

# 8. No Free Lunches

Wednesday, January 19, 2022 2:14 PM

## §5 "Bias Complexity Tradeoff" [SS]

- ① (today) Beyond selecting features, we need some prior knowledge still  
∅ universal learner, i.e. no learner can succeed at all tasks  
Do not interpret this as evidence to "give up". Very limited, technical scope

② (next class)

Split risk into

$$\underbrace{\text{Approx. Error / bias}}_{\text{if not expressive enough}} + \underbrace{\text{Estimation Error}}_{m \text{ not large enough} \\ (\text{small } \delta \text{ helps here})}$$

### §5.1 No Free Lunch theorems

Thm "No Free Lunch" (binary classification, 0-1 loss function)

let  $A$  be any learning algorithm, with output  $h_s = A(s)$ , domain  $X$ ,  
 $Y = \{0, 1\}$ ,  $|S| = m$  and  $|X| \geq 2m$  (e.g.  $|X| = \infty$ ). Then  $\exists$  distribution  $D$  over  $X \times Y$  such that

①  $\exists f: X \rightarrow Y$  w/  $L_D(f) = 0$  benchmark/comparison:  
"It's a fair problem"

② w.p.  $\geq Y_f$  (over  $S \sim D^m$ ),  $L_D(h_s) \geq Y_f$   
i.e., cannot take  $\delta \rightarrow 0$  nor  $\epsilon \rightarrow 0$

proof (a bit subtle but useful for later)  
set, not sequence, so no duplicates

Let  $C \subseteq X$  have size  $2m$

We observe  $S$  of size  $m$ , so we'll observe no more than  $Y_f$  of  $C$   
(and this gives room for our adversary to pick tricky  $D$ )

For  $D$ ,

first, can ignore  $D$  on  $X \setminus C$

Let  $f$  be an oracle function,  $f: C \rightarrow Y := \{0, 1\}$

Since  $|C| = 2m$ ,  $|Y| = 2$ , there are  $T := 2^{2m}$  such oracle functions,  
which we enumerate  $\{f_i\}_{i=1}^T$

Each  $f_i$  is compatible with only some distributions,  
i.e., need  $D(y|x) = f_i(x)$  cond.

and among all compatible distr., choose the one w/ a uniform  
marginal ( $\text{marginal } D_x(x) = \mathbb{E}_y D((x,y)/y)$ )  
 and for  $f_i$ , call this distr.  $D_i$

$$\text{i.e., } \forall i \in [T], \quad D_i(x, y) = \begin{cases} 1/c & \text{if } y = f_i(x) \\ 0 & \text{else} \end{cases}$$

so by design,

$$L_{D_i}(f_i) = 0$$

which  $i \in [T]$  to choose?

Our "adversary" goes second: we choose algo A

which returns  $h_s = A(s)$ , then adversary chooses  $i \in [T]$

such that  $\mathbb{E}_{S \sim D_i^m} L_{D_i}(h_s) \geq \gamma_4 \quad (\#)$

$$\text{So with } D = D_i, \text{ use } \mathbb{E}_{S \sim D^m} L_D(h_s) \geq \gamma_4 \quad \left. \begin{array}{l} \text{not completely trivial} \\ \text{Markov Ineq or similar} \\ \text{See Exerc. 5.1} \end{array} \right\}$$

$$\Rightarrow \Pr[L_D(h_s) \geq \gamma_8] \geq \gamma_7$$

To show (#), want

$$\max_{i \in [T]} \mathbb{E}_{S \sim D_i^m} L_{D_i}(h_s) \geq \gamma_4$$

Fixing  $i$ , note that since  $D_x$  is uniform, for  $m=1$ ,  $\{x_i\}$ ,  
every value of  $x_i$  is equally likely... and  $m>1$ ,  $\{x_1, \dots, x_m\}$ ,  
all sequences (duplicates OK) are equally likely.

Drawing  $m$  elements (w/o replacement) from  $2m$  elements of  $C$   
 gives  $K = (2m)^m$  such sequences (all equally likely).

From  $(x_1, \dots, x_m)$  we get one  $S = ((x_1, y_1), \dots, (x_m, y_m))$   
 since  $y = f_i(x)$  (via how we defined  $D_i$ )

Enumerate all these datasets as  $S_j^i$ ,  $j \in [K]$   $i \in [T]$

Still fixing  $i \in [T]$

$$\mathbb{E}_{S \sim D_i^m} L_{D_i}(h_s) = \frac{1}{K} \sum_{j=1}^K L_{D_i}(h_{S_j^i}) \text{ since } \underline{\text{uniform}}$$

Now deal with  $i$ . We'll use  $\max \geq \text{avg}$

$$\begin{aligned}
 \max_{i \in [T]} \mathbb{E}_{s \sim D_i^m} L_{D_i}(h_s) &\geq \frac{1}{T} \sum_{i=1}^T \left( \underbrace{\frac{1}{K} \sum_{j=1}^K L_{D_i}(h_{s_j})}_{\text{avg}} \right) \\
 &= \frac{1}{K} \sum_{j=1}^K \frac{1}{T} \sum_{i=1}^T L_{D_i}(h_{s_j}) \\
 \xrightarrow{\text{avg} \geq \min} &\geq \min_{j \in [K]} \frac{1}{T} \sum_{i=1}^T L_{D_i}(h_{s_j}) \quad (\star\star)
 \end{aligned}$$

Fix  $j \in [K]$  for the moment.  $\xrightarrow{\text{m or fewer distinct items}}$

Observed  $(x_1, \dots, x_m) \subseteq C$   $|C| = 2m$   
 so we didn't observe  $\{v_1, \dots, v_p\} \subseteq C$ ,  $p \geq m$

$$\begin{aligned}
 L_{D_j}(h) &:= \mathbb{E}_{(x,y) \sim D} l(h, (x, y)) \\
 &= \mathbb{E}_{x \sim D_x} l(h, (x, f(x))) \\
 &= \mathbb{E}_{x \sim D_x} \mathbb{I}_{h(x) \neq f(x)} \quad \text{since 0-1 loss} \\
 &= \frac{1}{2m} \sum_{x \in C} \mathbb{I}_{h(x) \neq f(x)} \\
 &\geq \frac{1}{2m} \sum_{r=1}^p \mathbb{I}_{h(v_r) \neq f(v_r)} \quad \text{since loss is non-negative} \\
 &\geq \frac{1}{2p} \sum_{r=1}^p \mathbb{I}_{h(v_r) \neq f(v_r)} \quad \text{since } p \geq m
 \end{aligned}$$

Looking back at  $(\star\star)$ ,

$$\begin{aligned}
 \frac{1}{T} \sum_{i=1}^T L_{D_i}(h_{s_j}) &\geq \frac{1}{T} \sum_{i=1}^T \frac{1}{2p} \sum_{r=1}^p \mathbb{I}_{h_j(v_r) \neq f(v_r)} \\
 &= \frac{1}{2p} \sum_{r=1}^p \underbrace{\frac{1}{T} \sum_{i=1}^T \mathbb{I}_{h_j(v_r) \neq f(v_r)}}_{\text{examine this } (\dagger\ddagger\ddagger)} \quad \text{swap sum}
 \end{aligned}$$

Fix  $\Gamma$  for the moment

Look at all  $f_i$ ,  $i \in [T]$ ,  $T = 2^{2m}$  (an even number)

All functions from  $C$  to  $\{0, 1\}$

Ex  $C = \{a, b, c\}$  where  $c = v_r$  (unobserved)  
 $y = \{0, 1\}$  so  $T = 2^3$  functions

	$f(a)$	$f(b)$	$f(c)$	
$i=1$	0	0	0	{ pair }
2	0	0	1	
3	0	1	0	{ pair }
4	0	1	1	
5	1	0	0	{ pair }
6	1	0	1	
7	1	1	0	{ pair }
8	1	1	1	

Write all  $T$  indices as  $T/2$  pairs  $(i, i')$

$$\text{where } f_i(a) = f_{i'}(a)$$

$$f_i(b) = f_{i'}(b)$$

$$f_i(c) \neq f_{i'}(c)$$

or more generally,  $f_i(x) = f_{i'}(x) \quad \forall x \in C \setminus \{v_r\}$   
 $\quad \quad \quad (\dagger) f_i(x) \neq f_{i'}(x) \quad x = v_r$

(think of binary notation,  $v_r$  index as least significant bit)

back to our term  $(\dagger\dagger\dagger)$

$$\frac{1}{T} \sum_i^T \mathbb{1}_{h(v_r) \neq f_i(v_r)} = \frac{1}{T} \left( \sum_{i=1}^{T/2} \mathbb{1}_{h(v_r) \neq f_i(v_r)} + \sum_{i'=1}^{T/2} \mathbb{1}_{h(v_r) \neq f_{i'}(v_r)} \right)$$

↑ paired up ↓

$(\dagger)$  implication  $f_i(v_r) = 0$  and  $f_{i'}(v_r) = 1$

or  $f_i(v_r) = 1$  and  $f_{i'}(v_r) = 0$ .

Now, is  $h_j^i = h_j^{i'}$ ?

At first, not necessarily since ones is

based on data  $S_j^i$ , the other on  $S_j^{i'}$

recall we've enumerated all of  
these, so not random.

where  $i$  controls the labels via  $f_i$

But, based on our pairs  $(i, i')$ ,

$f_i$  and  $f_{i'}$  agree everywhere but  $v_r$  and  
 $v_r$  is unobserved. So  $S_j^i = S_j^{i'}$  hence

$$h_j^i = h_j^{i'}$$

$$\begin{aligned}
 (\text{***}) &= \frac{1}{T} \left( \sum_{i=1}^{T/2} \mathbb{1}_{h_j^i(v_r) \neq f_i(v_r)} + \sum_{i=1}^{T/2} \mathbb{1}_{h_j^i(v_r) \neq f_i(v_r)} \right) \\
 &= \frac{1}{T} \left( \sum_{i=1}^{T/2} 1 \right) = \frac{1}{2}
 \end{aligned}$$

so exactly 1 term  
is 1, the other 0

Put it altogether:

$$\max_{i \in [T]} \mathbb{E}_{S \sim D_i^m} L_{D_i}(h_s) \geq \min_{j \in [K]} \frac{1}{T} \sum_{i=1}^T L_{D_i}(h_{S_j^i}) \quad (\text{**})$$

$$\begin{aligned}
 &\geq \min_{j \in [K]} \frac{1}{2P} \sum_{r=1}^P \underbrace{\frac{1}{T} \sum_{i=1}^T \mathbb{1}_{h_j^i(v_r) \neq f_i(v_r)}}_{= \frac{1}{2}} \quad (\text{***}) \\
 &= \frac{1}{2P} \sum_{r=1}^P \frac{1}{2} \\
 &= \frac{1}{4}
 \end{aligned}$$

$\square$

no more  $j$  dependence

## Implications of No Free Lunch

If  $|X| = \infty$  and  $\mathcal{H} = \{ \text{all functions from } X \text{ to } Y \}$   
then our reference  $f$  is in  $\mathcal{H}$ , so...  $\mathcal{H}$  isn't PAC learnable.

(Recall: PAC learnable,  $\exists A$  s.t.  $\forall \epsilon, \delta \in (0, 1)$ ,  $\exists m$  s.t. if  $|S| \geq m$  iid samples  
& distr.  $D$ , then w.p.  $\geq 1 - \delta$ ,  $L_D(h_s) \leq \epsilon + \min_{h \in \mathcal{H}} L_D(h_s)$ )

So we need  $\mathcal{H} \not\subseteq \mathcal{Y}^X$ , i.e., need some prior knowledge to pick  $\mathcal{H}$   
(can't avoid inductive bias)