

today: - structured hinge loss
- theory binary classification

III) structured hinge loss

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

"loss-augmented decoding"

carbon

$$\begin{array}{c} \overbrace{s(y)}^{> l(y, \tilde{y})} \\ \underline{s(\tilde{y}_{\text{next}})} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{array} \Rightarrow \sum s(x, y; w) = 0$$

b) $f(x, y; w) \geq l(y, h_w(x))$

why? $f(x, y; w) = \max_{\tilde{y}} [s(\tilde{y}) + l(y, \tilde{y})] - s(y)$

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

using $\hat{y} = \operatorname{argmax}_{\tilde{y} \in \mathcal{Y}(x)} s(\tilde{y}) = h_w(x)$

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

$$\geq l(\hat{y}) = l(y, h_w(x)) //$$

binary case:

$$u \in \{-1, +1\}$$

$$w = \begin{pmatrix} w_+ \\ w_- \end{pmatrix}$$

prediction: $h_w(x) = \operatorname{argmax} \{ \langle w_+, x \rangle, \langle w_-, x \rangle \}$

0-1 loss

predict +1 if $\langle w_+, x \rangle \geq \langle w_-, x \rangle$

$$\Leftrightarrow \langle w_+ - w_-, x \rangle \geq 0$$

$$\tilde{w} \triangleq w_+ - w_- \Rightarrow w_+ = \tilde{w} + w_-$$

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

$$J^{\text{svm}}(x, y; w) = \max_{\substack{\uparrow \\ \mathbb{I}\{y \neq +1\}}} \{ \langle w_+, x \rangle + \ell(y, +), \langle w_-, x \rangle + \ell(y, -) \} - \langle w_y, x \rangle$$

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

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

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

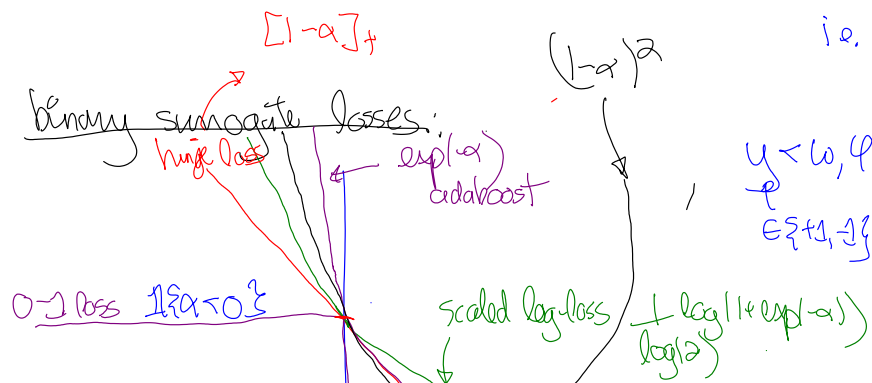
structured hinge

$$J(x, y; w) = [1 - y \langle \tilde{w}, x \rangle]_+$$

where $\tilde{w} = w_+ - w_-$

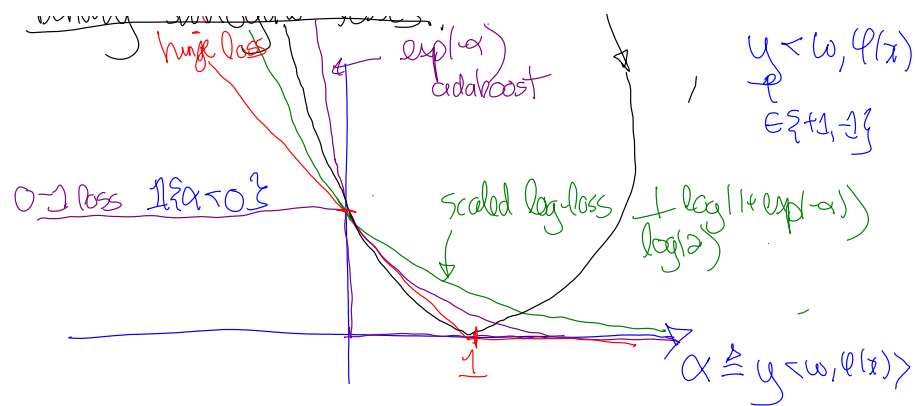
\Rightarrow that structure hinge loss reduces to standard binary SVM hinge loss when using $\ell(y, y') = \mathbb{I}\{y \neq y'\}$ and $\mathcal{Y} = \{-1, +1\}$

i.e. binary SVM is special case of structured SVM



$y \langle w, \phi(x) \rangle$ "margin" $\geq 0 \Rightarrow$ make no mistake
 $\mathcal{Y} = \{-1, +1\}$

$$\log\text{-loss: } \log(1 + \exp(-x))$$



$y \langle w, \phi(x) \rangle$ "margin" $\geq 0 \Rightarrow$ make no mistake
 $\in \{-1, 1\}$

log-loss: $\log(1 + \exp(-\alpha))$

[logistic regression
 $\log(1 + \exp(-y \langle w, \phi(x) \rangle))$]

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

Question on scaling the margin:

$$\frac{1}{n} \sum_{i=1}^n [1 - y_i \langle w, \phi(x_i) \rangle]_+ + \lambda \frac{\|w\|^2}{2}$$

$$\nu \left[\frac{1}{n} \sum_{i=1}^n [1 - y_i \langle \frac{w}{\nu}, \phi(x_i) \rangle]_+ + \frac{\lambda}{\nu} \frac{\|w\|^2}{2} + (\lambda \nu) \frac{\|\tilde{w}\|^2}{2} \right]$$

i.e. is equivalent to just rescale $\lambda \dots$

Theory basics

decision theory setup

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

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

task loss

ultimate goal is find $w^* = \arg \min_{w \in W} L_P(w)$

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

suppose $\underbrace{(x^{(i)}, y^{(i)})}_{\triangleq D_n}^n \stackrel{i.i.d.}{\sim} P$ training data

learning algorithm $\hat{w}_n = A(D_n)$
 \uparrow
 algorithm

\rightarrow we could look at $\hat{L}_n(w) = \frac{1}{n} \sum_{i=1}^n \ell(y^{(i)}, h_w(x^{(i)}))$

from statistics/prob

$$\hat{L}_n(w) \xrightarrow{a.s.} L_P(w) \quad \forall w$$

(LLN)

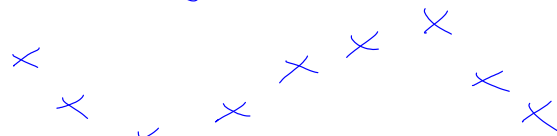
\rightarrow this is weaker than $\sup_w |\hat{L}_n(w) - L_P(w)| \xrightarrow{n \rightarrow \infty} 0$

*note: minimizing training error gives no guarantee in general

e.g. linear regression

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

\Rightarrow overfitting



in learning theory, want to study properties of learning algorithm

in particular, what can we say about $L_P(\underbrace{A(D_n)}_{\hat{w}_n})$?

different approaches: "frequentist risk" $R_P^F(A) \triangleq \mathbb{E}_{D_n \sim P^n} [L_P(A(D_n))]$

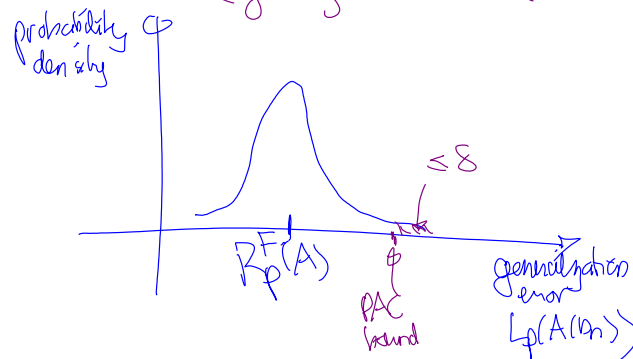
D_n is random

PAC framework
'probably approximately correct'

$$P\{L_P(A|D_n) > \text{some bound}\} \leq \delta$$

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

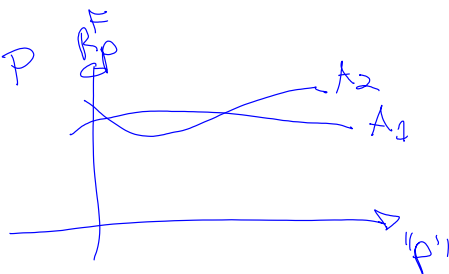
generalization error bound



$$EX = \underbrace{E[X \mathbb{1}\{X \leq B\}]}_{B \cdot (1-\delta)} + \underbrace{E[X \mathbb{1}\{X > B\}]}_{\leq B \cdot \delta}$$

↑
test error limit

issue with $R_P^F \rightarrow$ depends on P



weighted frequentist risk: $E_{G \sim p(G)} [R_{P_G}^F(A)]$

"Bayesian posterior risk"

$$R^{\text{post}}(w|D_n) = E_{G \sim p(G|D_n)} [L_{P_G}(w)]$$

prior $p(G)$ over distributions
 $p(D_n|G)$
 \Rightarrow posterior $p(G|D_n)$

Bayesian estimate $\hat{w}_n = \underset{w}{\text{argmin}} R^{\text{post}}(w|D_n)$

Bayesian is optimal for weighted frequentist risk using $p(G)$

No free lunch

frequentist risk analysis of learning algorithm A

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

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

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

we have
$$\sup_{P \in \mathcal{P}} [R_P^E(A) - L_P(h_P^*)] < \epsilon$$

"uniform result"

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

terminology: A is consistent for dist. P

if
$$\lim_{n \rightarrow \infty} \underbrace{R_P^E(A; n)}_{\mathbb{E}_n[L_P(A(n))]} - L_P(h_P^*) = 0$$

A is uniformly consistent for family \mathcal{P}

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

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

I) if X is finite; then the "voting procedure" (assign the most frequent label to a input x)

is uniformly and universally consistent

\hookrightarrow i.e. \mathcal{P} is all distributions on $X \times Y$

with (uniform) sample complexity $n(\mathcal{P} = \text{all}, \epsilon) \leq |X|$ (\dots $n(\mathcal{P} = \text{all}, \epsilon) \leq |X|$)

with (universal) sample complexity $n(\mathcal{D}, \varepsilon, A, \text{data}) \leq \frac{|X|}{\varepsilon^2}$

(free lunch? !!)

II) If X is infinite

no free lunch theorem:

(for binary classification with the 0-1 loss)

for any n and any learning algorithm A

$$\text{then } \sup_{\mathcal{P}} [R_{\mathcal{P}}^F(A; n) - L_{\mathcal{P}}(h_{\mathcal{P}}^*)] \geq \frac{1}{2}$$

i.e. always a distribution \mathcal{P} s.t. your algorithm A is not doing better than random prediction (\mathcal{P})

NFT II)

[thm. 7.2 in Devroye et al. 1996]

let ε_n be any non-increasing sequence converging to zero (could be arbitrarily slow) and any alg. A

$$\text{then } \exists \mathcal{P} \text{ s.t. } R_{\mathcal{P}}^F(A; n) - L_{\mathcal{P}}(h_{\mathcal{P}}^*) \geq \varepsilon_n \quad \forall n$$

⊗⊗ consequence: we need assumptions on \mathcal{P} to say anything

Occam's generalization error bound

- binary classification and 0-1 loss
- convex W to be a countable set

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

$|w|_{\pi}$ = "description length" of w

$$\leq \log_2 \frac{1}{\pi(w)}$$

Occam's bound

for any P_j with probability $\geq 1-\delta$ over training set $D_n \sim P^{\otimes n}$

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

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

⊛ bound is only useful for distribution P s.t. $|w_P^*|_{\pi}$ is small

$$w_P^* = \arg \min_{w \in W} L_P(w)$$

$$|w|_{\pi} = \log_2 \frac{1}{\pi(w)}$$

$$\text{if } \pi(w) \propto \exp(-\|w\|^2)$$

$$\text{then } |w|_{\pi} = \|w\|^2 + \text{const.}$$

proof: use 3 things

1) Chernoff bound
(concentration inequality)

$$P\{D_n: \hat{L}_n(w) \leq L(w) - \varepsilon\} \leq \exp(-2n\varepsilon^2) \quad \forall \varepsilon \geq 0$$

2) union bound

$$P\{\exists x \text{ s.t. } \text{prop}(x) \text{ is true}\} \leq \sum_x P\{\text{prop}(x) \text{ is true}\}$$

3) "Kraft's inequality" $\sum_w 2^{-|w|_{\pi}} \leq 1$

we say w is bad if bound fails

we say w is bad if bound fails

$$\text{bad}(w) = \mathbb{1} \left\{ L(w) > \hat{L}_n(w) + \underbrace{\frac{1}{\sqrt{n}} \sum_{i=1}^n \Omega_n(w; \delta)}_{\epsilon_n} \right\}$$

using Chernoff, $\hat{L}_n(w) \leq L(w) - \epsilon_n$

$$\begin{aligned} \mathbb{P}\{\text{bad}(w)\} &\leq \exp(-2n\epsilon_n^2) = \exp\left(-2n \frac{1}{2n} (\ln 2) |w|_{\pi} + \ln \frac{1}{\delta}\right) \\ &= \delta 2^{-|w|_{\pi}} \end{aligned}$$

using union bound

$$\mathbb{P}\{\exists w: \text{bad}(w)\} \leq \sum_w \mathbb{P}\{\text{bad}(w)\} \leq \sum_w \delta 2^{-|w|_{\pi}} \leq \delta //$$

example of complexity: $|w|_{\pi} = \|w\|_2^2 \rightarrow$ get ℓ_2 -regularization?

Surrogate loss:

Not hard to minimize $\hat{L}_n(w)$; replace with $\hat{J}_n(w)$ which is "surrogate" e.g. hinge loss