

Ch 21 Online Learning

Friday, March 27, 2020 3:05 PM

In Shalev-Shwartz + Ben-David, ch 21 Online Learning is first ch of "Part 3: Additional Learning Models"
(21: online learning, 22: clustering, 23: dimensionality reduction, 24: generative models (MLE, Bayes), 25: Feature Selection)
most similar to PAC... mostly orthogonal ideas, or (eg ch 24) you'd see in a standard Stat. class

Online Learning

Before, we had training data, then (phase 2) could apply classifier to test data

With online learning, no separate phases: train and test on same data

(many things naturally operate this way, e.g., SPAM filters, medical knowledge, ...)

Also, theory can be completely distribution free: instead, allow it to change, or even be adversarial

We'll have to change how we measure "learning" (don't use risk L_0 anymore)

Connections to game theory

Many similarities to PAC learning though (notions similar to VC-dimension, agnostic or not, online-to-batch conversion)

Popular recently ('05-'15 very hot)

Many online methods are cheap to implement, i.e., can deal w/ data streams

History:

Harnan '57, Rosenblatt '58, Novikoff '62

Modern start w/ Littlestone + Warmuth '89

Plan: ① Binary classification → realizable
 ↓
 agnostic

② Regression or Surrogate loss (need convexity)

1a) Online Classification, realizable (focus on learning, not computational complexity)

Setup: T rounds ($t \in [T]$), we assume $T \in \mathbb{N}$ but will take $\lim_{T \rightarrow \infty}$ sometimes

Each round, observe data/features X_t

you (or the algo.) predicts p_t ... then true answer y_t is revealed (for now, $y_t \in \{-1, 1\}$)

Goal: make as few mistakes as possible (i.e. 0-1 loss)

So... want a good prediction now and to learn so we make a good prediction in the future

(general examples: . SPAM filtering

. Restaurant problem (we'll see again in Reinforcement Learning)

example of "reveal": $y_t = \text{whether you liked the food } K_t$

)

How are x_t and y_t generated? Deterministic, Stochastic, adversarial ↗ our analysis, since distribution-free, worst-case

Of course, need some assumptions

(otherwise, choose $y_t = \neg p_t$ and it's hopeless)

Realizable Case Assume (x_t) is arbitrary/adversarial, but $\exists h^* \in \mathcal{H}$ st. $(\forall t) y_t = h^*(x_t)$

Def $M_A(\mathcal{H})$ is the maximum # of mistakes ($p_t \neq y_t$) made by algorithm A on a sequence of data $(x_t, y_t)_{t=1}^T$, over all (x_t) and all $h^* \in \mathcal{H}$

Def Given $S = ((x_1, h^*(x_1)), \dots, (x_T, h^*(x_T)))$, $h^* \in \mathcal{H}$, then $M_A(S)$ is # mistakes algo A makes on S, and $M_A(\mathcal{H})$ is supremum of $M_A(S)$ over all such S (arbitrary length), and a bound of the form $M_A(\mathcal{H}) \leq B$ ($\forall T$) is a mistake bound, and it is online learnable if \exists algo A w/ a mistake bound for \mathcal{H} .

Warmup: let $|\mathcal{H}| < \infty$ (like we did for PAC learning)

For PAC learning, our math tool was Empirical Risk Minimization (ERM)

We have something similar:

Algorithm: "Consistent" $|\mathcal{H}| < \infty$

Init. : $V_1 = \mathcal{H}$
 for $t=1, 2, \dots, T$ ↗ so this is an ERM for $S = (x_1, y_1), \dots, (x_{t-1}, y_{t-1})$
 receive x_t
 choose any $h \in V_t$, predict $p_t = h(x_t)$
 receive true label y_t
 update (prune) $V_{t+1} = \{h \in V_t : \underbrace{h(x_t) = y_t}_{\text{consistent}}\}$

Corollary 21.2 $A = \text{Consistent}$ then $M_A(\mathcal{H}) \leq |\mathcal{H}| - 1$

proof Every time we make a mistake, $|V_{t+1}| \leq |V_t| - 1$

So for M mistakes,

$$1 \leq |V_t| \leq |\mathcal{H}| - M$$

by realizability \square

Is that good? No. Easy fix:

Algorithm: "Halving" $|\mathcal{H}| < \infty$

Init. : $V_1 = \mathcal{H}$
 for $t=1, 2, \dots, T$
 receive x_t ↗ i.e. pick majority vote
 predict $p_t = \operatorname{argmax}_{r \in \{0,1\}} \left| \{h \in V_t : h(x_t) = r\} \right|$
 receive truth $y_t = h^*(x_t)$
 update $V_{t+1} = \{h \in V_t : h(x_t) = y_t\}$ ↗ same consistent update

Thm 21.3 $M_{\text{Hobby}}(H) \leq \log_2(|H|)$

Proof Everytime we make a mistake, we can reduce size of V_t a lot, since at least half get it wrong

$$|V_{t+1}| \leq \frac{1}{2}|V_t| \quad \text{so} \quad 1 \leq |V_{T_{t+1}}| \leq |H| 2^{-M}$$

by realizability \square

General case (H infinite is ok)

We want a measure of complexity of H , analogous to VC dimension

(Recall: $\text{VCdim}(H) \geq m \Rightarrow H$ can shatter m points, i.e., realize all possible dichotomies)

Analogy/example ("A theory that explains everything, explains nothing")

I create a new physics theory. Of course, it involves some unknown parameters, so in fact I actually have a class of theories $H = \{h_1, \dots, h_K\}$.

What's the fewest number of experiments needed for you to falsify my theory? (I am your adversary)
or, you are my adversary

Also, I am unscrupulous (or, "you can make the data say whatever you want")

i.e., you propose an experiment X_t with outcome $\neg P_t$, then I claim $y_t = P_t$

... but I have to be self-consistent (if you ask X_k again, I can't change my answer)

and I've claimed $y_t = h^*(X_t)$ for some $h^* \in \{h_1, \dots, h_K\}$

You need at least $\log_2(k)$ experiments before you could prove that I'm faking it (or that my theory is false)

(more generally)

Get in the mindset of an adversary.

Your annoying sibling has 1 piece of candy, puts it behind their back, and asks you to choose which hand

You choose "left hand". Probability you're correct? ~~50%~~? 0%.

Because they're your sibling, they cheat as much as possible without you being able to prove it

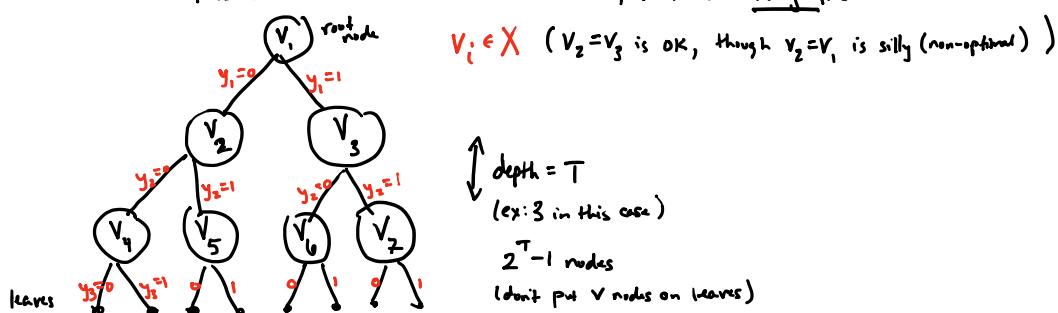
(Notes written during coronavirus quarantine w/ two daughters who just got Easter candy)

Our notion will be Ldim(H) (after Nick Littlestone)

Think of learner vs. environment: environment chooses x_t , learner chooses p_t , then environment chooses y_t

In fact, we only care about rounds where we make mistakes, so let $y_t = \neg p_t$ (i.e., $\neg p_t = 1 - p_t$) but must be realizable by some $h \in H$

To characterize what an adversarial environment can do, make a binary tree



Environment chooses X_t (i.e. V_i), but you choose P_t (so force $y_t = 1 - p_t$ if it's a mistake)
so you choose how to traverse the tree from root to leaf.

$$V_{i_2} = X_{i_2} \dots$$

Def A shattered tree of depth T is a tree w/ nodes $V_1, V_2, \dots, V_{2^T-1}$ such that any path

from root to leaf (via nodes $(V_{i_1} = V_1, V_{i_2}, \dots, V_{i_T})$ w/ labels (y_1, \dots, y_T))

can be realized by some $h \in H$ s.t. $h(X_{i_t}) = y_t \quad \forall t \in [T]$.

↳ i.e. defeats all learners,
or, all learners make at least
 T mistakes

Def $Ldim(H)$ is the maximal integer T s.t. \exists a shattered tree of depth T
environment's strategy

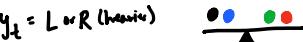
Immediately proves

Lemma 21.6: \forall algo A, $M_A(H) \geq Ldim(H)$

Ex 4 balls  3 identical, 1 is heavier i.e. $|H| = 4$

X_1 = weighty partition

$y_1 = L \text{ or } R$ (heavier)



you guess Left, so I choose R



$p_1 = L$

$y_1 = R$

$p_1 = R$

$y_1 = L$

you guess Right, so I choose Left (possible heavy balls: blue)



$p_1 = L$

$y_1 = R$

$p_1 = R$

$y_1 = L$

repeating X_t isn't going to help me cause you to make mistakes (nor will swapping sides)

I can cause one last mistake, but after this you won't make further mistakes since you know which ball is heavier

So $Ldim(H) \geq 2$

Ex $Ldim(H) = \log_2(|H|)$ if H finite

(either just see it, or prove via HALVING Algo)

Ex Unit vectors $X = \{1, \dots, d\}$, $H = \{h_1, \dots, h_d\}$ $h_j(x) = \begin{cases} 1 & x=j \\ 0 & x \neq j \end{cases}$

then $Ldim(H) = 1$

Why? take $V_1 = 1$, wlog (due to symmetry)

If algo says $p_1 = 1$, then no mistake and it's learned h

If algo says $p_1 = 0$, they make a mistake, but they see $y_1 = 1$ and learn h

Either way, the algo won't make any more mistakes

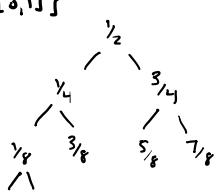
Since d can be arbitrarily large, this shows

$Ldim(H) \ll |H|$ is possible.

Ex $X = [0, 1] \subseteq \mathbb{R}$, $H = \left\{ \frac{1}{[x \leq a]} : a \in [0, 1] \right\}$

Recall $VCDim(H) = 1$

b.t $Ldim(H) = \infty$. Choose nodes



$x_1 = \frac{1}{2}$, you guess 1 ($x^* \leq \frac{1}{2}$)

so adversary decides $x^* > \frac{1}{2}$

$x_2 = 3x_1$, you guess 0 ($x^* > 3x_1$)

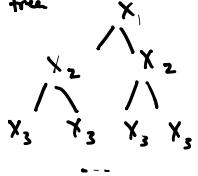
so adversary decides $x^* \leq 3x_1$

you'll make unbounded
of mistakes.

Ldim vs Vcdim: $\forall H$, $Vcdim(H) \leq Ldim(H)$, but $<$ possible

proof: let $d = Vcdim(H)$, and (x_1, \dots, x_d) a shattered set

make the tree



Now, we know $M_H(H) \geq Ldim(H)$. What about $\exists A$ s.t. $M_A(H) = Ldim(H)$?

Yes, constructive (generalize HALVING):

Alg: "Standard Optimal Algorithm" (SoA) for H ($|H| = \infty$ is ok)

Initialize $V_0 = H$

for $t=1, 2, \dots, T$

receive x_t

predict $p_t = \underset{r \in \{0, 1\}}{\operatorname{argmax}} Ldim(\{h \in V_t : h(x_t) = r\})$

receive true label y_t

update $V_{t+1} = \{h \in V_t : h(x_t) = y_t\}$

Lemma 21.7 $M_{SoA}(H) \leq Ldim(H)$

proof If we make a mistake, $Ldim(V_{t+1}) \leq Ldim(V_t) - 1$. Wlog, say we chose $p_t = 0$ (so actually $y_t = 1$)

If not, then $Ldim(V_{t+1}) = Ldim(V_t)$

$\Rightarrow V_t^{(0)} = V_t^{(1)}$ by update rule

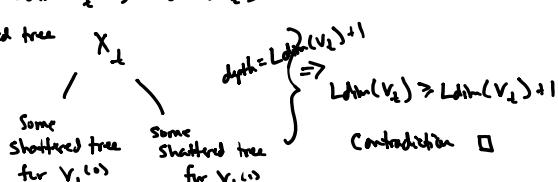
But also $Ldim(V_t) \geq Ldim(V_t^{(0)}) \geq Ldim(V_t^{(1)}) \geq Ldim(V_t)$

$Ldim(A) \geq Ldim(B)$ by how we choose r

if $A \neq B$

but then if $Ldim(V_t^{(0)}) = Ldim(V_t^{(1)}) = Ldim(V_t)$

can make a new shattered tree



Corollary H has a mistake bound iff $Ldim(H) < \infty$

(16) Online classification, unrealizable case

If we were realizable, it's possible to have a bounded # of mistakes even as $T \rightarrow \infty$

Not true if unrealizable, so need new measure of learning

Want our algo competitive with best fixed predictor, i.e., "regret"

$$\text{Def } \text{Regret}_A(h, T) = \sup_{(x_1, y_1), \dots, (x_T, y_T)} \sum_{t=1}^T |P_t - y_t| - \sum_{t=1}^T |h(x_t) - y_t|$$

this is for ℓ^1 loss,
but can be defined
more generally for
an arbitrary loss

$$\begin{aligned} \text{Def } \text{Regret}_A(h, T) &= \sup_{h \in H} \text{Regret}_A(h, T) = \sup_h \sup_{(x_i, y_i)} \sum |P_t - y_t| - \sum |h(x_t) - y_t| \\ &= \sup_{(x_i, y_i)} \left(\sum |P_t - y_t| + \sup_h - \sum |h(x_t) - y_t| \right) \\ &= \sup_{(x_i, y_i)} \left(\sum |P_t - y_t| - \inf_h \sum |h(x_t) - y_t| \right) \end{aligned}$$

i.e. best expert

We usually have $\inf_{h \in H} \sum_{t=1}^T |h(x_t) - y_t|$ grow sublinearly with T
(if not, then best expert is about as good as always choosing $P_t = 1$)

So if we want our algo to be useful, also need it to be sublinear in T

and want the regret to grow sublinearly in T

(\Rightarrow the difference between our learner and best hypothesis (w.r.t. hindsight) goes to 0 as $T \rightarrow \infty$)

Unfortunately, this isn't possible (Cover '65)

$$\text{Ex } H = \{h_0, h_1\} \quad h_0(x) = 0 \quad \forall x, \quad h_1(x) = 1 \quad \forall x$$

Not realizable, so adversary has no constraint. You guess (P_1, P_2, \dots, P_T)

and it sets $y_1 = \neg P_1, \dots, y_T = \neg P_T$, so you make T mistakes.

The best-in-class makes $\leq T/2$ mistakes, so $\text{regret} \geq T - T/2 = T/2$

which isn't sublinear.

So must limit adversary's power

Allow learner to randomize, and look at expectation

Assume randomness is independent of adversary, i.e., you flip a coin, and adversary

knows your general strategy but not the outcome of the coin flip.

In the binary case, output label $\hat{y} = \begin{cases} 1 & \text{w.p. } P_t \\ 0 & \text{w.p. } 1-P_t \end{cases}$

then expected loss is

$$\mathbb{E} |\hat{y} - y| = \mathbb{P}[\hat{y} \neq y] = \begin{cases} 1-P_t & y=1 \\ P_t & y=0 \end{cases} = |P_t - y_t|$$