

- Today:
- PAC-Bayes
  - probit loss
  - review of surrogate loss

PAC-Bayes :

Occam's band  $\rightarrow$  we linked  $\hat{L}_n(w)$  with  $L_p(w)$   
uniformly over all  $w \in W \rightarrow$  but countable  
 using complexity  $|w|_{\pi}$  "prior"

PAC-Bayes: generalize this to

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

concret: switch to a randomized predictor

ie. instead of learning  $\hat{w}$ , predicting  $y = h_{\hat{w}}(x)$   
 consider  $\hat{q}$  distribution over  $W$   
 predict: first  $w \sim \hat{q}(w)$ ;  $y = h_w(x)$

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

ie.  $\mathbb{E}_{(x,y) \sim p} \mathbb{E}_{w \sim \hat{q}} [\ell(y, h_w(x))]$   
 $\rightarrow$  empirical version

$\mathbb{E}_{\hat{q}} [\hat{L}_n(w)] \rightarrow$  an structured prediction will yield probit surrogate loss (see soon)  
 optimize over  $q$  to get  $\hat{q}$

PAC-Bayes theorem [McAllester 1999, 2003]

(let  $\ell(y, y') \in [0, 1]$ ) for any fixed prior  $\pi$  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  $W$

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

new complexity term

note: if  $W$  is countable: let  $a_{w_i} = \mathbb{1}_{\{w=w_i\}}$

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

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

probit loss for structured prediction [NeurIPS 2011, McAllester & Keshet]

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

then  $\mathbb{E}_{q_w} [L(w')] = \mathbb{E}_{w' \sim q_w} \mathbb{E}_{(x,y) \sim P} [l(y, h_{w'}(x))]$   
 $w' = w + \epsilon$  where  $\epsilon \sim N(0, I)$   
 $= \mathbb{E}_{(x,y) \sim P} [ \mathbb{E}_{\epsilon \sim N(0, I)} [l(y, h_{w+\epsilon}(x))] ]$   
 $\triangleq \mathcal{L}_{probit}(x, y; w)$

why name probit?

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

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

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

then  $\mathcal{L}_{probit}(x, y; w) = \mathbb{E}_{\epsilon \sim N(0, I)} \mathbb{1}_{\{y \neq h_{w+\epsilon}(x)\}}$   
 this equals one, when:

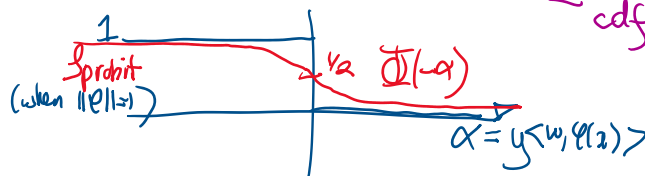
$y \langle w + \epsilon, \phi(x) \rangle < 0$

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

$-\alpha > \langle \epsilon, y \phi(x) \rangle$   
 Gaussian

RHS is a Gaussian with mean  $\mathbb{E} \langle \epsilon, y \phi(x) \rangle = \langle \mathbb{E} \epsilon, y \phi(x) \rangle = 0$   
 $\mathbb{E} [y^2 \phi^T \epsilon \epsilon^T \phi] = \phi^T \mathbb{E} [\epsilon \epsilon^T] \phi = \|\phi\|^2$   
 $\rightarrow N(0, \|\phi(x)\|^2)$

$\mathcal{L}_{probit} = P \{ \epsilon_1 \|\phi\| < -\alpha \} = \Phi \left( \frac{-\alpha}{\|\phi(x)\|} \right)$



$\hat{w}_n^{(probit)} = \arg \min_{w \in \mathcal{W}} \mathcal{L}_{probit}(w) + \frac{\lambda_n}{2n} \|w\|^2$

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

McAllester will use Catoni's PAC-Bayes thm. version

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

if we use  $\pi = N(0, I)$   
 $q_w = N(w, I)$

$$\frac{\hat{J}_{\text{prdbit}}(w) + \frac{\lambda_n}{n} \frac{\|w\|^2}{2}}{(*)}$$

motivates  $\hat{w}_n^{(\text{prdbit})}$

10/24

Thm. 1 in paper: Let  $\lambda_n \rightarrow 0$  slowly enough so that  $\frac{\lambda_n \ln n}{n} \rightarrow 0$   
 then  $\hat{J}_{\text{prdbit}}(\hat{w}_n) \xrightarrow[n \rightarrow \infty]{\text{a.s.}} L^* = \min_{w \in W} L(w)$

McAllester calls this "consistency"

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

[Lacoste-Julien unpublished result: ]

if  $L(w)$  is obs., then  $\hat{J}_{\text{prdbit}}(\hat{w}_n) \xrightarrow{\text{a.s.}} L^* \Rightarrow L(\hat{w}_n) \xrightarrow{\text{a.s.}} L^* = L(w^*)$

proof idea: use Catoni's PAC Bayes bound

with prob  $\geq 1 - \delta_n$

$$\hat{J}_{\text{prdbit}}(\hat{w}_n) \leq \left( \frac{1}{1-\delta_n} \right) \left( \hat{J}_{\text{prdbit}}(\hat{w}_n) + \frac{\lambda_n \|\hat{w}_n\|^2}{2n} + \frac{\lambda_n \ln \frac{1}{\delta_n}}{n} \right)$$

$$\leq \hat{J}_{\text{prdbit}}(\alpha w^*) + \frac{\lambda_n \alpha^2 \|w^*\|^2}{2n}$$

$$\leq \hat{J}_{\text{prdbit}}(\alpha w^*) + \sqrt{\frac{\lambda_n n}{n}}$$

Using Chernoff bound for  $\alpha w^*$  with prob  $\geq 1 - \frac{\delta_n}{n^2}$  (set  $\delta_n = \frac{1}{n^2}$ )

$\Rightarrow$  get  $\lim_{n \rightarrow \infty} \hat{J}_{\text{prdbit}}(\hat{w}_n) \leq \hat{J}_{\text{prdbit}}(\alpha w^*)$  with prob 1

$\otimes$  also use  $\lim_{\alpha \rightarrow 0} \hat{J}_{\text{prdbit}}(\alpha w^*) \leq L(w^*)$

$\Rightarrow \lim_{n \rightarrow \infty} \hat{J}_{\text{prdbit}}(\hat{w}_n) = L(w^*)$  [see paper for details]

problem:  $\hat{J}_{\text{prdbit}}(z, y; w)$  is non-convex in  $w \Rightarrow$  no optimization guarantees

now: convex surrogates on score

notation:  $S(\hat{y}) \triangleq S(x, y; w)$  i.e.  $x \neq w$  is implicit

Reminded convex-surrogate mentioned in Gen

now: convex surrogates on score

$$S(\tilde{y}) \equiv S(x, y^* | w) \text{ i.e. } x \ \& \ w \text{ is implicit}$$

Review of convex surrogates mentioned so far

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

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

assuming  $y \in \mathcal{Y}$

$$S_{\text{hinge}}(\quad) = \max_{\tilde{y}} [S(\tilde{y}) + \ell(y, \tilde{y})] - S(y)$$

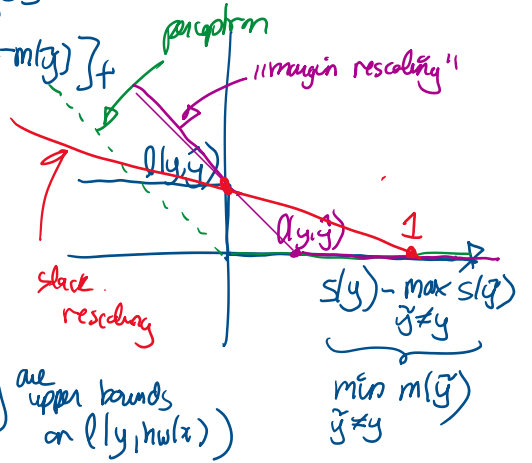
(structured SVM)

"margin rescaling"

vs. "slack rescaling"

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

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



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

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

suggests "smoothed hinge loss"

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

[e.g. Plefischer & 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 2 class]
  - ↳ relationship between  $L(w)$  &  $S(w)$

why structured score functions?

$$S(x, y) = \sum_{C \in \mathcal{C}} S_C(x, y_C)$$

motivation similar to graphical models

- 1) statistical efficiency: less # of parameters (simpler score fct.  $S_C$ )
  - ⇒ easier to learn [see Cortes & al. NIPS 2006] next class
  - (generalization guarantees)
- 2) computational "": compute  $\arg \max_{\tilde{y} \in \mathcal{Y}} S(\tilde{y})$