

$$\bar{T}_3 = \frac{1}{n} \sum_{i=1}^n L_0(x_i) - \mathbb{E}_x [L_0(x)]$$

concentration of measure meas.

$$\bar{T}_1 = \mathbb{E}_x [L_0(x)] - \frac{1}{n} \sum_{i=1}^n L_0(x_i)$$

challenging because $\hat{\theta}$ is a random quantity

$$\mathcal{G} = \{L_\theta(\cdot) | \theta \in \Theta\}$$

$$|\bar{T}_1| \leq \sup_{\theta \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n L_\theta(x_i) - \mathbb{E}_x [L_\theta(x)] \right| = \|P_n - P\|_{\mathcal{G}}$$

$$\mathbb{E}(\bar{\theta}_1 | \theta^*) \leq 2\|P_n - P\|_{\mathcal{G}}.$$

$$\|P_n - P\|_{\mathcal{G}} = \sup_{f \in \mathcal{G}} \left| \frac{1}{n} \sum_{i=1}^n f(x_i) - \mathbb{E}[f(x)] \right|$$

12.3. Uniform laws via Rademacher complexity

function class \mathcal{F}

$$x_1^n = (x_1, \dots, x_n)$$

$$\mathcal{F}(x_1^n) = \{f(x_1), f(x_2), \dots, f(x_n) | f \in \mathcal{F}\}$$

empirical Rademacher complexity

$$R(\mathcal{F}(x_1^n) | n) = \mathbb{E}_{\varepsilon} [\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(x_i) \right|]$$

$\{\varepsilon_i\}_{i=1}^n$... i.i.d. sequence of Rademacher RVs, taking ± 1 with equal prob.

$$\frac{1}{n} \sum_{i=1}^n \varepsilon_i f(x_i) = \frac{1}{n} (\varepsilon_1 \ \varepsilon_2 \dots \varepsilon_n) \underbrace{\begin{pmatrix} f(x_1) \\ f(x_2) \\ \vdots \\ f(x_n) \end{pmatrix}}_{\text{total of } 2^n \text{ realizations}}$$

$$\begin{array}{c} \cdot \\ + \\ \cdot \end{array}$$

$$X_1^n = \{x_i\}_{i=1}^n$$

$$R_n(\mathcal{F}) = R(\mathcal{F}(x_1^n) | n) = E_{\mathcal{X}} [R(\mathcal{F}(x_1^n) | n)]$$

$$= E_{\mathcal{X}} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(x_i) \right| \right]$$

Theorem 12.10. For any b-uniformly bounded function class, i.e.,

$\|f\|_\infty \leq b$, $\forall f \in \mathcal{F}$, any positive integer $n \geq 1$ and any $\delta \geq 0$, we have

$$\|P_n - P\|_{\mathcal{F}} \leq 2R_n(\mathcal{F}) + \delta$$

with prob. $\geq 1 - e^{-\frac{n\delta^2}{2b^2}}$. Hence, as long as $R_n(\mathcal{F}) = o(1)$, we have $\|P_n - P\|_{\mathcal{F}} \xrightarrow{\text{a.s.}} 0$.

sufficient condition for $\|P_n - P\|_{\mathcal{F}} \xrightarrow{\text{a.s.}} 0$ is: $R_n(\mathcal{F}) \xrightarrow{n \rightarrow \infty} 0$

Prop. 12.11. For every b-uniformly bounded function class \mathcal{F} , any integer $n \geq 1$ and any $\delta \geq 0$, we have

$$\|P_n - P\|_{\mathcal{F}} \geq \frac{1}{2} R_n(\mathcal{F}) - \frac{\sup_{f \in \mathcal{F}} |\mathbb{E}[Pf]|}{2\sqrt{n}} - \delta$$

with prob. $\geq 1 - e^{-\frac{n\delta^2}{2b^2}}$.

levels of escales in terms of upper-bounding Rademacher complexity

1. $|\mathcal{F}| < \infty$ union bounds won't consider

2. polynomial discrimination function classes ✓

3. VC (Vapnik-Chervonenkis) dimension ✓

4. metric entropy (chaining arguments) ✗

Empirical process theory

12.4. Function classes with polynomial discrimination

recall that $\mathcal{F}(x_1^n) = \{f(x_1), f(x_2), \dots, f(x_n) | f \in \mathcal{F}\}$

$$R_n(\mathcal{F}) = E_{\mathcal{E}X} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(x_i) \right| \right]$$

\mathcal{F} contains binary-valued functions

$$f \in \mathcal{F} : f(x) = \pm 1$$

$$|\mathcal{G}(x_1^n)| = \left| \underbrace{\{f(x_1), f(x_2), \dots, f(x_n)\}}_{\pm 1} \mid f \in \mathcal{F} \right| = 2^n$$

$|\mathcal{G}(x_1^n)|$ is polynomial in n

Def. 12.2 (Polynomial dimension). A class of functions \mathcal{F} with domain X has polynomial dimension of order $\nu \geq 1$, if for each $n \in \mathbb{N}$ and collection $x_1^n = (x_1, \dots, x_n)$ of n points in X , the set $\mathcal{G}(x_1^n)$ satisfies

$$|\mathcal{G}(x_1^n)| \leq (n+1)^\nu.$$

Lemma 12.13. Suppose that \mathcal{F} has polynomial dimension of order ν . Then, for all $n \in \mathbb{N}$ and any collection of points $x_1^n = (x_1, \dots, x_n)$, we have

$$E_{\mathcal{E}} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(x_i) \right| \right] \leq C D(x_1^n) \sqrt{\frac{\log(n+1)}{n}},$$

where $D(x_1^n) = \sup_{f \in \mathcal{F}} \sqrt{\frac{\sum_{i=1}^n f^2(x_i)}{n}}$ is the L_2 -radius of the set $\mathcal{G}(x_1^n)/\sqrt{n}$.

Proof.

$$E_{\mathcal{E}} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(x_i) \right| \right] = \frac{1}{n} E_{\mathcal{E}} \left[\max_{\theta \in \mathcal{G}(x_1^n)} |\langle \varepsilon, \theta \rangle| \right]$$

$$\varepsilon = (\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n)$$

The set $\mathcal{G}(x_1^n)$ has $\leq (n+1)^\nu$

will apply

Lemma 2.4. For $n \geq 2$, let $\{X_i\}_{i=1}^n$ be a set of zero-mean RVs, each sub-gaussian with parameter σ . Then,

$$E \left[\max_{i=1,\dots,n} |X_i| \right] \leq 2\sigma \sqrt{\log(n)}.$$

sub-gaussian $E [e^{\lambda(X-\mu)}] \leq e^{\frac{\sigma^2 \lambda^2}{2}}, \forall \lambda \in \mathbb{R}$

$\langle \varepsilon, \theta \rangle$ is sub-gaussian (with parameter $\sigma = \sup_{\theta \in \mathcal{F}(x_1^n)} \sqrt{\sum_{i=1}^n \theta_i^2}$)
 (used independence of ε_i)

$$\begin{aligned} E_{\varepsilon} [e^{\lambda \langle \varepsilon, \theta \rangle}] &= \prod_{i=1}^n E_{\varepsilon_i} [e^{\lambda \varepsilon_i \theta_i}] \\ &= \prod_{i=1}^n \frac{e^{\lambda \theta_i} + e^{-\lambda \theta_i}}{2} \end{aligned}$$

$$\begin{aligned} e^x + e^{-x} &\leq 2e^{\frac{x^2}{2}} \\ &\leq e^{\frac{\lambda^2}{2} \sum_{i=1}^n \theta_i^2} \\ &\leq e^{\frac{\lambda^2}{2} \underbrace{\left\{ \sup_{\theta \in \mathcal{F}(x_1^n)} \sum_{i=1}^n \theta_i^2 \right\}}_{\sigma^2}} \end{aligned}$$

$$R(\mathcal{F}(x_1^n) | n) \leq 2 D(x_1^n) \sqrt{\frac{\log(n+1)}{n}}. \quad \square$$

$$D(x_1^n) = \sup_{f \in \mathcal{F}} \sqrt{\frac{\sum_{i=1}^n f(x_i)}{n}}$$

$$R_n(\mathcal{F}) = E_x [R(\mathcal{F}(x_1^n) | n)] \leq 4 E_{x_1^n} [D(x_1^n)] \sqrt{\frac{\log(n+1)}{n}}.$$

b-uniformly bounded function class

$$R_n(\mathcal{F}) \leq 4b \sqrt{\frac{\log(n+1)}{n}}, \text{ for all } n \in \mathbb{N}$$

Corollary 12.4. (Classical Glivenko-Cantelli). Let $\bar{F}(t) = P[X \leq t]$ be the CDF of a random variable $X \sim P$, and let \hat{F}_n be the empirical CDF based on n i.i.d. samples $X_i \sim P$. Then,

$$P \left[|\hat{f}_n - f|_{\infty} \geq \sqrt{\frac{\log(n+1)}{n}} + \delta \right] \leq e^{-\delta^2/2}, \delta \geq 0$$

and hence $|\hat{f}_n - f|_{\infty} \xrightarrow{\text{as.}} 0$.

Proof $\mathcal{F}(x_1^n) = \{f(x_1), f(x_2), \dots, f(x_n) \mid f \in \mathcal{F}\}$

$x_1^n = (x_1, \dots, x_n) \in \mathbb{R}^n$, consider the set $\mathcal{F}(x_1^n)$ of $\{0,1\}$ -valued indicator functions of the half-intervals $(-\infty, t]$, $t \in \mathbb{R}$.

$$x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)},$$

$$\overbrace{\quad \quad \quad}^{x_{(1)}} \overbrace{\quad \quad \quad}^{x_{(2)}} \overbrace{\quad \quad \quad}^{x_{(n)}} \quad \quad \quad$$

divides \mathbb{R} up into $n+1$ intervals

for a given $t \in \mathbb{R}$, $1_{(-\infty, t]}(x) = \begin{cases} 1, & x \leq t \\ 0, & \text{else} \end{cases}$

$$t \leq x_{(1)} : (f(x_1), f(x_2), \dots, f(x_n)) = (0, 0, \dots, 0)$$

$$x_{(1)} \leq t \leq x_{(2)} : \quad \overbrace{\quad \quad \quad}^{1, 1, \dots} = (1, 0, \dots, 0)$$

$$x_{(2)} \leq t \leq x_{(3)} : \quad \overbrace{\quad \quad \quad}^{1, 1, \dots} = (1, 1, 0, \dots, 0)$$

!

$$(1, 1, \dots, 1)$$

$$|\mathcal{F}(x_1^n)| \leq (n+1)^n = (n+1)^n, n=1$$

with $n=1$, Lemma 12.13.

$$E \mathbb{E} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \mathbb{E}[f(x_i)] \right| \right] \leq \sqrt{\frac{\log(n+1)}{n}}$$

$$E \text{ w.r.t. } X \Rightarrow R(\mathcal{F}) \leq \sqrt{\frac{\log(n+1)}{n}}$$

apply Thm. 12.10. □

12.5 Vapnik-Chervonenkis dimension

$$\mathcal{F}(x_1^n) = \{f(x_1), \dots, f(x_n) \mid f \in \mathcal{F}\}$$

$$|\mathcal{F}(x_1^n)| \leq 2^n$$

just had an example with $|F(x_1^n)| \leq (n+1)$

Def 12.15. (Shattering and VC dimension). Given a class \mathcal{F} of binary-valued functions, we say that the set $x_1^n = (x_1, \dots, x_n)$ is shattered by \mathcal{F} if $|F(x_1^n)| = 2^n$.

The VC dimension $V(\mathcal{F})$ is the largest integer n for which there is a collection $x_1^n = (x_1, \dots, x_n)$ of n points that is shattered by \mathcal{F} .

When $V(\mathcal{F})$ is finite, \mathcal{F} is said to be a VC class.

Example 12.16. $S_{\text{left}} = \{(-\infty, a] \mid a \in \mathbb{R}\}$, \mathcal{F} = indicator functions on S_{left}

$$n=1: |\mathcal{F}| = 2 \text{, with } \mathcal{F} = \{0, 1\}$$

$$= 2^1$$

$$n=2:$$

$$\begin{array}{c} \times \quad \times \\ \hline x_{(1)} \quad x_{(2)} \end{array}$$

$$\mathcal{F}(x_1^2) = \{f(x_{(1)}), f(x_{(2)}) \mid f \in \mathcal{F}\}$$

$$\mathcal{F}(x_1^2) = \{(0,0), (1,0), (1,1)\}$$

\Rightarrow 2-point sets are not being shattered

$$\Rightarrow V(\mathcal{F}) = 1.$$

$S_{\text{two}} = \{[a,b] \mid a, b \in \mathbb{R} \text{ s.t. } a < b\}$

$$\begin{array}{c} \times \quad \times \\ \hline x_{(1)} \quad x_{(2)} \end{array}$$

$$\mathcal{F} = \{(0,0), (1,0), (0,1), (1,1)\}$$

$$|\mathcal{F}| = 4 = 2^2$$

$$\begin{array}{c} \times \quad \times \quad \times \\ \hline x_{(1)} \quad x_{(2)} \quad x_{(3)} \end{array}$$

$$\mathcal{F} = \{(1,0,0), (0,1,0), (0,0,1), (1,1,0), (0,1,1), (1,1,1), (0,0,0), (\cancel{1,0,1})\}$$

$$\Rightarrow |\mathcal{F}| = 7 < 2^3$$

$$\Rightarrow V(\mathcal{F}) = 2.$$

$$|\mathcal{F}| \leq (n+1)^2 \Rightarrow V(\mathcal{F}) = 2.$$

all examples we considered are of finite VC dim. and also have polynomial discernment.

Q: Does every finite VC class have polynomial discernment?

$$n=1$$

$$2, |\mathcal{F}| = 2^1$$

$$n=2$$

$$2^2 = 4, |\mathcal{F}| = 2^2$$

$$n=3$$

$$2^3 = 8, |\mathcal{F}| = 2^3$$

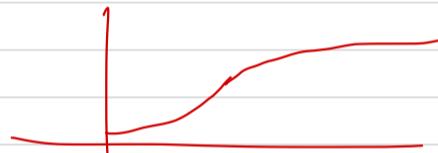
$$n=4$$

$$\text{breaks} < 16$$

$$|\mathcal{F}| \leq (n+1)^2$$

:

:



Theorem 12.17. (Vapnik-Chervonenkis, Sauer-Shelah). Consider a set class S with $w(S) < \infty$. Then, for any collection of points $P = (x_1, \dots, x_n)$ with $n \geq w(S)$, we have

$$|S(P)| \leq \sum_{i=0}^{w(S)} \binom{n}{i} \leq (n+1)^{w(S)}.$$