

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

$$\forall w \in \mathcal{W} \quad L_{01}(w) \leq \hat{L}_n(w) + \frac{1}{\sqrt{n}} \sqrt{d \log \frac{d}{\delta} + \log \frac{2}{\delta}}$$

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

VC-dimension of $\mathcal{H} \triangleq \max \{m : \exists \text{ a set of } m \text{ pts. s.t. } \forall \text{ labelings of these pts. } \exists w \text{ s.t. } h_w \text{ gives the correct label on those points}\}$
"shattering the set of pts"

of prediction functions on m pts. is 2^m

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

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

\Rightarrow motivates going to data distribution dependent measure of complexity

example: empirical Rademacher complexity

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

"correlations with random noise"

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

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

$$\forall w \quad L_{01}(w) \leq \hat{L}_n(w) + \hat{R}_{D_n}(\mathcal{H}) + \frac{1}{\sqrt{n}} 3 \sqrt{\frac{\log 2/\delta}{2}}$$

$$\forall w \quad |L(w) - \hat{L}(w)| \leq L_n(w) + \frac{1}{\sqrt{n}} \sum \sqrt{\log \frac{1}{\delta}} + \frac{1}{\sqrt{n}} \sum \sqrt{\log \frac{1}{\delta}}$$

complexity depends on D_n (implicitly on P)

high level idea to prove bound:

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

"symmetrization trick" \rightarrow bound sup of differences between $\mathbb{E}_{D_n'} [L_n'(w)]$
 $\hat{L}_n(w)$

+ 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 $y \neq y'$

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

\hookrightarrow set of cliques of a graph, model G / factor graphs

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

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

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

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

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

and consider $\mathcal{W}_R \triangleq \{w : \|w\|_2 \leq R\}$; let $R = \max_{i,c,y} \|\psi_c(x_i, y)\|_2$

$$\text{then } \hat{R}_n^G(\mathcal{H}_{\mathcal{W}_R}) \leq \frac{R \sum |\mathcal{C}| \sqrt{\max |\mathcal{Y}_c|}}{\sqrt{n}}$$

\rightarrow so want small degrees!

* plug Thm. 2 back in Thm. 7:

$$L(w) \leq L_{\text{hinge}}(w) + \underbrace{\left(\frac{R \sum |\mathcal{C}| \sqrt{\max |\mathcal{Y}_c|}}{\sqrt{n}} \right)}_{\frac{\lambda_n}{2}} \underbrace{\|w\|_2}_{\lambda_n} + \text{cst.}$$

min of RHS suggests

SVM struct alg. $\hat{w}_n = \underset{w}{\text{argmin}} \sum_{i=1}^n \text{hinge}(w) + \frac{\lambda}{2} \|w\|^2$

Missing link: ① min $f(w)$ s.t. $\|w\| \leq R$

(if f is convex) use Lagrangian duality for $\lambda(L)$ st.

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

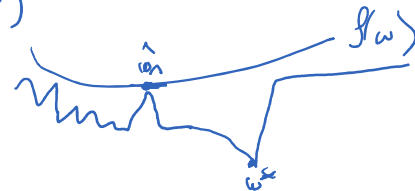
sol'n of ② gives same solution as ①

[side note: 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

caution: also note here: no consistency guarantee

next: consistency

4/4/1

consistency & calibration

need to relate $f(w)$ to $L(w)$: tool "calibration function" [Steinwart]

relationship is usually very complicated

⇒ current results look mainly at non-parametric setting (∞ # of parameters)

all functions $h: X \rightarrow Y$ are considered ⇒ this evaluates 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)

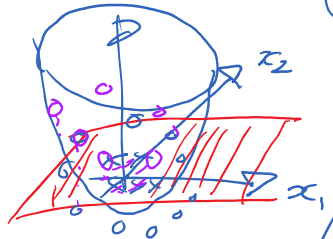
→ can do this using a universal kernel

$$S(\cdot, \cdot; \omega) \in \mathcal{H}_{x \times x}$$

RKHS

motivation:

illustration:
(not to scale!)



generalize linear structure

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

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

formal example of kernel trick

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

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

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

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

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

even have ∞ dim, e.g. $k(x, x') = \exp(-\frac{1}{2}\|x-x'\|^2)$ (RBF kernel)

RKHS (reproducing kernel Hilbert space)

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

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

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

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

"pre-Hilbert" space [inner product space]

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

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

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

Then RKHS H is = completion (\tilde{H}) using $\|\cdot\|_H$ as you norm

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

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

* "reproducing" property of H : for $f \in H$

$$\langle f, k(x, \cdot) \rangle_H = f(x)$$

also
natural $\Phi(x) = k(x, \cdot)$

nice property of RKHS, fct. evaluation is acts operation

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

$$|f(x) - g(x)| = |\langle f - g, k(x, \cdot) \rangle_H|$$

$$\stackrel{CS}{\leq} \|f - g\|_H \|k(x, \cdot)\|_H \quad \text{i.e. } E_x \text{ is Lipschitz } f_{\mathbb{R}} \\ \text{with } L = \|k(x, \cdot)\|_H$$

⊗ this property is
important to do statistics

$$f: X \rightarrow \mathbb{R}$$

$$S: X \times Y \rightarrow \mathbb{R}$$

or

$$x \rightarrow \mathbb{R}^k$$