

# 21. Boosting

Tuesday, February 27, 2024 2:48 PM

(Throughout, we're assuming binary classification and 0-1 loss)

An example of an ensemble method

- ① We know there's a tradeoff (bias v. variance) based on the complexity of  $h_l$ , but it's not always clear how to change the complexity
- ② Finding the ERM for most  $h_l$  is difficult

Boosting helps w/ both the above issues

- pick a small, simple class  $h_l$  (or  $\mathcal{B}$  for "base" class) for which we can solve ERM. It should be "good enough" of a learner (aka a "weak learner")
- we'll boost  $\mathcal{B}$  (solving ERM $_{\mathcal{B}}$  several times) to get a "strong" learner.

History:

1990 MIT grad. student Robert Schapire, turned practical in '95

Schapire + Yoav Freund's AdaBoost

(for  $h_l$ )

PAC-learner aka **strong learner**:  $m \geq m_{H_l}(\epsilon, \delta)$

an algo A s.t. if  $m$  is sufficiently large,  $|S| = m$  has iid samples, then  $L_{D,f}(A(S)) \leq \epsilon$  with probability at least  $1 - \delta$

we're in the realizable (non-agnostic) setting

Def An algo A is a  $\gamma$ -weak learner if  $\exists m_{H_l}: (0, 1) \rightarrow \mathbb{N}$  s.t.  $(\gamma > 0)$

$\forall$  dist. D,  $\forall f$ , if S has  $m \geq m_{H_l}(\delta)$  iid samples then

$L_{D,f}(A(S)) \leq \frac{1}{2} - \gamma$  with prob. at least  $1 - \delta$

$= P[\text{misclassification}]$  since 0-1 loss

↑  
the "edge"  
like the  
Casino's  
"edge"

... and  $h_l$  is  $\gamma$ -weak-learnable if  $\exists$  a  $\gamma$ -weak-learner for  $h_l$ .

Theoretically, for binary classif.,  $h_l$  is PAC learnable iff  $\text{VCdim}(h_l) < \infty$

But a particular algo. A might be a weak (but not strong) learner.

And there may be computational issues w/ ERM

/ iff  $h_l$  is  $\gamma$ -weak-learnable.  
iff ERM is a PAC learner

## 21a. Boosting (example)

Tuesday, February 27, 2024 3:05 PM

**Ex. 10.1**  $\mathcal{H}$  = "3-piece classifiers", algo A is "decision stump"

Setup:  $X = \mathbb{R}$ ,  $Y = \{\pm 1\}$ ,  $\mathcal{H} = \{h_{\theta_1, \theta_2, b} : \theta_1, \theta_2 \in \mathbb{R}, w, \theta_1 < \theta_2, b \in \{\pm 1\}\}$

tiny decision trees  
bit flip

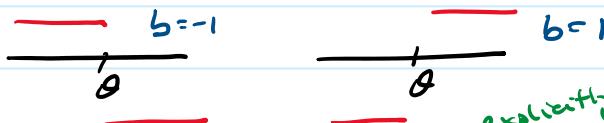
$$h_{\theta_1, \theta_2, b}(x) = \begin{cases} +b & x \notin [\theta_1, \theta_2] \\ -b & x \in [\theta_1, \theta_2] \end{cases}$$

$\text{VCdim}(\mathcal{H}) = 3$ , so it's PAC learnable



Let  $\mathcal{B} = \{x \mapsto \text{sign}(x - \theta) \cdot b, \theta \in \mathbb{R}, b \in \{\pm 1\}\}$

↑ decision stumps  
aka depth-1 decision trees  
(A common choice in boosting)



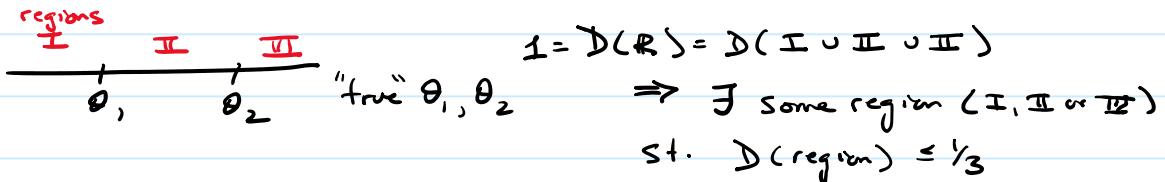
$\mathcal{B}$  not as expressive as  $\mathcal{H}$ , and  $\text{VCdim}(\mathcal{B}) = 2$  (easy to show)

explicitly or  
bound as  
union of  
left/right  
thresholding

Claim:  $A = \text{ERM}_{\mathcal{B}}$  is a  $\gamma = 1/2$  weak learner for  $\mathcal{H}$

proof Realizability: arbitrary distr.  $D$  over  $x$ , and  $\exists f \in \mathcal{H}$   
s.t.  $y = f(x)$ .

We claim  $\exists h \in \mathcal{B}$  s.t.  $L_{D,f}(h) \leq 1/3$ .



Then ignore this region!

e.g.  $D(I) \leq 1/3$ , so set  $\theta = \theta_2$ , get 100% correct on regions II + III.

Next, since  $\text{VCdim}(\mathcal{B}) = 2 < \infty$ ,  $\mathcal{B}$  is PAC learnable AGNOSTIC VERSIONS!  
Subtle!

so apply (quantitative) Fund. Thm. of ML (Thm 6.8) to  $\mathcal{B}$ ,

if  $m \geq m_g(\epsilon, \delta)$ ,  $m_g(\epsilon, \delta) \leq \text{const} \frac{2 + \log(\gamma/\delta)}{\epsilon^2}$  then

$\text{ERM}_{\mathcal{B}}$  returns  $h \in \mathcal{B}$  s.t.  $L_{D,f}(h) \leq \min_{h' \in \mathcal{B}} L_{D,f}(h') + \epsilon$

choose  $\epsilon = 1/2$   
so loss  $\leq \gamma/2 = 1/2 - 1/2$  ✓  $\leq 1/3 + \epsilon$

## 21b. Boosting (example)

Tuesday, February 27, 2024 3:22 PM

(example continued)

Claim:  $\text{ERM}_B$  is easy to implement

Let  $\mathcal{H}_{ht} = \{x \mapsto \text{Sign}(x - \theta) : \theta \in \mathbb{R}\}$  hard-thresholding  
 $\mathcal{H}_{-ht} = \{x \mapsto -\text{Sign}(x - \theta) : \theta \in \mathbb{R}\}$  aka 1D hyperplanes!

$\mathcal{B} = \mathcal{H}_{ht} \cup \mathcal{H}_{-ht}$  so  $\rightarrow$  solve  $h \in \text{ERM}_{ht}$

$\min_{x \in A \cup B} f(x) = \min(\min_{x \in A} f(x), \min_{x \in B} f(x))$

2) solve  $h_2 \in \text{ERM}_{-ht}$

3) choose whichever is better,  $h_1$  or  $h_2$

$$\text{so } \min_{h \in \mathcal{H}_{ht}} \hat{\mathcal{L}}_S(h) := \sum_{i=1}^m \underbrace{l(h, (x_i, y_i))}_{= \begin{cases} 0 & \text{if } h(x_i) = y_i \\ 1 & \text{else} \end{cases}}$$

or, since it'll be useful later, consider a slight generalization

$$\min_{h \in \mathcal{H}_{ht}} \sum_{i=1}^m D_i l(h, (x_i, y_i)) \quad \text{for non-neg. weights } D_i$$

don't confuse with  $D$

$$\text{eg } D_i = \frac{1}{m}$$

$$\overline{D} = (D_i)_{i=1}^m$$

while we're at it,

let's also generalize to  $x \in \mathbb{R}^d$ , not just  $d=1$

$$\mathcal{H}_{ht} = \{x \in \mathbb{R}^d \mapsto \text{Sign}(\theta - x^{(j)}), \theta \in \mathbb{R}, j \in [d]\}$$

$$\mathcal{H}_{ds} = \mathcal{H}_{ht} \cup \mathcal{H}_{-ht} \quad ds = \text{decision tree}$$



need to find  $\theta \in \mathbb{R}$  and  $j \in [d]$

$$\min_{j \in [d]} \min_{\theta \in \mathbb{R}} \sum_{i=1}^m D_i \cdot \mathbb{1}_{h_{j,\theta}(x_i) \neq y_i}$$

loop over  $j$

$$\sum_{i: y_i = +1} D_i \cdot \mathbb{1}_{x_{i,j} > 0} + \sum_{i: y_i = -1} D_i \cdot \mathbb{1}_{x_{i,j} \leq 0}$$

$x_{i,j} = j^{\text{th}} \text{ coordinate of } i^{\text{th}} \text{ example } \vec{x}_i$

Fixing  $j \in [d]$ : back to 1D problem with  $\{x_{i,j}\}$ . Let's sort them,

$$\tilde{x}_1 \leq \dots \leq \tilde{x}_m \quad (\text{all in 1D now})$$

Objective is piecewise constant in  $\theta$ , look for "break points" / "turning pts"

e.g. restrict search to  $\theta \in \bigoplus := \{\tilde{x}_1 - 1, \frac{1}{2}(\tilde{x}_1 + \tilde{x}_2), \frac{1}{2}(\tilde{x}_2 + \tilde{x}_3), \dots, \frac{1}{2}(\tilde{x}_{m-1} + \tilde{x}_m), \tilde{x}_m + 1\}$

$m+1$  things to try,  $O(m)$  cost per try, but can reduce that

if clever (cumulative sum), so total  $O(m)$  cost plus cost of sort

so altogether  $O(d \cdot m \cdot \log(m))$  complexity

## 21c. Boosting / sorting / shuffling / bagging

Tuesday, February 27, 2024 3:42 PM

Still on decision stump example:

$\text{ERM}_S$  costs  $O(d m \log(m))$  flops

How much would  $\text{ERM}_{\text{FC}}$  cost?

Can still do similar tricks w/ bitflip b and coordinate j

But now  $\theta \in \mathbb{R}^2$  not  $\mathbb{R}$ , much trickier. It's still

piecewise constant in regions but not that helpful

(analogy: in 1D rootfinding, we can do the bisection method)

w/  $\log(\frac{1}{\epsilon})$  steps... in 2D, no equivalent  
"ellipsoid method" is closest equiv. )

At best it'd be  $O(m^2)$ ,

so  $\text{ERM}_{\text{FC}}$  at least  $\Omega(d m^2)$  cost.

Aside on complexity of sorting, shuffling, etc.

Given a list of d elements:

- Sorting is  $O(d \log d)$  for generic comparison sorts,

or  $O(d)$  for radix sorts (eg. integers) but complicated due

- finding max or min is  $O(d)$

- finding top-k is  $O(k \log(k) d)$ , maybe surprising

(naively, it's  $\min(k, \log(d)) \cdot d$ )

- finding median is  $O(d)$ , surprisingly! (naively its  $\log(d) \cdot d$ )

(but complicated... often better to use a well-implemented sort in practice)

- shuffling data (applying random permutation) is  $O(d)$ , also surprisingly

(naively it's done via sorting in  $O(d \log d)$ )

Use Fisher-Yates / Knuth shuffle

Aside on other well-known ensemble methods

like GCV...

Bootstrap resampling (Bradley Efron '79), extends jackknife, a bit like cross-validation

Resample as much as you want from a fixed sample

with replacement

... to estimate errors, confidence intervals

Refs: Larry Wasserman's "All of Statistics", or work of Victor Chernozhukov (econ at MIT)

Bagging = Bootstrap Aggregating reduces variance, helps w/ overfitting, often used on decision trees.

Resample dataset S into k new datasets (via bootstrap), train k learners, combine learners via averaging (regression) or voting (classification)