

## Last time:

- analysis of 3-stage boosting, justification of why it gives 35.2%-error  $h$  wrt  $\mathcal{D}$
- discuss ext. to "full" booster ( $1-\epsilon$  accuracy)
- " boosting over a fixed sample
- AdaBoost: simple, practical booster over a fixed sample;  
described alg, main ideas
  - reweight  $\mathcal{D}_t \rightarrow \mathcal{D}_{t+1}$  s.t.  
 $h_t$  is 50% acc. under  $\mathcal{D}_{t+1}$
  - clever weighted MAJ vote of  $h_1, \dots, h_T$  is final  $h$ .

Today: • analysis of AdaBoost  $\rightarrow$  state & prove thm about its performance  
• start unit on PAC learning in the presence of noise

• general framework,  
partic. noise models.

## Questions?

---

Let's consider reweighting rule:  $\left( \alpha_t = \frac{1}{2} \ln \left( \frac{1-\epsilon_t}{\epsilon_t} \right) \right)$   
 $(h_t, y_i \pm 1 \text{ valued})$

$$\mathcal{D}_{t+1}(i) = \frac{\mathcal{D}_t(i) \cdot \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

$\rightarrow$  • If  $h_t(x_i) = y_i$  (right):

$$\begin{aligned} \exp(-\alpha_t \overbrace{y_i h_t(x_i)}^{=1}) &= \exp(-\alpha_t) \\ &= \sqrt{\frac{\epsilon_t}{1-\epsilon_t}} \end{aligned}$$

Tot wt on  $i$  s.t.  $h_t(x_i) = y_i$  is  $1 - \epsilon_t$  :  
 mult by  $\exp(-\alpha_t x_i h_t(x_i))$

$$\hookrightarrow (1 - \epsilon_t) \sqrt{\frac{\epsilon_t}{1 - \epsilon_t}} = \sqrt{\epsilon_t (1 - \epsilon_t)}$$


---

• If  $h_t(x_i) \neq y_i$  (wrong):

$$\begin{aligned} \exp(-\alpha_t x_i \overset{=-1}{h_t(x_i)}) &= \exp(\alpha_t) \\ &= \sqrt{\frac{1 - \epsilon_t}{\epsilon_t}} \end{aligned}$$

Tot wt on  $i$  s.t.  $h_t(x_i) \neq y_i$  is  $\epsilon_t$  :  
 mult by  $\exp(-\alpha_t x_i h_t(x_i))$

$$\hookrightarrow (\epsilon_t) \sqrt{\frac{1 - \epsilon_t}{\epsilon_t}} = \sqrt{\epsilon_t (1 - \epsilon_t)}$$


---

Thm about AdaBoost:

Suppose AdaBoost is run for  $T$  stages.

Its final hyp  $H: \{x_1, \dots, x_m\} \rightarrow \{+1, -1\}$  makes errors on  $\leq \alpha$

$$\prod_{t=1}^T \sqrt{1 - 4\delta_t^2} \leq \exp\left(-2 \sum_{t=1}^T \delta_t^2\right)$$

frac of  $m$  pts in  $\{x_1, \dots, x_m\}$ .

$$(\delta_t = \frac{1}{2} - \epsilon_t).$$


---

Cor: If each  $\gamma_t \geq \gamma$ , can run Ada Boost for

$$T = \frac{1}{2\gamma^2} \ln\left(\frac{1}{\epsilon}\right) \text{ stages, } \epsilon \leq \epsilon.$$

( $\epsilon < \frac{1}{m} \rightarrow$  zero error on ex  $x_1, \dots, x_m$ )

optimal!

Pf: From 3 lemmas.

$$\begin{aligned} H(x) &= \text{sign}(f(x)) \\ f(x) &= \sum_{t=1}^T \alpha_t h_t(x) \end{aligned}$$

$$(L1): \frac{1}{m} |\{i \in [m] : H(x_i) \neq y_i\}| \leq \frac{1}{m} \cdot \sum_{i=1}^m \exp(-y_i f(x_i))$$

$$(L2): \frac{1}{m} \cdot \sum_{i=1}^m \exp(-y_i f(x_i)) \leq \prod_{t=1}^T z_t$$

$$(L3): \prod_{t=1}^T z_t \leq \prod_{t=1}^T \sqrt{1 - 4\gamma_t^2}$$

Pf of (L1):

Suppose  $i$  is s.t.  $H_t(x_i) \neq y_i$ .

$\text{sign}(f(x_i)) \neq y_i$ , so  $f(x_i) \cdot y_i < 0$ .

so

$-y_i f(x_i) > 0$ , so  $\exp(-y_i f(x_i)) > 1$ .

And if  $i$  is s.t.  $H_t(x_i) = y_i$ , get 0

on LHS of         $\forall > 0$  on RHS of       

So  $\forall i$ , LHS  $<$  RHS  
contrib. from  $i$       contrib. from  $i$ .

PF of L2: Recall def of  $\mathcal{D}_{T+1}(i)$ :

$$\begin{aligned} \mathcal{D}_{T+1}(i) &= \mathcal{D}_T(i) \cdot \frac{\exp(-\alpha_T y_i h_T(x_i))}{z_T} \\ &= \frac{\exp(-\alpha_T y_i h_T(x_i))}{z_T} \cdot \mathcal{D}_T(i) \\ &= \frac{\exp(-\alpha_T y_i h_T(x_i))}{z_T} \cdot \frac{\exp(-\alpha_{T-1} y_i h_{T-1}(x_i))}{z_{T-1}} \cdot \mathcal{D}_{T-1}(i) \\ &= \dots \\ &= \frac{\exp(-\alpha_T y_i h_T(x_i)) \cdot \exp(-\alpha_{T-1} y_i h_{T-1}(x_i)) \cdot \dots \cdot \exp(-\alpha_1 y_i h_1(x_i))}{z_T \cdot z_{T-1} \cdot \dots \cdot z_1} \mathcal{D}_1(i) \\ &= \frac{1}{m} \cdot \frac{\exp(-\sum_{t=1}^T \alpha_t y_i h_t(x_i))}{\prod_{t=1}^T z_t} \end{aligned}$$

Sum  $i=1, \dots, m$  both sides:

$$1 = \frac{1}{m} \sum_{i=1}^m \frac{\exp(-\sum_{t=1}^T \alpha_t y_i h_t(x_i))}{\prod_{t=1}^T z_t}, \text{ i.e.}$$

$$\prod_{t=1}^T Z_t = \frac{1}{m} \cdot \sum_{i=1}^m \exp(-y_i \sum_{t=1}^T \alpha_t h_t(x_i))$$

$\downarrow$   
 $f(x_i)$

Note: didn't yet use, in L1 or L2, setting of  $\alpha_t$

PF of L3: We'll show  $Z_t \leq \sqrt{1 - 4\epsilon_t^2}$ .

$$Z_t = \sum_{i=1}^m \mathcal{D}_t(i) \exp(-\alpha_t y_i h_t(x_i)) = A + B$$

$$\sum_{\substack{i: \\ h_t(x_i) \neq y_i}} \mathcal{D}_t(i) \exp(-\alpha_t y_i h_t(x_i)) = A$$

+

$$\sum_{\substack{i: \\ h_t(x_i) = y_i}} \mathcal{D}_t(i) \exp(-\alpha_t y_i h_t(x_i)) = B$$

$$\sum_{\substack{i: \\ h_t(x_i) \neq y_i}} \mathcal{D}_t(i) = \epsilon_t$$

We saw  $\epsilon_t$ , & saw that for  $i \in \text{sum}$ ,  $A$

$$\exp(-\alpha_t y_i h_t(x_i)) = \sqrt{\frac{1 - \epsilon_t}{\epsilon_t}}, \text{ so}$$

$$A = \epsilon_t \sqrt{\frac{1 - \epsilon_t}{\epsilon_t}} = \sqrt{\epsilon_t (1 - \epsilon_t)}$$


Similarly saw  $B = (1 - \epsilon_t) \sqrt{\frac{\epsilon_t}{1 - \epsilon_t}} = \sqrt{\epsilon_t(1 - \epsilon_t)}$

So  $Z_t = 2 \sqrt{\epsilon_t(1 - \epsilon_t)} = \sqrt{4\epsilon_t(1 - \epsilon_t)}$

Since  $\gamma_t = \frac{1}{2} - \epsilon_t$ , have  $\epsilon_t = \frac{1}{2} - \gamma_t$

$$2\epsilon_t = 1 - 2\gamma_t$$

$$2(1 - \epsilon_t) = 1 + 2\gamma_t$$

$\rightarrow = \sqrt{(1 - 2\gamma_t)(1 + 2\gamma_t)} = \sqrt{1 - 4\gamma_t^2}$  . 

$$1 - x \leq e^{-x} \quad \text{so} \quad \sqrt{1 - x} \leq e^{-x/2}$$

$$x = 4\gamma_t^2: \text{ done.}$$

---

---

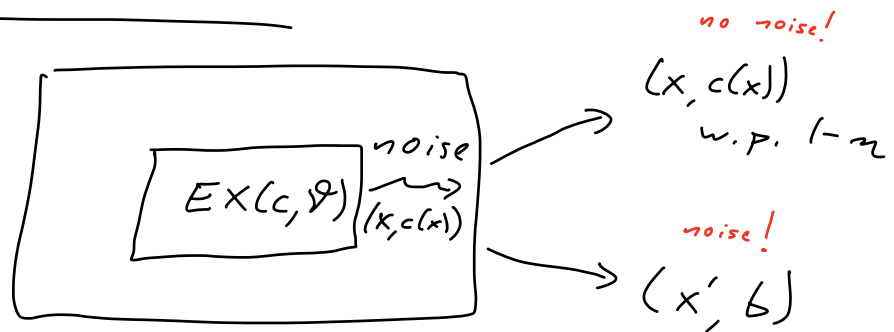
## PAC Learning in the Presence of Noise

General noise framework:

still have  $\epsilon$ , dist  $\mathcal{D}$ , still want

$\epsilon$ -acc by w.r.t  $c$  under  $\mathcal{D}$  w.p.  $1-\delta$ .

But now learner accesses a noisy example oracle: noise rate  $\eta$



2 diff noise models in this framework:

RCN

① Random (mis)classification noise:

$$x' = x, \quad b = \overline{c(x)}.$$

② Malicious noise:  $(x', b)$  is arbitrary.

(Should view as gen. by omniscient malevolent adv. who "knows"  $\mathcal{D}$ ,  $c$ ,  $\eta$ , state of learning alg, etc.)

Hi-level takeaway:

② Mal. noise very challenging: • if  $\eta \approx 2\epsilon$ ,

can't achieve error  $\leq \epsilon$ .



- Best known methods to deal w/ mal. noise: weak.

---

① RCN generally easy to deal with.

General method lets us convert many PAC alg's into variants that can handle RCN, even at  $\eta = 0.49$ .

( $\eta = \frac{1}{2}$ : learning impossible.)

---

Simple example of above:

Consider 2 worlds: (noise-free data sources)

---

$W1$	$W2$
$\frac{3}{8} + \text{pos ex}$	$\frac{5}{8} + \text{pos}$
$\frac{5}{8} - \text{neg ex}$	$\frac{3}{8} - \text{neg.}$

---

$W1$  or  $W2$  ?



Mal noise at rate  $\frac{1}{5} = n$ :

in both worlds, can make data look

50-50  $\overset{''}{\underset{''}{\updownarrow}}$

---

RCN at  $n < \frac{1}{2}$ :

W1: see + w.p.  
(true pos) (no noise) (true neg) (noise)  
 $\frac{3}{8} \cdot (1-n) + \frac{5}{8} \cdot n$

$$= \frac{3}{8} + \frac{n}{4}$$

W2: see + w.p.  $\frac{5}{8} - \frac{n}{4}$

IF  $n < \frac{1}{2}$ :  $\downarrow < \frac{1}{2}$ ,  $\downarrow > \frac{1}{2}$ .

So with enough  $\approx \left( \frac{1}{(1-2n)^2} \right)$  ex, can

tell whether in W1 or W2.

---

---

---

---