# CIS 5200: MACHINE LEARNING LEARNINGTHEORY

Content here draws from material by Rob Schapire (Princeton), Hamed Hassani (UPenn) and Michael Kearns (UPenn)



2 March 2023

### Surbhi Goel

#### Spring 2023





#### OUTLINE - TODAY

#### Survey Overview

- \*What about generalization?
- \* Probably Approximately Correct (PAC) learning
- \* Finite Function Classes are PAC learnable

#### (PAC) learning earnable

## SUPERVISED LEARNING - SO FAR \* Training dataset $\mathcal{S} = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$ Function class $\mathcal{F}$ , loss function $\ell$ \* Empirical Risk Minimizer:

We have looked at various methods to find the ERM Is this good enough for learning?



 $\hat{f} = \arg\min_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^{m} \ell(f(x_i), y_i)$ 

 $\hat{R}(f)$ 

#### MEMORIZATION

# Memorizer predictor $f_{\text{mem}}(\cdot)$ $f_{\text{mem}}(x) = \begin{cases} y_i & \text{if } \exists (x_i, y_i) \in \mathcal{S}, x = x_i, \\ 0 & \text{otherwise.} \end{cases}$

This gets 0 training loss  $\hat{R}(f_{mem}) = 0$ , so it is an ERM. But is it a good predictor?

### GENERALIZATION

We want the predictor to perform well not just on the training data but on examples it will see in the future.

#### **Recall how we formalized this:**

Training dataset is drawn independently and identically from some unknown but fixed distribution  $\mathcal{D}$ 

> loss on future examples = loss over the distribution  $R(\hat{f}) = \mathbb{E}_{(x,y)\sim\mathcal{D}} \left| \ell(\hat{f}(x), y) \right|$

We ideally want to minimize true risk R and not just empirical risk  $\hat{R}$ 

#### GENERALIZATION



We don't have access to the true risk, we can only see a training set

$$R(\hat{f}) = (R(\hat{f}) - \hat{R}(\hat{f})) + \hat{R}(\hat{f})$$

generalization gap

In this lecture, we will bound this generalization gap

# loss on future examples = loss over the distribution $R(\hat{f}) = \mathbb{E}_{(x,y)\sim \mathcal{D}} \left| \ell(\hat{f}(x), y) \right|$

ERM can guarantee that this is small

This will depend on the size of the training set and the complexity of the function class  ${\mathscr F}$ 



### LET US FORMALIZE THIS!

(realizable learning model)

We want to show that  $R(\hat{f})$  is small for any empirical risk minimizer  $\hat{f}$ 

**Challenge I:** Can we find exactly  $f_*$  or get exactly 0 error? **Challenge 2:** Can we find a good predictor for all datasets?

 $\exists f_* \in \mathscr{F}$  such that  $R(f_*) = 0$ Let us work in the classification setting and assume that there is a perfect classifier

 $\hat{R}(\hat{f}) = 0$  since  $\hat{R}(f_*) = 0$ 



#### EXAMPLE - THRESHOLDS

-1



#### One dimensional half space or thresholds

#### EXAMPLE - ZERO RISK?

Suppose the data is uniformly distributed on this line, and you observe the following:



 $f_a(x) = \begin{cases} 1 & \text{if } x \ge a \\ -1 & \text{otherwise.} \end{cases}$ 

From finite samples, it is hard to exactly find  $f_*$  to get 0 error





 $f_a(x) = \begin{cases} 1 & \text{if } x \ge a \\ -1 & \text{otherwise.} \end{cases}$ 

Suppose the data is uniformly distributed on this line, and you observe the following:



Would not know where to put the threshold, however this is a very unlikely sample

### PROBABLY APPROXIMATELY CORRECT (PAC) LEARNING

Introduced by Leslie Valiant in 1984, captures the notion of finding approximately good Error parameter  $\epsilon$ predictors with high probability Confidence parameter  $\delta$ **Definition:** 

A function class  $\mathcal{F}$  is PAC learnable if there exists an algorithm  $\mathcal{A}$  and a function  $m_{\mathcal{F}}: (0,1)^2 \to \mathbb{N}$  with the following property: for every labelling function  $f \in \mathcal{F}$ , for every distribution  $\mathcal{D}$  on feature space  $\mathcal{X}$ , and for all  $\epsilon, \delta \in (0,1)$ , if  $\mathscr{A}$  is given access to a training dataset S of size  $m \geq m_{\mathscr{F}}(\epsilon, \delta)$  where the features are drawn from  $\mathcal{D}$  and labels are according to f, then with probability  $1 - \delta$  (over the choice of the training dataset),  $\mathscr{A}$  outputs a predictor  $\hat{f}$  such that  $\Pr |\hat{f}(x) \neq f(x)| \leq \epsilon$ .  $x \sim \mathcal{D} L$ 



### PROBABLY APPROXIMATELY CORRECT (PAC) LEARNING

#### **Definition:**

A function class  $\mathcal{F}$  is PAC learnable if there exists an algorithm  $\mathscr{A}$  and a function  $m_{\mathcal{F}}: (0,1)^2 \to \mathbb{N}$  with the following property: for every labelling function  $f \in \mathcal{F}$ , for every distribution  $\mathcal{D}$  on feature space  $\mathcal{X}$ , and for all  $\epsilon, \delta \in (0,1)$ , if  $\mathscr{A}$  is given access to a training dataset S of size  $m \geq m_{\mathscr{F}}(\epsilon, \delta)$  where the features are drawn from  $\mathcal{D}$  and labels are according to f, then with probability  $1 - \delta$  (over the choice of the training dataset),  $\mathscr{A}$  outputs a predictor  $\hat{f}$  such that  $\Pr_{x \to \mathscr{O}} |\hat{f}(x) \neq f(x)| \leq \epsilon$ .

Function  $m_{\mathcal{F}}: (0,1)^2 \to \mathbb{N}$  captures the sample complexity of learning Depends on complexity of  ${\mathcal F}$ 



#### EXAMPLE - NOT PAC LEARNABLE

dataset of size *m* we would have only seen the labels on those *m* points.

The true function could take any value outside!

We cannot possibly guarantee small generalization error!



# Class of all possible predictors from $\mathcal{X} \to \{-1,1\}$ is not PAC learnable, for any

#### EXAMPLE - THRESHOLDS

-1



One dimensional half space or thresholds

#### EXAMPLE - PAC LEARNABLE





As we see more and more samples, the mass of the region of error will shrink In the next lecture we will quantify how many samples we will need for this

### GENERAL - FINITE CLASSES ARE PAC LEARNABLE

Consider a finite function class  $\mathcal{F} = \{f_1, \dots, f_{|\mathcal{F}|}\}$ 

#### **Theorem:**

Every finite function class  $\mathcal{F}$  is PAC learnable with sample complexity

Observe that it depends on the size of  ${\mathscr F}$  which is a natural notion of complexity of  ${\mathscr F}$ 

 $m_{\mathcal{F}}(\epsilon,\delta) \leq \left| \frac{\log(|\mathcal{F}|/\delta)}{\epsilon} \right|.$ 



### GENERAL - FINITE CLASSES ARE PAC LEARNABLE BY ERM

Consider a finite function class  $\mathcal{F} = \{f_1, \dots, f_{|\mathcal{F}|}\}$ 

#### **Theorem:**

Every finite function class  $\mathcal{F}$  is PAC learnable with sample complexity

where the algorithm  $\mathscr{A}$  is any empirical risk minimization algorithm.

# $m_{\mathcal{F}}(\epsilon, \delta) \leq \left| \frac{\log(|\mathcal{F}|/\delta)}{\epsilon} \right|$

Proof on the iPad

### GENERAL - FINITE CLASSES ARE PAC LEARNABLE BY ERM

#### Another way to state this is:

#### **Theorem:**

# For any ERM $\hat{f}$ evaluated over training set of size *m*, with probability $1 - \delta$ , $R(\hat{f}_S) \leq \frac{\log(|\mathcal{F}|/\delta)}{\delta}$ .



M



SUMMARY

### We studied the notion of PAC learning where we allowed approximately correct learning with high probability

#### We proved that finite classes are PAC learnable using ERM

**Next class:** How do we handle infinite classes?