

# Lecture 4 - scribbles - PAC Bayes

Friday, February 2, 2018

14:31

today: concentration inequality  
PAC-Bayes

## Concentration inequality

Markov inequality: R.V.  $X \geq 0$ ,  $a > 0$   
 (ie.  $X(\omega) \geq 0$   
 $\forall \omega \in \Omega$ )

$$\mathbb{P}\{X \geq a\} \leq \frac{\mathbb{E}X}{a}$$

proof:

$$a \mathbb{1}\{X \geq a\} \leq X$$

$$\mathbb{E}[ \quad ] \leq \mathbb{E}X$$

$$a \mathbb{P}\{X \geq a\} \leq \mathbb{E}X$$

Corollary: if  $\varphi: \mathbb{R} \rightarrow \mathbb{R}^+$  is an increasing fct.  $x \geq a \Rightarrow \varphi(x) \geq \varphi(a)$

then  $\varphi(x) \geq \varphi(a) \Leftrightarrow x \geq a$

$$\mathbb{P}\{X \geq a\} \leq \mathbb{P}\{\varphi(X) \geq \varphi(a)\} \leq \frac{\mathbb{E}(\varphi(X))}{\varphi(a)}$$

is true even if  $X$  is negative

Chernoff trick: use  $\varphi(x) \triangleq e^{tx}$  for  $t > 0$

thus  $\mathbb{P}\{X \geq a\} \leq \frac{\mathbb{E}e^{tX}}{e^{ta}} \forall t > 0$  by Markov

ie.  $\mathbb{P}\{X \geq a\} \leq \inf_{t > 0} \frac{\mathbb{E}e^{tX}}{e^{ta}}$

is not Chernoff bound

$$\hat{L}_n(\omega) = \prod_{i=1}^n \mathbb{1}\{y_i \neq h_{\omega}(x_i)\}$$

$\triangleq \text{Bin} \sim \text{Bernoulli}(L(\omega))$

$$L(\omega) = \mathbb{E} \mathbb{1}\{y_i \neq h_{\omega}(x_i)\}$$

$$= \mathbb{P}\{y_i \neq h_{\omega}(x_i)\} = P$$

want to bound the prob that:

$$L - \hat{L}_n \geq \epsilon$$

so get Chernoff bound,

$$\text{let } X \triangleq L - \hat{L}_n = \frac{1}{n} \sum_{i=1}^n (p - B_i) \quad \mathbb{E} e^{tX} = \mathbb{E} \exp\left(\frac{t}{n} \sum_{i=1}^n (p - B_i)\right)$$

[Chernoff bound on wikipedia](#)

independence of  $B_i$ 's

$$= \prod_i \mathbb{E} \exp\left(\frac{t}{n} (p - B_i)\right)$$

Hoeffding's Lemma:

for r.v.  $Y$  s.t.  $a \leq Y \leq b$  and  $\mathbb{E}Y = 0$

$$\text{then } \mathbb{E} e^{tY} \leq \exp\left(\frac{t^2}{8} (b-a)^2\right)$$

(use Taylor's expansion of  $\exp$ .)

[wikipedia proof](#)

$$\text{let } Y = \frac{p - B_i}{n} \in \left[\frac{p-1}{n}, \frac{p}{n}\right] \Rightarrow (b-a) = \frac{1}{n}$$

$$\mathbb{E} e^{tY} \leq \exp\left(\frac{t^2}{8n^2}\right)$$

$$\mathbb{P}\{L - \hat{L}_n \geq \epsilon\} \leq \frac{\left(\mathbb{E} \exp(tY)\right)^n}{e^{t\epsilon}} \leq \exp\left(\frac{t^2}{8n} - t\epsilon\right) \quad \forall t > 0$$

min with respect to  $t \Rightarrow t^* = 4n\epsilon$

$$\text{i.e. } \mathbb{P}\{L - \hat{L}_n \geq \epsilon\} \leq \exp(-2n\epsilon^2)$$

## PAC-bayes

Occam's bound  $\rightarrow$  we linked  $\hat{L}_n(w)$  with  $L(w)$

uniformly over all  $w \in W$  (countable)

using complexity  $|w|_{\text{pr}}^*$  ("prior")

PAC-bayes: generalize this to

- arbitrary  $W$
- general  $l(y, y') \in [0, 1]$

by using a randomized predictor

ie. instead of  $\hat{w}$   $y = h_{\hat{w}}(x)$

considers  $\hat{q}$  distribution over  $W$

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

use  $\mathbb{E}_{\hat{q}}[L(w)]$  as the generalization error for  $\hat{q}$  ie.  $\mathbb{E}_{(x,y) \sim p} \mathbb{E}_{w \sim \hat{q}} l(y, h_w(x))$

$\mathbb{E}_{\hat{q}}[\hat{L}_n(w)]$

↓ empirical version  
 ↘ structured prediction  
 this will yield probit surrogate loss (see soon)

PAC-Bayes Thm. [McAllester 1999, 2003]

$(l(y, y') \in [0, 1])$  for any fixed prior  $\pi$  over  $W$   
 and any distribution  $p$  on  $X \times Y$

then with  $\geq 1 - \delta$  over  $D_n \sim p^{\otimes n}$

it holds that  $\forall$  distribution  $q$   $\mathbb{E}_q[L_p(w)] \leq \mathbb{E}_q[\hat{L}_n(w)] + \frac{L}{\sqrt{2(n-1)}} \sqrt{KL(q||\pi) + \ln \frac{n}{\delta}}$

note: if  $W$  is countable; let  $q_{w_0} = \mathbb{1}\{w=w_0\}$

then  $KL(q_{w_0} || \pi) = \sum_w q(w) \ln \frac{q(w)}{\pi(w)} = \ln \frac{1}{\pi(w_0)} = (Q_n \mathcal{E}) / w_0$

probit loss for structured prediction: (NIPS 2011 McAllester & Keshet)

$$\text{if } q(w) \cong N(w | w, I)$$

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

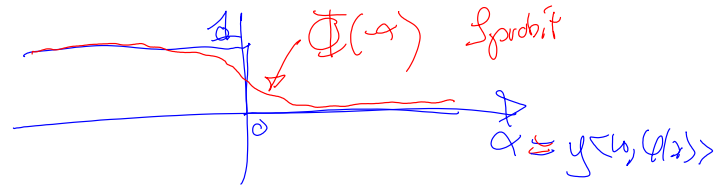
name probit: binary classification  $y \in \{-1, +1\}$   
0-1 loss

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

let margin  $\alpha = y \langle w, \phi(x) \rangle$

$$\text{then } \text{probit}(x, y; w) = \mathbb{E}_{\epsilon \sim N(0, I)} \mathbb{1}\{y \neq h_{w+\epsilon}(x)\} \quad \text{cdf of a Gaussian}$$

$$= \mathbb{P}_{\epsilon \sim N(0, I)} \{\epsilon \geq \alpha\} = \Phi(-\alpha)$$



McAllester still uses Cohen's PAC-Bayes version:

$$\forall q, \mathbb{E}_q[L(w)] \leq \left( \frac{1}{1 - \frac{1}{2\lambda_n}} \right) \left[ \mathbb{E}_q[\hat{L}_n(w)] + \frac{1}{\lambda_n} \left( \underbrace{KL(q || \pi)}_{\frac{1}{2} \|w\|^2} + \ln \frac{1}{\delta} \right) \right]$$

if we use  $\pi = N(0, I)$   
 $q_w = N(w, I)$  } define  $\hat{w}_n^{(\text{probit})} = \underset{w \in W}{\text{argmin}} \hat{L}_{\text{probit}}(w) + \frac{\lambda_n}{2n} \|w\|^2$

thm. 1: let  $\lambda_n \rightarrow \infty$  slowly enough so that  $\frac{\lambda_n}{n \ln n} \rightarrow 0$

$$\text{then } \text{probit}(\hat{w}_n) \xrightarrow{\text{a.s.}} L^* = \min_w L(w)$$

McAllester calls this "consistency"

for consistency would be  $L(\hat{w}_n) \xrightarrow{a.s.} L^*$

(Lacoste-Surken unpublished fix: if  $L(w)$  is cts<sub>0</sub>  
 then  $\mathbb{E} \text{profit}(\hat{w}_n) \xrightarrow{a.s.} L^*$   
 $\Rightarrow L(\hat{w}_n) \xrightarrow{a.s.} L^*$ )

proof idea: use Catoni's PAC-Bayes bound

$$\begin{aligned} \mathbb{E} \text{profit}(\hat{w}_n) &\leq \frac{1}{1 - \frac{1}{2\lambda n}} \left( \underbrace{\mathbb{E} \text{profit}(\hat{w}_n) + \frac{\lambda n \|\hat{w}_n\|^2}{2n}}_{\text{pick } \alpha} + \ln \frac{1}{\delta_n} \right) \quad \text{with prob. } 1 - \delta_n \\ &\leq \mathbb{E} \text{profit}(\alpha w^*) + \frac{\lambda n^2 \|w^*\|^2}{2n} \\ &\leq \mathbb{E} \text{profit}(\alpha w^*) + \sqrt{\frac{\ln n}{n}} \quad \text{using Chernoff bound for } \alpha w^* \end{aligned}$$

use  $\lim_{\alpha \rightarrow 1} \mathbb{E} \text{profit}(\alpha w^*) \leq L(w^*)$

$\therefore \lim_{n \rightarrow \infty} \mathbb{E} \text{profit}(\hat{w}_n) = L(w^*) //$

(see paper for details)