

No Free Lunch Theorem

①

Recap from last time:

- X = domain set of "feature vectors"; Y = "label" set
- $Z = X \times Y$, $z = (x, y) \leftarrow$ sample.
- $S = \{ (x_1, y_1), \dots, (x_m, y_m) \}$ "Training set"
- $\mathcal{D}(x, y) = \mathcal{D}(y|x) \mathcal{D}(x)$ unknown distn generating iid samples.
- $h: X \rightarrow Y$ a rule or an hypothesis.
- $\mathcal{H} \ni h$ a class of rules or "hypothesis class"
- $l: \mathcal{H} \times Z \rightarrow \mathbb{R}_+$ loss fct
 $(h, z) \mapsto l(h, z)$
 - 0-1 loss is
 $l(h, z) = \mathbb{1}(h(x) \neq y)$
 - square loss is
 $l(h, z) = (y - h(x))^2$
- True loss, true error, gen error
 $L_{\mathcal{D}}(h) = \mathbb{E}_{z \sim \mathcal{D}} [l(h, z)]$
- Empirical loss
 $L_S(h) = \frac{1}{m} \sum_{i=1}^m l(h, z_i)$

ERM rule : $h_S = \text{ERM}(S) = \text{argmin}_h L_S(h)$

Definition : Agnostic PAC learning

\mathcal{H} is agnostic PAC learnable w.r.t \mathcal{Z}, ℓ if

$\exists m_{\epsilon, \delta} : [0, 1]^2 \rightarrow \mathbb{N} ; (\epsilon, \delta) \mapsto m_{\epsilon, \delta}(\epsilon, \delta)$ and a

learning rule A such that :

$\forall (\epsilon, \delta)$ and $\forall \mathcal{D}$ (on \mathcal{Z}), when running

A on $m > m_{\epsilon, \delta}(\epsilon, \delta)$ iid samples $S \sim \mathcal{D}^m$

$$\text{Prob}_S \left\{ L_{\mathcal{D}}(A(S)) \leq \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \epsilon \right\} \geq 1 - \delta$$

Terminology :

• PAC = Probably Approximately Correct

• Realizable case is when $\exists h \in \mathcal{H}$ s.t $L_{\mathcal{D}}(h) = 0$ so $\min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) = 0$. We say that \mathcal{H} is PAC learnable

if $\exists m_{\epsilon, \delta}$ and A s.t $\forall (\epsilon, \delta), \forall \mathcal{D}$ realizable (i.e $\exists h \in \mathcal{H}$ with $L_{\mathcal{D}}(h) = 0$) for $m > m_{\epsilon, \delta}(\epsilon, \delta)$ we have

$$\text{Prob}_S \left\{ L_{\mathcal{D}}(A(S)) \leq \epsilon \right\} \geq 1 - \delta.$$

• Agnostic PAC learnable is more general and considers any \mathcal{D} .

③

Last time we saw that all finite hypothesis classes $|H| < +\infty$ are agnostic PAC learnable

with the ERM rule for $m > m_{\text{ERM}}^{\text{UC}}(\epsilon, \delta) = \left\lceil \frac{1}{2\epsilon^2} \log \frac{2|H|}{\delta} \right\rceil$

Of course this implies that finite hypothesis classes are also PAC learnable.

In exercises on the other hand you saw that infinite classes $|H| = +\infty$ may be PAC learnable.

For example $H =$ Threshold functions (and the loss ℓ is the 0-1 classifier).

This brings the questions:

* Are all infinite hyp classes PAC learnable?

* We may also ask a more naive and basic question: Do we need

to restrict ourselves to rules belonging to some

hypothesis class? In other words could we find a universal Algo that learns for any \mathcal{D} irrespective of how we restrict $H \in \mathcal{H}$?

③

The answer to first question will be given in next chapters and turns out to be subtle. Some infinite classes are PAC learnable when they are in some sense of effective finite size as captured by the notion of VC (Vapnik-Chernovenkis) dimension.

In this chapter we answer the second question:

From now on we assume that we are in the realizable case. So basically we consider \mathcal{H} and \mathcal{D} such that $\exists h \in \mathcal{H} : L_{\mathcal{D}}(h) = 0$.

Moreover we address the question only for $y \in \{0, 1\}$

classification problems and the loss function

$l(h; x, y) = \mathbb{1}(h(x) \neq y)$. In particular

$$L_{\mathcal{D}}(h) = \mathbb{E}_{\mathcal{D}}(\mathbb{1}(h(x) \neq y)) = \mathbb{P}(h(x) \neq y).$$



Proposition.

(e.g. $X = \mathbb{R}^N$)

Let X be an infinite domain set and let \mathcal{H} be the set of all functions $h: X \rightarrow \{0, 1\}$.

Then \mathcal{H} is not PAC learnable.

Discussion:

This proposition basically

says that some form of prior knowledge is necessary for PAC learnability. The hypothesis class should be suitably restricted (in particular the set of all functions $h: X \rightarrow \{0, 1\}$ is much too big when X is infinite).

Proof ← Read after NFL theorem

Since X is infinite we can take $C \subset X$ with $|C| = 2m$ whatever integer m size of S is. By the construction of No Free Lunch Theorem (see later) $\exists \mathcal{D}$ which is realizable (for some f) and at the same time $\mathbb{P}(L_{\mathcal{D}}(f) \geq \frac{1}{8}) \geq \frac{1}{7}$.

Take any $\epsilon < \frac{1}{8}$, $\delta < \frac{1}{7}$. We have:

$$\delta < \frac{1}{7} \leq \mathbb{P}(L_{\mathcal{D}}(f) \geq \frac{1}{8}) \leq \mathbb{P}(L_{\mathcal{D}}(f) \geq \epsilon) \Rightarrow \mathbb{P}(L_{\mathcal{D}}(f) \leq \epsilon) \leq 1 - \delta$$

which contradicts PAC learnability.

5

Theorem: No Free Lunch.

Let A be any learning rule for the task of binary classification over a finite domain X .

(0-1 loss assumed here). Let $m < \frac{|X|}{2}$ be the size of the training set $S = (x_1, \dots, x_m)$.

Then there exist a distribution \mathcal{D} over

$X \times \{0, 1\} = \mathcal{Z}$ such that:

a) There exist $f: X \rightarrow \{0, 1\}$ with $L_{\mathcal{D}}(f) = 0$

b) $\mathbb{P}(L_{\mathcal{D}}(A(S)) \geq \frac{1}{8}) \geq \frac{1}{7}$.

for some \mathcal{D} .

The theorem says that a) there exist a learner

which succeeds since we are in realizable case it suffices

to take ERM and "minimize" it over the set $\{f\}$ trivially.

b) at the same time the learner A will fail on \mathcal{D} .

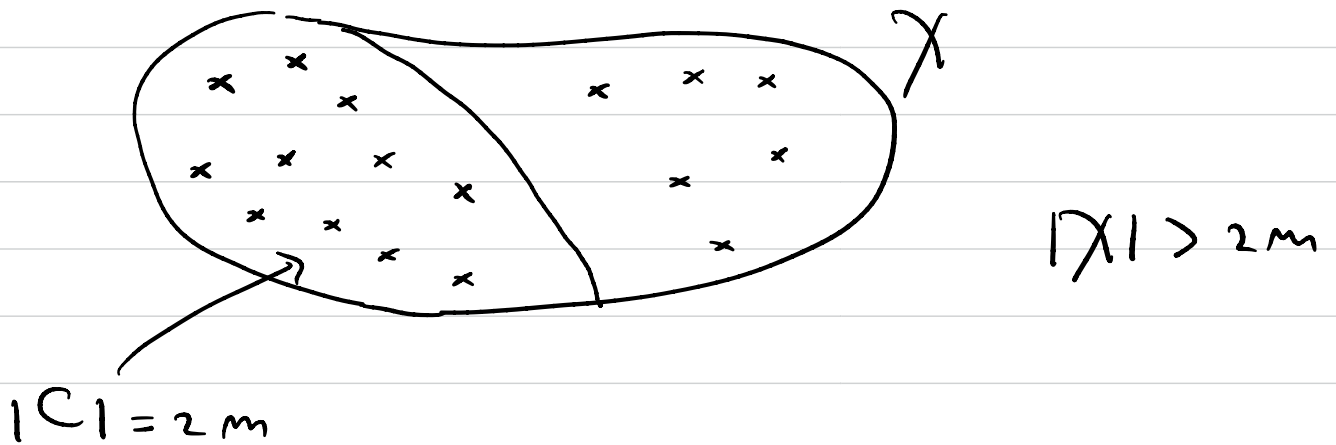
So there is no universally successful A and "No free lunch".

6

Proof.

a) we first show $\exists \mathcal{D}$ s.t. $L_{\mathcal{D}}(f) = 0$.

let $C \subset \mathcal{X}$ with $|C| = 2m$



→ Set of all functions $f_i : C \rightarrow \{0, 1\}$

has cardinality $T = 2^{2m}$

(for each input $x_1 \dots x_{2m}$ we have $2 \cdot 2 \cdot \dots \cdot 2 = 2^{2m}$ possible values that define the set).

→ Define $\mathcal{D}_i(x, y) = \begin{cases} \frac{1}{|C|} & \text{if } y = f_i(x) \\ 0 & \text{if } y \neq f_i(x) \end{cases}$

prob to choose (x, y) is $\frac{1}{|C|}$ if label is correct according to f_i .

$$\mathbb{P}_{\mathcal{D}_i} (\gamma \neq f_i(x)) = 0$$

so equivalently $\mathbb{E}_{\mathcal{D}_i} (\mathbb{1} (\gamma \neq f_i(x)))$
 $= L_{\mathcal{D}_i} (f_i) = 0.$

We see that all \mathcal{D}_i 's $i = 1 \dots 2^{2m}$
 are realizable.

This achieves to prove part a).

Now we turn to b).

We will show that

$$\star \max_{i=1 \dots 2^{2m}} \mathbb{E}_{S \sim \mathcal{D}_i^m} [L_{\mathcal{D}_i} (A(S_i))] \geq \frac{1}{4}$$

So there is some i_0 for which

$$\mathbb{E}_{S \sim \mathcal{D}_{i_0}^m} (L_{\mathcal{D}_{i_0}} (A(S_{i_0}))) \geq \frac{1}{4}$$

We will prove later (see lemma) that this implies

$$\mathbb{P}_{S \sim \mathcal{D}_{i_0}^m} (L_{\mathcal{D}_{i_0}} (A(S)) \geq \frac{1}{8}) \geq \frac{1}{7} \quad +$$

To show * we proceed as follows:

Consider all possible sequences $S = (x_1, \dots, x_m)$ with samples in C . There are $K = |C|^m$ of them and recall $|C| = 2^m$ so $K = (2^m)^m$.

Denote them S_1, S_2, \dots, S_K .

For $S_j = (x_1, \dots, x_m)$ we denote by $S_j^i = ((x_1, f_i(x_1)), \dots, (x_m, f_i(x_m)))$ the training sets with the labels corresponding to f_i .

If samples are drawn $\sim \mathcal{D}_i$ then the algorithm A receives the possible training sets $S_1^i, S_2^i, \dots, S_K^i$ with uniform prob (since the samples are drawn in C with unif prob)

$$\Rightarrow \mathbb{E}_{S \sim \mathcal{D}_i^m} [L_{\mathcal{D}_i}(A(S))] = \frac{1}{K} \sum_{j=1}^K L_{\mathcal{D}_i}(A(S_j^i))$$

(9)

$$\Rightarrow \max_{i=1 \dots 2^{2m}} \mathbb{E}_{S \sim \mathcal{D}_i} [L_{\mathcal{D}_i}(A(S))]$$

$$\geq \frac{1}{2^{2m}} \sum_{i=1}^{2^{2m}} \mathbb{E}_{S \sim \mathcal{D}_i} [L_{\mathcal{D}_i}(A(S))]$$

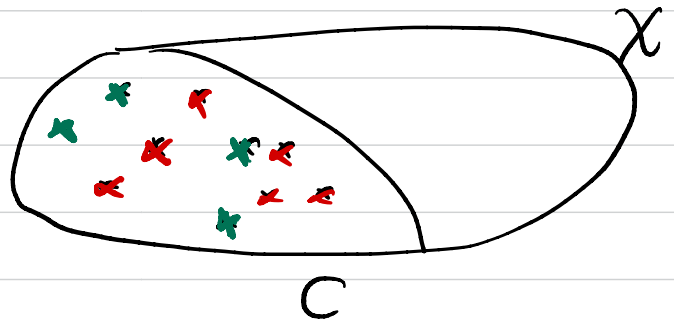
$$= \frac{1}{K} \sum_{j=1}^K \frac{1}{2^{2m}} \sum_{i=1}^{2^{2m}} L_{\mathcal{D}_i}(A(S_j^i))$$

$$\geq \min_{j=1 \dots K} \frac{1}{2^{2m}} \sum_{i=1}^{2^{2m}} L_{\mathcal{D}_i}(A(S_j^i))$$

Now fix some sequence $j \in \{1 \dots K\}$ (recall $K = |C|^m$)

$S_j = (\underline{x_1 \dots x_m}) \subset C$. Let $\underline{r_1 \dots r_p}$ be $x \in C$

that are not in S_j .



$\left\{ \begin{array}{l} p \geq m \text{ since } |C| = 2^m \\ \text{and there can be repetitions in } S_j. \end{array} \right.$

for any rule h by definition :

$$L_{\mathcal{D}_i}(h) = \frac{1}{|C|} \sum_{x \in C} \mathbb{1}(f_i(x) \neq h(x))$$

since \mathcal{D}_i has weight $\frac{1}{|C|}$ on $(x, y = f_i(x))$.

$$\Rightarrow L_{\mathcal{D}_i}(h) \geq \frac{1}{|C|} \sum_{r=1}^P \mathbb{1}(f_i(x_r) \neq h(x_r))$$

↑
take only x_r that do not appear
in S_j

$$\geq \frac{1}{2^p} \sum_{r=1}^P \mathbb{1}(f_i(x_r) \neq h(x_r))$$

since $|C| = 2^m \leq 2^p$.

$$\Rightarrow L_{\mathcal{D}_i}(A(S_j^i)) \geq \frac{1}{2^p} \sum_{r=1}^P \mathbb{1}(f_i(x_r) \neq A(S_j^i)(x_r))$$

$$\Rightarrow \frac{1}{2^{2^m}} \sum_{i=1}^{2^{2^m}} L_{\mathcal{D}_i}(A(S_j^i)) \geq \frac{1}{2^p} \sum_{r=1}^P \underbrace{\frac{1}{2^{2^m}} \sum_{i=1}^{2^{2^m}} \mathbb{1}(f_i(x_r) \neq A(S_j^i)(x_r))}_{\geq \frac{1}{2}}$$

we show below that

$$\geq \frac{1}{2}$$

(11)

$$\text{Thus } \frac{1}{2^{2m}} \sum_{i=1}^{2^{2m}} L_{\mathcal{D}_i} (A(S_j^{i'})) \geq \frac{1}{4} \quad \forall j.$$

and combining with inequ above on page 9 ;

$$\begin{aligned} & \max_{i=1, \dots, 2^{2m}} \bar{E}_{S \sim \mathcal{D}_i} [L_{\mathcal{D}_i} (A(S))] \\ & \geq \min_{j=1, \dots, K} \frac{1}{2^{2m}} \sum_{i=1}^{2^{2m}} L_{\mathcal{D}_i} (A(S_j^{i'})) \geq \frac{1}{4}. \end{aligned}$$

Now it remains to show :

$$\frac{1}{2^{2m}} \sum_{i=1}^{2^{2m}} \mathbb{1} (f_i(\nu_r) \neq A(S_j^{i'})(\nu_r)) \geq \frac{1}{2}.$$

This sum is over set of functions $f_1, f_2, \dots, f_{2^{2m}}$

There are $\frac{2^{2m}}{2}$ disjoint pairs $(f_i, f_{i'})$

such that ;

$$\forall x \in C \quad f_i(x) \neq f_{i'}(x) \quad \text{iff} \quad x = \nu_r.$$

Recall ν_1, \dots, ν_p are the $x \in C$ not seen in

S_j . The pairs are as follows;

	x	$f_i(x)$	$f_{i'}(x)$	
S_j	x_1	1	1	all possible equal sequences.
	x_2	0	0	
	\vdots	\vdots	\vdots	
	\vdots	\vdots	\vdots	
	x_m	1	1	
Not in S_j	ν_1	1	0	complementary sequences.
	ν_2	0	1	
	\vdots	\vdots	\vdots	
	\vdots	\vdots	\vdots	
	ν_p	1	0	
<hr style="width: 20%; margin: 0 auto;"/> C				

$$\frac{1}{2^{2m}} \sum_{i=1}^{2^m} \mathbb{1} \left(f_i(\nu_r) \neq A(S_j^i)(\nu_r) \right) = \frac{1}{2^{2m}} \frac{2^m}{2}$$

$$= 1 \text{ for half of } f_i$$

$$= 0 \text{ for other half of } f_i'$$

$$= \frac{1}{2}$$

This ends proof of No free lunch.



Lemma (used previously).

$$\mathbb{E}(L_{2^i}(A(s_i))) \geq \frac{1}{4} \Rightarrow \text{Prob}(L_{2^i}(A(s_i)) \geq \frac{1}{8}) \geq \frac{1}{7}$$

Proof.

It suffices to show that for a r.v. $Z \in [0, 1]$

we must have:

$$\mathbb{P}(Z \geq a) \geq \frac{\mathbb{E}(Z) - a}{1 - a}.$$

Choose $Z = L_{2^i}(A(s_i))$ & $a = \frac{1}{8}$.

The proof is an application of Markov's inequality

for a positive r.v. $X \geq 0$:

$$\mathbb{P}(X \geq b) \leq \frac{\mathbb{E}(X)}{b}.$$

Now let $X = 1 - Z$, This is ≥ 0 , Thus

(14)

$$\mathbb{P}(1 - Z \geq b) \leq \frac{1 - E(Z)}{b}$$

$$\begin{aligned} \Rightarrow \mathbb{P}(1 - Z \leq b) &\geq 1 - \frac{1 - E(Z)}{b} \\ &= \frac{b - 1 + E(Z)}{b} \end{aligned}$$

Take $b = 1 - a$. We get

$$\mathbb{P}\left(\underbrace{1 - Z \leq 1 - a}_{Z \geq a}\right) \geq \frac{E(Z) - a}{1 - a}$$

□

Remarks on Bias-Complexity Tradeoff.

Let $h_S = \text{ERM}(S) = \underset{h}{\text{argmin}} L_S(h)$

the hypothesis class given by the ERM rule.

The test error of an ERM predictor is $L_D(h_S)$.

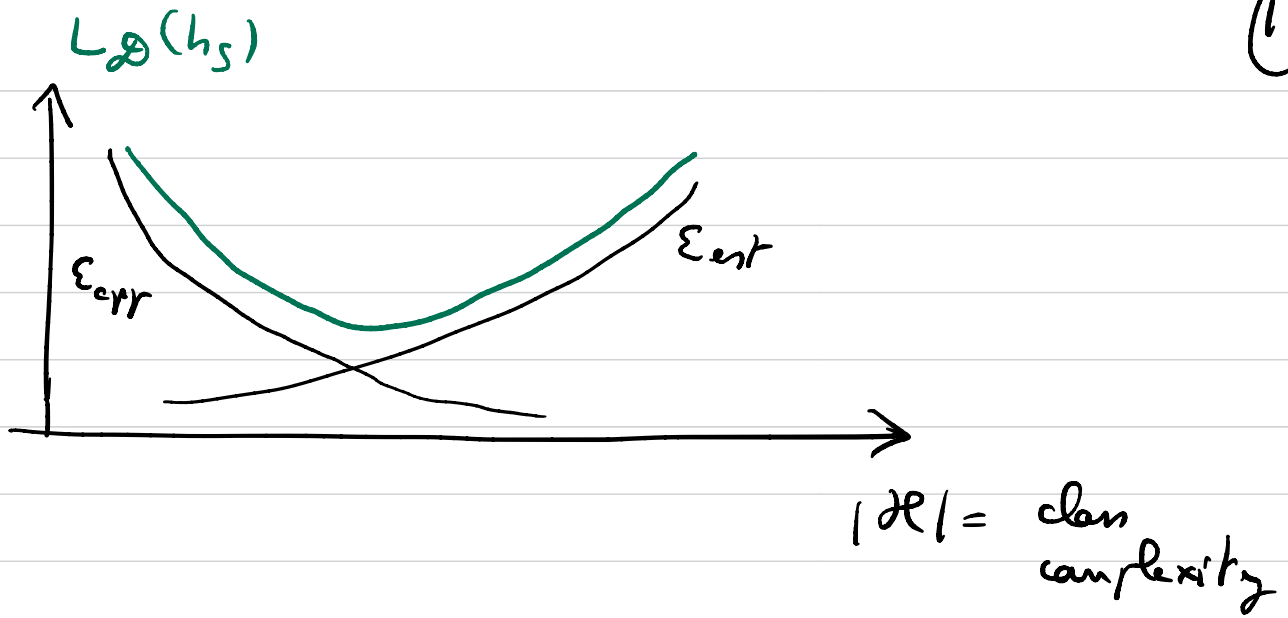
We can decompose it as:

$$L_D(h_S) = \underbrace{\min_{h \in \mathcal{H}} L_D(h)}_{\substack{\text{approx error} \\ \text{or "bias" of} \\ \text{hyp class} = \epsilon_{\text{app}}}} + \underbrace{\left(L_D(h_S) - \min_{h \in \mathcal{H}} L_D(h) \right)}_{\substack{\text{estimation error} \\ \text{of ERM.} = \epsilon_{\text{est}}.}}$$

- We have obviously that $\epsilon_{\text{app}} \rightarrow$ as $|\mathcal{H}| \nearrow$
- We also expect usually $\epsilon_{\text{est}} \rightarrow$ as $|\mathcal{H}| \nearrow$.

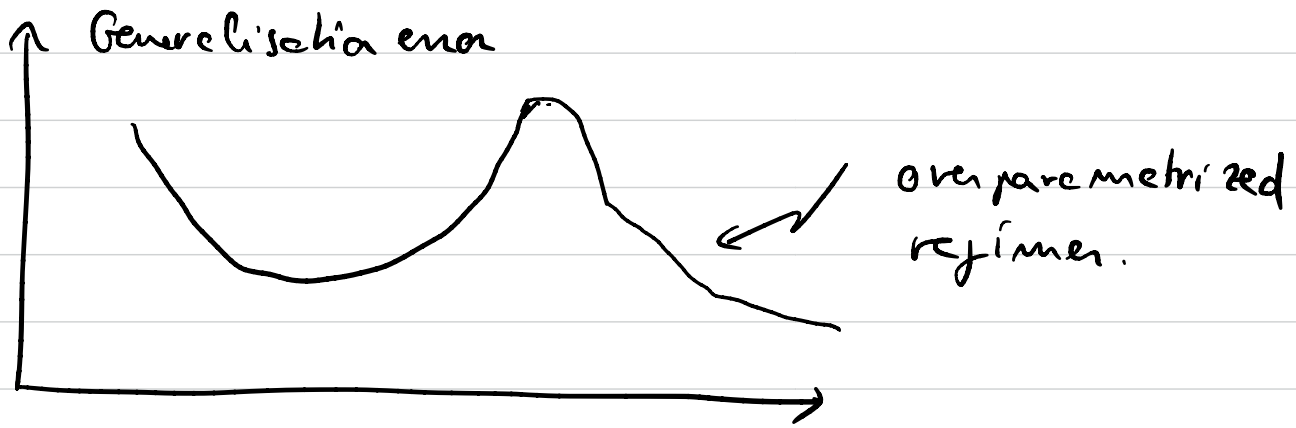
For example the generalization bound shown for finite classes

suggests $\epsilon_{\text{est}} \sim \sqrt{\frac{1}{2m} \log \frac{2|\mathcal{H}|}{\delta}} \rightarrow$ as $|\mathcal{H}| \nearrow$



This is the traditional bias-complexity trade-off picture. Increasing or decreasing our bias but at some point may lead to overfitting.

This picture has been challenged in recent years. Indeed there are situations displaying a "double descent" phenomenon



and it is recognized that the picture is much richer.