

# Fundamental Theorem of Statistical Learning

STATS 303 Statistical Machine Learning

Spring 2022

Lecture 22

# Fundamental Theorem of Statistical Learning

Let  $\mathcal{H}$  be a hypothesis class of functions from a domain  $\mathcal{X}$  to  $\{0,1\}$  and let the loss function be 0-1 loss. Then the following statements are equivalent:

- 
1.  $\mathcal{H}$  has UCP
  2. Any ERM is a successful agnostic PAC learner for  $\mathcal{H}$
  3.  $\mathcal{H}$  is agnostic PAC learnable
  4.  $\mathcal{H}$  is PAC learnable
  5. Any ERM is a successful PAC learner for  $\mathcal{H}$
  6.  $\mathcal{H}$  has a finite VC-dimension

Idea of proof:

- If  $\text{VCdim}(\mathcal{H}) = d$ , then the "effective size" of  $\mathcal{H}$  is small.
- Small "effective size"  $\Rightarrow$  UCP

### Sauer's Lemma

Def. (growth function) Let  $\mathcal{H}$  be a hypothesis class. The growth function of  $\mathcal{H}$ , denoted by

$$T_{\mathcal{H}}: \mathbb{N} \rightarrow \mathbb{N},$$

is defined by

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

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

$$m \leq d \Rightarrow T_{\mathcal{H}}(m) = 2^m.$$

Lemma (Sauer-Shelah-Pearles) Let  $\mathcal{H}$  be a hypothesis class with  $\text{VCdim}(\mathcal{H}) \leq d < \infty$ . Then for all  $m \in \mathbb{N}$ ,

$$\tau_{\mathcal{H}}(m) \leq \sum_{i=0}^d \binom{m}{i}$$


Remark: In particular, if  $m \geq d+1$ , then this implies


$$\tau_{\mathcal{H}}(m) \leq \left(\frac{em}{d}\right)^d$$

following Stirling's approximation.

Proof of Sauer's Lemma:

We prove a stronger claim:

  $\left\{ \begin{array}{l} \text{For any } C = \{c_1, \dots, c_m\} \subset \mathcal{X}, \text{ for any } \mathcal{H}, \\ |\mathcal{H}_C| \leq |\{B \subset C : \mathcal{H} \text{ shatters } B\}| \end{array} \right\}$

If we can prove , then if  $\text{VCdim}(\mathcal{H}) \leq d$ , then no set with size larger than  $d$  can be shattered by  $\mathcal{H}$ .

Therefore,  $|\{B \subset C : \mathcal{H} \text{ shatters } B\}| \leq \underbrace{\sum_{i=0}^d \binom{m}{i}}_{\text{total \# of subsets with size } \leq d}$ .

To prove  $\textcircled{AA}$ , we use mathematical induction.

- For  $m=1$ ,  $C = \{c_1\}$ . There are two cases:
  - $|H_C| = 1$ . Only the empty set  $\emptyset$  is shattered by  $H$ .
  - $|H_C| = 2$ . Both  $\emptyset$  and  $C$  are shattered by  $H$ .

In both cases,  $|H_C| = |\{B \subset C : B \text{ is shattered by } H\}|$ .

Therefore,  $\textcircled{AA}$  holds for  $m=1$ .

- Assume  $\textcircled{AA}$  holds for all sets of size  $< m$ . We are going to prove that  $\textcircled{AA}$  also holds for  $m$ .


Fix  $H$ . Let  $C = \{c_1, \dots, c_m\}$ .


Denote  $C' = \{c_2, \dots, c_m\}$

In addition, define the following two sets:

$$\gamma_0 = \{ (\overset{\text{"labels"}}{y_1}, y_2, \dots, y_m) : (0, y_2, \dots, y_m) \in H_C \text{ or } (1, y_2, \dots, y_m) \in H_C \}$$

$$\gamma_1 = \{ (y_1, y_2, \dots, y_m) : (0, y_2, \dots, y_m) \in H_C \text{ and } (1, y_2, \dots, y_m) \in H_C \}$$

 Claim:  $|H_C| = |\gamma_0| + |\gamma_1|$

The reason  holds is the following:

If  $(y_1, y_2, \dots, y_m) \in H_C$ , then either

$$(y_2, \dots, y_m) \in \gamma_0 \text{ but } (y_2, \dots, y_m) \notin \gamma_1$$

or

$$(y_2, \dots, y_m) \in \gamma_0 \cap \gamma_1.$$

In the former case,  $(\overset{\uparrow}{\bullet}, y_2, \dots, y_m) \in H_C$ ; in the  
either 0 or 1.


latter case,  $(0, y_2, \dots, y_m) \in H_C$  and  $(1, y_2, \dots, y_m) \in H_C$ .

$$\text{Therefore, } |\gamma_0| + |\gamma_1| = |\gamma_0 - \gamma_1| + 2|\gamma_1| = |H_C|$$

Moreover,  $Y_0 = H_{C'}$ .

From induction,

$$|Y_0| = |H_{C'}| \leq |\{B \subset C' : H \text{ shatters } B\}|$$


$$= |\{B \subset C : C_1 \neq B \text{ and } H \text{ shatters } B\}|$$

Next, define  $H' \subset H$  to be

$$H' := \left\{ h \in H : \text{there exists } h' \in H \text{ s.t.} \right. \\ \left. (1-h'(C_1), h'(C_2), \dots, h'(C_m)) \right. \\ \left. = (h(C_1), h(C_2), \dots, h(C_m)) \right\}$$

That is,  $H'$  is the hypothesis class whose members come in pairs that differ at  $C_1$  only.

Obviously,  $H'$  shatters  $B \subset C'$

$$\Leftrightarrow H' \text{ shatters } B \cup \{C_1\}.$$

Also,  $Y_1 = H'_{C'}$


By induction,

$$|\gamma_1| = |\mathcal{H}'_{C'}| \leq |\{B \subset C' : \mathcal{H}' \text{ shatters } B\}|$$



$$= |\{B \subset C : \mathcal{H}' \text{ shatters } B \cup \{c_i\}\}|$$

$$\leq |\{B \subset C : \mathcal{H} \text{ shatters } B \text{ and } c_i \in B\}|$$


Bringing  and  together, we have

$$|\gamma_0| + |\gamma_1| \leq |\{B \subset C : c_i \neq B \text{ and } \mathcal{H} \text{ shatters } B\}|$$

+

$$|\{B \subset C : \mathcal{H} \text{ shatters } B \text{ and } c_i \in B\}|$$

$$= |\{B \subset C : \mathcal{H} \text{ shatters } B\}|$$

Noting that  $|\mathcal{H}_C| = |\gamma_0| + |\gamma_1|$  (by ) , we are done with the proof.





## UCP for classes with small "effective size"

Thm\* Let  $H$  be a hypothesis class and  $T_H$  be its growth function. Then for any distribution  $\mathbb{D}$  and any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ ,

$$|L_{\mathbb{D}}(h) - L_S(h)| \leq \frac{4 + \sqrt{\log T_H(2m)}}{\delta \sqrt{2m}}.$$

Suppose  $\text{Thm}^*$  holds. For  $m > d$ ,  $T_H(2m) \leq \left(\frac{2em}{d}\right)^d$ .

Therefore, with probability at least  $1 - \delta$ ,

$$|L_{\mathbb{D}}(h) - L_S(h)| \leq \frac{4 + \sqrt{d \log(2em/d)}}{\delta \sqrt{2m}}$$

For simplicity, let's assume  $4 \leq \sqrt{d \log(2em/d)}$ , which can be done by choosing a large  $m$ .

Then

$$|L_{\mathbb{D}}(h) - L_S(h)| \leq \frac{\sqrt{2d \log(2em/d)}}{\delta \sqrt{m}}$$

Let's choose  $m$  s.t.

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

That is,

$$\sqrt{m} \geq \frac{\sqrt{2 d \log m + 2 d \log (2e/d)}}{\varepsilon \delta}$$

or

$$m \geq \frac{2d \log m}{(\varepsilon \delta)^2} + \frac{2d \log (2e/d)}{(\varepsilon \delta)^2}$$

Lemma. Let  $a \geq 1$  and  $b > 0$ . Then

$$x \geq 4a \log(2a) + 2b \Rightarrow x \geq a \log(x) + b$$

According to this lemma, we only need to take

$$m \geq 4 \frac{2d}{(\varepsilon \delta)^4} \log \left( \frac{4d}{(\varepsilon \delta)^4} \right) + 2 \frac{2d \log (2e/d)}{(\varepsilon \delta)^2}.$$

Suppose we sample this many data points. Then with

probability at least  $1 - \delta$ ,  $|L_S(h) - L_S(h)| \leq \varepsilon$ .

Therefore, we have the ucp.

## Proof of Thm\*

First note that it suffices to prove

$$\infty \quad \mathbb{E}_{\mathcal{D}^m} \left[ \sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| \right] \leq \frac{4 + \sqrt{\log T_{\mathcal{H}}(2m)}}{\sqrt{2m}}$$

If we have proved  $\infty$ , then Markov Inequality implies

$$\begin{aligned} & \mathbb{P}_{\mathcal{D}^m} \left( \sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| > \frac{4 + \sqrt{\log T_{\mathcal{H}}(2m)}}{\delta \sqrt{2m}} \right) \\ & \leq \frac{\mathbb{E}_{\mathcal{D}^m} \left[ \sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| \right]}{\frac{4 + \sqrt{\log T_{\mathcal{H}}(2m)}}{\delta \sqrt{2m}}} \leq \frac{\frac{4 + \sqrt{\log T_{\mathcal{H}}(2m)}}{\sqrt{2m}}}{\frac{4 + \sqrt{\log T_{\mathcal{H}}(2m)}}{\delta \sqrt{2m}}} \\ & = \delta \end{aligned}$$

That is, with probability at least  $1 - \delta$ ,

$$\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| \leq \frac{4 + \sqrt{\log T_{\mathcal{H}}(2m)}}{\delta \sqrt{2m}}$$

To prove  $\circ$ , write

$$L_D(h) = \mathbb{E}_{s' \sim \mathcal{D}^m} [L_{s'}(h)]$$

Then the LHS of  $\circ$ ,  $\mathbb{E}_{\mathcal{D}^m} \left[ \sup_{h \in \mathcal{H}} |L_D(h) - L_S(h)| \right]$ ,

becomes

$$\mathbb{E}_{s \sim \mathcal{D}^m} \left[ \sup_{h \in \mathcal{H}} | \mathbb{E}_{s' \sim \mathcal{D}^m} [L_{s'}(h)] - L_S(h) | \right],$$

Note that  $| \mathbb{E}_{s' \sim \mathcal{D}^m} [L_{s'}(h)] - L_S(h) |$

$$\leq \mathbb{E}_{s' \sim \mathcal{D}^m} |L_{s'}(h) - L_S(h)|$$

$$\text{Also, } \sup_{h \in \mathcal{H}} \mathbb{E}_{s' \sim \mathcal{D}^m} |L_{s'}(h) - L_S(h)|$$

$$\leq \sup_{h \in \mathcal{H}} \mathbb{E}_{s' \sim \mathcal{D}^m} \sup_{h \in \mathcal{H}} |L_{s'}(h) - L_S(h)|$$

$$= \mathbb{E}_{s' \sim \mathcal{D}^m} \sup_{h \in \mathcal{H}} |L_{s'}(h) - L_S(h)|$$

Therefore,  $\mathbb{E}_{\mathcal{D}^m} \left[ \sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| \right]$

$$\leq \mathbb{E}_{S \sim \mathcal{D}^m} \mathbb{E}_{S' \sim \mathcal{D}^m} \sup_{h \in \mathcal{H}} |L_{S'}(h) - L_S(h)|$$

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

Let  $z_i$  be the samples in  $S$   
and  $z'_i$  be the samples in  $S'$

Due to symmetry, we can change

$$\sum_{i=1}^m [\ell(h, z'_i) - \ell(h, z_i)]$$

$$\text{to } \sum_{i=1}^m \sigma_i (\ell(h, z'_i) - \ell(h, z_i))$$

where each  $\sigma_i \in \{\pm 1\}$ .

$$= \mathbb{E}_{\sigma \sim \mathcal{U}_{\pm}^m} \mathbb{E}_{S, S' \sim \mathcal{D}^m} \sup_{h \in \mathcal{H}} \frac{1}{m} \left| \sum_{i=1}^m \sigma_i [\ell(h, z'_i) - \ell(h, z_i)] \right|$$

$$= \mathbb{E}_{S, S' \sim \mathcal{D}^m} \mathbb{E}_{\sigma \sim \mathcal{U}_{\pm}^m} \sup_{h \in \mathcal{H}} \frac{1}{m} \left| \sum_{i=1}^m \sigma_i [\ell(h, z'_i) - \ell(h, z_i)] \right|$$

Fix  $S, S'$ . Let  $C$  be the set of instances in  $S$  and  $S'$ .

Then

$$\begin{aligned} & \mathbb{E}_{\sigma \sim \mathcal{U}_{\pm}^m} \sup_{h \in \mathcal{H}} \frac{1}{m} \left| \sum_{i=1}^m \sigma_i [\ell(h, z'_i) - \ell(h, z_i)] \right| \\ &= \mathbb{E}_{\sigma \sim \mathcal{U}_{\pm}^m} \sup_{h \in \mathcal{H}_C} \frac{1}{m} \left| \sum_{i=1}^m \sigma_i [\ell(h, z'_i) - \ell(h, z_i)] \right| \end{aligned}$$

Fix  $h \in \mathcal{H}_C$ . Denote  $\Theta_h = \frac{1}{m} \sum_{i=1}^m \sigma_i [\ell(h, z'_i) - \ell(h, z_i)]$

Since  $\mathbb{E}_{\sigma \sim \mathcal{U}_{\pm}^m} [\Theta_h] = 0$ , by Hoeffding's Inequality,

$$\mathbb{P}(|\Theta_h| > \rho) \leq 2 \exp(-2m\rho^2).$$

By a union bound,

$$\mathbb{P}\left(\max_{h \in \mathcal{H}_C} |\Theta_h| > \rho\right) \leq 2 |\mathcal{H}_C| \exp(-2m\rho^2).$$

Therefore, following some Calculus,

$$\begin{aligned} \mathbb{E}_{\sigma \sim \mathcal{U}_{\pm}^m} \left[ \max_{h \in \mathcal{H}_C} |\Theta_h| \right] &\leq \frac{4 + \sqrt{\log(\mathcal{H}_C)}}{\sqrt{2m}} \\ &\leq \frac{4 + \sqrt{\log \mathcal{H}(2m)}}{\sqrt{2m}} \end{aligned}$$

Therefore,  $\mathbb{E}_{\mathcal{D}^m} \left[ \sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| \right]$

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



# Questions?

## *Reference*

- *FTSL*
  - *[S-S] Ch 6.4*

