

# Valiant's Theory

## The PAC Model

Theory of Machine Learning – Spring 22

January 13, 2022

# Last class

## Logistics

- ▶ **Course webpage:** Canvas, can find lecture schedule, slides, scribe template, . . .
- ▶ TA: Chris Harker
- ▶ Scribe for today?

# ML until the 1980s

- ▶ Many informal notions of learning: Rosenblatt and the “perceptron” algorithm, neural networks, ...
- ▶ Limitations of perceptrons
- ▶ No formal theory to reason about, no clear definitions

**Question:** Can we formally define “learning”?

# Theory of the Learnable

- ▶ Leslie Valiant 1983 – Theory of the Learnable (CACM)
- ▶ Drawing the boundaries of learnability – how to define it? what is possible?
- ▶ Really a theory of *supervised* learning
- ▶ I.e., deals with classification or *prediction* problems
  - given some description of a "scenario".  
what to do next.?
  - given an input, <sup>make</sup> prediction – label.

# Theory of the Learnable

image: Set of all pixel values

- ▶ **Input**: “features” of input
- ▶ **Hypothesis/model**: function from input to prediction/label  
$$\text{hypothesis } h: \underset{\text{(all inputs)}}{I} \rightarrow \underset{\text{(all labels)}}{L}$$
- ▶ **Definition of a learning algorithm.** an algorithm that can find a (good hypothesis) without explicitly being told what it is!
  - what all does a learning alg. need?

Most natural way. Give examples of inputs and ~~outputs~~ their labels

Good hypothesis:

"Low error": agreement with "true" label.

Qn: should it agree on all inputs?

Input:  $x \rightsquigarrow$  collection of feature values.

$$\left( \begin{array}{c} \uparrow \\ \text{pixel 1-value} \end{array} , \text{pixel 2 value} , \dots \right) \in \mathbb{R}^m.$$

Ans: No, but we must have agreement on all  $m$  pixel.  
"inputs of interest".

# Good hypothesis?

$x$	$y$
$x_0$	

(given examples.)

- ▶ Must do well on given inputs (hopefully perfectly)
- ▶ Must also do well on "unseen" inputs (generalization)
- ▶ **How to formalize this?**

**Valiant's key assumption.** Assume an "input distribution"  
(unknown to the learner)

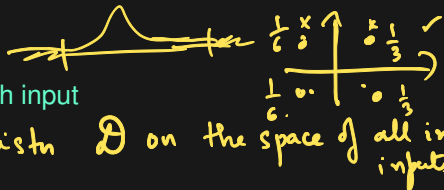
✓  
we care about error  
"wrt" this distribution.

↙  
probability distribution on the space  
of all inputs.

# Good hypothesis

assuming there is a true label for each input

∃ some (unknown to learner) distn  $\mathcal{D}$  on the space of all inputs



Risk minimization is the goal:

Given a hypothesis  $h$ , a true label function  $l$ , the risk of  $h$  w.r.t. a dist  $\mathcal{D} := R_{\mathcal{D}}(h) = \Pr_{\underline{x} \sim \mathcal{D}} [h(\underline{x}) \neq \underline{l(x)}]$ .

Definition of learnability

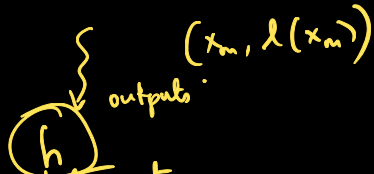
→ We say that a hypothesis ( $h$ ) is learnable, if  $\forall \mathcal{D}$  for any  $\epsilon > 0$ , there exists an  $m(\epsilon)$  (training size) such that given  $m$  iid examples  $x_1, \dots, x_m$  from  $\mathcal{D}$  and  $l(x_1), \dots, l(x_m)$ , we can produce a  $h$  s.t.  $R_{\mathcal{D}}(h) \leq \epsilon$ , with  $\text{prob}^{\mathcal{D}} > 99\%$ .



"We can produce"?



- $\exists$  an efficient algorithm  $A$ .  
(poly in  $m$  - # training examples)  
that takes  $(x_1, l(x_1)), (x_2, l(x_2)) \dots$



- \* Inherently a probabilistic statement.
- \* Training samples come from same dist as test

# Complexity of ground truth label

Importance of hypothesis class

$(\mathbb{R}^m)$  (labels,  $\pm 1$ )

Label function  $l : \mathcal{I} \rightarrow L$

How "rich" can this fun be?



- Sample complexity & training time depend on "how complex"  $l$  is.

$$\text{sign}(x_1^2 + 3x_2 + x_4)$$

- Assume: label function  $l$  is in a certain "hypothesis class"  $\mathcal{H}$  (which is known to also ~~be~~  $A$ ).

$\mathcal{H}$ : set of all polynomials of degree  $d$  in input features  $x_1, \dots, x_m$ .

# Learnability with finite hypothesis classes

Theorem: (informal): Any finite hypothesis class  $\mathcal{H}$  is learnable with  $\sim \log |\mathcal{H}|$  training examples.

$\forall \epsilon$ , you can produce  $h$  such that  
 $R_D(h) \leq \epsilon$  w.p.  $\geq 90\%$   
using  $\frac{\log |\mathcal{H}|}{\epsilon^2}$  samples.