

# CS492(D) Computational Learning Theory. 08 / 29 / 2023

§ 1.

## Objectives.

- Study basic ML algorithms.
- Study math techniques for analysing those algorithms.

§ 2

## Preview.

$\mathcal{X}$  ... input space.

$\mathcal{Y}$  ... output space.

Unknown distribution  $D \sim \Pr(x,y)$

①

Given loss function  $L$ ,  $R(h) = \underbrace{\mathbb{E}[L(ha), y]}_{(a,y) \sim D}$

✓  
e.g. quadratic loss

↳ Generalization error(risk)  
Population risk(?)

We want  $\arg\min_{h \in \mathcal{H}} R(h)$ , but we don't know  $D$ .

$\mathcal{H}$   
 $[x \rightarrow y]$

②

Then, do what?

Sample from  $D$  is available.

-  $S = ((x_1, y_1), \dots, (x_m, y_m)) \in (\mathcal{X} \times \mathcal{Y})^m$

- Proxy objective of  $R(h)$ , using  $S$ .

Then solve optimization problem with this objective.

- Eg.  $\hat{R}_S(h) := \frac{1}{m} \sum_{i=1}^m (h(x_i) - y_i)^2 = \mathbb{E}_{(x,y) \sim \text{Unif}(S)} [(ha) - y]^2$

- Empirical risk optimization  $\arg\min_{h \in \mathcal{H}} \hat{R}_S(h)$

What is relationship btwn  $\hat{R}_S(h)$  &  $R(h)$ ?

Result for regression:

**Theorem** Under some conditions,  $\exists c_1, c_2, c_3$  s.t.

$$\forall \delta > 0, P[\forall h \in H \quad R(h) \leq \hat{R}_S(h) + c_1 \sqrt{\frac{\log C_m}{m}} + c_2 \sqrt{\frac{\log(1/\delta)}{m}}] \geq 1 - \delta.$$

③

How do we solve  $\hat{R}_S(h)$  quickly?

- Design a good hypo. set  $H$ .

$$- H_1 = \{f \in [x \rightarrow y] \mid f(w) = \langle w, \Phi(\alpha) \rangle \text{ for } w \in \mathbb{R}^F\}$$

$$\Phi: \mathbb{R}^I \rightarrow \mathbb{R}^F.$$

→ Traditional regression

↑ Strong duality.

$$- H_2 = \{f \in [x \rightarrow y] \mid f(x) = \langle \alpha, (y_i K(x, x_i))_{i \in [m]} \rangle, \alpha \in \mathbb{R}^m\}$$

$$K: \mathbb{R}^I \times \mathbb{R}^I \rightarrow \mathbb{R}\}$$

→ Kernel (Nonparametric approach).

§3

## Logistics

\*

Ch 2-6 of the textbook.

\*

Final 40%.

Group Project 40%.

HW (2-3 Problem sheet). 20%.

• (i) 2-4 people

| (ii) Pick a paper COLT 2019-2023

| (iii) Report ≤ 4 pages

• (iv) 35-min talk with slides

(only 4 teams present with extra +5 marks).

# Ch2 PAC learning framework

08 / 31 / 2023

§1.

## Classification Problem

$\mathcal{X}$ : Input Space ( $\mathbb{R}^n$  most of the time)

$\mathcal{Y}$ : Output Space ( $\{0, 1\}$ )

$$c \in [\mathcal{X} \rightarrow \mathcal{Y}]$$

$$D \in \mathcal{P}_r(\mathcal{X}), S = (x_1, y_1), \dots, (x_m, y_m) \in (\mathcal{X} \times \mathcal{Y})^m$$

$$H \subseteq [\mathcal{X} \rightarrow \mathcal{Y}]$$

$$x_i \stackrel{\text{iid}}{\sim} D, y_i = c(x_i)$$

Goal: Find a good  $h \in H$  s.t.  $h \approx c$  using  $\mathcal{X}, \mathcal{Y}, S$

§2

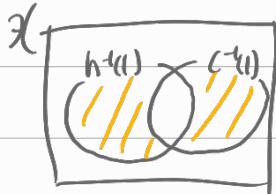
## Generalization Error, Empirical Error

$$R(h)$$

$$\hat{R}_S(h)$$

Generalization Err (i) For  $h \in H$ ,  $R(h) := \mathbb{E}_{x \sim D} [\mathbb{1}_{\{h(x) \neq c(x)\}}] = \mathbb{P}_{x \sim D}(h(x) \neq c(x))$

Can understand  $h, c$  as subset of input:  $h^{-1}(1), c^{-1}(1)$



Empirical Err

$$(ii) \hat{R}_S(h) = \frac{1}{m} \sum_{i=1}^m [\mathbb{1}_{\{h(x_i) \neq c(x_i)\}}]$$
$$= \mathbb{E}_{x_i \sim \text{Unif}(S)} [h(x_i) \neq c(x_i)]$$

$$(iii) \text{Prove } \mathbb{E}_{S \sim D^m} [\hat{R}_S(h)] = R(h).$$

By linearity of expectation, it suffices to show  $(*) = \mathbb{E}_{S \sim D^m} [\mathbb{1}_{\{h(x_i) \neq c(x_i)\}}] = R(h)$

$$(*) = \underset{(x_i, y_i) \sim D}{\mathbb{E}} [I_{h(x_i) \neq c(x_i)}] = R(h) \text{ by definition.}$$

§3

## PAC Learnability

PAC: Probably Approximately Correct

① Setup

$$x, y, H, C \subseteq \mathbb{R}^n \rightarrow \mathbb{R}$$

- $C$  is PAC-learnable by  $H$  if  $\exists$  an algorithm  $\mathcal{A}$  and a polynomial  $p(\cdot, \cdot)$  such that

$$\forall C \in C, \forall D \in \Pr(x), \forall \varepsilon, \delta > 0, \exists m \in \mathbb{N}$$

$$m \geq p\left(\frac{1}{\varepsilon}, \frac{1}{\delta}, n\right) \Rightarrow \underset{S \sim D^m}{\mathbb{P}} [R(h_{A,S}) \leq \varepsilon] \geq 1 - \delta.$$

(in polynomial scale)

i.e. If we see enough samples ( $m$ ), with high probability, the algorithm produce good answers.

- $C$  is efficiently PAC-learnable if  $\exists$  polynomial-complexity algorithm  $\mathcal{A}$  that makes  $C$  PAC-learnable.

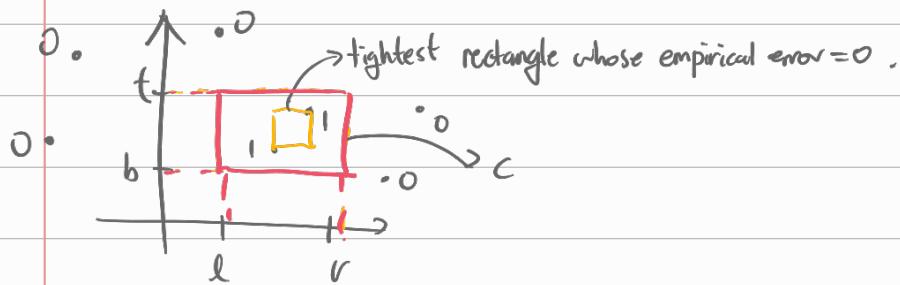
Remarks

- This should work for all distribution  $D$
- If  $H \not\subseteq C$ , then  $C$  is PAC is not PAC-learnable.  
For any algorithm  $\mathcal{A}$ ,  
Take  $c \notin H$ .  $\forall h \in H$ , we have  $c \neq h$ .  
Then,  $\exists x \in X$  s.t.  $c(x) \neq h(x)$ .  
Let

## §4. Examples - Learning Axis Aligned Rectangles

### ① Setup

$$\mathcal{X} = \mathbb{R}^2, \mathcal{Y} = \{0, 1\}, \mathcal{H} = \mathcal{C} = \{(a, b) \mapsto \mathbb{1}_{(a, b) \in [l, r] \times [b, t]} \mid l \leq r, b \leq t\}$$



### ② Polynomial p

$$P(\frac{l}{t}, \frac{b}{t}) = \frac{4}{t} \log \frac{4}{\delta}$$

### ③ Why PAC

Let  $c \in \mathcal{C}, D \in \mathcal{P}_1(\mathbb{R}), \varepsilon, \delta > 0$  be given,  $m = \frac{4}{\varepsilon} \log \frac{4}{\delta}$ .

Goal:  $\underset{s \sim D^m}{P}[R(h_s) \leq \varepsilon] \geq 1 - \delta$ .

$$\Leftrightarrow \underset{s \sim D^m}{P}[R(h_s) > \varepsilon] < \delta.$$

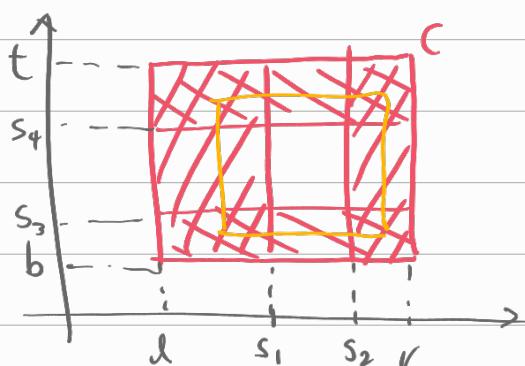
(pf)

(Case 1)  $\underset{s \sim D^m}{P}[x \in c] \leq \varepsilon$ . As  $\underset{\substack{h \in \mathcal{C} \\ \text{By algorithm}}}{h \subset c}$  (In sense  $h^{-1}(1) \subseteq h^{-1}(c)$ ),

$$h \setminus c (= h \setminus c) \cup c \setminus h = (h^{-1}(1) \setminus c^{-1}(1)) \cup (c^{-1}(1) \setminus h^{-1}(1)) \text{ has}$$

$$P(h \setminus c) = R(h) < \varepsilon.$$

(Case 2)  $\underset{s \sim D^m}{P}[x \in c] > \varepsilon$ .



Take  $s_1, s_2$  s.t.

$$P(x \in [l, s_1] \times [b, t]) \geq \frac{\varepsilon}{4}$$

$$P(x \in [l, s_2] \times [b, t]) \leq \frac{\varepsilon}{4}.$$

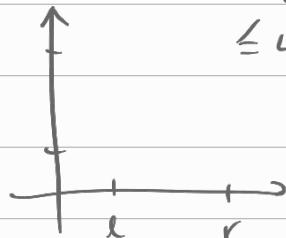
Find  $s_2, s_3, s_4$  similarly.  
 If  $h$  lives in  $\mathcal{H}$ , then Boundary of  $h$  is included in  $\mathcal{H}$   
 $(\Leftrightarrow h \cap r_i = \emptyset \forall i \in [4])$

$$R(h) \leq \Pr(\mathcal{H}) \leq 4 \cdot (\varepsilon/4) = \varepsilon$$

(Union bound.)

Hence, it suffices to choose  $\delta$  produce  $h_s$  in  $\mathcal{H}$  with prob  $\geq 1 - \delta$ .

$$\begin{aligned} \Pr(R(h_s) > \varepsilon) &\leq \Pr(\exists i \in [4] \text{ s.t. } h_s \cap r_i \neq \emptyset) \\ &\leq \sum_{i=1}^4 \Pr(h_s \cap r_i \neq \emptyset) \leq \sum_{i=1}^4 (1 - \varepsilon/4)^m \\ &\leq 4e^{-\varepsilon m/4} = 4e^{-\log \frac{1}{\delta}} = \delta \end{aligned}$$



### Theorem

Let  $\mathcal{H} \subseteq \mathcal{X} \rightarrow \mathcal{Y}$  be finite with  $C \subseteq \mathcal{H}$ .

Let  $\mathcal{A}$  be an algorithm s.t.  $\forall c \in C$  and  $S = ((x_1, y_1), \dots, (x_m, y_m))$  for  $c$ ,  $\hat{R}_S(h_c) = 0$ .

Then,  $\forall \varepsilon, \delta > 0$ ,  $\Pr_{S \sim D^m}[R(h_c) \leq \varepsilon] \geq 1 - \delta$  if  $m \geq \frac{1}{\varepsilon}(\log |\mathcal{H}| + \log \frac{1}{\delta})$ .

### Example

① Conjunction of Boolean literals.

$\mathcal{X} = \{0, 1\}^n$ ,  $C = \mathcal{H} = \{x_1 \wedge x_2, x_2 \wedge x_3, \dots, \text{?0? false?}\}$

$|\mathcal{H}| = 3^n \Rightarrow \log |\mathcal{H}| = n \log 3 \Rightarrow$  PAC learnable.

Efficiently PAC learnable. ( $O(n)$  algorithm).

② Universal Concept Class

$C = \mathcal{H} = \mathcal{X} \rightarrow \mathcal{Y}$ . Cannot classify using them.

If  $C$  is PAC learnable, every  $C' \subseteq C$  is also PAC learnable.

$|\mathcal{H}| = |\mathcal{Y}^{\mathcal{X}}|$ .  $\log |\mathcal{H}| = |\mathcal{X}| \log |\mathcal{Y}|$ .

### ③ k-term DNF formulas

k-disjunction of conjunction  
 $|H| = (3^n)^k = 3^{nk}$  ( $|H| = \sum_{i=1}^k 3^{n_i}$ )

$$( ) \vee ( \underbrace{ }_{\substack{s \\ \vdots \\ k}} ) \vee ( \underbrace{ }_{\substack{s \\ \vdots \\ j}} ) \vee ( )$$

$\log |H| = nk \log 3 \Rightarrow$  PAC learnable using them.  
 Not eff. pac learnable unless RP = NP

### ④ k-term CNF formulas.

$$( ) \wedge ( \underbrace{ }_{\substack{k \\ \vdots \\ 1}} ) \wedge ( ) \dots$$

CNF  $\Leftrightarrow$  DNF ?

$$2^{\dots} \Rightarrow (2n+1)^k \text{ variables}$$

$$|H| = 3^{\# \text{ of var.}} \Rightarrow$$
 PAC learnable.

Efficiently PAC learnable.

(pf of theorem) Goal:  $\underset{S \sim D^m}{P} [R(h_s) > \varepsilon] \leq \delta$ .

$$\begin{aligned} & \underset{S \sim D^m}{P} [\exists h \text{ s.t. } \hat{R}_S(h) = 0 \wedge R(h) > \varepsilon] \\ & \leq P_S \left[ \bigvee_{\substack{h \in H \\ R(h) > \varepsilon}} \hat{R}_S(h) = 0 \right] \end{aligned}$$

$$\leq \sum_{\substack{h \in H \\ R(h) > \varepsilon}} P_{S \sim D^m} [R_S(h) = 0]$$

$$\leq \sum_{\substack{h \in H \\ R(h) > \varepsilon}} (1 - \varepsilon)^m \leq |H| (1 - \varepsilon)^m \leq |H| e^{-m\varepsilon}$$



Thus, for  $m \geq \frac{1}{\varepsilon} (\log |H| + \log \frac{1}{\delta})$ ,

$$\delta \leq |H| e^{-(\log |H| + \log \frac{1}{\delta})} = \delta.$$

Definition

[Agnostic PAC Learnable]

"

if



$$m \geq P\left(\frac{1}{\epsilon}, \frac{1}{\delta}, n\right) \Rightarrow \Pr_{S \sim D^n} [R(h_S) \leq \inf_{h \in H} R(h) + \epsilon] \geq 1 - \delta.$$

Theorem

$\forall$  algorithm s.t.  $H \subseteq \mathcal{C}$  and  $\forall S = ((x_1, y_1), \dots, (x_m, y_m)) \in \mathcal{C}$ ,

$$\hat{R}_S(h_S) = \min_{h \in H} \hat{R}_S(h). \quad (\text{ERM})$$

$\Rightarrow \forall C \subseteq \mathcal{C} \ \forall D \in \Pr(S) \ \forall \epsilon, \delta > 0, \ \forall m \in \mathbb{N}$ ,

if  $m \geq \frac{2}{\epsilon^2} (\log |H| + \log \frac{2}{\delta})$ , then  $\Pr_{S \sim D^n} [R(h_S) \leq \min_{h \in H} R(h) + \epsilon] \geq 1 - \delta$

Lemma

(c.f. proposition 4.1)

Denote  $h_S$  above as  $h^{ERM}$ .

$$\Pr [ \exists_{h \in H} |R(h) - \hat{R}_S(h)| > \frac{\epsilon}{2} ]$$

$$\Pr_{S \sim D^n} [R(h^{ERM}) \geq \min_{h \in H} R(h) + \epsilon] \leq \Pr_{S \sim D^n} \left[ \max_{h \in H} |R(h) - \hat{R}_S(h)| > \frac{\epsilon}{2} \right]$$

Thm 2.13'

$$m \geq \frac{1}{2\epsilon^2} (\log |H| + \log \frac{2}{\delta}). \Rightarrow \Pr_{S \sim D^n} [\forall h \in H \quad |R(h) - \hat{R}_S(h)| \leq \epsilon] \geq 1 - \delta$$

$\uparrow$

$$\Pr [ \exists_{h \in H} |R(h) - \hat{R}_S(h)| > \epsilon ] \leq \delta$$

(pf)

$$m \geq \frac{1}{2\epsilon^2} ( \quad ) \Rightarrow \Pr [ \exists_{h \in H} \text{ s.t. } |R(h) - \hat{R}_S(h)| > \frac{\epsilon}{2} ] \leq \delta$$
$$\Rightarrow \Pr [R(h^{ERM}) \geq \min_{h \in H} R(h) + \epsilon] \leq \delta. \blacksquare$$

Definition (weakly PAC learnable)  
 ≡ Agnostic PAC learnable.

(pt of proposition 4.1) Let  $h_0 \in \arg \min_{h \in H} R(h)$ .

$$\begin{aligned}
 \text{Then, } R(h_s^{\text{ERM}}) - \min_{h \in H} R(h) &= R(h_s^{\text{ERM}}) - R(h_0) \\
 &\leq R(h_s^{\text{ERM}}) - R(h_0) \\
 &\leq R(h_s^{\text{ERM}}) - \hat{R}(h_s^{\text{ERM}}) + \hat{R}(h_s^{\text{ERM}}) - \hat{R}(h_0) + \hat{R}(h_0) - R(h_0) \\
 &\leq 2 \max_{h \in H} |R(h) - \hat{R}(h)| + \hat{R}(h_s^{\text{ERM}}) - \hat{R}(h_0) \\
 &\quad \text{by def of } h_s^{\text{ERM}}
 \end{aligned}$$

$$\text{Thus, } \mathbb{P}[R(h_s^{\text{ERM}}) - \min_{h \in H} R(h) > \varepsilon] \leq \mathbb{P}\left[\max_{h \in H} |R(h) - \hat{R}_s(h)| > \frac{\varepsilon}{2}\right]$$

(pt of theorem 2.13(b)) Hoeffding's Inequality (theorem D.2)

If  $X_1, \dots, X_m$ : independent R.V. s.t.  $X_i \in [a_i, b_i]$ .  $\forall i$ .

$$\begin{aligned}
 \text{For, } A := \frac{1}{m} \sum_{i=1}^m X_i, \varepsilon > 0, \quad &2m\varepsilon^2 \\
 \mathbb{P}[A - \mathbb{E}[A] \geq \varepsilon] &\leq e^{-\frac{2m\varepsilon^2}{\sum_{i=1}^m (b_i - a_i)^2}} \\
 \mathbb{P}[A - \mathbb{E}[A] \leq -\varepsilon] &\leq e^{-\frac{2m\varepsilon^2}{\sum_{i=1}^m (b_i - a_i)^2}} \\
 \Rightarrow \mathbb{P}[|A - \mathbb{E}[A]| \geq \varepsilon] &\leq 2e^{-\frac{2m\varepsilon^2}{\sum_{i=1}^m (b_i - a_i)^2}}.
 \end{aligned}$$

$$\mathbb{E}_{S \sim D^m} [\hat{R}_s(h)] = R(h).$$

$$\hat{R}_s(h) = \frac{1}{m} \sum_{i=1}^m \underbrace{\mathbf{1}_{h(x_i) \neq c(x_i)}}_{X_i} \quad X_i \in \{0, 1\} \text{ is a RV.}$$

$$\begin{aligned}
 \text{Thus by Hoeffding, } \mathbb{P}[|\hat{R}_s(h) - R(h)| \geq \varepsilon] &\leq 2e^{-\frac{2m\varepsilon^2}{\sum_{i=1}^m 1^2}} \\
 &= 2e^{-2m\varepsilon^2}.
 \end{aligned}$$

$$\text{for } m \geq \frac{1}{2\varepsilon^2} \log(H + \log \frac{2}{\delta}), \quad \varepsilon \leq \delta \text{ under union bd in HFD.}$$

$$\begin{aligned}
& \Pr[\max_{h \in H} |R(h) - \hat{R}_S(h)| > \varepsilon] \\
&= \Pr[\exists h \in H \text{ s.t. } |R(h) - \hat{R}_S(h)| > \varepsilon] \\
&\leq \sum_{h \in H} \Pr[|R(h) - \hat{R}_S(h)| > \varepsilon] \\
&= |H| 2e^{-2m\varepsilon^2} \leq 8.
\end{aligned}$$

## §5 Stochastic Scenarios

- $x \sim D \quad y = c(x) \quad \dots \text{stochastic scenario}$
- $D \in \Pr(x \sim y) \quad (x, y) \sim D \rightarrow y \text{ also given by probability.}$
- $S = ((x_1, y_1), \dots, (x_m, y_m)) \sim D^m.$
- $\hat{R}_S(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\{h(x_i) \neq y_i\}}$
- $R(h) = \mathbb{E}[\hat{R}_S(h)] = \Pr_{(x, y) \sim D} [h(x) \neq y].$

Definition

(Agnostic PAC-learning algo.)      Stochastic Cast

$\exists$  a poly p s.t.  $\forall D \in \Pr(x \sim y) \quad \forall \varepsilon, \delta > 0 \quad \exists m \in \mathbb{N}$

$m \geq p(\frac{1}{\varepsilon}, \frac{1}{\delta}, n) \Rightarrow \Pr_{S \sim D^m} [R(h_S) \leq \inf_{h \in H} R(h) + \varepsilon] \geq 1 - \delta.$

## §5.2 Bayes Classifier

- $D \in \Pr(x \sim y)$
- $h \in [x \rightarrow y], h \underset{h_o \in [x \rightarrow y]}{\operatorname{argmin}} R(h_o)$

$$h(x) = \underset{i \in \{0, 1\}}{\operatorname{argmax}} (\Pr[y=i|x]).$$

Proposition

$$R(h_{BC}) \leq R(h) \quad \forall h \in H.$$

conditional prob

$$\begin{aligned}
\mathbb{E}_{(x, y) \sim D} [\mathbb{1}_{\{h(x) \neq y\}}] &= \mathbb{E}_x [\mathbb{E}_y [\mathbb{1}_{\{h(x) \neq y\}} | x]] \\
&= \mathbb{E}_x [\mathbb{E}_y [\mathbb{1}_{\{h(x)=0\}} \mathbb{1}_{y=1} + \mathbb{1}_{\{h(x)=1\}} \mathbb{1}_{y=0} | x]] \\
&= \mathbb{E}_x [\mathbb{1}_{\{h(x)=0\}} \mathbb{E}_y [\mathbb{1}_{y=1} | x] + \mathbb{1}_{\{h(x)=1\}} \mathbb{E}_y [\mathbb{1}_{y=0} | x]]
\end{aligned}$$

$$\begin{aligned}
 &\geq \mathbb{E}_x [\min(\mathbb{P}[y=0|x], \mathbb{P}[y=1|x])] \\
 &= 1 - \mathbb{E}_x [\max(\mathbb{P}[y=0|x], \mathbb{P}[y=1|x]))] \\
 &= 1 - \mathbb{E}_x [h(x) = y]
 \end{aligned}$$

need completion

- ① Stochastic  $\Leftrightarrow$  deterministic.
- ② If infinite PAC learnable theorems?
- ③  $C \subseteq H$  necessary. Why.

- Consistent / ERM 아닌 A로도 PAC-learnability 이야기 할 수 있는지
- Non PAC-learnable 의사?
- Theorem 2.13의 strictly strong statement 살피기.
  - (1) for  $h_S$  vs  $H$ . (2) Bound  $\frac{2}{\epsilon^2} (\log |H| + \log \frac{2}{\delta})$  vs  $\frac{1}{\epsilon^2} (\dots)$
  - (3) A need not be ERM.