

Lecture 6 - generalization error bounds

Tuesday, February 2, 2021 14:32

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. $\geq 1-\delta$ on D_n

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

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

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

\exists W s.t. h_w gives the correct label on these points
"shattering the set of points"

of prediction functions on m pts
is 2^m

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

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

\Rightarrow motivates going to data distribution dependent measure of complexity

example: empirical Rademacher complexity

$$\hat{R}_{D_n}(H) \triangleq \mathbb{E}_P \left[\sup_{h_w \in H} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i \prod_{j \neq i} \mathbb{I}\{y_j \neq h_w(x_j)\} \right| \right]$$

"correlation with random noise"

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

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

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

complexity depends on D_n (implicitly on P)

complexity depends on D_3 (implicitly on P)

structured prediction generalization bounds [Cortes et al. NIPS 2016]

general loss $l(y, y')$ s.t. $l(y, y') \neq 0$ $\forall y \neq y'$

suppose $s(x,y) = \sum_{c \in \mathcal{C}} s_c(x, y_c)$

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

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

$$\text{true } L(w) \leq \hat{L}_{\text{hinge}}(w) + 4\sqrt{R_{D_n}^G(\hat{L}(w))} + 3L_{\max} \sqrt{\log \frac{1}{\delta}}$$

where $\hat{R}_{D_n}^G \triangleq \frac{1}{n} \mathbb{E}_G \left[\sup_{w \in W} \sum_{i=1}^n \sqrt{P_{\theta_i}} \sum_{c \in \mathcal{C}_i} \sum_{y_j \in \mathcal{Y}_j} S_c(x_i, y_j; w) \right]$

actually only depends on $(x_i^{(j)})_{i=1}^n$

"empirical feature graph complexity"

indep Rademacher R.V.

Thm 2: if $s_c(x, y_c; \omega) = \langle \omega, \varphi_c(x, y_c) \rangle$

and consider $W_n \triangleq \{w : \|w\|_2 \leq n\}$; let $R = \max_{i, c, y} \|(\ell_c(x_i, y))\|_2$

$$\text{then } \boxed{R_{Dn}^G(H_{Wn}) \leq \frac{R\Delta}{\sqrt{n}} \|f\| \sqrt{\max_c \|f_c\|}}$$

so want small clynes!

→ plug them back in them. 7

$$L(w) \leq \underbrace{J_{\text{change}}(w)}_{\frac{\partial J}{\partial w}} + \left(\underbrace{R|\ell| \sqrt{\max_{\zeta} |\zeta|^2}}_{\frac{\partial R}{\partial w}} \right) \underbrace{\|w\|_2}_{\text{||}w\text{||}_2} + \text{const.}$$

Min of RHS suggests

$$\text{SVM soft alg.} \quad \vec{w}_n = \arg \min_{\vec{w}} \frac{\text{Shage}(\vec{w}) + \lambda_n}{2} \|\vec{w}\|^2$$

missing link: (1) min $f(w)$ st. $\|w\|_2 \leq L$ (if f is convex), use Lagrangian duality + a $\lambda(L)$ s.t.

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

\leftarrow sol'n to (2) gives same
solution as (1)

$$\leftarrow \min_{w \in \mathbb{R}^n} \|w\|_2 + \frac{\lambda}{2} \|w\|^2$$

min to same solution as ①

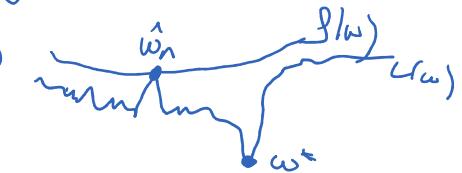
[sidenote: constrained formulation can have solutions not achievable for ② when f is non-convex
but penalized/reg.-formulation is less sensitive to choice of Λ vs. constrained formulation]

can think of SVM struct as minimizing an upper bound on gen. error

properties:

- minimize upper bound, hope that $\min L(w)$

but no general guarantees



• Can evaluate bound to get guarantees

next: consistency + convex surrogate

15h 40

☞ cautious:

also note here: no consistency guarantees

consistency & calibration

need to relate $L(w)$ to $\hat{L}(w)$ = tool "calibration fct." [Steinwart]

relationship is usually very complicated

⇒ current result look mainly at non-parametric setting (so # of parameters)

all functions $h: X \rightarrow \mathcal{Y}$ are considered ⇒ this evacuates the dependence on x of the analysis
 ↗ pointwise analysis!

i.e. we suppose that $s(x_i; w)$ can be arbitrary for any x
 (i.e. w is n -dim)

→ can do this using a universal kernel

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

BkHS:

motivation:

generalize linear structure

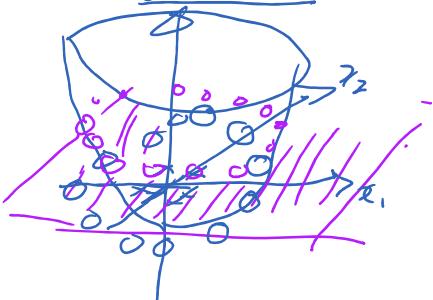
Illustration

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

→ kernel trick $\langle \phi(x), \phi(z) \rangle = k(x, z)$

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

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



~~so far~~

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

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

polynomial kernel e.g. $(\langle x, x' \rangle + 1)^p = k(x, x')$

equivalent to mapping data to a space of dimension exponential in p , $\langle \Phi(x), \Phi(x') \rangle$

$$\text{even have RBF, } k(x, x') = \exp\left(-\frac{\|x-x'\|_2^2}{2}\right)$$

"RBF kernel"

RkHS (reproducing kernel Hilbert space)

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

\mathcal{H} is a space of fct. $X \rightarrow \mathbb{R}$

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

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

"pre-Hilbert" space [inner product space]

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

$$\|f\|_{\tilde{\mathcal{H}}} \triangleq \sqrt{\langle f, f \rangle_{\tilde{\mathcal{H}}}}$$

$$\langle k(x_i^f, \cdot), k(x_j^g, \cdot) \rangle_{\tilde{\mathcal{H}}}$$

α_i^f means α_i in $f = \sum_i \alpha_i k(x_i, \cdot)$

then RkHS \mathcal{H} is \cong completion of $\tilde{\mathcal{H}}$ using $\|\cdot\|_{\tilde{\mathcal{H}}}$ as norm

i.e. add all limit points of $\tilde{\mathcal{H}}$ -Cauchy sequences to get \mathcal{H}

you could think $f = \sum_{i=1}^{\infty} \alpha_i k(x_i, \cdot)$



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

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

$$\Phi(x)$$

Nice property of RkHS, fct. evaluation is a cls. operation

mapping $E_x: \mathcal{H} \rightarrow \mathbb{R}$

$$E_x(f) = f(x)$$

$$|f(x) - g(x)| \leq |(f - g, k(\cdot, \cdot))_H|$$

⊗ this property is important to do statistics

$\leq \|f - g\|_H \|k(x, \cdot)\|_H$ i.e. E_x is Lipschitz with $L = \|k(x, \cdot)\|_H$