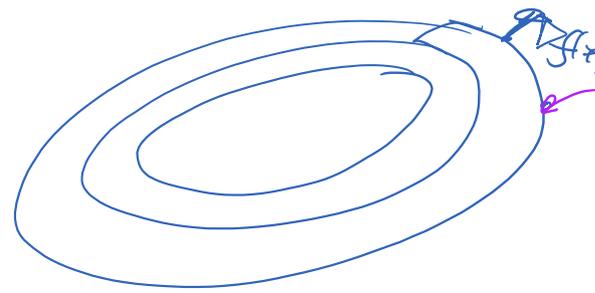
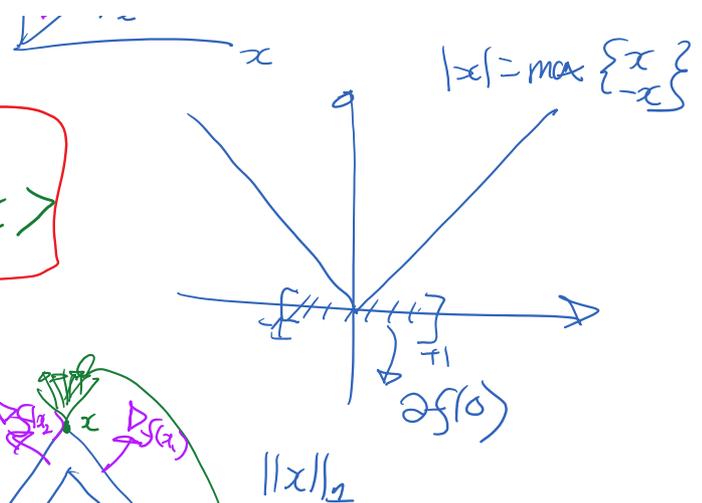
$|x| = \max \{ x, -x \}$

(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
 $\partial f(x) = \text{conv}\{\nabla f_i(x) : i \in \text{argmax}_j f_j(x)\}$

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

some standard assumptions:

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

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

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

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

Strong convexity constant

f is L -smooth i.e. f has L -Lipschitz gradient $\forall x$
 $\Leftrightarrow \|\nabla f(x) - \nabla f(y)\|_* \leq L \|x - y\|$

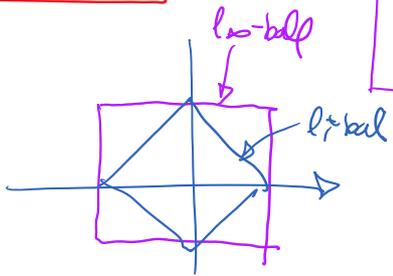
$$\|w\|_* \triangleq \sup_{\|v\| \leq 1} \langle w, v \rangle$$

generalized C.S.
 $\langle w, v \rangle \leq \|w\|_* \|v\|$

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

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

$$p=1 \Leftrightarrow q=\infty$$



Fundamental descent lemma:

when ∇f is L -Lipschitz (even if f is not nec. convex)

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

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

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

→ minimize RHS with respect to δ

gives $\delta^* = 1$

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

gives

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

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

[Think of 2nd order Taylor expansion]

$$f(y) = f(x) + \langle \nabla f(x), y-x \rangle + \frac{1}{2} \int_0^1 \langle y-x, H(x+\delta(y-x)) y-x \rangle d\delta$$

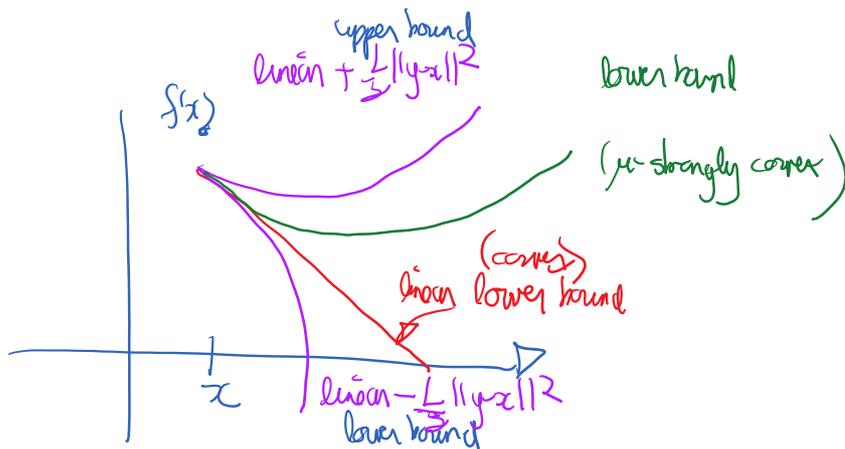
top-evalue of $H \leq L$
in absolute value

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

or
"pro-way" is to use fundamental thm. of calculus

$$f(y) = f(x) + \int_0^1 \underbrace{\frac{d}{d\delta} f(x+\delta(y-x))}_{\langle \nabla f(x+\delta(y-x)), y-x \rangle} d\delta$$

(see Nesterov)



when f is twice differentiable
 $L = \lambda_{\max}(H)$ Hessian

$\frac{1}{L}$ | linear - $\frac{1}{L} \|y - x\|^2$
 lower bound

$$L = \lambda_{\max}(H)$$

$$\mu = \lambda_{\min}(H)$$

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

gradient descent: $x_{t+1} = x_t - \gamma \nabla f(x_t)$ $\gamma = \frac{1}{L}$

a) when f is convex & L -smooth

$$f(x_t) - \underbrace{\min_x f(x)}_{f^*} \leq O\left(\frac{L r_0^2}{t}\right) \quad \text{where } r_0 \geq \text{dist}(x_0, x^*)$$

"sublinear rate"

origin $f(x)$

\hookrightarrow [see Nesterov book for $O(\frac{1}{t})$ rate]²

note: no guarantee on $\text{dist}(x_t, x^*)$ (in general)

\rightarrow Nesterov lower bounds

b) if f is μ -strongly convex & L -smooth "linear rate"

$$f(x_t) - f(x^*) \leq O\left(\exp\left(-\frac{\mu}{L} t\right)\right)$$

$\frac{L}{\mu} \triangleq$ condition # of f



~~~~~

Newton's method:  $x_{k+1} = x_k - \delta H(x_k)^{-1} \nabla f(x_k)$

↓

$\nabla \nabla f(x_k)$