

Lecture 6 - generalization error bounds

Thursday, February 2, 2023 1:24 PM

- today:
 - review surrogate losses
 - generalization error bounds
 - structured SVM

Review of convex surrogate losses mentioned so far:

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

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

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

$$l_{\text{hinge}}(\quad) = \max_{\tilde{y} \in \mathcal{Y}} [s(\tilde{y}) + l(y, \tilde{y})] - s(y)$$

$$\text{"margin rescaling"} \rightarrow = \max_{\tilde{y} \in \mathcal{Y}} [l(y, \tilde{y}) - m(\tilde{y})]$$

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

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

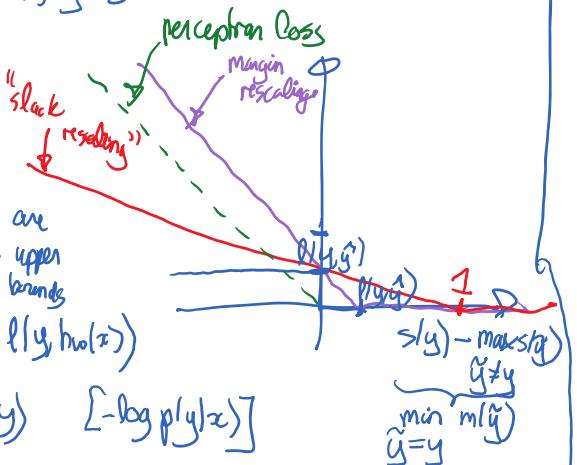
$\beta \rightarrow \infty \Rightarrow \text{perceptron loss}$

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

suggests
"smoothed
hinge loss"

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

[e.g. Flitscher et al. 2018]



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

but more computationally costly

What are theoretical properties could we look at?

a) generalization error bounds [today]

b) consistency properties { calibration function [next class]

↳ relationship between $L(\omega)$ & $\ell(\omega)$

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 : loss # of parameters (sample score fct. Σ)

\Rightarrow easier to learn
(generalization guarantees) [see Cortes et al. NeurIPS 2016]
(today)

2) computational " : compute argmax _{$g \in \mathcal{G}$} $S(g)$

generalization error bounds:

for binary classification

a classical PAC bound is :

for any fixed dist. p on data
with prob $\geq 1-\delta$ on D_n

$$\forall w \in W \quad L_{\text{err}}(w) \leq \hat{L}_n(w) + \frac{1}{\sqrt{n}} \sqrt{d \log d + \log \frac{1}{\delta}}$$

where d is VC-dimension of $H = \{h_w : w \in W\}$

VC-dimension of $H \triangleq \max \{m : \exists \text{ a set of } m \text{ points s.t.}$

\forall labelings of these points

$\exists w \in W$ that gives the

correct label on these points

"shattering the set of points"

\hookrightarrow # of binary functions on m points

$\hookrightarrow 2^m$
for $H = \{\text{linear classifiers of } p \text{ parameters}\} \quad \text{VC-dim}(H) = p+1$

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

\Rightarrow motivates going to data distribution dependent measure of complexity

example : empirical Rademacher complexity

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

"correlation with random noise"

$\sigma_i = \begin{cases} +1 & \text{uniformly "Rademacher R.V."} \\ -1 \end{cases}$

bound : with prob $\geq 1-\delta$

$$\forall w \quad L_{\text{err}}(w) \leq \hat{L}_n(w) + \hat{R}_{D_n}(H) + \frac{1}{\sqrt{n}} \sqrt{3 \log \frac{2}{\delta}}$$

$$\text{Hw } L_0(\omega) \leq \hat{L}_n(\omega) + \hat{R}_{D_n}^G(H) + \frac{1}{\sqrt{n}} 3 \sqrt{\log \frac{2}{\delta}}$$

complexity depends on D_n (implicitly on P)

| 4th BO

structured prediction generalization bounds [Cortes & al. NeurIPS 2016]

general loss fct. $l(y, y')$ s.t. $l(y, y') \geq 0 \quad \forall y, y'$

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

↳ set of cliques of a graph model G / factor graph

thm. 7 with prob $\geq 1 - \delta$

$$\text{Hw } L(\omega) \leq \hat{L}_{\text{hinge}}(\omega) + 4\sqrt{2} \hat{R}_{D_n}^G(H) + 3 \frac{\max_{y, y'} l(y, y')}{\sqrt{2n}} \sqrt{\log \frac{1}{\delta}}$$

$$\text{where } \hat{R}_{D_n}^G \triangleq \frac{1}{n} \mathbb{E}_\Omega \left[\sup_{\omega \in W} \sum_{i=1}^n \sqrt{|b_i|} \sum_{c \in C: y_c \in Y_c} \sum_{\tilde{o}_i, c, \tilde{y}_c} \delta_{\tilde{o}_i, c, \tilde{y}_c} S_c(x_i, \tilde{y}_c; \omega) \right]$$

actually
only depends on $(x^{(i)})_{i=1}^n$ "empirical factor graph complexity" under Rademacher R.V.

thm. 2 : If $S_c(x, y_c; \omega) = \langle \omega, \varphi_c(x, y_c) \rangle$

and consider $W_R \triangleq \{ \omega : \|\omega\|_2 \leq R \}$; let $R = \max_{i, c, y_c} \|\varphi_c(x_i, y_c)\|_2$

$$\text{then } \hat{R}_{D_n}^G(H_{W_R}) \leq \frac{R \cdot \sqrt{n} \max_i |b_i| \sqrt{\max_{c \in C} |c|}}{\sqrt{n}}$$

so want small cliques?

* plug thm. 2 in thm. 7 :

$$L(\omega) \leq \hat{L}_{\text{hinge}}(\omega) + \left(\frac{R |\mathcal{C}|}{\sqrt{n}} \sqrt{\max_{c \in C} |c|} \right) \frac{1}{\|\omega\|_2} + \text{const.}$$

min of RHS suggests

$$\text{SVM struct d.g. } \hat{w} = \arg \min_{\omega} \hat{L}_{\text{hinge}}(\omega) + \frac{\lambda n}{2} \|\omega\|_2^2$$

missing link : ① $\min_{\text{st. } \|\omega\|_2 \leq R} f(\omega)$ (if f is convex), use Lagrangian

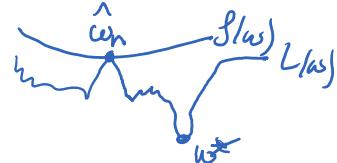
missing link: ① $\min f(w)$
s.t. $\|w\|_2 \leq L$

② $\min f(w) + \frac{\lambda}{2} \|w\|^2$

(if f is convex)
use Lagrangian
for a $\lambda(L)$ s.t.
reln to ② same
solution b ①

[sidenote: constrained formulation can have solutions not achievable for ② when f is non-convex
but penalized/reg. formulation is less sensitive to choice of λ vs. constrained formulation]

properties: - minimize upper bound, hope that $\min L(w)$
but no general guarantees



• Can evaluate bound to get guarantees

caveat:
also note here: no consistency guarantee

next: consistency + convex surrogate

consistency & calibration

need to relate $f(w)$ to $L(w)$: bad "calibration fact." [Steinwart]

relationship is usually very complicated

\Rightarrow usual results lack mainly at non-parametric setting (no # of parameters)

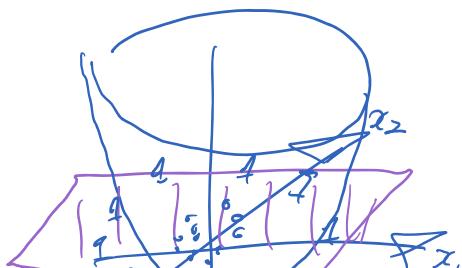
all functions $h: X \rightarrow \mathbb{R}$ are considered \Rightarrow this erases the dependence
on x of the analysis
(pointwise analysis)
i.e. we suppose that $s(x, y; w)$ can be arbitrary for any x
(i.e. w is ∞ -dim.)

\rightarrow can do this using a universal kernel

$$s(\cdot, \cdot; w) \in \mathcal{H}_{X \times \mathbb{R}}$$

RkHS:

motivation:



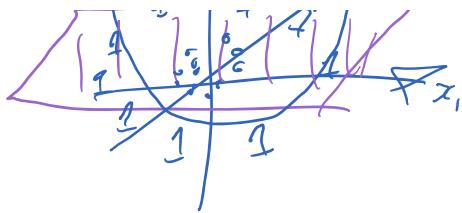
generalize linear structure

$\langle w, \varphi(x) \rangle$ to higher dim.

+ kernel trick $\langle \varphi(x), \varphi(x') \rangle = k(x, x')$

$$\Phi: X \rightarrow \mathbb{R}^3$$

$\xrightarrow{x-1} \quad \xrightarrow{x-2} \quad \xrightarrow{x}$



$$\phi : X \rightarrow \mathbb{R}^3$$

$$\phi(x) = \begin{pmatrix} x_1^2 \\ x_1 x_2 \\ x_2^2 \end{pmatrix}$$

$$\langle \phi(x), \phi(x') \rangle_{\mathbb{R}^3} = (\langle x, x' \rangle_{\mathbb{R}^2})^2 = k(x, x')$$

polynomial kernel e.g. $(\langle x, x' \rangle + 1)^p = k_{\text{poly}}(x, x')$

equivalent to mapping data to a space of dimension exponential in p

$$\langle \phi(x), \phi(x') \rangle$$

$$\text{even have to-dim e.g. } k(x, x') = \exp(-\frac{\|x - x'\|^2}{2})$$

"RBF kernel"