

Lecture 5 - scribbles - generalization error bounds

Tuesday, February 6, 2018
14:33

today: • surrogate losses
• generalization error bounds

last time: $\mathcal{L}_{\text{profit}}(x, y; w) = \mathbb{E}_{\mathcal{E} \sim \mathcal{P}(0, \pi)} [\ell(y, h_{w+\mathcal{E}}(x))]$ $h_w(x) = \operatorname{argmax}_{y \in \mathcal{Y}} s(x, y; w)$

issue: non-convex in $w \Rightarrow$ no optimization guarantee

today: convex surrogate $\mathcal{L}(\tilde{y}) \triangleq \mathcal{L}(x, \tilde{y}; w)$ i.e. x & w are implicit

recap:

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

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

$$\mathcal{L}_{\text{hinge}}(x, y; w) = \max_{\tilde{y} \in \mathcal{Y}} [s(\tilde{y}) + \ell(y, \tilde{y})] - s(y)$$

(structured sum)
"margin rescaling"

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

"slack rescaling"

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

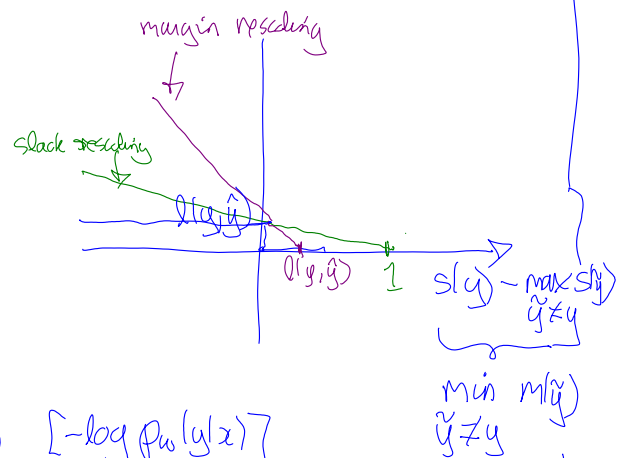
\rightarrow they are both upper bounds on $\ell(y, h_w(x))$

log-loss (CRF) (soft-max) $\mathcal{L}_{\text{CRF}}(x, y; w) = \left[\frac{1}{\beta} \log \left(\sum_{\tilde{y}} \exp(\beta s(\tilde{y})) \right) \right] - s(y)$

$\beta \rightarrow \infty \Rightarrow$ perceptron loss

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

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



$$[-\log p_w(y|z)]$$

"smoothed hinge loss"

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

[e.g. Pletcher & al. 2010]

what are the...? \rightarrow

what are theoretical properties?

- a) generalization error bounds [today]
- b) consistency properties & calibration [cf. [next class]]

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 [see later today
(generalization guarantees) Cortes & El NIPS 2016]

2) computational || : compute $\arg\max_{\tilde{y} \in \mathcal{Y}} s(\tilde{y})$

BUT compare to what happens for Hamming loss:

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

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

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

for Hamming loss:

$$l_{q_x}(\tilde{y}) = \mathbb{E}_{q_x(y)} \left[\sum_p \underbrace{1\{y_p \neq \tilde{y}_p\}}_{1 - 1\{y_p = \tilde{y}_p\}} \right] = \sum_p (1 - q_x(\tilde{y}_p)) \stackrel{\text{marginal of } p(y|x)}{=} \sum_{y: y_p = \tilde{y}_p} q_x(y)$$

\Rightarrow best decision $y^* = \arg\min_{\tilde{y} \in \mathcal{Y}} l_x(\tilde{y})$ is just $y_p^* = \arg\max_{\tilde{y}_p} p(\tilde{y}_p | x)$

'marginal decoding'

* thus, if no constraint, i.e. no need "consistency" between parts,

could just train independently models for each part marginal $p(y|x)$

i.e. $S_p(y|x; w_p)$

model $p(y|x) \propto \exp(S_p(y|x; w))$

but a) this function might be too complicated

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

Generalization error bounds

for binary classification, a classical PAC bound is
for any fixed distribution on data
with prob $\geq 1-\delta$ over D_n

$$\forall w \in W, L_0(w) \leq \hat{L}_n(w) + \frac{1}{\sqrt{n}} \sqrt{d \log \frac{d}{n} + \log \frac{2}{\delta}}$$

where d is VC-dimension

of $\mathcal{H} = \{h_w : w \in W\}$

VC-dimension of $\mathcal{H} \equiv \max \left\{ m : \begin{array}{l} \text{exists a set of } m \text{ points s.t.} \\ \text{for any labeling of } m \text{ points,} \\ \exists w \text{ s.t. } h_w \text{ gives the correct} \\ \text{labels on those points} \end{array} \right\}$

"scattering the set of points"

also means that # of prediction functions on fixed m points is 2^m

for linear classifiers of p parameters, VC-dim = $p+1$

* one issue for this bound is that true for all distributions \Rightarrow too loose bound

\Rightarrow motivates going to data distribution dependent measure of complexity

example: empirical Rademacher complexity

$$\hat{R}_n(\mathcal{H}) \triangleq \mathbb{E}_{\sigma} \left[\sup_{h \in \mathcal{H}} \left| \frac{2}{n} \sum_{i=1}^n \sigma_i \mathbb{1}\{y_i \neq h(x_i)\} \right| \right]$$

'correlations with random noise'

$\sigma_i = \sum_{j=1}^{+1} -1$ uniformly "Rademacher" R.V.

$$\forall w \quad L_0(w) \leq \hat{L}_n(w) + \hat{R}_n(\mathcal{H}) + \frac{1}{\sqrt{n}} 3 \sqrt{\frac{\log 2/s}{2}}$$

complexity depends on \mathcal{D}_n (implicitly on p)

high level idea to prove bound:

"double sample trick" \rightarrow use a second sample \mathcal{D}_n' for generalization error $L(w) = \mathbb{E}_{\mathcal{D}_n'} [\hat{L}_n(w)]$

+

"symmetrization trick" \rightarrow bound the sup of differences between $L(w)$ & $\hat{L}_n(w)$

+

union bound as usual + concentration inequality

structured prediction generalization bounds [Kortis & al. NIPS 2016]

general loss $\ell(y, y')$ s.t. $\ell(y, y') \neq 0$ if $y \neq y'$

$$\text{suppose } S(x, y) = \sum_{c \in \mathcal{C}} S_c(x, y_c)$$

\rightarrow set of cliques of a graphical model / factor graph

thm. 7

with prob. $\geq 1-\delta$

$$\forall w \in \mathcal{W} \quad L(w) \leq \text{Shinge}(w) + 4\sqrt{2} \hat{R}_n(\mathcal{H}_w) + 3L_{\max} \sqrt{\frac{\log 1/\delta}{2n}}$$

depends on $\ell(y, y')$

(truncated)

$\max_{y, y'} \ell(y, y')$

$$\forall w \in W \quad L(w) \leq \text{Shinge}(w) + 4\sqrt{2} \hat{R}_n^G(Hw) + 3L_{\max} \sqrt{\frac{\log 1/s}{2n}}$$

where $\hat{R}_n^G \triangleq \frac{1}{n} \mathbb{E}_{\sigma} \left[\sup_{w \in W} \sum_{i=1}^n \sqrt{|\phi_i|} \sum_{c \in \mathcal{C}} y_c \in \mathcal{Y}_c \sigma_{i,c,y_c} s_c(x_i, y_c; w) \right]$

↓ "empirical factor graph complexity"

only uses $(x_i^{(i)})_{i=1}^n$

set of labels for y_c

Rademacher R.V.

Thm. 2: if $s_c(x_i, y_c; w) = \langle w, \phi_c(x_i, y_c) \rangle$

and consider $W_{\Lambda} \triangleq \{w : \|w\|_2 \leq \Lambda\}$; let $R = \max_{i,c,y} \|\phi_c(x_i, y_c)\|_2$

then $\hat{R}_n^G(Hw_{\Lambda}) \leq \frac{R \Lambda}{\sqrt{n}} |\phi| \sqrt{\max_c |\mathcal{Y}_c|}$

so want small degrees

$$L(w) \leq \text{Shinge}(w) + \underbrace{\left(\frac{R |\phi| \sqrt{\max_c |\mathcal{Y}_c|}}{\sqrt{n}} \right)}_{\lambda_n} \underbrace{\|w\|_2}_{\frac{\|w\|_2^2}{2}}$$

min of RHS gives SVMstruct $\hat{w}_n = \arg \min_w \text{Shinge}(w) + \lambda_n \frac{\|w\|_2^2}{2}$
(if f is convex)

missing links: (1) $\min_{\|w\| \leq \Lambda} f(w)$ $\leadsto \exists \lambda(\Lambda)$ s.t. $\min f(w) + \lambda \frac{\|w\|_2^2}{2}$ (2)
gives same solution as (1)

SVMstruct can be interpreted as minimizing upper bound on generalization error

properties: • minimize upper bound, hope that minimize $L(w)$

* can evaluate bound to get guarantees

caveat:

minimizing an upper bound is not same thing
as minimizing $L(w)$
(i.e. here, no consistency
guarantees)

next: consistency

