

# Lecture 16 - scribbles - PAC-Bayes

Friday, March 10, 2017  
14:58

theory: PAC-Bayes  $\rightarrow$  probit loss

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

$$\text{(usually): } h_w(x) = \arg \max_{y \in \mathcal{Y}} s(x, y; w)$$

$$\langle w, \phi(x, y) \rangle \quad (\text{linear score})$$

theory for structured prediction:

Last lecture: Linked  $\hat{L}_n(w)$  with  $L(w)$

uniformly over all  $w \in \mathcal{W}$  (countable)

using complexity  $|w|_{\pi}$  prior

PAC-Bayes: generalize this to

- arbitrary  $\mathcal{W}$
- general  $\ell(y, y') \in [0, 1]$

by using randomized predictor  
i.e. instead of  $\hat{w}$   $y = h_w(x)$

consider  $\hat{q}$  distribution over  $\mathcal{W}$

predict, first  $w \sim \hat{q}(w)$ ;  $y = h_w(x)$

then we work with  $\mathbb{E}_q[L(w)]$  as the generalization  
of this process  
 $\mathbb{E}_q[\hat{L}_n(w)]$  empirical version

PAC-Bayes thm. [McAllester 1999, 2003]

$(\ell(y, y') \in [0, 1])$  for any fixed prior  $\pi$  over  $\mathcal{W}$   
and any dist.  $P$  on  $\mathcal{X} \times \mathcal{Y}$

then with prob  $\geq 1 - \delta$  over  $D_n \sim P^n$  it holds

aspects of structured prediction

1) constraints on  $\mathcal{Y}$  ] theory open?

2)  $\ell(y, y')$

3) structured score functions:

$$S(x, y; w) = \sum_C S_C(x, y_C; w)$$

$$L_P(w) = \mathbb{E}_{(x, y) \sim P} [\ell(y, h_w(x))] \quad \text{generalization error V-risk}$$

$(L(w) \rightarrow P \text{ from context})$

$$\hat{L}_n(w) = \hat{L}_{D_n}(w) \triangleq \frac{1}{n} \sum_{(x^{(i)}, y^{(i)}) \in D_n} \ell(y^{(i)}, h_w(x^{(i)}))$$

$$\forall \text{ distributions } q, \quad \mathbb{E}_q[L(\omega)] \leq \mathbb{E}_q[\hat{L}_n(\omega)] + \frac{1}{\sqrt{2(n-1)}} \sqrt{KL(q||\pi) + \ln \frac{n}{\delta}}$$

if  $\omega$  is countable; let  $q(\omega) = \mathbb{I}\{\omega = \omega_0\}$

$$\text{then } KL(q||\pi) = \sum_{\omega} q(\omega) \ln \frac{q(\omega)}{\pi(\omega)} = \ln \frac{1}{\pi(\omega_0)} = \ln 2 / \pi(\omega_0)$$

Probit loss for structured prediction → NIPS 2011 McAllester & Keshet

if  $q(\omega') = N(\omega' | \omega, I)$

$$\text{then } \mathbb{E}_q[L(\omega)] = \mathbb{E}_{\omega \sim q(\omega)} [\mathbb{E}_{(x,y) \sim \pi} \ell(y, h_{\omega}(x))] = \mathbb{E}_{(x,y) \sim \pi} \left[ \underbrace{\mathbb{E}_{\epsilon \sim N(0,I)} \ell(y, h_{\omega+\epsilon}(x))}_{\text{probit}(x,y;\omega)} \right]$$

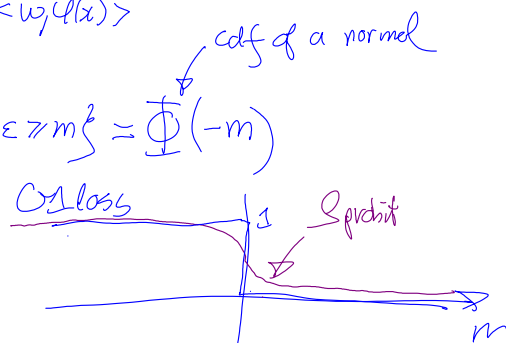
name probit: binary classification  $\mathcal{Y} = \{-1, +1\}$   
0-1 loss

$$\ell(x,y) = \frac{1}{2} (1 - y \phi(x))$$

$$h_{\omega}(x) = \text{sgn}(\langle \omega, \phi(x) \rangle)$$

let margin  
 $m = y \langle \omega, \phi(x) \rangle$

$$\text{then } \text{probit}(x,y;\omega) = \mathbb{P}_{\epsilon \sim N(0,1)} \{ \epsilon \geq m \} = \Phi(-m)$$



McAllester 2011 uses Catoni's PAC-Bayes version:  
for fixed  $\frac{1}{2}$   
(...)

$$\forall q \quad \mathbb{E}_q[L(\omega)] \leq \left( \frac{1}{1 - \frac{1}{2\alpha}} \right) \left( \mathbb{E}_q[\hat{L}_n(\omega)] + \frac{1}{n} (KL(q||\pi) + \ln \frac{1}{\delta}) \right)$$

if we use  $\pi = N(0, I)$   
 $q = N(\omega, I)$

↓  
, /

$$S_{\text{prbit}}(w) \leq \left( \frac{1}{1 - \frac{1}{2\lambda_n}} \right) \left( S_{\text{prbit}}(w) + \frac{\lambda_n}{n} \left( \frac{1}{2} \|w\|^2 + \ln \frac{1}{\delta} \right) \right) \quad \forall w$$

define  $\hat{w}_n = \underset{w \in W}{\text{argmin}} \left( \hat{S}_{\text{prbit}}(w) + \frac{\lambda_n}{2n} \|w\|^2 \right)$  ( $S_{\text{prbit}}(w) \triangleq \mathbb{E}_{(x,y) \sim p} [S_{\text{prbit}}(x,y,w)]$ )

thm 1: Let  $\lambda_n \nearrow \infty$  slowly enough so that  $\frac{\lambda_n}{n/\ln n} \rightarrow 0$

then  $S_{\text{prbit}}(\hat{w}_n) \xrightarrow{\text{a.s.}} L^* = \min_{w \in W} L(w)$

McAllester calls this 'consistency'

[Lacoste-Sulén unpublished fix: if  $L(w)$  is cts,

then  $L(\hat{w}_n) \xrightarrow{\text{a.s.}} L^* = \min_{w \in W} L(w)$

consistency even in the misspecified setting

well-specified  $\rightarrow$  let  $h^* = \underset{h \in \mathcal{H}}{\text{argmin}} L(h)$ ;

$\hookrightarrow \exists \tilde{w} \in W$  s.t.  $L(h\tilde{w}) = L(h^*)$

$$L(\hat{w}_n) - L(h^*) = \underbrace{L(\hat{w}_n) - L(w^*)}_{\text{estimation 'regret'}} + \underbrace{L(w^*) - L(h^*)}_{\text{approximation 'regret'}}$$

proof idea: use Catoni's PAC-Bayes bound

$$S_{\text{prbit}}(\hat{w}_n) \leq \left( \frac{1}{1 - \frac{1}{2\lambda_n}} \right) \left( S_{\text{prbit}}(\hat{w}_n) + \frac{\lambda_n}{n} \left( \frac{1}{2} \|\hat{w}_n\|^2 + \ln \frac{1}{\delta_n} \right) \right) \quad \text{with prob} \geq 1 - \delta_n$$

