



# THEORY OF MACHINE LEARNING

## LECTURE 5

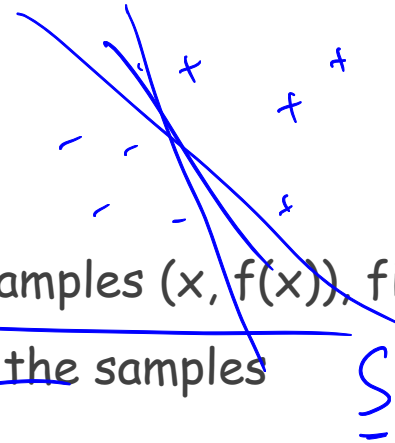
PAC MODEL, VC DIMENSION

## LAST WEEK

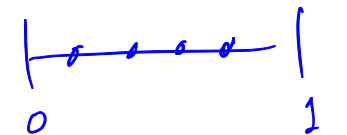
Goal of statistical learning: given some iid examples from  $\mathcal{D}$  + labels, learn a good  $h$ .  
Unknown distn  $\mathcal{D}$   
 $h$  ;  $\text{risk}_{\mathcal{D}}(h) = \Pr_{x \sim \mathcal{D}}[h(x) \neq \text{label}(x)]$ .

- Motivation: do we really need to restrict the hypothesis/concept class before starting learning? - yes! (No free lunch theorem)
- (PAC Learning): Learnability of a concept class  $H$  over domain  $X$ 
  - Informally, for any  $f \in H$  and any distribution  $\mathcal{D}$  over  $X$ , given examples of the form  $(x, f(x))$ , <sup>with  $n \sim \mathcal{D}$</sup>  we can learn a hypothesis 'h' such that  $\text{Risk}_{\mathcal{D}}(h)$  is  $< \epsilon$ , with high prob.  $((1 - \delta), \text{ for some parameter } \delta)$
  - Sample size only function of  $H, \epsilon, \delta$  (not distribution)  <sup>$\mathcal{D}$</sup>  (realizable Setting).
  - Learned hypothesis need not belong to  $H$  (improper learning)
- (Agnostic):  $f$  need not belong to  $H$   $\rightarrow$  want the risk to be  
$$\leq \min_{h' \in H} \text{risk}_{\mathcal{D}}(h') + \epsilon.$$

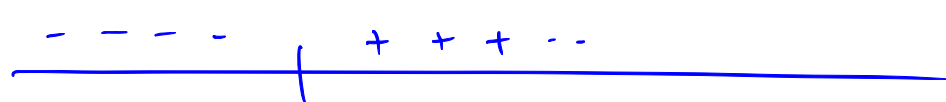
# GENERIC ALGORITHM



- Empirical risk minimization: given the samples  $(x, f(x))$ , find hypothesis  $h \in H$  that minimizes the total error on the samples
  - most natural algorithm == minimize training error
- How to do it efficiently? (learning half spaces - can do it in the realizable case...)
- When does it work? (outputs  $h$  with small gen. error).
  - If sample is "representative" of distribution --- for every hypothesis in class, error on sample  $\approx$  error on distribution (i.e. risk)

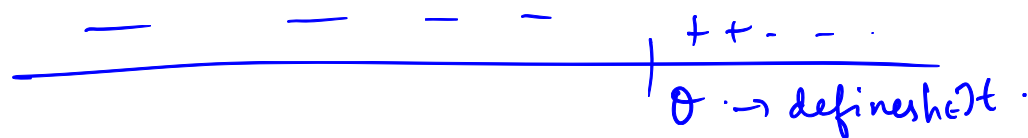


H: ~~set~~ all linear threshold functions in  $\mathbb{R}^1$ -D.



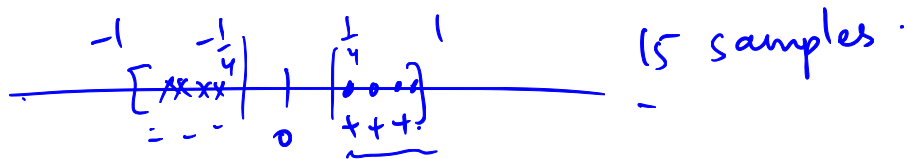
~~set~~

$X = \mathbb{R}$ .



$\mathcal{D}$ :  $[-1, 1]$  Unif $(-1, 1)$ .

goal: learn with error (risk bound)  $\leq 0.1$



\* Informal statement: #samples needed depends on both hypothesis class & distribution (also on  $\epsilon, \delta \dots$ ).

# REPRESENTATIVE SAMPLE

[applies to realizable & agnostic cases]

- Let  $H$  be a hypothesis class and  $X$  be an input space with a distribution  $D$  on it, and let  $f$  be a target function. Sample  $S \subseteq X$  is said to be  $\epsilon$  - "representative" if **for all**  $h$  in  $H$ , we have:

$$\left| \frac{1}{|S|} \text{error}(S, h) - \text{risk}_D(h) \right| < \epsilon$$

- If we happen to get a representative sample, ~~we have~~ <sup>ERM gives the</sup> desired bound on risk! (why?)
- Is a sample  <sup>$\epsilon$</sup>  representative "with high probability"?

Claim: If we perform ERM on a  $\epsilon$ -representative sample  $S$ , we obtain a hyp. with risk  $\leq$  (best-risk-in- $H$ )  $+\epsilon$ .


 $\rightarrow$  ERM  $h \rightarrow \frac{1}{|S|} \text{error}(S, h) \approx_{\varepsilon} \text{risk}_{\mathcal{D}}(h) \rightarrow \textcircled{\text{I}}$

opt hypothesis  $= h^*$ , say.  
 (hypo<sup>in S</sup> with the least risk w.r.t.  $\mathcal{D}$ .)

$\leadsto \frac{1}{|S|} \text{error}(S, h^*) \approx_{\varepsilon} \text{risk}_{\mathcal{D}}(h^*) \leadsto \textcircled{\text{II}}$

$$\text{error}(S, h) := \sum_{x \in S} \mathbb{1}[h(x) \neq \text{label}(x)].$$

indicator.

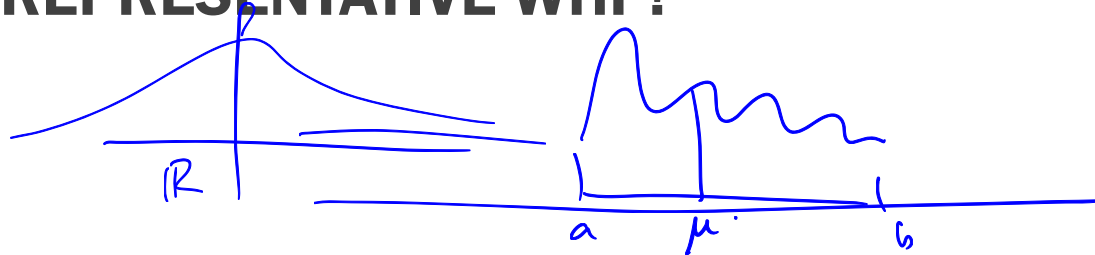
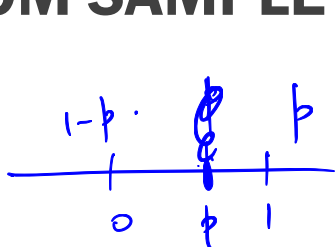
$$\underbrace{\text{risk}_{\mathcal{D}}(h)}_{\substack{\downarrow \\ \text{from } \textcircled{\text{I}}}} \leq \frac{1}{|S|} \cdot \text{error}(S, h) + \varepsilon \leq \frac{1}{|S|} \cdot \text{error}(S, h^*) + \varepsilon$$

because of  
ERM

$$\leq \underbrace{\text{risk}_{\mathcal{D}}(h^*)}_{\substack{\downarrow \\ \text{from } \textcircled{\text{II}}}} + \underbrace{(\varepsilon + \varepsilon)}_{2\varepsilon}$$

(in the realizable case, you only get  $\varepsilon$ .)

# RANDOM SAMPLE IS REPRESENTATIVE WHP!



- Chernoff bound (Hoeffding). Suppose  $X_1, X_2, \dots, X_n$  are  $n$  iid samples from a distribution with mean  $\mu$  and support  $[a, b]$ . Then we have

(Very useful.)

$$\Pr \left[ \left| \frac{1}{n} (X_1 + \dots + X_n) - \mu \right| > \underline{\epsilon} \right] \leq 2 \exp \left( - \frac{\epsilon^2 n}{(a-b)^2} \right) \rightarrow \frac{1}{n^2 \epsilon^4} < \delta$$

- Note: exponential dependence on  $n$

$$\Pr \left[ \left| \text{sample avg} - \mu \right| > \epsilon \right] \leq \text{exponentially small in \# samples.}$$

$$\text{Suppose } (a-b)^2 = 1$$

$$n \approx \frac{1}{\epsilon^2} \cdot \log \left( \frac{2}{\delta} \right)$$

# FINITE CLASSES ARE LEARNABLE

$\epsilon$ : accuracy.

$\delta$ : prob. parameter or confidence parameter.

$H$ : finite hypothesis class.

- Claim: for any  $X$  and distribution  $D$  over it, a sample of size  $O\left(\frac{1}{\epsilon^2} \log \frac{|H|}{\delta}\right)$  is representative with prob. at least  $1 - \delta$
- Proof idea: first start with a single hypothesis  $h \in H$ ; what is the probability that error on sample  $\approx$  error on  $D$ ?

Want to show:

$$m = \left( 4 \cdot \frac{1}{\epsilon^2} \cdot \log \left( \frac{2|H|}{\delta} \right) \right)$$

A random sample with  $m$  examples is  $\epsilon$ -representative with prob.  $\geq 1 - \delta$ .  $\rightarrow \forall h \in H \quad \left| \frac{1}{|S|} \text{error}(S, h) - \text{risk}_D(h) \right| \leq \epsilon$ .

Fix some  $h$ ; Claim: w.p.  $\geq 1 - \frac{\delta}{|H|}$   $\left| \frac{1}{|S|} \text{error}(S, h) - \text{risk}_D(h) \right| \leq \epsilon$ .



$$\text{i.e., } \Pr \left[ \left| \underbrace{\frac{1}{|S|} \text{error}(S, h)}_{\text{sample-based}} - \underbrace{\text{risk}_{\mathcal{D}}(h)}_{\text{distr.}} \right| > \varepsilon \right] \leq \frac{\delta}{|H|}$$

$$\text{Let } Y_i = \begin{cases} 1 & \text{if sample } h(x_i) \neq \text{label}(x_i) \\ 0 & \text{if } h(x_i) = \text{label}(x_i) \end{cases}$$

$h$   
 $x_i \sim \mathcal{D}$

$$\Pr_{x \sim \mathcal{D}} \left[ h(x) \neq \text{label}(x) \right] = \text{risk}(h) = \mathcal{R}(h)$$

If we sample  $x_i \sim \mathcal{D}$ , how is  $Y_i$  distributed?

$$Y_i = \begin{cases} 1 & \text{w.p. } \mathcal{R}(h) \\ 0 & \text{w.p. } 1 - \mathcal{R}(h) \end{cases} \Rightarrow \mu \text{ of this distr.} = \mathcal{R}(h)$$

$$\text{Chernoff [Hoeffding]: } \Pr \left[ \left| \frac{Y_1 + \dots + Y_m}{m} - \mu \right| > \varepsilon \right] \leq 2 \cdot e^{-\frac{\varepsilon^2 m}{1}}$$

$$\leq \frac{\delta}{|H|}$$

(plugging in  $m = \frac{4 \log \left( \frac{2|H|}{\delta} \right)}{\varepsilon^2}$ )

Union bound:

Want: some event  $\mathcal{E}$  <sup>occurs</sup> ~~holds~~ with prob. 99%.

$$\Pr[\text{event does not occur}] \leq 1\%.$$

Say event  $\mathcal{E}$  = want A and B to hold.  $\leq \frac{1}{200}$ .

$\mathcal{E}$  not occurring  $\equiv$  either A doesn't occur or B doesn't occur.  
 $\leq \frac{1}{200}$ .



## WHAT ABOUT INFINITE CLASSES?

- Note: if sample is representative, we are good!  
(modulo inefficiency of ERM)
- What if we can divide hypotheses into finitely many “classes”?
- Example of threshold functions on a line

---

## GROWTH FUNCTION OF A CLASS

- For a class  $H$  and an input space  $X$ , we can define a notion of “growth function”

# LEARNABILITY IN TERMS OF THE GROWTH FUNCTION

- Theorem: Suppose  $\tau_H(m)$  is an upper bound on the total number of “distinct sign patterns” possible for any sample of size  $m$ . Then for any  $X$ ,  $D$ , if we take a sample  $S$  of size  $m$ , we have, with prob.  $1-\delta$ ,

$$\sup_{h \in H} |err(h, S) - err(h, D)| \leq \frac{4 + \sqrt{\log \tau_H(2m)}}{\delta \sqrt{2m}}$$

---

## HOW TO BOUND GROWTH FUNCTION?

- Shattering.
- VC dimension.

## SAUER-SHELAH LEMMA (VAPNIK-CHERVONENKIS)

- **Lemma.** Let  $H$  be a hypothesis class of finite VC dimension  $d$ . Then for every  $m$ , we have: