

# TTIC 31250 An Introduction to the Theory of Machine Learning

Avrim Blum

04/08/20

## Lecture 2: Online learning

### Mistake-bound model:

- Basic results, relation to PAC, halving algorithm
- Connections to information theory

### Combining "expert advice":

- (Randomized) Weighted Majority algorithm
- Regret-bounds, connections to game theory

1

## Recap from last time

- Last time: PAC model and Occam's razor.
  - If data set has  $m \geq (1/\epsilon)[s \ln(2) + \ln(1/\delta)]$  examples, then whp any consistent hypothesis with  $\text{size}(h) < s$  has  $\text{err}(h) < \epsilon$ .
  - Equivalently,  $\text{size}(h) \leq (\epsilon m - \ln(1/\delta))/\ln(2)$  suffices.
  - "compression  $\Rightarrow$  learning"
- Occam bounds  $\Rightarrow$  any class is learnable in PAC model if computation time is no object.

2

## Online learning

- What if we don't want to make assumption that data is coming from some fixed distribution? Or any assumptions at all?
- Can no longer talk about past performance predicting future results.
- Can we hope to say anything interesting??

Idea: mistake bounds & regret bounds.

3

## Mistake-bound model

- View learning as a sequence of stages.
- In each stage, algorithm is given  $x$ , asked to predict  $f(x)$ , and then is told correct value.
- Make no assumptions about order of examples.
- Goal is to bound total number of mistakes.

Alg A learns class C with mistake bound M if A makes  $\leq M$  mistakes on any sequence of examples consistent with some  $f \in C$ .

4

## Mistake-bound model

Alg A learns class C with mistake bound M if A makes  $\leq M$  mistakes on any sequence of examples consistent with some  $f \in C$ .

- Note: can no longer talk about "how much data do I need to converge?" Maybe see same examples over again and learn nothing new. But that's OK if don't make mistakes either...
- Try to bound in terms of size of examples  $n$  and complexity of target  $s$ .
- C is learnable in MB model if exists alg with mistake bound and running time per stage  $\text{poly}(n, s)$ .

5

## Simple example: disjunctions

- Suppose features are boolean:  $X = \{0, 1\}^n$ .
- Target is an OR function, like  $x_3 \vee x_9 \vee x_{12}$ .
- Can we find an on-line strategy that makes at most  $n$  mistakes?
- Sure.
  - Start with  $h(x) = x_1 \vee x_2 \vee \dots \vee x_n$
  - Invariant:  $\{\text{features in } h\} \supseteq \{\text{features in } f\}$
  - Mistake on negative: discard features in  $h$  set to 1 in  $x$ . Maintains invariant & decreases  $|h|$  by 1.
  - No mistakes on positives. So at most  $n$  mistakes total.

6

### Simple example: disjunctions

- Algorithm makes at most  $n$  mistakes.
- No deterministic alg can do better:

1 0 0 0 0 0 + or - ?

0 1 0 0 0 0 + or - ?

0 0 1 0 0 0 + or - ?

0 0 0 1 0 0 + or - ?

...

7

### MB model properties

An alg  $A$  is "conservative" if it only changes its state when it makes a mistake.

Claim: if  $C$  is learnable with mistake-bound  $M$ , then it is learnable by a conservative alg.

Why?

- Take generic alg  $A$ . Create new conservative  $A'$  by running  $A$ , but rewinding state if no mistake is made.
- Still  $\leq M$  mistakes because  $A$  still sees a legal sequence of examples.

8

### MB learnable $\Rightarrow$ PAC learnable

Say alg  $A$  learns  $C$  with mistake-bound  $M$ .

Transformation 1:

- Run (conservative)  $A$  until it produces a hyp  $h$  that survives  $\geq (1/\epsilon)\ln(M/\delta)$  examples.
- $\Pr(\text{fooled by any given } h) \leq \delta/M$ .
- $\Pr(\text{fooled ever}) \leq \delta$ .

Uses at most  $(M/\epsilon)\ln(M/\delta)$  examples total.

- Fancier method gets  $O(\epsilon^{-1}[M + \ln(1/\delta)])$

9

### One more example...

- Say we view each example as an integer between 0 and  $2^n-1$ .
- $C = \{[0, a] : a < 2^n\}$ . (device fails if it gets too hot)
- In PAC model we could just pick any consistent hypothesis. Does this work in MB model?
- What would work?

10

### What can we do with unbounded computation time?

- "Halving algorithm": take majority vote over all consistent  $h \in C$ . Makes at most  $\lg(|C|)$  mistakes.  $\frac{1}{2} \leq p_f$
- What if we had a "prior"  $p$  over fns in  $C$ ?
  - Weight the vote according to  $p$ . Make at most  $\lg(1/p_f)$  mistakes, where  $f$  is target fn.
- What if  $f$  was really chosen according to  $p$ ?
  - Expected number of mistakes  $\leq \sum_h [p_h \lg(1/p_h)]$   
= entropy of distribution  $p$ .

11

### What can we do with unbounded computation time?

- "Halving algorithm": take majority vote over all consistent  $h \in C$ . Makes at most  $\lg(|C|)$  mistakes.
- What if  $C$  has functions of different sizes?
- For any (prefix-free) representation, can make at most 1 mistake per bit of target.
  - Think of writing random 0s and 1s until hit a legal hypothesis or no longer a prefix of one.
  - $p_f = \Pr(\text{reach } f) = 1/2^{\text{size}(f)}$
  - $\lg(1/p_f) = \text{size}(f)$ .

12

## Is halving alg optimal?

- Not necessarily
- Can think of MB model as 2-player game between alg and adversary.
  - Adversary picks  $x$  to split  $C$  into  $C_-(x)$  and  $C_+(x)$ . [fns that label  $x$  as - or + respectively]
  - Alg gets to pick one to throw out.
  - Game ends when all fns left are equivalent.
  - Adversary wants to make game last as long as possible.
- $OPT(C) = MB$  when both play optimally.

13

## Is halving alg optimal?

- Halving algorithm: throw out larger set.
- Optimal algorithm: throw out set with larger mistake bound.

14

## What if there is no perfect function?

Think of as  $h \in C$  as "experts" giving advice to you. Want to do nearly as well as best of them in hindsight.

These are called "regret bounds".

➤ Show that our algorithm does nearly as well as best predictor in some class.

We'll look at a strategy whose running time is  $O(|C|)$ . So, only computationally efficient when  $C$  is small.

15

## Using "expert" advice

Say we want to predict the stock market.

- We solicit  $n$  "experts" for their advice. (Will the market go up or down?)
- We then want to use their advice somehow to make our prediction. E.g.,

Expt 1	Expt 2	Expt 3	neighbor's dog	truth
down	up	up	up	up
down	up	up	down	down
...	...	...	...	...

Can we do nearly as well as best in hindsight?

["expert" = someone with an opinion. Not necessarily someone who knows anything.]

[note: would be trivial in PAC (i.i.d.) setting]

16

## Using "expert" advice

If one expert is perfect, can get  $\leq \lg(n)$  mistakes with halving alg.

But what if none is perfect? Can we do nearly as well as the best one in hindsight?

Strategy #1:

- Iterated halving algorithm. Same as before, but once we've crossed off all the experts, restart from the beginning.
- Makes at most  $\lg(n)[OPT+1]$  mistakes, where  $OPT$  is #mistakes of the best expert in hindsight.

Seems wasteful. Constantly forgetting what we've "learned". Can we do better?

17

## Weighted Majority Algorithm

Intuition: Making a mistake doesn't completely disqualify an expert. So, instead of crossing off, just lower its weight.

Weighted Majority Alg:

- Start with all experts having weight 1.
- Predict based on weighted majority vote.
- Penalize mistakes by cutting weight in half.

Weights:	1	1	1	1	
Predictions:	U	U	U	D	We predict: U Truth: D
Weights:	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{2}$	1	

18

### Analysis: do nearly as well as best expert in hindsight

- $M$  = # mistakes we've made so far.
- $m$  = # mistakes best expert has made so far.
- $W$  = total weight (starts at  $n$ ).
- After each mistake,  $W$  drops by at least 25%. So, after  $M$  mistakes,  $W$  is at most  $n(3/4)^M$ .
- Weight of best expert is  $(1/2)^m$ . So,

$$\begin{aligned} (1/2)^m &\leq n(3/4)^M \\ (4/3)^M &\leq n2^m \\ M &\leq 2.4(m + \lg n) \end{aligned}$$

constant ratio

19

### Randomized Weighted Majority

- $2.4(m + \lg n)$  not so good if the best expert makes a mistake 20% of the time. Can we do better? **Yes.**
- Instead of taking majority vote, use weights as probabilities. (e.g., if 70% on up, 30% on down, then pick 70:30) **Idea:** smooth out the worst case.
  - Also, generalize  $\frac{1}{2}$  to  $1 - \epsilon$ .

$$\text{Solves to: } M \leq \frac{-m \ln(1 - \epsilon) + \ln(n)}{\epsilon} \approx (1 + \epsilon/2)m + \frac{1}{\epsilon} \ln(n)$$

$$M \leq 1.39m + 2 \ln n \quad \leftarrow \epsilon = 1/2$$

$$M \leq 1.15m + 4 \ln n \quad \leftarrow \epsilon = 1/4$$

$$M \leq 1.07m + 8 \ln n \quad \leftarrow \epsilon = 1/8$$

unlike most worst-case bounds, numbers are pretty good.

20

### Analysis



- Say at time  $t$  we have fraction  $F_t$  of weight on experts that made mistake.
- So, we have probability  $F_t$  of making a mistake, and we remove an  $\epsilon F_t$  fraction of the total weight.
  - $W_{\text{final}} = n(1 - \epsilon F_1)(1 - \epsilon F_2) \dots$
  - $\ln(W_{\text{final}}) = \ln(n) + \sum_t [\ln(1 - \epsilon F_t)] \leq \ln(n) - \epsilon \sum_t F_t$ 

(using  $\ln(1-x) \leq -x$ )  
( $\sum F_t = E[\text{\# mistakes}]$ )
- If best expert makes  $m$  mistakes, then  $\ln(W_{\text{final}}) > \ln((1 - \epsilon)^m)$ .
- Now solve:  $\ln(n) - \epsilon M > m \ln(1 - \epsilon)$ .

$$M \leq \frac{-m \ln(1 - \epsilon) + \ln(n)}{\epsilon} \approx (1 + \epsilon/2)m + \frac{1}{\epsilon} \log(n)$$

21

### Summarizing

- $E[\text{\# mistakes}] \leq (1 + \epsilon)OPT + \epsilon^{-1} \log(n)$   
 $= OPT + (\epsilon OPT + \epsilon^{-1} \log(n))$ 

Assuming here that  $OPT \geq \log(n)$
- If set  $\epsilon = (\log(n)/OPT)^{1/2}$  to balance the two terms out (or use guess-and-double), get bound of  $M \leq OPT + 2(OPT \log n)^{1/2} \leq OPT + 2(T \log n)^{1/2}$
- Define **average regret** in  $T$  time steps as:  
 (avg per-day cost of alg) - (avg per-day cost of best fixed expert in hindsight).  
 Goes to 0 or better as  $T \rightarrow \infty$  = "no-regret" algorithm].

22

### Extensions

- What if experts are actions? (**rows in a matrix game, ways to drive to work,...**)
- At each time  $t$ , each has a loss (cost) in  $\{0,1\}$ .
- Can still run the algorithm
  - Rather than viewing as "pick a prediction with prob proportional to its weight",
  - View as "pick an expert with probability proportional to its weight"
  - Alg pays expected cost  $\vec{p}_t \cdot \vec{c}_t = F_t$ .
- Same analysis applies.  
**Do nearly as well as best action in hindsight!**

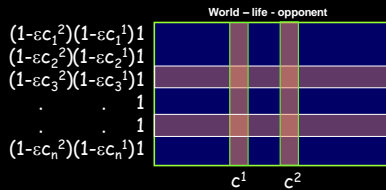
23

### Extensions

- What if losses (costs) in  $[0,1]$ ?
- Just modify alg update rule:  $w_i \leftarrow w_i(1 - \epsilon c_i)$ .
- Fraction of wt removed from system is:  
 $(\sum_i w_i \epsilon c_i) / (\sum_j w_j) = \epsilon \sum_i p_i c_i = \epsilon [\text{our expected cost}]$
- Analysis very similar to case of  $\{0,1\}$ .

24

## RWM (multiplicative weights alg)



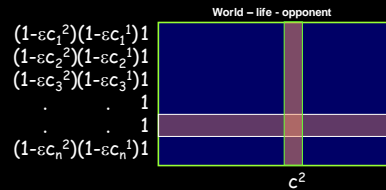
scaling  
so costs  
in  $[0,1]$

Guarantee: do nearly as well as fixed row in hindsight

Which implies doing nearly as well (or better)  
than minimax optimal

25

## Connections to minimax optimality



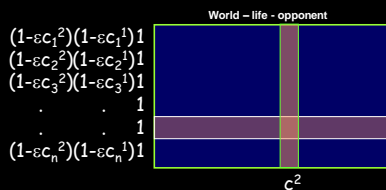
scaling  
so costs  
in  $[0,1]$

If play RWM against a best-response oracle,  $\vec{p}$  will  
approach minimax optimality.

(If it didn't, wouldn't be getting promised guarantee)

26

## Connections to minimax optimality



scaling  
so costs  
in  $[0,1]$

If play two RWM against each other, then empirical  
distributions must be near-minimax-optimal.

(Else, one or the other could & would take advantage)

27