

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

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

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

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

$\delta: 1/m \rightarrow \{1/2, 1/3\}$

note: # of prediction functions on m points is 2^m

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

⊗ one issue for this bound 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

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

"correlation with random noise"

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

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

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

complexity depends on D_n (implicitly depends on P)

structured prediction generalization bounds [Cortes et al NeurIPS 2016]

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

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

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

$\mathcal{C} \subseteq \mathcal{E}$ set of edges of a graph model G / factor graph

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

$$\forall w \in W \quad L(w) \leq \hat{J}_{\text{hinge}}(w) + 4\sqrt{\lambda} \hat{R}_{D_n}^G(\mathcal{H}(w)) + \frac{3\sqrt{\lambda} \max_{y, y' \in \mathcal{Y}} \ell(y, y')}{\sqrt{2n}}$$

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

"empirical factor graph complexity"

indep Rademacher R.V.

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

thm 2: if $\mathcal{L}(x, y_c; w) = \langle w, \varphi_c(x, y_c) \rangle$

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

$$\text{then } \hat{R}_{D_n}^G(\mathcal{H}_{W_\lambda}) \leq \frac{R \sqrt{\lambda} |\mathcal{C}| \sqrt{\max_c |\mathcal{Y}_c|}}{\sqrt{n}}$$

so want small degrees!

* plug thm. 2 back in thm. 7

$$L(w) \leq \hat{J}_{\text{hinge}}(w) + \underbrace{\left(\frac{R \sqrt{\lambda} |\mathcal{C}| \sqrt{\max_c |\mathcal{Y}_c|}}{\sqrt{n}} \right)}_{\frac{\lambda n}{2}} \underbrace{\|w\|_2}_{\|w\|_2} + \text{cst.}$$

min of R.H.S suggests

SVM struct. alg. $\hat{w}_\lambda = \underset{w}{\text{argmin}} \hat{J}_{\text{hinge}}(w) + \frac{\lambda n}{2} \|w\|_2^2$

missing link:

- ① $\min f(w)$ s.t. $\|w\|_2 \leq \lambda$ (if f is convex) use Lagrangian duality $\exists \alpha(\lambda)$ s.t.
- ② $\min f(w) + \frac{\lambda}{2} \|w\|_2^2$ sol'n to ② gives same solution to ①

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

10h36

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 + convex surrogate

consistency & calibration

need to relate $g(w)$ to $L(w)$: let "calibration fct." [Steinwart]

relationship is usually very complicated

⇒ classical results lack mainly at non-parametric setting (no # of parameters)

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

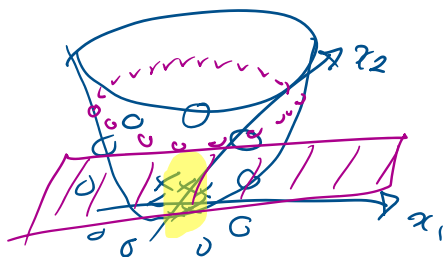
i.e. we suppose that $s(x, y; w)$ can be arbitrary for any x (i.e. w is h_0 -dim)

→ can do this using a universal kernel

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

RKHS:

motivation for "kernel trick"



generalize linear structure

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

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

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

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

$$\langle \Phi(x), \Phi(x') \rangle_{\mathbb{R}^3} = (\langle x, x' \rangle_{\mathbb{R}^2})^2 = (x_1 x'_1 + x_2 x'_2)^2 = k(x, x')$$

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

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

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

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

even have to do with e.g. $k(x, x') = \exp\left(-\frac{\|x-x'\|^2}{2}\right)$

"RBF kernel"

RKHS (reproducing kernel Hilbert space)

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

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

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

(here $\Phi(x) = k(x, \cdot)$)

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

'pre-Hilbert' space [inner product space]

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

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

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

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

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

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

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

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

$\Phi(x) \leftarrow \langle w, \Phi(x) \rangle$

nice property of RKHS, fct. evaluation is a cb. operation

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

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

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

$$\leq \|f - g\|_{\mathcal{H}} \|k(x, \cdot)\|_{\mathcal{H}} \quad \text{i.e. } E_x \text{ is Lipschitz w.r.t. } L = \|k(x, \cdot)\|_{\mathcal{H}}$$

⊗ this property is important to do statistics