

Lecture 5 - scribbles - PAC-Bayes

Tuesday, January 22, 2019 14:33

today: PAC-Bayes
probit Loss
review Surrogate Losses

PAC-Bayes

Occam's bound \rightarrow we linked $\hat{L}_n(w)$ with $L_p(w)$
uniformly over all $w \in W$ (countable)
using complexity $I_w q$
 q ("prior")

PAC-Bayes: generalizes this to:
• arbitrary W
• general $l(y, y') \in [0, 1]$

Caveat: switch to a randomized predictor

i.e. instead of \hat{w} $y = \hat{h}_{\hat{w}}(x)$
consider q distribution over W

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

use $E_{\hat{q}}[\hat{L}(w)]$ as the generalization error for \hat{q} i.e. $E_{(x,y) \sim p} E_{w \sim \hat{q}} l(y, h_w(x))$
 \hat{q} empirical version

$$\mathbb{E}_q [L_{\pi}(w)]$$

structured prediction,
this will yield probabilistic surrogate loss (see soon)

PAC-Bayes Thm. [McAllester 1999, 2003]

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

and any dist. p on $X \times Y$

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

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

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

new complexity term.

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

probabilistic loss for structured prediction: [NIPS 2011 McAllester & Keshet]

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

$$\text{then } \mathbb{E}_{q_w} [L(w)] = \mathbb{E}_{w \sim q_w} \mathbb{E}_{(x,y) \sim p} [\ell(y, h_w(x))]$$

$$= \mathbb{E}_{(x_1, y_1), \dots, (x_n, y_n) \sim p} [\ell(y_i, h_{w_i}(x_i))]$$

$$= \mathbb{E}_{(x,y) \sim P} \left[\mathbb{E}_{\epsilon \sim N(0,1)} l(y, h_{w+\epsilon}(x)) \right]$$

\triangleq

$$\text{Sprob}(x, y; w)$$

why name probit?

binary classification $Y = \{-1, +1\}$
with $C=1$ loss

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

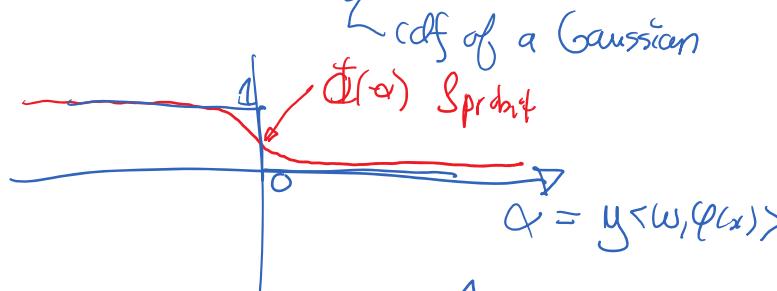
$$\text{Then } \text{Sprob}(x, y; w) = \mathbb{E}_{\epsilon \sim N(0,1)} \mathbb{I}\{\epsilon \neq h_{w+\epsilon}(x)\}$$

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

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

$$(\text{supposing } \|\varphi(x)\| = 1) \Rightarrow \alpha < -y \langle \epsilon, \varphi(x) \rangle$$

$$= P\{\epsilon \geq \alpha\} = \Phi(-\alpha)$$



McAllester showed the consistency of the \hat{w}^{probit} which minimizes a ℓ_2 -reg. Sprob

McAllester 2011 uses Catoni's PAC-Bayes version:

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

$$\text{if we use } \pi = N(0, I) \quad \Rightarrow \quad \hat{S}_{\text{prob}}(w) = \frac{1}{2} \|w\|^2$$

$$q_w = N(w, I) \quad \text{define } \hat{w}_n^{(\text{prob})} = \underset{w \in W}{\operatorname{arg\min}} \hat{S}_{\text{prob}}(w) + \frac{\lambda_n}{n} \|w\|^2$$

Thm 1

in paper: Let $\lambda_n \nearrow \infty$ slowly enough so that $\frac{\lambda_n}{n \ln n} \rightarrow 0$

McAllester calls this "consistency"

$$\text{then } \hat{S}_{\text{prob}}(\hat{w}_n) \xrightarrow[n \rightarrow \infty]{a.s.} L^* = \min_{w \in W} L_p(w)$$

but

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

(Lacoste-Jutten unpublished fix : if $L(w)$ is cts.

$$\text{then } \hat{S}_{\text{prob}}(\hat{w}_n) \xrightarrow{a.s.} L^* \\ \Rightarrow L(\hat{w}_n) \xrightarrow{a.s.} L^*$$

proof idea! use Catoni's PAC-Bayes bound

with prob $\geq 1 - \delta_n$

$$\begin{aligned} \hat{S}_{\text{prob}}(\hat{w}_n) &\leq \left(\frac{1}{1-\frac{1}{2\lambda_n}}\right) \left(\hat{S}_{\text{prob}}(\hat{w}_n) + \frac{\lambda_n}{n} \|\hat{w}_n\|^2 + \ln \frac{1}{\delta_n} \right) \\ &\leq \hat{S}_{\text{prob}}(g_w) + \lambda_n \alpha^2 \|w\|^2 \end{aligned}$$

$$\leq \hat{S}_{\text{probit}}(\alpha w^*) + \frac{\lambda \alpha^2 \|w^*\|^2}{n}$$

$$\leq S_{\text{probit}}(\alpha w^*) + \sqrt{\frac{\lambda n}{n}} \quad \text{using Chernoff bound for } \alpha w^*$$

use $\lim_{\alpha \rightarrow \infty} \hat{S}_{\text{probit}}(\alpha w^*) \leq L(w^*)$

$\therefore \lim_{n \rightarrow \infty} \hat{S}_{\text{probit}}(\hat{w}_n) = L(w^*) \quad [\text{see paper for details}]$

15/13)

problem: $\hat{S}_{\text{probit}}(x, y; w)$ is non-convex \Rightarrow no optimization guarantee

now: convex surrogates $s(\tilde{y}) \triangleq s(x, \tilde{y}; w)$ ie. $x \not\perp w$ are implicit

Review of convex surrogates mentioned so far:

$$S_{\text{perceptron}}(x, y; w) = \max_{\tilde{y} \in \mathcal{Y}} s(\tilde{y}) - s(y)$$

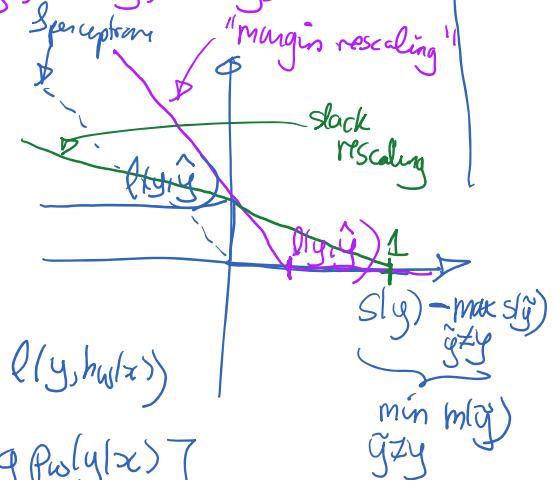
$$= \max_{\tilde{y} \in \mathcal{Y}} [-m(\tilde{y})] = [\max_{\tilde{y} \in \mathcal{Y}} m(\tilde{y})]_+$$

$$\text{f hinge} (__) = \max_{\tilde{y} \in \mathcal{Y}} [s(\tilde{y}) + \ell(y, \tilde{y})] - s(y)$$

$$\text{"margin rescaling"} \quad \text{vs} \quad \text{f hinge} (__) = \max_{\tilde{y} \in \mathcal{Y}} [s(\tilde{y}) - m(\tilde{y})]$$

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

$$\text{let } m(\tilde{y}) \triangleq s(y) - s(\tilde{y})$$



$$S_{\text{CRF}}(_) = \lfloor \log(\frac{1}{2} \exp(B \delta(_))) \rfloor - s(_), \quad \text{f-log Prob}(u|x)$$

$$\text{S-CRF}(\underline{\quad}) = \frac{1}{\beta} \log \left(\sum_{\tilde{y}} \exp(\beta s(\tilde{y})) \right) - s(y) \quad [-\log p_{\theta}(y|x)]$$

(Log-loss for CRF)

$$\beta \rightarrow 0 \Rightarrow \text{perceptron loss} \\ \frac{1}{\beta} \log \left(\sum_{\tilde{y}} \exp(\beta m(\tilde{y})) \right)$$

"smoothed hinge loss" \rightarrow

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

[e.g. Pletscher et al. 2010]

$\min_{\tilde{y}} m(\tilde{y})$

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

what theoretical properties could we look at?

- generalization error bounds [today]
- consistency properties & calibration [next class]

why structured score functions? $s(x, y) = \sum_{c \in C} s_c(x, y_c)$

motivations similar to graphical models

1) statistical efficiency: less # of parameters (simpler score functions s_c)

\Rightarrow easier to learn [see Cartes et al. NIPS 2016]
(generalization guarantees) [beg. of next class]

2) computational: compute argmax _{$\tilde{y} \in \mathcal{Y}$} $s(\tilde{y})$

BUT compare to what happens for Hamming Loss:

BUT compare to what happens for Hamming Loss:

given true conditional $Q_x(y) \triangleq p(y|x)$ generating data

expected error when using \tilde{y} as prediction is $\mathbb{E}_{y \sim Q_x(y)} [l(y, \tilde{y})] \triangleq l_{Q_x}(\tilde{y})$

$$\text{for Hamming Loss: } l_{Q_x}(\tilde{y}) = \mathbb{E}_{Q_x(y)} \left[\sum_p \underbrace{\mathbb{I}\{y_p \neq \tilde{y}_p\}}_{1 - \mathbb{P}\{y_p = \tilde{y}_p\}} \right] = \sum_p (1 - Q_x(\tilde{y}_p))$$

Marginal on y_p
 i.e. $\sum_{y: y_p = \tilde{y}_p} Q_x(y)$

\Rightarrow best decision $\tilde{y}^* = \arg \min_{\tilde{y} \in \mathcal{Y}} l_x(\tilde{y})$

is just independent predictions $\tilde{y}_p = \arg \max_{y_p} p(\tilde{y}_p|x)$

max marginal decoding

marginal of
 $p(y|x)$

thus; If there is no constraints, then can just train independently models for each part marginal $p(y_p|x)$

i.e. $s_p(y_p|x; w_p)$

e.g. $p(y_p|x) \propto \exp(s_p(y_p; w))$

but a) this function might be too complicated

b) statistically, could be beneficial to share learning together "transfer learning" between parts