

Recap.

Probably Approximately Correct (PAC) learning

Learner

- does not know distribution  $D$  or true hypothesis  $f$
- chooses parameters for accuracy  $\epsilon$  and confidence  $\delta$ .
- training data  $S = \{(X_i, Y_i = f(X_i))\}$ , of size  $m(\epsilon, \delta)$
- outputs a hypothesis  $h$  s.t.

$$\mathbb{P}(\mathcal{L}_{D,f}(h) \leq \epsilon) \geq 1 - \delta$$

then, this learner is Probably Approximately Correct.

u.p.  $\geq 1 - \delta$       error  $\leq \epsilon$

Empirical Risk Minimization (ERM)

- train error:  $\mathcal{L}_S(h) \triangleq \frac{1}{m} |\{i \in [m] : h(X_i) \neq Y_i\}|$

- ERM:  $h_S \in \arg \min_{h \in \mathcal{H}} \mathcal{L}_S(h)$

↑ hypothesis class of some complexity, like  $|\mathcal{H}|$

- claim: If  $|\mathcal{H}| < \infty$ ,  $f, D$  are realizable, and

sample complexity:  $m \geq \frac{\log\left(\frac{|\mathcal{H}|}{\delta}\right)}{\epsilon}$ , then for any  $(\epsilon, \delta)$

PAC learnable :  $\mathbb{P}_S(\mathcal{L}_{D,f}(h_S) > \epsilon) \leq \delta$

↑ ERM

Lecture 11. Key words: Uniform Convergence, Agnostically PAC learnable, Hoeffding's Inequality.

Relaxing the realizability assumption:  $y_i = f(x_i)$ ,  $f \in \mathcal{H}$ .  
 $\uparrow$  deterministic

Data generating distribution  $\mathcal{D}$  is over  $\mathcal{X} \times \mathcal{Y}$ :  $\mathcal{Z} = (x, y) \sim \mathcal{D}$

True error:  $L_{\mathcal{D}}(h) \triangleq \mathbb{P}_{(x,y) \sim \mathcal{D}}(h(x) \neq y)$   $\uparrow$  stochastic = random

Empirical error:  $L_S(h) \triangleq \frac{|\{i \in [m] : h(x_i) \neq y_i\}|}{m}$

Goal: find  $h$  minimizing  $L_{\mathcal{D}}(h)$

if given oracle access to  $\mathcal{D}$ , then Bayes Optimal Predictor achieves lower error than any other predictor [Theory HW 1, Q2] due 5/14

$$f_{\mathcal{D}}(x) = \begin{cases} 1 & \text{if } \mathbb{P}_{\mathcal{D}}(y=1|x) \geq \frac{1}{2} \\ 0 & \text{else} \end{cases}$$

\* Discuss.

### Agnostic PAC learnability.

A class of hypotheses,  $\mathcal{H}$ , is Agnostic PAC learnable

if  $\exists m_{\mathcal{H}}$  and a learning algorithm s.t. for  $\epsilon, \delta \in (0, 1)$ ,

running the Algorithm with  $m \geq m_{\mathcal{H}}(\epsilon, \delta)$  samples returns

$h$  s.t.

$$L_{\mathcal{D}}(h) \leq \min_{h' \in \mathcal{H}} L_{\mathcal{D}}(h') + \epsilon$$

Related accuracy requirement.

Relaxing Binary Classification: what needs to change in our framework?

- Multi-class classification  $\mathcal{Y} = [K]$

loss:  $\ell_h(x_i, y_i) = \mathbb{I}(h(x_i) \neq y_i)$

- Regression  $\mathcal{Y} = \mathbb{R}$

loss:  $\ell_h(x_i, y_i) = (h(x_i) - y_i)^2$

same

Risk function  $L_D(h) \triangleq \mathbb{E}_{(x,y) \sim D} [\ell(h, (x,y))]$

Empirical Risk function  $L_S(h) \triangleq \frac{1}{m} \sum_{i=1}^m \ell(h, (x_i, y_i))$

[Q 3.3 from book.] Problem Definition:

$\mathcal{X} = \mathbb{R}^2$

$\mathcal{Y} = \{0, 1\}$

$\mathcal{H} = \{h_r : r \in \mathbb{R}_+\}$ ,  $h_r(x) = \mathbb{I}_{\{\|x\| \leq r\}}$   
 $\uparrow$   $L_2$  norm

$|\mathcal{H}| = \infty$

Prove that  $\mathcal{H}$  is PAC learnable (assuming realizability)

with  $m_{\mathcal{H}}(\epsilon, \delta) \leq \left\lceil \frac{\log \frac{1}{\delta}}{\epsilon} \right\rceil$

Solution: ① Under realizability assumption  $\longrightarrow$  fix  $x \sim D$ , and  $f_{r^*}$

② Define : Algorithm : return smallest circle that includes all +.



$\mathcal{H}_{\text{BAD}} = \{h_r : r \leq \tilde{r} < r^*, \text{ where } \mathbb{P}_{x \sim D} (h_{\tilde{r}}(x) \neq h_{r^*}(x)) = \epsilon\}$   
 $= \{h_r : r \leq \tilde{r} < r^*, \text{ where } \mathbb{P}_{x \sim D} (\tilde{r} \leq \|x\| \leq r^*) = \epsilon\}$

③ Analysis : claim:  $\text{Algo} \leq r^*$

$\mathbb{P}_S(\text{Algo} < \tilde{r}) = \mathbb{P}(\|x\| < \tilde{r})^m \stackrel{m \sum \frac{1}{\epsilon} \log \frac{1}{\delta}}{\leq} (1-\epsilon)^m \leq e^{-\epsilon m} \stackrel{\downarrow}{\leq} e^{-\log \frac{1}{\delta}} \leq \delta$

# Chapter 4

Uniform Convergence is sufficient for learnability.

Definition: ( $\epsilon$ -representative sample)

Training set  $S$  is  $\epsilon$ -representative if  $\forall h \in \mathcal{H}$

$$|L_S(h) - L_D(h)| \leq \epsilon$$

Intuitively,  $S$  is large enough and good representation of  $D$ .

$$L_S(h) = \frac{1}{m} \sum_{i=1}^m l(h, z_i) : \text{Empirical Risk} \longleftrightarrow \text{population Risk} : L_D(h) = \mathbb{E}_D[l(h, z)]$$

arg min  $L_S(h)$  over  $h \in \mathcal{H}$       ERM  $\longleftrightarrow$  best in  $\mathcal{H}$ .      arg min  $L_D(h)$  over  $h \in \mathcal{H}$

Lemma. Assume we are given a training set  $S$  that is  $\frac{\epsilon}{2}$ -representative, then any output from ERM satisfies

$$L_D(h_S) \leq \min_{h \in \mathcal{H}} L_D(h) + \epsilon$$

How do we apply this lemma? Assume  $m_{\mathcal{H}} \geq \square$   $\xrightarrow{\text{Uniform Convergence}}$   $S$  is  $\frac{\epsilon}{2}$ -representative for  $\mathcal{H}$   $\xrightarrow{\text{Lemma}}$  ERM is agnostic PAC

Proof of Lemma by triangular inequality:  $a \leq b + |a-b|$

$$\begin{aligned}
 L_D(h_S) &\stackrel{\Delta \text{ ineq.}}{\leq} |L_D(h_S) - L_S(h_S)| + L_S(h_S) \\
 &\leq \frac{\epsilon}{2} + L_S(h_S), \quad \forall h \in \mathcal{H} \\
 &\stackrel{\Delta \text{ ineq.}}{\leq} \frac{\epsilon}{2} + |L_S(h) - L_D(h)| + L_D(h) \\
 &\leq \frac{\epsilon}{2} + \frac{\epsilon}{2} + L_D(h).
 \end{aligned}$$

$\uparrow$  representative  $S$        $\uparrow$  def ERM

To show ERM is agnostic PAC learner, we are left to show that, with probability  $1-\delta$ ,  $S$  is  $\frac{\epsilon}{2}$ -representative.

Definition (Uniform Convergence)

A hypothesis class  $\mathcal{H}$  has the uniform convergence property, if  $\exists m_{\mathcal{H}}^{UC}$  s.t. for every  $\epsilon, \delta$  and every  $D$ , the following hold:  
 if  $m \geq m_{\mathcal{H}}^{UC}(\epsilon, \delta)$  then  $S$  is  $\epsilon$ -representative w.p  $\geq 1-\delta$ .

Corollary 1.  $\mathcal{H}$  is agnostically PAC learnable with sample complexity

$$m_{\mathcal{H}}(\epsilon, \delta) \leq m_{\mathcal{H}}^{UC}\left(\frac{\epsilon}{2}, \delta\right)$$

and ERM is the agnostic PAC learner.

Next, we want to show that finite  $\mathcal{H}$ , i.e.  $|\mathcal{H}| < \infty$ , is agnostic PAC learnable  $\leftarrow$  by showing uniform convergence.

Claim. Any finite  $\mathcal{H}$ , i.e.,  $|\mathcal{H}| < \infty$ , is uniformly convergent with

$$m_{\mathcal{H}}^{UC}(\epsilon, \delta) \leq \left\lceil \frac{\log \frac{2|\mathcal{H}|}{\delta}}{\epsilon} \right\rceil$$

Corollary 2. ERM is agnostic PAC learner with sample complexity

$$m_{\mathcal{H}}(\epsilon, \delta) \leq \underset{\text{Corollary 1}}{m_{\mathcal{H}}^{UC}\left(\frac{\epsilon}{2}, \delta\right)} \leq \underset{\text{Claim}}{\left\lceil \frac{2 \log \frac{2|\mathcal{H}|}{\delta}}{\epsilon} \right\rceil}$$

we are left to prove the claim.

← Q. assuming  $0 \leq l(h, z_i) \leq 1$ , what is the sample size we need to guarantee  $|L_S(h) - L_D(h)| \leq \epsilon$  w.p  $1-\delta$ ?

n.t.s (need to show)  $P_S(\{S: \forall h \in H, |L_S(h) - L_D(h)| \leq \epsilon\}) \geq 1-\delta$

n.t.s  $\leftrightarrow P_S(\{S: \exists h \in H, |L_S(h) - L_D(h)| > \epsilon\}) < \delta$

$$\text{LHS} \stackrel{\text{Union bound}}{\leq} \sum_{h \in H} P_S(|L_S(h) - L_D(h)| > \epsilon) = \sum_{h \in H} P\left(\left|\frac{1}{m} \sum_{i=1}^m l(h, z_i) - \mathbb{E}[l(h, z)]\right| > \epsilon\right) (*)$$

Recall:  $L_S(h) = \frac{1}{m} \sum_{i=1}^m l(h, z_i)$  ← Average of iid r.v.s,  $z_i = (X_i, Y_i)$

$L_D(h) = \mathbb{E}_{z \sim D}[l(h, z)] = \mathbb{E}[L_S(h)]$  ← expectation of each term

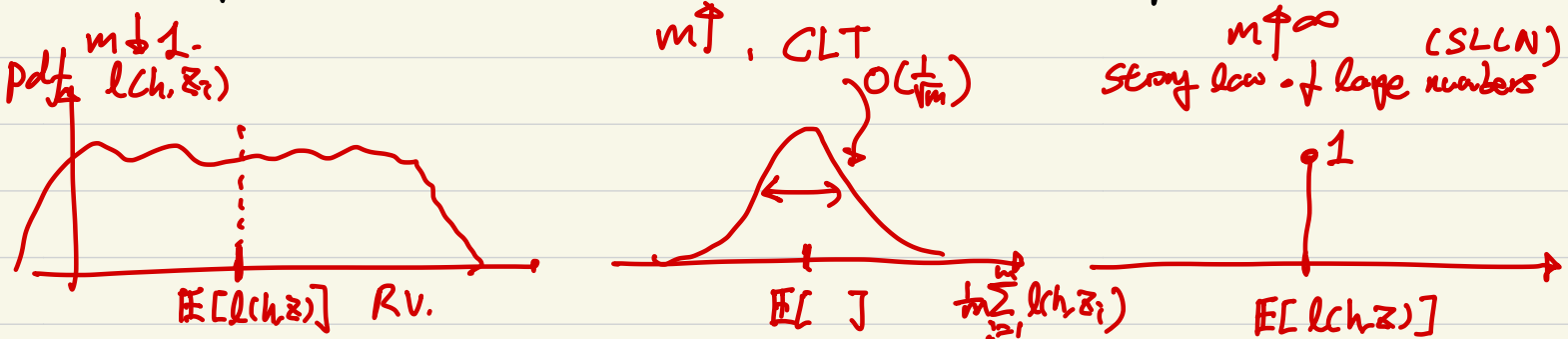
\* We need to show that the measure of R.V.  $L_S(h)$  distribution

concentrates around its mean.

\* the most important tool in statistical analysis for learning.

[ concentration of measure  
tail bound

Q. How fast (in terms of # of samples) does the measure of empirical mean concentrate around its expected value



End of lecture

Tail bound:  $P\left(\left|\frac{1}{m} \sum_{i=1}^m l(h, z_i) - \mathbb{E}[J]\right| > \epsilon\right) \leq F(m, \epsilon)$

# \* Hoeffding's Inequality

Let  $\theta_1, \theta_2, \dots, \theta_m$  be a set of i.i.d. random variables.  
with  $\mu \equiv \mathbb{E}[\theta_i]$ , satisfying.

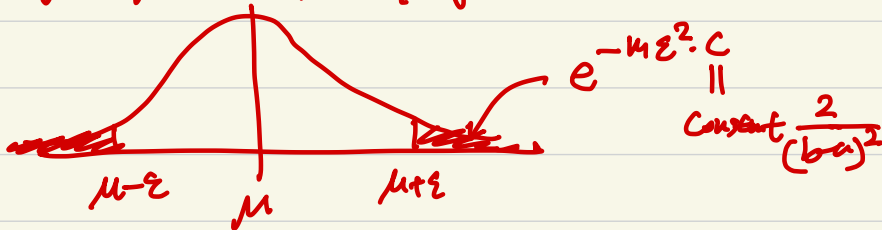
$$\mathbb{P}(a \leq \theta \leq b) = 1, \quad [\text{Boundedness}]$$

then for any  $\epsilon > 0$ ,

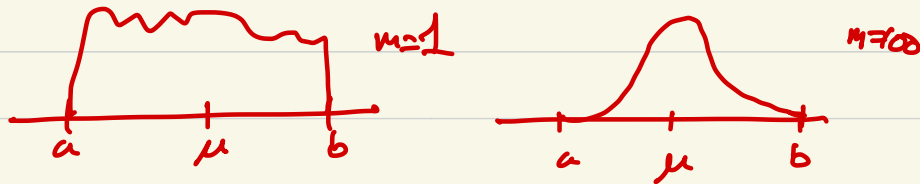
$$\mathbb{P}\left(\left|\frac{1}{m} \sum_{i=1}^m \theta_i - \mu\right| > \epsilon\right) \leq 2 \cdot e^{-\frac{2m\epsilon^2}{(b-a)^2}}$$

\* first example of concentration of measure.

for fixed  $m$ , varying  $\epsilon$



for fixed  $\epsilon$ , varying  $m$ .



letting  $\theta_i = \mathbb{1}(h, z_i)$ ,  $0 \leq \mathbb{1}(h, z_i) \leq 1$ ,  $\mu = \mathbb{1}_D(h)$

$$\text{Hoeffding's} \Rightarrow \mathbb{P}_S\left(\left|\mathbb{1}_S(h) - \mathbb{1}_D(h)\right| > \epsilon\right) = \mathbb{P}\left(\left|\frac{1}{m} \sum_{i=1}^m \theta_i - \mu\right| > \epsilon\right)$$

$$\leq 2 \cdot e^{-2m\epsilon^2}$$

$$(*) \leq \sum_{h \in \mathcal{H}} 2e^{-2m\epsilon^2} \leq 2|\mathcal{H}| \cdot e^{-2m\epsilon^2} \leq \delta$$

$$\text{let } m(\epsilon, \delta) \geq \frac{\ln\left(\frac{2|\mathcal{H}|}{\delta}\right)}{2\epsilon^2}$$

\* Summary.

$$\text{if } m \geq \frac{\log \frac{2M}{\delta}}{2\varepsilon^2}$$

then  $\mathcal{H}$  has uniform convergence

Lemma:  $(\frac{\varepsilon}{2}, \delta)$ -Uniform Convergence

↓  
 $(\varepsilon, \delta)$  PAC learnable with ERM

$$\text{if } m \geq \left\lceil \frac{2 \cdot \log \left( \frac{2M}{\delta} \right)}{\varepsilon} \right\rceil \geq m_{\mathcal{H}}^{\text{U.C.}} \left( \frac{\varepsilon}{2}, \delta \right) \geq m_{\mathcal{H}}(\varepsilon, \delta).$$