

Lecture 6 - scribbles - generalization error bounds

Thursday, January 24, 2019 13:36

today: · generalization error bounds
 & structured SVM

generalization error bounds

for binary classification,

a classical PAC bound is:

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

$$\forall w \in \mathcal{W} \quad L_{\text{PAC}}(w) \leq \hat{L}_n(w) + \frac{1}{\sqrt{n}} \sqrt{d \log d + \log \frac{2}{\delta}}$$

where d is VC-dimension

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

VC-dimension of $\mathcal{H} \triangleq \max \{ m : \exists \text{ a set of } m \text{ pts, } \}$

st. \forall labelings of these points,
 $\exists w$ st. h_w gives the
 correct label on these points
 "Shattering the set of pts"

of prediction functions
 on m pts. is 2^m

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

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

\Rightarrow motivates going to data distribution dependent

→ motivates going to data distribution dependent
measure of complexity

Example: empirical Rademacher complexity

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

"correlations
with random noise"

$\sigma_i = \sum_{j=1}^{+1}$ uniformly "Rademacher" RV

bound:

with prob $\geq 1-\delta$

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

complexity depends on D_n (implicitly on P)

high level idea to prove bound:

"double sample trick" → use a second sample D_n for gen. error. $L(w) = \mathbb{E}_{D_n} [\hat{L}_n(w)]$

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

+ union bound as usual + concentration inequality

Structured prediction generalization bounds [Cortes & al. NIPS 2016]

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

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

set of cliques of
a graph. model / factor graph

Thm. 7 with prob $\geq 1 - \delta$ depends on $\ell(y, y')$

$$\forall w \in W \quad L(w) \leq \hat{L}_{\text{hinge}}(w) + 4\sqrt{2} \hat{R}_{D_n}^G(\mathcal{H}_w) + 3 \frac{L_{\max}}{\sqrt{n}} \sqrt{\log \frac{1}{\delta}}$$

where $\hat{R}_{D_n}^G \triangleq \frac{1}{n} \mathbb{E}_0 \left[\sup_{w \in W} \sum_{i=1}^n \sqrt{16d_i} \sum_{c \in C_i} \sum_{y_c \in \mathcal{Y}_c} S_c(x_i, y_c; w) \right]$

only depends $(x_i^{(t)})_{i=1}^n$ "empirical factor graph complexity"

set of labels for y_c Ferdinand K.V.

Thm. 2: if $S_c(x, y_c; w) = \langle w, \varphi_c(x, y_c) \rangle$

and consider $W_L \triangleq \{w : \|w\|_2 \leq L\}$; let $R = \max_{c, \tilde{y}} \|\varphi_c(x_i, \tilde{y})\|_2$

then $\hat{R}_{D_n}^G(\mathcal{H}_{W_L}) \leq \frac{R}{\sqrt{n}} |C| \sqrt{\max_{c, \tilde{y}} 16d_i}$

so want small cliques ?

back to bound: $L(w) \leq \hat{L}_{\text{hinge}}(w) + \left(R \frac{|C| \sqrt{\max_{c, \tilde{y}} 16d_i}}{\sqrt{n}} \right) L$

back to bound: $L(w) \leq \text{Singe}(w) + \left(\underbrace{\frac{K \|g\|_1 + \max |f_i|\}}_{\Delta_n} \right) \|w\|_2$

min of RMS suggests

SVM struct alg. $\hat{w}_n = \arg \min \text{Singe}(w) + \Delta_n \|w\|^2$

missing link: (1) $\min f(w)$ st. $\|w\| \leq L$ (if f is convex) $\xrightarrow{\text{use Lagrangian duality}} \exists \alpha \Delta(L) \text{ s.t.}$

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

gives same solution as (1)

sideline: penalized formulation less sensitive b choice of Δ vs. constrained formulation

can think of SVM struct as minimizing upper bound on generalization error

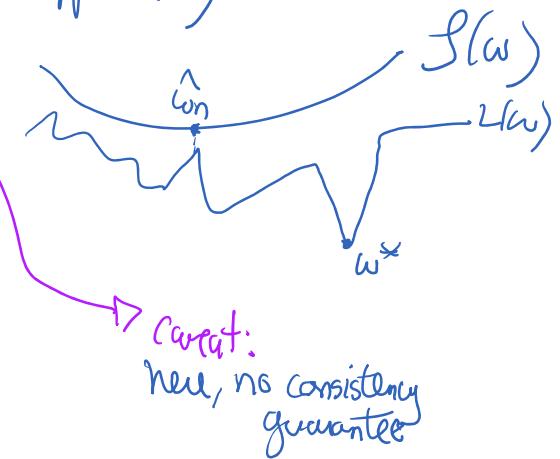
Properties:

- minimize upper bound, hope that minimizes app. $L(w)$

but no general guarantees?

- can evaluate bound to get guarantees

next: margin...



next: consistency
14th Mar

"guarantee"

consistency \nsubseteq calibration fct.:

Need to relate $\mathcal{L}(w)$ to $L(w)$ \Rightarrow tool "calibration function" [Steinwart]

relationship is usually very complicated

\Rightarrow current results look mainly at non-parametric setting (\propto # of parameters)

all functions $h: \mathcal{X} \rightarrow \mathcal{Y}$ are considered \Rightarrow this erases the dependence on x in the analysis

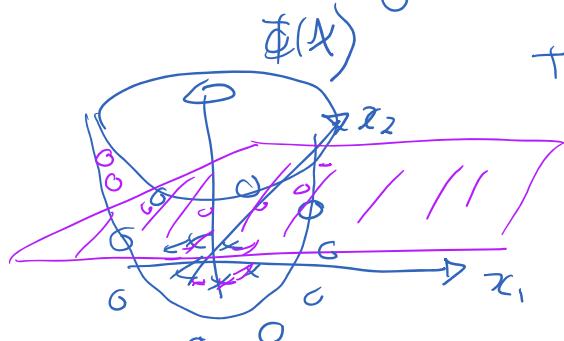
i.e. we suppose that $S(x, y; w)$ can be arbitrary for any x (i.e. w is ~~no dim~~) \rightarrow could use a universal kernel

$$S(\cdot, \cdot; w) \in \mathcal{H}_{\mathcal{X} \times \mathcal{Y}}$$

RKHS

motivation: generalize $\langle w, \varphi(x) \rangle$ to higher dim. space

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



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

$$\Phi(x) = \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{pmatrix}$$

ϕ

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

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

polynomial kernel e.g. $(\langle x, x' \rangle + 1)^d = \tilde{k}(x, x')$

equivalent to mapping data to space of dimension exponential on d

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

even have ∞ -dim., e.g. $k(x, x') = \exp(-\|x - x'\|^2)$ (RBF kernel)

RKHS (Reproducing Kernel Hilbert Space)

$\Phi: X \rightarrow \mathcal{H}$ s.t. $\langle \Phi(x), \Phi(x') \rangle_{\mathcal{H}} = k(x, x')$ (uniqueness property of RKHS)

Let $\tilde{\mathcal{H}} = \text{span}\{k(x, \cdot) : x \in X\}$

e.g. $f \in \tilde{\mathcal{H}} \Rightarrow f = \sum_i \alpha_i k(x_i, \cdot)$ for some finite $\{x_i\}_{i=1}^{n(f)}$
 $\alpha_i \in \mathbb{R}^{n(f)}$

↪ "pre-Hilbert" space [inner product space]

$$\text{with } \langle f, g \rangle_{\tilde{\mathcal{H}}} \triangleq \underbrace{\sum_{i,j} \alpha_i^f \alpha_j^g k(x_i^f, x_j^g)}_{\langle k(x_i^f, \cdot), k(x_j^g, \cdot) \rangle_{\tilde{\mathcal{H}}}} \quad \|f\|_{\tilde{\mathcal{H}}} \triangleq \sqrt{\langle f, f \rangle_{\tilde{\mathcal{H}}}}$$

Then RKHS \mathcal{H} is $= \text{completion}(\tilde{\mathcal{H}})$ using $\|\cdot\|_{\tilde{\mathcal{H}}}$ as your norm

Then RKHS \mathcal{H} is = completion($\tilde{\mathcal{H}}$) using $\|\cdot\|_{\mathcal{H}}$ as your norm
 i.e. add all limit points of $\tilde{\mathcal{H}}$ -Cauchy sequences to get \mathcal{H}
 you could think of $f = \sum_{i=1}^{\infty} \alpha_i k(x_i, \cdot)$

"reproducing" property of \mathcal{H} : for $f \in \mathcal{H}$

$$\langle f, k(x, \cdot) \rangle_{\mathcal{H}} = f(x)$$

nice property of RKHS, fct. evaluation is cb.

mapping $E_x: \mathcal{H} \rightarrow \mathbb{R}$
 $E_x(f) = f(x)$

$$|f(x) - g(x)| = |\langle f - g, k(x, \cdot) \rangle_{\mathcal{H}}|$$

≤ $\|f - g\|_{\mathcal{H}} \|k(x, \cdot)\|_{\mathcal{H}}$ i.e. E_x is Lipschitz c.f.
 with $L = \|k(x, \cdot)\|_{\mathcal{H}}$

★ this property
 important to do statistics