

today:

- finish Occam's bound
- PAC Bayes
- probit loss

Occam's bound

for any fixed P ; with prob. $\geq 1 - \delta$ over training set $D_n \sim P^{\otimes n}$

$$\forall w \in W \quad L_P(w) \leq \hat{L}_P(w) + \frac{1}{\sqrt{2n}} \mathcal{D}_{\pi}(w; \delta)$$

$$\text{where } \mathcal{D}_{\pi}(w; \delta) \triangleq \sqrt{(\ln 2)} \underbrace{\|w\|_{\pi} + b_1}_{8} \underbrace{\delta}_{\text{Complexity measure}}$$

* bound is useful only for dist. P s.t. $\|w\|_{\pi}$ is small
 \hookrightarrow argmin $L_P(w)$
 $w \in W$

example of learning alg: $\hat{w}_n = \underset{w \in W}{\operatorname{argmin}}$ RHS of bound
 \Rightarrow universally consistent!

$$\|w\|_{\pi} = \log \sum_{\pi(w)} \quad \text{if } \pi(w) \propto \exp(-\|w\|^2) \quad \rightarrow \text{RHS is } l_2 \text{-regularized ERM}$$

then $\|w\|_{\pi} = \|w\|^2 + \text{const.}$

proof: use 3 things:

1) Chernoff bound
 (concentration inequality)

$$P\{\exists D_n : \hat{L}_n(w) \leq L(w) - \varepsilon\} \leq \exp(-2n\varepsilon^2) \quad \forall \varepsilon \geq 0$$

2) union bound

$$P\{\exists x \text{ s.t. prop}(x) \text{ is true}\} \leq \sum_x P\{\text{prop}(x) \text{ is true}\}$$

$$\sum_w 2^{-l(w)} \leq 1$$

3) "kraft's ineq."

note: 0-1 loss assumption
 appears in constant of Chernoff

we say that w is haughty if bound fails

$$\text{bad}(w) = \mathbb{1}\{\hat{L}_n(w) > L(w) + \frac{1}{\sqrt{2n}} \mathcal{D}_{\pi}(w; \delta)\}$$

$$L - \varepsilon > \hat{L}_n$$

$$\varepsilon_n(w)$$

using Chernoff, $\hat{L}_n(w) \leq L(w) - \varepsilon_n(w)$ with small prob...

$$P\{\text{bad}(w)\} \leq \exp(-2n\varepsilon_n(w)^2) = \exp\left(-2n \frac{1}{2n} (\ln 2) \|w\|_{\pi} + \ln \frac{1}{\delta}\right)$$

$$\propto \sim \|w\|_{\pi}$$

$$\int_2 \exp(-\alpha \ln(w)) = \exp\left(-\alpha \frac{1}{w} (\ln(w) + \ln(\frac{1}{w}))\right)$$

$$= \delta 2^{-\ln(w)}$$

using union bound:

$$\Pr_{\omega} \{ \text{bad}(w) \} \leq \sum_{\omega} \Pr_{\omega} \{ \text{bad}(w) \} \stackrel{\text{kraft's inequality}}{\leq} \delta 2^{-\ln(w)} \leq \delta //$$

Surrogate loss: NP hard to minimize $\hat{L}_n(w)$; replace with $\hat{S}_n(w)$ which is "surrogate"

e.g. hinge loss
log loss

Next: countable \rightarrow uncountable
"PAC-Bayes"

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 $\|w\|_{\pi, \text{"prior"}}$

PAC-Bayes: generalize this to

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

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

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

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

idea: use empirical version

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

PAC-Bayes thm. [McAllester 1999, 2003]

PAC-Bayes thm. [McAllester 1999, 2003]

(let $l(y, \hat{y}) \in [0, 1]$) ; for any fixed prior Π over W

and any dist. p over $X \times \mathcal{Y}$

Then with $\geq 1 - \delta$ over $D_n \sim p \otimes n$

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

$$\mathbb{E}_q[L_p(\omega)] \leq \mathbb{E}_q[\hat{L}_n(\omega)] + \frac{1}{\sqrt{2^{(n-1)}}} \underbrace{\text{KL}(q || \pi) + \ln \frac{n}{8}}_{\uparrow}$$

new complexity term

Note: if ω is countable; let $g_{\omega_0} = \mathbb{1}_{\{\omega = \omega_0\}}$

$$\text{Then } \text{KL}(q || \pi) = \sum_w q(w) \log \frac{q(w)}{\pi(w)} = \log \frac{1}{\pi(w)} = (\log 2) |w| / \pi$$

15h33

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

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

$$\text{then } \mathbb{E}_{q_w} [L(w)] = \mathbb{E}_{w \sim q_w} \mathbb{E}_{(x,y) \sim p} [l(y, h_w(x))] \quad w' = w + \varepsilon \text{ where } \varepsilon \sim N(0, I)$$

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

\downarrow

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

Why called probit?

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

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

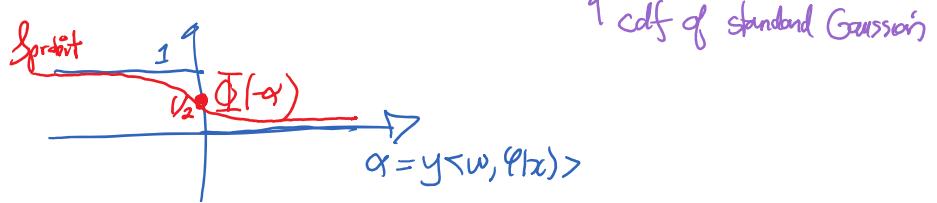
then $\text{Sprob}(\mathbf{z}, y; \omega) = \mathbb{E}_{\varepsilon \sim N(0, I)} \left[\mathbb{1}_{\{\mathbf{y} \neq h_{\omega+\varepsilon}(\mathbf{z})\}} \right]$
 i.e. $y < \omega + \varepsilon, \varrho(\mathbf{z}) < 0$

$$\underbrace{y \in w, \varphi(z)}_{\vartheta} \leftarrow y \varepsilon, \varphi(x)$$

$$1 \text{ when } -\alpha > \langle \varepsilon_y, \varphi(x) \rangle$$

$$N(0, \|\varphi(\tau)\|^2) \quad (\text{assume } \|\varphi(\tau)\| = 1)$$

$$\text{Sprobbit} = \Pr\{\varepsilon_1 < -\alpha\} = \Phi(-\alpha)$$



④ define $\hat{w}_n^{(\text{probbit})} = \underset{w \in W}{\operatorname{arg\min}} \text{Sprobbit}(w) + \frac{\lambda_n}{2n} \|w\|^2$ (*)

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

McAllester 2011 uses Catoni's PAC-Bayes' version

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

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

↓
 $\text{Sprobbit}(w) + \frac{\lambda_n}{n} \frac{\|w\|^2}{2}$
 Motivates $\hat{w}_n^{(\text{probbit})}$ (*)

thm 1
in paper Let $\lambda_n \nearrow \infty$ slowly enough so that $\frac{\lambda_n \ln n}{n} \rightarrow 0$

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

McAllester
call this
"consistency"

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

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

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

$$\Rightarrow L(\hat{w}_n) \xrightarrow{\text{a.s.}} L^* - L(w^*)$$

proof idea: use Catoni's PAC Bayes bound

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

set $\delta_n = \frac{1}{n^2}$

$\xrightarrow{\text{by def. of } \hat{w}_n} \text{Sprobbit}(w^*) + \frac{\lambda_n}{2n} \alpha^2 \|w^*\|^2$

$\mathcal{L}_{\text{probit}}(\alpha w^*) + \sqrt{\frac{O_n}{n}}$ using Chernoff bound for αw^*

with prob $\geq 1 - \frac{1}{n^2}$

get $\lim_{n \rightarrow \infty} \mathcal{L}_{\text{probit}}(\hat{w}_n) \leq \mathcal{L}_{\text{probit}}(\alpha w^*)$ with prob 1

* also show that $\lim_{\alpha \rightarrow \infty} \mathcal{L}_{\text{probit}}(\alpha w^*) \leq L(w^*)$ [see paper for details]

problem: $\mathcal{L}_{\text{probit}}(x, y; w)$ is non-convex in $w \Rightarrow$ no optimization guarantees

now: convex surrogates \mathcal{L}