

Lecture 5 - PAC-Bayes

Thursday, January 28, 2021 13:33

- today:
- PAC-Bayes
 - probit loss
 - review surrogate loss

PAC-Bayes:

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

uniformly over all $\omega \in W \rightarrow$ but
(countable)

using complexity $\|w\|_p$
"prior"

PAC-Bayes \Rightarrow generalize this to

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

caveat: switch to a randomized predictor

i.e. instead of learning \hat{w} , predicting $y = h_{\hat{w}}(x)$

consider \hat{q} distribution over W

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

\Rightarrow use $\mathbb{E}_{\hat{q}}[L(w)]$ as the generalization error for \hat{q}

{ i.e. $\mathbb{E}_{(x,y) \sim P} \mathbb{E}_{w \sim \hat{q}}[l(y, h_w(x))]$

\downarrow empirical version

$\mathbb{E}_{\hat{q}}[\hat{L}_n(\omega)] \rightsquigarrow$ structured prediction
on will yield probit surrogate loss (see soon)
optimize over q

PAC-Bayes theorem [McAllester 1999, 2003]

(let $l(y, y') \in [0, 1]$) for any fixed prior π over W

and any dist. P over $X \times Y$

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

it holds that $\forall q$ dist. over W

$$\mathbb{E}_{\hat{q}}[l_p(\omega)] \leq \mathbb{E}_{\hat{q}}[\hat{L}_n(\omega)]$$

it holds that $\forall q \text{ dist. over } W$

$$\mathbb{E}_q [L_p(\omega)] \leq \mathbb{E}_q [L_n(\omega)]$$

$$+ \frac{1}{\sqrt{2(n-1)}} \sqrt{\text{KL}(q||\pi) + \ln \frac{n}{8}}$$

Note: if W is countable; let $Q_{W_0} = \{w_i\}_{w=w_0}$ new complexity term

$$\text{then } \text{KL}(q||\pi) = \sum_w q(w) \log \frac{q(w)}{\pi(w)} = \log \frac{1}{\pi(w_0)} = (\ln 2) \|w_0\|_\pi$$

Probit loss for structured prediction [NIPS 2011 McAllester & Keshet]

$$\text{if } q_w(w) \triangleq N(w|w, I)$$

$$\begin{aligned} \text{then } \mathbb{E}_{Q_W} [L(w^*)] &= \mathbb{E}_{w \sim N(w^*, I)} \mathbb{E}_{(x,y) \sim P} [L(y, h_{w^*}(x))] \\ &= \mathbb{E}_{(x,y) \sim P} \left[\mathbb{E}_{\varepsilon \sim N(0, I)} [L(y, h_{w+\varepsilon}(x))] \right] \\ &\boxed{\text{Sprob} (x, y; w)} \end{aligned}$$

Why name probit?

binary classi.: $\mathcal{Y} = \{-1, +1\}$ with 0-1 loss

$$h_{w^*}(x) = \text{sgn}(\langle w^*, \varphi(x) \rangle)$$

Let margin
 $\alpha = y \langle w, \varphi(x) \rangle$

$$\text{then } \text{Sprob} (x, y; w) = \mathbb{E}_{\varepsilon \sim N(0, I)} \mathbb{1}\{\varepsilon \neq h_{w+\varepsilon}(x)\}$$

$$y \langle w + \varepsilon, \varphi(x) \rangle < 0$$

$$\underbrace{y \langle w, \varphi(x) \rangle}_{\alpha} < -y \langle \varepsilon, \varphi(x) \rangle$$

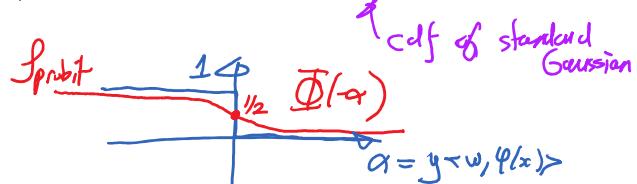
$$-\alpha > \langle \varepsilon, \varphi(x) \rangle$$

(suppose $|\varphi(x)| = 1$)

same prob. as

$$-\alpha > \langle \varepsilon, \varphi(x) \rangle$$

$$\text{Sprob} = P\{\varepsilon < -\alpha\} = \Phi(-\alpha)$$



⊕ define $\hat{w}_n^{(\text{probit})} = \arg \min_{w \in W} \text{Sprob} (w) + \frac{\lambda_n}{2n} \|w\|^2$

$$\textcircled{2} \text{ define } \hat{w}_n^{(\text{prob})} = \underset{w \in W}{\operatorname{arg\min}} \text{ Sprobbit}(w) + \frac{\lambda_n}{2n} \|w\|^2 \quad (*)$$

McAllester showed the consistency of $\hat{w}_n^{(\text{prob})}$

14h23

McAllester 2011 uses Caboni's PAC-Bayes version

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

↓
 $\hat{L}_n(w)$ ↓
 $\frac{\lambda_n}{n}$ $\|w\|^2$
 if we use $\pi = N(0, I)$
 $q_w = N(w, I)$

↑
 $\text{Sprobbit}(w)$

Motivates $\hat{w}_n^{(\text{prob})}$

thm. 1

in paper : let $\lambda_n \gg n$ slowly enough so that $\frac{\lambda_n}{n} \rightarrow 0$

McAllester
calls this
"consistency"

then $\text{Sprobbit}(\hat{w}_n) \xrightarrow{n \rightarrow \infty} L^* = \min_{w \in W} L(w)$

but true consistency would be $L(\hat{w}_n) \xrightarrow{n \rightarrow \infty} L^*$

[Lacoste-Julien unpublished result :
if $L(w)$ is ctg.]

then $\text{Sprobbit}(\hat{w}_n) \xrightarrow{n \rightarrow \infty} L^*$

$\Rightarrow L(\hat{w}_n) \xrightarrow{n \rightarrow \infty} L^* = L(w^*)$

proof idea : use Caboni's PAC Bayes bound

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

↑
 $\text{Sprobbit}(\alpha w^*) + \frac{\lambda_n}{2n} \alpha^2 \|w^*\|^2$

by def. of \hat{w}_n

pick α

$$\leq \text{Sprobbit}(\alpha w^*) + \sqrt{\frac{\lambda_n}{n}} \text{ using Chernoff bound for } w$$

* also use $\lim_{\alpha \rightarrow \infty} \text{Sprobbit}(\alpha w^*) \leq L(w^*)$

$$\lim_{\|\omega\| \rightarrow 0} \text{Sprobbit}(\hat{\omega}) = L/\|\omega\| \quad [\text{see paper for details}]$$

problem: $\text{Sprobbit}(x, y; \omega)$ is non-convex in $\omega \Rightarrow$ no optimization guarantees

how: convex surrogates $s(\tilde{y}) \triangleq s(x, \tilde{y}; \omega)$ i.e. $x \nparallel \omega$ is implicit

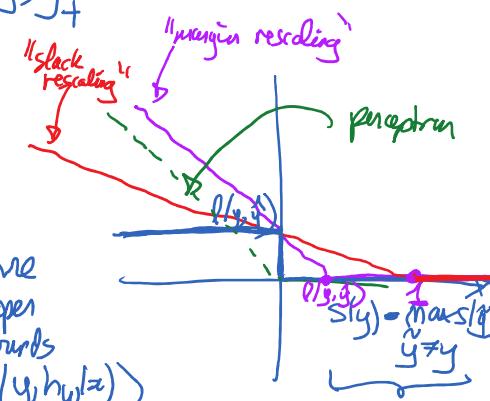
Review of convex surrogates mentioned so far:

$$\text{Sperceptron}(x, y; \omega) = \max_{\tilde{y} \in \mathcal{Y}} s(\tilde{y}) - s(y) \quad \text{let } m(\tilde{y}) = s(\tilde{y}) - s(y)$$

$$= \max_{\tilde{y} \in \mathcal{Y}} [-m(\tilde{y})] = \max_{\tilde{y} \neq y} -m(\tilde{y})_+$$

hinge (structured sum)
"margin rescaling"

$$= \max_{\tilde{y}} [s(\tilde{y}) + l(y, \tilde{y})] - s(y)$$



"margin rescaling"
vs.
"slack rescaling"

$$= \max_{\tilde{y}} [l(y, \tilde{y}) - m(\tilde{y})]$$

$$= \max_{\tilde{y}} l(y, \tilde{y}) [1 - m(\tilde{y})]$$

one upper bounds
 $l(y, h_w(x))$

$$\mathcal{L}_{\text{CRF}}(\cdot) = \frac{1}{\beta} \log \left(\sum_{\tilde{y}} \exp(\beta s(\tilde{y})) \right) - s(y) \quad [-\log p_{\omega}(y|x)]$$

$$\frac{1}{\beta} \log \left(\sum_{\tilde{y}} \exp(-\beta m(\tilde{y})) \right)$$

suggests
"smooth hinge",

$$\frac{1}{\beta} \log \left(\sum_{\tilde{y}} \exp(\beta [s(\tilde{y}) - m(\tilde{y})]) \right)$$

[e.g. Pletscher et al. 2010]

note: slack rescaling more robust when have small $l(y, \tilde{y})$ [e.g. 0]

but more computationally costly

what theoretical properties could we look at?

a) generalization error bounds [next class]

b) consistency properties & calibration fact [next 2 classes]
 \hookrightarrow relationship between $L(\omega) \triangleq \mathcal{L}(\omega)$

why structural score functions?

$$s(x, y) = \sum_{c \in C} s_c(x, y_c)$$

Motivation similar to graphical models

- 1) statistical efficiency : less # of parameters (simpler score fn. Σ)
 \Rightarrow easier to learn [see (Orlits & al. NIPS 2006 next class)
(generalization guarantees)
- 2) computational " : compute cugmate $s(\vec{y})$