

Theorem [Fundamental Theorem of Statistical Learning]

Let H be a hypothesis class, from \mathcal{X} to $\{0, 1\}$,
and the loss function is a 0-1 loss, then the
following are equivalent.

1. H has the uniform convergence property
2. Any ERM rule is a successful agnostic PAC learner of H .
3. H is agnostic PAC learnable.
4. H is PAC learnable
5. Any ERM rule is a successful PAC learner for H
6. H has a finite VC-dimension.

We already learned that $1 \rightarrow 2$ (lecture 3)

$$\begin{array}{ll} 2 \rightarrow 3 & \text{trivial} \\ 3 \rightarrow 4 & \text{trivial} \\ 2 \rightarrow 5 & \text{trivial} \\ 2 \rightarrow 6 & : n \rightarrow n^2 \text{ no free lunch theorem.} \end{array}$$

We are left to prove $6 \rightarrow 1$

Theorem [Fundamental theorem, Quantitative version]

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

1. H has uniform convergence property with

$$C_1 \cdot \frac{dt \log \frac{1}{\delta}}{\epsilon^2} \leq m_H^{\text{VC}}(\epsilon, \delta) \leq C_2 \cdot \frac{dt \log \frac{1}{\delta}}{\epsilon^2}$$

2. H is agnostic PAC learnable with

$$C_1 \cdot \frac{dt \log \frac{1}{\delta}}{\epsilon^2} \leq m_H(\epsilon, \delta) \leq C_2 \cdot \frac{dt \log \frac{1}{\delta}}{\epsilon^2}$$

3. H is PAC learnable with,

$$C_1 \cdot dt \log \frac{1}{\delta} \leq m_H(\epsilon, \delta) \leq C_2 \cdot \frac{dt \log \epsilon + \log \frac{1}{\delta}}{\epsilon^2}$$

proof sketch.

proof strategy for fundamental theorem of statistical learning.

① [Sauer's lemma]

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

the effective size of \mathcal{H} restricted to C , i.e., $|\mathcal{H}_C|$ is only $O(|C|^d) \ll 2^{|C|}$

polynomial exponential

$2^{|C|}$

worst-case

bif-O notation

$|\mathcal{H}_C| = O(f_{\mathcal{C}}(c, d))$ means that $\exists c > 0$

$\leq c f_{\mathcal{C}}(c, d)$ for large enough c .

Definition (Restriction of \mathcal{H} to C)

\mathcal{H} is class of functions $\mathcal{X} \rightarrow \{0, 1\}$, and $C = (c_1, \dots, c_m) \subseteq \mathcal{X}$,

the Restriction of \mathcal{H} to C

$$\mathcal{H}_C \triangleq \left\{ h: C \rightarrow \{0, 1\} \mid h \in \mathcal{H} \right\}$$

$c_i \mapsto h(c_i)$

equivalently we can represent h by a vector $(h(c_1), \dots, h(c_m))$.

② If $|\mathcal{H}_C|$ grows polynomially in $|C|$, then we have uniform Convergence Property.

Definition (Uniform Convergence)

A hypothesis class \mathcal{H} has the uniform convergence property,

if $\exists M_{\mathcal{H}}^{UC}$ s.t. for every ϵ, δ and every D , if $m \geq M_{\mathcal{H}}^{UC}(\epsilon, \delta)$

the S is ϵ -representative with probability at least $1 - \delta$.

$$\frac{1}{m} \sum_{i=1}^m |L_D(h_i) - L_S(h_i)| < \epsilon. \quad \forall h \in \mathcal{H}.$$

Definition [Growth Function]

The Growth Function $T_{\mathcal{H}}: \mathbb{N} \rightarrow \mathbb{N}$, is defined as

$$T_{\mathcal{H}}(m) = \max_{C \subseteq \mathcal{X}: |C|=m} |\mathcal{H}_C|$$

which is the number of different functions from C of size m to $\{0, 1\}$ that can be obtained by restricting \mathcal{H} to C .

* If $m \leq \text{VCdim}(\mathcal{H}) = d$, then by definition of VCdim , $T_{\mathcal{H}}(m) = 2^m$.
large $m \leq T_{\mathcal{H}}(m) = 2^m$.

* For $m > d$, m grows polynomially in m .

① Lemma [Sauer-Shelah-Parkes]

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

$$\rightarrow d \leq m, T_{\mathcal{H}}(m) \leq \sum_{i=0}^m \binom{m}{i} = 2^m \leftarrow \text{exp in } m \quad m \text{ choose } i, \frac{m!}{i!(m-i)!}$$

$$\rightarrow d \leq m, T_{\mathcal{H}}(m) \leq \binom{em}{d} \leftarrow \text{polynomial in } m$$

↑ Stirling's formula.

② Theorem [for uniform convergence] for every $D, \delta, h \in \mathbb{N}$

$$|\hat{L}_D(h) - L_S(h)| \leq \frac{4 + \sqrt{\log(\gamma_1(2m))}}{\delta \cdot \sqrt{m/2}} \quad \text{w.p. } 1 - \delta.$$

Proof of fundamental theorem of statistical learning.

using ① Sauer's Lemma, we get from ②, for $m > d$

$$|\hat{L}_D(h) - L_S(h)| \leq \frac{4 + \sqrt{d \log(\frac{2em}{\delta})}}{\delta \cdot \sqrt{m/2}}$$

for simplicity let $\frac{d \log(\frac{2em}{\delta})}{\delta} \geq 4$,

$$\leq \frac{2}{\delta} \sqrt{\frac{2d \log(\frac{2em}{\delta})}{m}} \stackrel{\text{want}}{\leq} \varepsilon$$

$$\rightarrow \frac{8d \log(\frac{2em}{\delta})}{(\delta\varepsilon)^2} \leq m$$

$$\rightarrow m \geq 4 \cdot \frac{8d}{(\delta\varepsilon)^2} \log\left(\frac{2d}{(\delta\varepsilon)^2}\right) + \frac{8d \log(\frac{2e}{\delta})}{(\delta\varepsilon)^2} \text{ is sufficient.}$$

this gives a loose bound on $m_{\mathcal{H}}(\delta, \varepsilon)$, but

sufficient to prove PAC learnability.

Agnostic

Proof of ② claim: $\mathbb{E}_{\delta} \left[\sup_{h \in \mathcal{H}} |\hat{L}_D(h) - L_S(h)| \right] \leq \frac{4 + \sqrt{\log(\gamma_1(2m))}}{\sqrt{m/2}}$

$$\xrightarrow{\text{Markov's Ineq: }} \mathbb{P} \left(\sup_{h \in \mathcal{H}} |\hat{L}_D(h) - L_S(h)| > \varepsilon \right) \leq \frac{4 + \sqrt{\log(\gamma_1(2m))}}{\varepsilon \sqrt{m/2}}$$

$$\mathbb{E}_{\delta} \left[\sup_{h \in \mathcal{H}} |\hat{L}_D(h) - L_S(h)| \right] = \mathbb{E}_{S=\{(x_i, y_i)\}_{i=1}^m} \left[\sup_{h \in \mathcal{H}} \left| \mathbb{E}_{S'} \left[\hat{L}_S'(h) - L_S(h) \right] \right| \right]$$

$$\begin{aligned} \text{Jensen's Ineq.} &\Rightarrow \mathbb{E}_{\delta} \left[\sup_{h \in \mathcal{H}} \left| \hat{L}_S'(h) - L_S(h) \right| \right] \\ \text{Convex f.} &\Rightarrow \mathbb{E}_{\delta} \left[\mathbb{E}_{S'} \left[\sup_{h \in \mathcal{H}} \left| \hat{L}_S'(h) - L_S(h) \right| \right] \right] \\ \mathbb{E}[f(z)] \geq f(\mathbb{E}[z]) &\leq \mathbb{E}_{\delta} \left[\mathbb{E}_{S'} \left[\sup_{h \in \mathcal{H}} \left| \hat{L}_S'(h) - L_S(h) \right| \right] \right] \\ \checkmark |z| \text{ max linear f.g.} &= \mathbb{E}_{S, S'} \left[\sup_{h \in \mathcal{H}} \left| \frac{1}{m} \sum_{i=1}^m (\hat{L}_h(z'_i) - \hat{L}_h(z_i)) \right| \right] \quad (*) \end{aligned}$$

$$\hat{L}(h(z'_i)) - \hat{L}(h(z_i)) \stackrel{\text{def.}}{=} b_i \cdot (\hat{L}_h(z'_i) - \hat{L}_h(z_i)), \quad b_i \in \{-1, 1\}$$

$$\stackrel{z'_i, z_i \text{ i.i.d.}}{\stackrel{\text{def.}}{=}} b_i \cdot (\hat{L}_h(z'_i) - \hat{L}_h(z_i)), \quad b_i \sim \{-1, 1\}$$

$$\Rightarrow (*) = \mathbb{E}_{C=(a_1, \dots, a_m)} \left[\mathbb{E}_{S, S'} \left[\sup_{h \in \mathcal{H}} \left| \frac{1}{m} \sum_{i=1}^m b_i (\hat{L}_h(z'_i) - \hat{L}_h(z_i)) \right| \right] \right]$$

$$\text{Linearity of expectation} \Rightarrow \mathbb{E}_{S, S'} \left[\mathbb{E}_{C=(a_1, \dots, a_m)} \left[\sup_{h \in \mathcal{H}} \left| \frac{1}{m} \sum_{i=1}^m b_i (\hat{L}_h(z'_i) - \hat{L}_h(z_i)) \right| \middle| S, S' \right] \right]$$

$$= \mathbb{E}_S \left[\max_{h \in \mathcal{H}_C} \left| \frac{1}{m} \sum_{i=1}^m b_i (\hat{L}_h(z'_i) - \hat{L}_h(z_i)) \right| \middle| S, S' \right]$$

$$\theta_i, \mathbb{E}\theta_i = 0, -1 \leq \theta_i \leq 1$$

$$\text{Hoeffding's inequality} \Rightarrow \mathbb{P} \left(\left| \frac{1}{m} \sum_{i=1}^m \theta_i \right| > p \right) \leq 2 \cdot e^{-\frac{2mp^2}{m}}$$

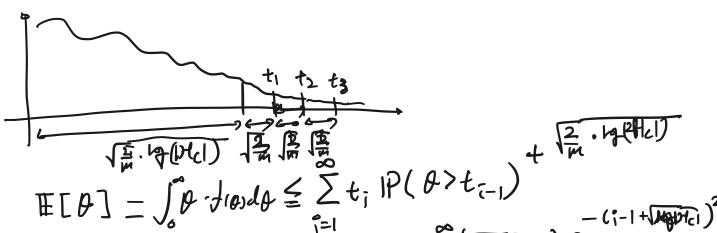
$$\text{Union Bound} \Rightarrow \mathbb{P} \left(\max_{h \in \mathcal{H}_C} \left| \frac{1}{m} \sum_{i=1}^m \theta_i \right| > p \right) \leq |\mathcal{H}_C| \cdot \mathbb{P} \left(\left| \frac{1}{m} \sum_{i=1}^m \theta_i \right| > p \right) \leq 2|\mathcal{H}_C| \cdot e^{-\frac{2mp^2}{m}}$$

$$(*) \Rightarrow \mathbb{E} \left[\max_{h \in \mathcal{H}_C} \left| \frac{1}{m} \sum_{i=1}^m (\hat{L}_h(z'_i) - \hat{L}_h(z_i)) \right| \right] \leq \frac{4 + \sqrt{\log(|\mathcal{H}_C|)}}{\sqrt{m/2}} \leq \frac{4 + \sqrt{\log(\gamma_1(2^{2m}))}}{\sqrt{m/2}}$$

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

Lemma A.4 back

Proof P.D.F



$$\mathbb{E}[\theta] = \int_0^\infty \theta \cdot f(\theta) d\theta \leq \sum_{i=1}^{\infty} t_i \cdot \mathbb{P}(\theta > t_{i-1}) + \frac{\sqrt{2} \cdot \sqrt{\log(|\mathcal{H}_C|)}}{\sqrt{m/2}}$$

$$\leq 2\sqrt{\frac{2}{m}} \cdot |\mathcal{H}_C| \cdot \sum_{i=1}^d \left(\sqrt{4\eta(\mathcal{H}_C)} + 1 \right) e$$

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

Proof of Sauer's Lemma >

We will show that $\forall d$

def of VC dimension of H

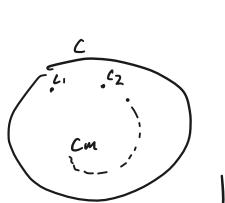
$$|\mathcal{H}_C| \stackrel{(*)}{\leq} |\{B \subseteq C : H \text{ shatters } B\}| \leq \sum_{i=0}^d \binom{m}{i}$$

because largest set that can be shattered is d .

Proof of $(*)$ by induction

- $m=1$: n.t.s $d=0 \rightarrow |\mathcal{H}_C| \leq 1$. trivial since $d < 1 \rightarrow$ does not shatter when $m=1$.
 $d=1 \rightarrow |\mathcal{H}_C| \leq 2 = 2^m$

- Suppose $(*)$ holds for sets $C' \subset C$, n.t.s. $(*)$ holds for set size m .



$$\mathcal{H}_C = \underbrace{\dots}_{0} \underbrace{\dots}_{1} \underbrace{\dots}_{2} \dots$$

$\forall B \subseteq C$ that H shatters B .

$$|\mathcal{H}_C| = |\mathcal{Y}_0| + |\mathcal{Y}_1|$$

$\underbrace{\quad}_{\text{ass } H}, \quad \underbrace{\quad}_{\{ c_i \in B, H \text{ shatters } B \}}$