

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

PAC-Bayes :

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

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

using complexity $|w|_\pi$
 \uparrow "prior"

PAC-Bayes : generalizes this to

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

caveat : switch to a randomized predictor

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

consider \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. $\mathbb{E}_{(x,y) \sim p} \mathbb{E}_{w \sim \hat{q}} l(y, h_w(x))$
 empirical version

$\mathbb{E}_{\hat{q}} [L_n(w)] \rightsquigarrow$ structural prediction
 this will yield probit surrogate loss (see soon)

PAC-Bayes Thm [McClellan 1999, 2003]

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

and any dist p on $X \times Y$

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

it holds that \forall dist. q on W $\mathbb{E}_q [L_p(w)] \leq \mathbb{E}_q [L_n(w)] +$

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

+ ...

non-constant

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

new complexity term

$$\text{then } k(L(q|\pi)) = \sum_w q(w) \ln \frac{q(w)}{\pi(w)} = \ln \frac{1}{\pi(w_0)} = (\ln 2) |w|$$

probit loss for structured prediction [NIPS 2011 McAllester & Keshtor]

if $q_w(w') \triangleq N(w'|w, I)$

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

$$= \mathbb{E}_{(x,y) \sim p} \left[\mathbb{E}_{\epsilon \sim N(0, I)} [l(y, h_{w+\epsilon}(x))] \right] \quad \text{where } \epsilon \in N(0, I)$$

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

why name probit?

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

let margin

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

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

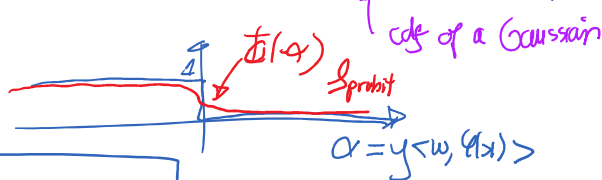
$$\text{then } \text{probit}(x, y; w) = \mathbb{E}_{\epsilon \sim N(0, I)} \mathbb{1}\{y \neq h_{w+\epsilon}(x)\}$$

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

(supposing $\|\phi(x)\|=1$)

$$\frac{y \langle w, \phi(x) \rangle}{\alpha} < -\frac{y \langle \epsilon, \phi(x) \rangle}{\alpha}$$

$$\text{probit} = \mathbb{P}\{\sum \epsilon_i < -\alpha\} = \Phi(-\alpha) \quad \text{where } \Phi \text{ is the CDF of a Gaussian}$$



$$\text{define } \hat{w}_n^{(\text{probit})} = \underset{w \in W}{\text{argmin}} \left[\hat{\text{probit}}(w) + \frac{\lambda_n}{2n} \|w\|^2 \right]$$

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

McAllester 2011 vs Caloni's PAC-Bayes version:

$$\forall q, \mathbb{E}_q [L(w)] \leq \frac{1}{2} \left[\mathbb{E}_q [L_n(w)] + \lambda_n [k(L(q|\pi)) + \ln \frac{1}{\pi(w_0)}] \right]$$

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

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

minimizes $\hat{w}_n^{(\text{probit})}$

Thm 1

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

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

McAllester call this "consistency"

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

[Lacoste-Julian unpublished fixed: if $L(w)$ is cts.]

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

proof idea: use Cramer's PAC Bayes bound

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

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

$$\leq \hat{S}_{\text{probit}}(\alpha w^*) + \sqrt{\frac{2 \ln n}{n}} \quad \text{using Chernoff bound for } \alpha w^*$$

$$\times \text{ also use } \lim_{\alpha \rightarrow 0} \hat{S}_{\text{probit}}(\alpha w^*) \leq L(w^*)$$

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

15h37

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

now: convex surrogates $s(\tilde{y}) \triangleq s(x, \tilde{y}; w)$ i.e. x & w are implicit

Review of convex surrogates mentioned so far:

$\mathcal{L}_{\text{perceptron}}(x, y; w) = \max_{\tilde{y} \neq y} s(\tilde{y}) - s(y)$
let $m(\tilde{y}) \triangleq s(y) - s(\tilde{y})$

$= \max_{\tilde{y} \neq y} [-m(\tilde{y})] = \left[\max_{\tilde{y} \neq y} -m(\tilde{y}) \right]_+$

$\mathcal{L}_{\text{hinge}}(\text{structured sum}) = \max_{\tilde{y}} [s(\tilde{y}) + \ell(y, \tilde{y})] - s(y)$

"margin rescaling" $\rightarrow \max_{\tilde{y}} [\ell(y, \tilde{y}) - m(\tilde{y})]$

"slack rescaling" $\rightarrow \max_{\tilde{y}} \ell(y, \tilde{y}) [1 - m(\tilde{y})]$

"are with upper bounds $\ell(y, \tilde{y})$ "

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

$\frac{1}{\beta} \log \left(\sum_{\tilde{y}} \exp(\beta m(\tilde{y})) \right)$ $\xrightarrow{\text{suggested "smoothed hinge loss"}}$ $\frac{1}{\beta} \log \left(\sum_{\tilde{y}} \exp(\beta [\ell(y, \tilde{y}) - m(\tilde{y})]) \right)$

[e.g. Plefcher et al. 2010]

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

what theoretical properties could we look at?

- a) generalization error bounds [next class]
- b) consistency properties & calibration fct. [next class]
 - \hookrightarrow relationship between $L(w)$ & $\mathcal{L}(w)$

why structured score functions?

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

motivations similar to graphical models

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

\Rightarrow easier to learn (generalization guarantees)

[see Cortes & el. NIPS 2016] next class

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

but consider what happens for Hamming loss

$$y = (y_1, \dots, y_p, \dots, y_L)$$

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

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

$$\text{for Hamming loss: } l_{q_x}(\tilde{y}) = \mathbb{E}_{q_x(y)} \left[\sum_p \frac{\mathbb{1}\{y_p \neq \tilde{y}_p\}}{1 - \mathbb{1}\{y_p = \tilde{y}_p\}} \right] = \sum_p (1 - q_x(\tilde{y}_p))$$

\Rightarrow best decision $y^* = \operatorname{argmin}_y l_{q_x}(y)$

is just independent predictions

marginal on \tilde{y}_p

$$\text{ie. } \sum_{y_p} q_x(y_p)$$

Max marginal decoding	$y_p^* = \operatorname{argmax}_{y_p} p(\tilde{y}_p x)$
	↑ marginal of $p(y x)$

⊛ thus: if there is no constraints, then can just train independently models for each

part marginal $p(y_p|x)$

ie. $S_p(y_p|x; w_p)$

e.g. $p(y_p|x)$ or $\exp(S_p(x))$

but a) this function might be too complicated

b) statistically, could be inferred to share learning together

"transfer learning" between parts