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

* Recap: * Probably Approximately Correct (PAC) learning * Finite Function Classes are PAC learnable * What about infinite classes? * VC Dimension * VC Classes are PAC Learnable

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:

fixed distribution \mathcal{D}

- Training dataset is drawn independently and identically from some unknown but
 - 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 want to minimize true risk R but only have access to a training set

PROBABLY APPROXIMATELY CORRECT (PAC) LEARNING

Introduced by Leslie Valiant in 1984, captures the notion of finding approximately good predictors with high probability Error parameter ϵ Confidence parameter δ

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:

$$R(\hat{f}) = \Pr_{x \sim \mathcal{D}}$$

- 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

$$\left[\hat{f}(x) \neq f(x)\right] \leq \epsilon$$
.



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 \mathcal{A} is any empirical risk minimization algorithm.

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

FINITE CLASSES ARE PAC LEARNABLE BY ERM

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

Theorem:

For any ERM \hat{f}_S over training set S of size m, with probability $1 - \delta$,



$R(\hat{f}_S) \leq \frac{\log(|\mathcal{F}|/\delta)}{\delta}$

M

EXAMPLE - INTEGER WEIGHT LINEAR CLASSIFIER

w is such that each $w_i \in \{-10, ..., 0, ..., 10\}$ for $i \in [d]$

How many possible functions are there?

What is the sample complexity for PAC learning?

- Consider the class \mathcal{F} of integer weight linear classifiers where the parameter
 - 21^d

 $m_{\mathcal{F}}(\epsilon, \delta) \lesssim \frac{d \log 21 + \log(1/\delta)}{d \log 21 + \log(1/\delta)}$



WHAT ABOUT INFINITE CLASSES?

If the size of the class is infinite, then the previous bound is useless.

What quantity should replace $\log |\mathcal{F}|$?

Let us think about the behavior of the function on our training dataset: $\Pi_{\mathscr{F}}(S) = \{ (f(x_1), \dots, f(x_m)) : f \in \mathscr{F} \}$



- This is all possible labeling that the training points could have according to ${\mathscr F}$



WHAT ABOUT INFINITE CLASSES?

Let us think about the behavior of the function on our training dataset:

Define the maximum possible labelings over all training sets of size m

Growth function

 2^m What is an upper bound on $\Pi_{\mathcal{F}}(m)$?

But for many \mathcal{F} , this is actually much smaller!



$\Pi_{\mathscr{F}}(S) = \{ (f(x_1), \dots, f(x_m)) : f \in \mathscr{F} \}$

$\Pi_{\mathcal{F}}(m) = \max_{S; |S|=m} |\Pi_{\mathcal{F}}(S)|$

(SOME) INFINITE CLASSES ARE PAC LEARNABLE BY ERM

Consider an infinite function class \mathcal{F}

Theorem:

For any ERM \hat{f}_{S} over training set S of size m, with probability $1 - \delta$,

Note that if $\Pi_{\mathscr{F}}(m) = 2^m$ then this is vacuous!

We will not cover the proof in class since it is a bit involved

$R(\hat{f}_S) \leq \left\lceil \frac{\log(|\Pi_{\mathcal{F}}(2m)|/\delta)}{m} \right\rceil.$

EXAMPLE - THRESHOLDS

—1

Consider a dataset of three points $x_1 < x_2 < x_3$

How are the possible labelings with the class of thresholds?



EXAMPLE - INTERVALS



Consider a dataset of three points $x_1 < x_2 < x_3$

How are the possible labelings with the class of intervals?

 $f_{a,b}(x) = \begin{cases} 1 & \text{if } a \le x \le b \\ -1 & \text{otherwise.} \end{cases}$



GENERAL BOUNDS - VC DIMENSION

bounds $\Pi_{\mathcal{F}}(m)$

Definition (shattering):

is, \mathcal{F} can realize all possible labelings for the set of points in S.

Definition (VC dimension):

VC dimension of a function class $\mathcal{F}(VC(\mathcal{F}))$ is the size of the largest set S that can be shattered by \mathcal{F} .



There is a notion of complexity called Vapnik-Chervonenkis (VC) dimension that

A set S of inputs is said to be shattered by function class \mathscr{F} if $|\Pi_{\mathscr{F}}(S)| = 2^{|S|}$, that



GENERAL BOUNDS - VC DIMENSION

Definition (VC dimension):

can be shattered by \mathcal{F} .

To show that a function class has $VC(\mathcal{F}) = d$, we must show that,

- There is a set S of d points that is shattered by \mathcal{F}
- There is no set S of d+1 points that is shattered by \mathcal{F}

VC dimension of a function class $\mathcal{F}(VC(\mathcal{F}))$ is the size of the largest set S that

EXAMPLE - THRESHOLDS

What is the VC dimension of the class of thresholds?

-1



EXAMPLE - INTERVALS



What is the VC dimension of the class of intervals?

 $f_{a,b}(x) = \begin{cases} 1 & \text{if } a \le x \le b \\ -1 & \text{otherwise.} \end{cases}$



CONNECTION - VC DIMENSION & GROWTH FUNCTION **Theorem (Sauer's Lemma):** Let $d = VC(\mathcal{F})$, then • $\Pi_{\mathscr{F}}(m) = 2^m$ for $m \leq d$ • $\Pi_{\mathscr{F}}(m) = O(m^d)$ for m > d

Theorem:



For any ERM \hat{f}_{S} over training set S of size m > d, with probability $1 - \delta$,

M

EXAMPLE - RECTANGLES



What is the VC dimension of the class of rectangles?





EXAMPLE - LINEAR CLASSIFIERS

What is the VC dimension of the class of linear classifiers?



 $f_w(x) = \operatorname{sgn}(w^{\top}x)$



CONCIUSION

Theorem:

 $R(\hat{f}_S) \lesssim \frac{d + \log(1/\delta)}{m}.$

To quantify how many samples we need to learn a particular function class, we can use the VC dimension as a complexity measure

For any ERM \hat{f}_{S} over training set S of size m > d, with probability $1 - \delta$,

M

