

# RECITATION 7

## LEARNING THEORY

10-701: INTRODUCTION TO MACHINE LEARNING

April 3 2026

## 1 Learning Theory

### 1.1 PAC Learning

1. Basic notation:

- Probability distribution (unknown):  $X \sim p^*$
- **True function** (unknown):  $c^* : X \rightarrow Y$
- **Hypothesis space**  $\mathcal{H}$  and **hypothesis**  $h \in \mathcal{H} : X \rightarrow Y$
- Training dataset  $D = \{x^{(1)}, \dots, x^{(N)}\}$

2. **True Error (expected risk)**

$$R(h) = P_{x \sim p^*(x)}(c^*(x) \neq h(x))$$

3. **Train Error (empirical risk)**

$$\begin{aligned}\hat{R}(h) &= P_{x \sim D}(c^*(x) \neq h(x)) \\ &= \frac{1}{N} \sum_{i=1}^N \mathbb{1}(c^*(x^{(i)}) \neq h(x^{(i)})) \\ &= \frac{1}{N} \sum_{i=1}^N \mathbb{1}(y^{(i)} \neq h(x^{(i)}))\end{aligned}$$

The **PAC criterion** is that we produce a high accuracy hypothesis with high probability. More formally,

$$P(\forall h \in \mathcal{H}, \text{_____} \leq \text{_____}) \geq \text{_____}$$

**Sample Complexity** is the minimum number of training examples  $N$  such that the PAC criterion is satisfied for a given  $\epsilon$  and  $\delta$ .

Sample Complexity for 4 Cases: See Figure 1. Note that

- **Realizable** means  $c^* \in \mathcal{H}$
- **Agnostic** means  $c^*$  may or may not be in  $\mathcal{H}$

	Realizable	Agnostic
Finite $ \mathcal{H} $	<b>Thm. 1</b> $N \geq \frac{1}{\epsilon} [\log( \mathcal{H} ) + \log(\frac{1}{\delta})]$ labeled examples are sufficient so that with probability $(1 - \delta)$ all $h \in \mathcal{H}$ with $\hat{R}(h) = 0$ have $R(h) \leq \epsilon$ .	<b>Thm. 2</b> $N \geq \frac{1}{2\epsilon^2} [\log( \mathcal{H} ) + \log(\frac{2}{\delta})]$ labeled examples are sufficient so that with probability $(1 - \delta)$ for all $h \in \mathcal{H}$ we have that $ R(h) - \hat{R}(h)  \leq \epsilon$ .
Infinite $ \mathcal{H} $	<b>Thm. 3</b> $N = O(\frac{1}{\epsilon} [\text{VC}(\mathcal{H}) \log(\frac{1}{\epsilon}) + \log(\frac{1}{\delta})])$ labeled examples are sufficient so that with probability $(1 - \delta)$ all $h \in \mathcal{H}$ with $\hat{R}(h) = 0$ have $R(h) \leq \epsilon$ .	<b>Thm. 4</b> $N = O(\frac{1}{\epsilon^2} [\text{VC}(\mathcal{H}) + \log(\frac{1}{\delta})])$ labeled examples are sufficient so that with probability $(1 - \delta)$ for all $h \in \mathcal{H}$ we have that $ R(h) - \hat{R}(h)  \leq \epsilon$ .

12

Figure 1: Sample Complexity for 4 Cases

The **VC dimension** of a hypothesis space  $\mathcal{H}$ , denoted  $\text{VC}(\mathcal{H})$  or  $d_{\text{VC}}(\mathcal{H})$ , is the maximum number of points such that there exists at least one arrangement of these points and a hypothesis  $h \in \mathcal{H}$  that is consistent with any labelling of this arrangement of points.

To show that  $\text{VC}(\mathcal{H}) = n$ :

- 
- 

## Questions

1. For the following examples, write whether or not there exists a dataset with the given properties that can be shattered by a linear classifier.
  - 2 points in 1D
  - 3 points in 1D
  - 3 points in 2D
  - 4 points in 2D

How many points can a linear boundary (with bias) classify exactly for  $d$ -Dimensions?

2. Consider a rectangle classifier (i.e. the classifier is uniquely defined 3 points  $x_1, x_2, x_3 \in \mathbb{R}^2$  that specify 3 out of the four corners), where all points within the rectangle must equal 1 and all points outside must equal  $-1$ .

(a) Which of the configurations of 4 points in Figure 2 can a rectangle shatter?

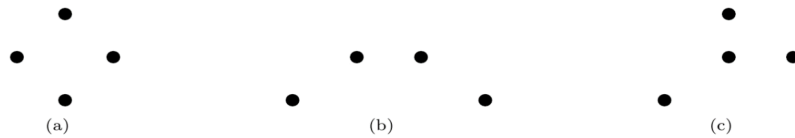


Figure 2

(b) What about the configurations of 5 points in Figure 3?



Figure 3

3. Let  $x_1, x_2, \dots, x_n$  be  $n$  random variables that represent binary literals ( $x \in \{0, 1\}^n$ ). Let the hypothesis class  $\mathcal{H}_n$  denote the conjunctions of no more than  $n$  literals in which each variable occurs at most once. Assume that  $c^* \in \mathcal{H}_n$ .

Example: For  $n = 4$ ,  $(x_1 \wedge x_2 \wedge x_4), (x_1 \wedge \neg x_3) \in \mathcal{H}_4$

Find the minimum number of examples required to learn  $h \in \mathcal{H}_{10}$  which guarantees at least 99% accuracy with at least 98% confidence.