

* we focus on binary classification.

Also holds for other tasks like regression with quadratic loss.

But exists example where PAC learnable but not uniformly convergent.

chapter 13 exercise 2. $\left\{ \begin{array}{l} \text{Uniform Convergence} \\ \text{learnability} \end{array} \right. \begin{array}{l} m \propto \log d \\ m \propto \frac{1}{\epsilon} : \text{indep}(d) \end{array}$

Theorem [Fundamental Theorem, Quantitative Version]

If $VCdim(\mathcal{H}) = d$, then

1. \mathcal{H} has **Uniform Convergence** property with

$$C_1 \cdot \frac{d + \log \frac{1}{\delta}}{\epsilon^2} \leq m_{\mathcal{H}}^{U.C.}(\epsilon, \delta) \leq C_2 \cdot \frac{d + \log \frac{1}{\delta}}{\epsilon^2}$$

2. \mathcal{H} is **Agnostic PAC learnable** with

$$C_1 \cdot \frac{d + \log \frac{1}{\delta}}{\epsilon^2} \leq m_{\mathcal{H}}(\epsilon, \delta) \leq C_2 \cdot \frac{d + \log \frac{1}{\delta}}{\epsilon^2}$$

3. \mathcal{H} is **PAC learnable** with

$$C_1 \cdot \frac{d + \log \frac{1}{\delta}}{\epsilon} \leq m_{\mathcal{H}}(\epsilon, \delta) \leq C_2 \cdot \frac{d \log \frac{1}{\epsilon} + \log \frac{1}{\delta}}{\epsilon}$$

proof later.

proof sketch of fundamental theorem of learning.

① [Sauer's Lemma]

If $VCdim(H) = d$, then for any $C \subseteq \mathcal{X}$

the effective size of H restricted to C ,

i.e., $|H_C|$ is only $O(|C|^d) \ll 2^{|C|}$

polynomial

worst-case
exponential in $|C|$

*Recap [Restriction of H to $C = \{c_1, \dots, c_m\} \subseteq \mathcal{X}$]

$$H_C \triangleq \{ \tilde{h} : C \rightarrow \{0, 1\} \mid h \in H \}$$

$$\tilde{h} = (\tilde{h}(c_1), \tilde{h}(c_2), \dots, \tilde{h}(c_m))$$

$$= (0 \quad 1 \quad 0 \quad \dots \quad 0)$$

② If $|H_C|$ grows polynomially in $|C|$,

then we have uniform convergence: $\forall m \geq m_{\mathcal{H}}^{U.C.}(\epsilon, \delta)$

$$\mathbb{P}_{\mathcal{S}} \left(|L_D(h) - L_S(h)| < \epsilon, \forall h \in H \right) \geq 1 - \delta.$$

In lecture 11, we used Hoeffding's to show $|H| < \infty \rightarrow$ Uniform Convergence.

Definition [Growth function]

The growth function $\tau_H : \mathbb{N} \rightarrow \mathbb{N}$ is defined as

$$\tau_H(m) = \max_{C \subseteq \mathcal{X} : |C|=m} |H_C|$$

* the number of different functions from C to $\{0, 1\}$, restricting H to C .

* If $m \leq VCdim(H) = d$, then by definition $\tau_H(m) = 2^m$
largest m s.t. $\tau_H(m) = 2^m$

* If $m > d$, τ_H grows polynomially in m .

Precisely,

① Lemma [Sauer-Shelah-Perles]

If $\text{VCdim}(\mathcal{H}) = d < \infty$, then for all m , $T_{\mathcal{H}}(m) \leq \sum_{i=0}^d \binom{m}{i}$

$\rightarrow m \leq d$, $T_{\mathcal{H}}(m) \leq \sum_{i=0}^m \binom{m}{i} = 2^m$ ← exponential

$\rightarrow m \geq d+1$, $T_{\mathcal{H}}(m) \leq \left(\frac{em}{d}\right)^d$ ← polynomial

↑
Stirling's formula

m choose i
 $\frac{m!}{i!(m-i)!}$

② Theorem [for Uniform Convergence] for every $D, \delta, h \in \mathcal{H}$

$$|L_0(h) - L_S(h)| \leq \frac{4 + \sqrt{4d \lg(T_{\mathcal{H}}(2m))}}{\delta \cdot \sqrt{m/2}} \quad \text{w.p } 1-\delta.$$

Proof of fundamental theorem

Apply ① Sauer's lemma to ②, for $m \geq d$

$$|L_0(h) - L_S(h)| \leq \frac{4 + \sqrt{d \lg \frac{2em}{d}}}{\delta \sqrt{m/2}}$$

for simplicity, let $\sqrt{d \lg \frac{2em}{d}} \geq 4$

$$\leq \frac{2 \sqrt{d \lg \frac{2em}{d}}}{\delta \sqrt{m/2}} \stackrel{\text{want}}{\leq} \epsilon$$

$$\rightarrow \frac{8d \lg \left(\frac{2em}{d}\right)}{\delta^2 \epsilon^2} \leq m$$

$$\rightarrow m \geq 4 \cdot \frac{8d \lg \frac{2e}{\delta^2 \epsilon^2}}{\delta^2 \epsilon^2} + \frac{16 \cdot d \cdot \lg \left(\frac{2e}{d}\right)}{\delta^2 \epsilon^2} \quad \text{is sufficient}$$

this is loose bound but enough for PAC learnability

Proof of ② > Claim: $\mathbb{E}_S \left[\sup_{h \in \mathcal{H}} |L_D(h) - L_S(h)| \right] \leq \frac{4 + \sqrt{\log(\mathcal{N}_H(2m))}}{\sqrt{m/2}}$

Markov's ineq. $\implies \mathbb{P}_S \left(\sup_{h \in \mathcal{H}} |L_D(h) - L_S(h)| > \varepsilon \right) \leq \frac{4 + \sqrt{\log(\mathcal{N}_H(2m))}}{\varepsilon \cdot \sqrt{m/2}}$

How to handle when $|X| = \infty \rightarrow$ symmetrization

$$\mathbb{E}_S \left[\sup_{h \in \mathcal{H}} |L_D(h) - L_S(h)| \right] = \mathbb{E}_S \left[\sup_{h \in \mathcal{H}} \left| \mathbb{E}_{S'} [L_S(h)] - L_S(h) \right| \right]$$

\parallel
 $\mathbb{E}_{S'} [L_S(h)]$

$\{z_i\}_{i=1}^m$
 \parallel
 $\{z'_i\}_{i=1}^m$

Jensen's ineq. \rightarrow

$$\leq \mathbb{E}_S \left[\sup_{h \in \mathcal{H}} \mathbb{E}_{S'} \left[|L_{S'}(h) - L_S(h)| \right] \right]$$

any convex $f(\cdot)$



$$\leq \mathbb{E}_S \left[\mathbb{E}_{S'} \left[\sup_{h \in \mathcal{H}} |L_{S'}(h) - L_S(h)| \right] \right]$$

$$\mathbb{E}[f(z)] \geq f(\mathbb{E}[z]) \quad = \mathbb{E}_{S,S'} \left[\sup_{h \in \mathcal{H}} \left| \frac{1}{m} \sum_{i=1}^m (l_h(z_i) - l_h(z'_i)) \right| \right] (*)$$

Symmetrization trick:

to retain randomness even for fixed S, S'

$$l_h(z'_i) - l_h(z_i) \stackrel{\text{dist}}{=} \beta_i \cdot (l_h(z'_i) - l_h(z_i))$$

$\beta_i \in \{+1, -1\}$

$$\stackrel{\text{dist}}{=} \beta_i \cdot (l_h(z'_i) - l_h(z_i)) \quad , \beta_i \sim \begin{cases} +1 & 1/2 \\ -1 & 1/2 \end{cases} \text{ b.p.}$$

$$\Rightarrow (*) = \mathbb{E}_{\beta=(\beta_1 \dots \beta_m)} \left[\mathbb{E}_{S,S'} \left[\sup_{h \in \mathcal{H}} \left| \frac{1}{m} \sum_{i=1}^m \beta_i (l_h(z'_i) - l_h(z_i)) \right| \mid \beta \right] \right]$$

$$= \mathbb{E}_{S,S'} \left[\mathbb{E}_{\beta} \left[\sup_{h \in \mathcal{H}} \left| \frac{1}{m} \sum_{i=1}^m \beta_i (l_h(z'_i) - l_h(z_i)) \right| \mid S, S' \right] \right]$$

S, S' fixed, let $C = S \cup S'$, $|C| \leq 2m$, only \mathcal{H}_C matters in the supremum.

$$= \mathbb{E}_G \left[\max_{h \in \mathcal{H}_C} \left| \frac{1}{m} \sum_{i=1}^m \theta_i (l_h(z_i') - l_h(z_i)) \right| \middle| S, S' \right]$$

θ_i is random, $\mathbb{E}\theta_i = 0, -1 \leq \theta_i \leq 1$.

Hoeffding's ineq. $\rightarrow \mathbb{P} \left(\left| \frac{1}{m} \sum_{i=1}^m \theta_i \right| > \rho \right) \leq 2 \cdot e^{-\frac{2m\rho^2}{4}}$

Union bound $\rightarrow \mathbb{P} \left(\max_{h \in \mathcal{H}_C} \left| \frac{1}{m} \sum_{i=1}^m \theta_i \right| > \rho \right) \leq |\mathcal{H}_C| \cdot \mathbb{P} \left(\left| \frac{1}{m} \sum_{i=1}^m \theta_i \right| > \rho \right)$

$$\leq 2 \cdot |\mathcal{H}_C| \cdot e^{-\frac{2m\rho^2}{4}}$$

Sub-Gaussian tail (**)
 \downarrow
 bounded expectation

$$\mathbb{E} \left[\max_{h \in \mathcal{H}_C} \left| \frac{1}{m} \sum_{i=1}^m (l_h(z_i') - l_h(z_i)) \right| \right]$$

$$\leq \frac{4 + \sqrt{\log(|\mathcal{H}_C|)}}{\sqrt{m/2}}$$

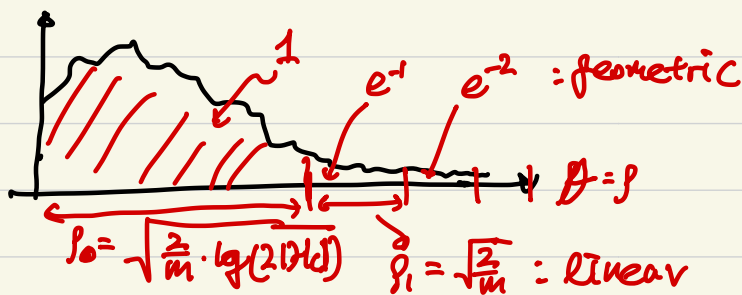
$|C| \leq 2m \rightarrow$

$$\leq \frac{4 + \sqrt{\log(2m)}}{\sqrt{m/2}}$$

(**) [lemma A.4 in the book]

$$\mathbb{P}(\theta > \rho) \leq 2|\mathcal{H}_C| \cdot e^{-\frac{m\rho^2}{2}} \Rightarrow \mathbb{E}[\theta] \leq \frac{4 + \sqrt{\log(2|\mathcal{H}_C|)}}{\sqrt{m/2}}$$

Proof sketch \rangle sub-Gaussian pdf



$$\mathbb{E}[\theta] = \int_0^\infty \theta \cdot f(\theta) d\theta \leq \rho_0 + \sum_{i=1}^{\infty} e^{-i} (\rho_0 + i \cdot \rho_1) \leq \dots \leq 2\sqrt{\frac{2}{m}} + \sqrt{\frac{2}{m} \cdot \log(2|\mathcal{H}_C|)}$$

