

today: finish surrogate losses
theory basics

III) structured hinge loss

$$f^{\text{SVM}}(x, y; \omega) = \max_{\tilde{y} \in \mathcal{Y}(x)} [s(x, \tilde{y}; \omega) + l(y, \tilde{y})] - s(x, y; \omega)$$

loss-augmented decoding

a) cartoon:

$$\Rightarrow f^{\text{SVM}}(x, y; \omega) = 0$$

b) $f^{\text{SVM}}(x, y; \omega) \geq l(y, h_{\omega}(x))$

why? $f(x, y; \omega) = \max_{\tilde{y}} [s(\tilde{y}) + l(y, \tilde{y})] - s(y)$ Let $\hat{y} = \underset{\tilde{y}}{\text{argmax}} s(\tilde{y})$

$$\geq s(\hat{y}) + l(y, \hat{y}) - s(y)$$

if $y \in \mathcal{Y}(x) \Rightarrow s(\hat{y}) \geq s(y)$

$$\Rightarrow l(\hat{y}) = l(y, h_{\omega}(x)) //$$

binary case: structured hinge loss reduces to binary SVM loss

$$y \in \{-1, +1\} \quad \omega = \begin{pmatrix} \omega_+ \\ \omega_- \end{pmatrix} \quad \varphi(x, +1) = \begin{pmatrix} \varphi(x) \\ 0 \end{pmatrix}$$

$$h_{\omega}(x) = \underset{y \in \{-1, +1\}}{\text{argmax}} \{ \langle \omega_+, \varphi(x) \rangle, \langle \omega_-, \varphi(x) \rangle \}$$

predict +1 if $\langle \omega_+, \varphi(x) \rangle \geq \langle \omega_-, \varphi(x) \rangle$

$$\Leftrightarrow \langle \underbrace{\omega_+ - \omega_-}_{\tilde{\omega}}, \varphi(x) \rangle \geq 0$$

$$\tilde{\omega} \triangleq \omega_+ - \omega_- \quad h_{\omega}(x) = \text{sgn}(\langle \tilde{\omega}, \varphi(x) \rangle)$$

(see notes to show)

structured hinge loss

$$\tilde{w} = w_+ - w_- \quad h_w(x) = \text{sgn}(\langle \tilde{w}, x \rangle)$$

$$J^{\text{SUM}}(x, y; w) = \max \left\{ \langle w_+, x \rangle + \ell(y, +), \langle w_-, x \rangle + \ell(y, -) \right\} - \langle w_y, x \rangle$$

$w_+ = \tilde{w} + w_-$ $\mathbb{1}\{y \neq +1\}$ $1 - \mathbb{1}\{y \neq +1\}$

$$= \max \left\{ \langle \tilde{w}, x \rangle + \langle w_-, x \rangle + \mathbb{1}\{y \neq +1\}, \langle w_-, x \rangle + 1 - \mathbb{1}\{y \neq +1\} \right\} - \langle w_y, x \rangle$$

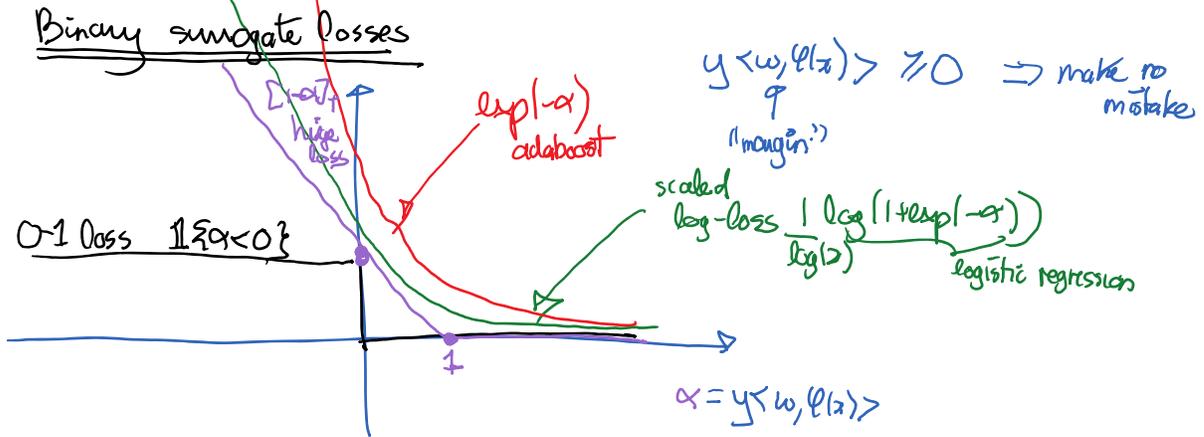
$$= \max \left\{ \langle \tilde{w}, x \rangle + 1, 1 - 1 \right\} + \langle w_-, x \rangle - \langle w_y, x \rangle$$

case $y = +1$: $\max \{ \langle \tilde{w}, x \rangle, 1 \} - \langle \tilde{w}, x \rangle = [1 - y \langle \tilde{w}, x \rangle]_+$
 case $y = -1$: $\max \{ \langle \tilde{w}, x \rangle + 1, 0 \} + 0 = [1 - y \langle \tilde{w}, x \rangle]_+$

overall: $[1 - y \langle \tilde{w}, x \rangle]_+$

$$J^{\text{SUM}}(x, y; w) = [1 - y \langle \tilde{w}, \phi(x) \rangle]_+ \quad \text{where } \tilde{w} \triangleq w_+ - w_-$$

ie. structured hinge loss reduces to binary SUM hinge loss when using $\ell(y, y') = \mathbb{1}\{y \neq y'\}$ and $\mathcal{Y} = \{-1, +1\}$



[Bartlett et al. 2006] \rightarrow showed all these methods are consistent

Theory basics

decision theory setup

estimate $h_w: X \rightarrow \mathcal{Y}$

generalization error = $L_P(w) \triangleq \mathbb{E}_{(x,y) \sim P} [\ell(y, h_w(x))]$

ultimate goal is to find $w^* = \underset{w \in \mathcal{H}}{\text{argmin}} L_P(w)$

ultimate goal is to find $w^* = \underset{w \in \mathcal{W}}{\operatorname{argmin}} L_P(w)$

problem: do not know P ("true" distribution on (x, y))

suppose $(x^{(i)}, y^{(i)})_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} P$ \rightarrow we could look at $\hat{L}_n(w) = \frac{1}{n} \sum_{i=1}^n \ell(y_i^{(i)}; w(x^{(i)}))$

$\cong D_n$
training dataset

from statistics / prob. theory

$\hat{L}_n(w) \xrightarrow[n \rightarrow \infty]{\text{a.s.}} L_P(w)$ for each fixed w (pointwise)
(LLN)

this is weaker than $\sup_w |\hat{L}_n(w) - L_P(w)| \xrightarrow[n \rightarrow \infty]{} 0$

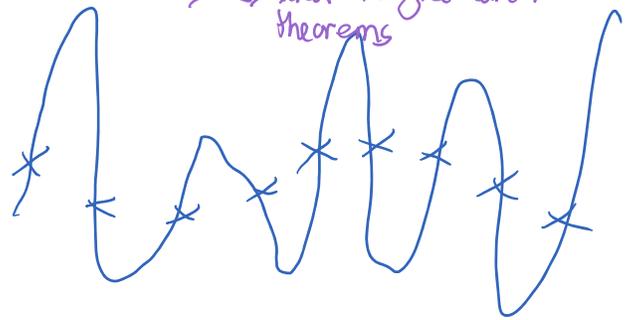
note: minimizing the training error gives no guarantee in general?
on $L_P(w)$

\rightarrow later no free lunch theorems

e.g. polynomial regression

for n points, can get zero training error with poly of degree $n-1$

\Rightarrow overfitting



in learning theory: study properties of learning alg. $A: D_n \rightarrow \mathcal{W}$

in particular, what can we say about $L_P(A(D_n))$

different approaches:

a) "frequentist risk" $R_{P,n}^F(A) \triangleq \mathbb{E}_{D_n \sim P^n} [L_P(A(D_n))]$

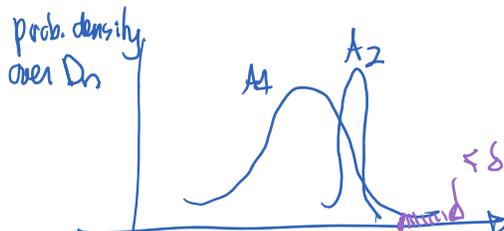
D_n is random

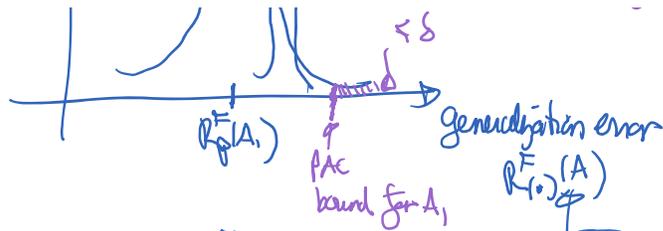
b) PAC framework "probably approximately correct" $P \sum L_P(A(D_n)) > \text{some bound } \leq \delta$

"tail bound"

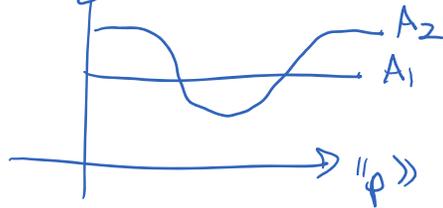
ie. $L_P(A(D_n)) \leq \text{some bound}$ with prob. $\geq 1 - \delta$

generalization error bound





issue with $R_p^F \rightarrow$ depends on P
"risk profiles"



weighted frequentist risk $\mathbb{E}_{\theta \sim \pi(\theta)} [R_{P_\theta}^F(A)]$

c) "Bayesian posterior risk"

$$R^{Post}(\omega | D_n) \triangleq \mathbb{E}_{\theta \sim p(\theta | D_n)} [L_{P_\theta}(\omega)]$$

Bayesian estimate $\hat{\omega}_n^{Bayes} = \underset{\omega}{\operatorname{argmin}} R^{Post}(\omega | D_n)$

$A^{Bayesian}$ is optimal for weighted frequentist risk using $\pi(\theta) = p(\theta)$

- prior $p(\theta)$ over dist.
- observation model $p(D_n | \theta)$
- \Rightarrow posterior $p(\theta | D_n)$

14h53

no free lunch?

frequentist risk analysis learning deg. A

let \mathcal{P} be a set of distributions on $X \times \mathcal{Y}$

sample complexity of A with respect to \mathcal{P}

is the smallest $n(\mathcal{P}, A, \epsilon) \leq \infty$ s.t. $\forall n \geq n(\mathcal{P}, A, \epsilon)$

we have $\sup_{P \in \mathcal{P}} [R_P^F(A; n) - L_P(h_P^*)] \leq \epsilon$

"uniform result"

$h_P^* = \underset{h: X \rightarrow \mathcal{Y}}{\operatorname{argmin}} L_P(h)$

terminology: • A is consistent for dist. p

is $\lim_{n \rightarrow \infty} R_p^F(A; n) - L_p(h_p^*) = 0$

• A is uniformly consistent for a family \mathcal{P}

is $\lim_{n \rightarrow \infty} \left[\sup_{P \in \mathcal{P}} [R_P^F(A; n) - L_P(h_P^*)] \right] = 0$

Binary classification $\mathcal{Y} = \{-1, +1\}$

I) if X is finite; then the "voting procedure" (assign the most frequent label to an input x)
 is uniformly and universally consistent
 \hookrightarrow i.e. \mathcal{P} is all dist. on $X \times Y$
 with (universal) sample complexity

$$n(\mathcal{P}, \varepsilon, \text{Averaging}) \leq \frac{|X|}{\varepsilon^2} \quad (\text{free lunch?})$$

II) if X is infinite

no free lunch theorem \exists (for binary with ℓ the 0-1 loss)

for any n and any learning alg. A

$$\text{then } \sup_{P \text{ all dist.}} [R_P^F(A; n) - L_P(h_P^*)] \geq \frac{1}{2}$$

i.e. \exists always a dist. $P_{A,n}$ s.t. A is worse than random predictor for $P_{A,n}$

NFL II: [Thm. 7.2 in Devroye & de 1996]
 let ε_n be any non-increasing seq. converging to 0 + ($\varepsilon_n \leq 1$)
 for A , then $\exists P_A$ s.t.
 $[R_{P_A}^F(A; n) - L_{P_A}(h_{P_A}^*)] \geq \varepsilon_n \ln n$
 (could be arbitrarily slow e.g. $\frac{1}{\lg(\lg(\lg(\dots \ln)))}$)

~~consequence~~ consequence is we need assumptions on \mathcal{P} to say anything useful

Occam's generalization error bound

- binary class. & 0-1 loss
- consider W to be a countable set

let's define a prior over W : $\pi(w)$ i.e. $\sum_{w \in W} \pi(w) = 1$ $\pi(w) \geq 0 \forall w$

$$|w|_\pi = \text{"description length" of } w \triangleq \log_2 \frac{1}{\pi(w)} \implies \sum_w 2^{-|w|_\pi} \leq 1$$

"Kraft's inequality"

Occam's bound

"(110) Kraft's inequality"

Occam's bound

for any fixed P ; with prob. $\geq 1-\delta$ over training set $D_n \sim P^{\otimes n}$

$$\forall w \in W \quad L_P(w) \leq \hat{L}_P(w) + \frac{1}{\sqrt{2n}} \Omega_{\pi}(w; \delta)$$

$$\text{where } \Omega_{\pi}(w; \delta) \stackrel{\text{def}}{=} \sqrt{(\ln 2) \underbrace{|w|}_{\text{complexity measure}} + \ln \frac{1}{\delta}}$$