

Last time:

- finish proof of concept: 3-stage booster to go from 40% error to 35.2% error
- discuss boosting-by-filtering vs boosting over a fixed sample
- start discussion + analysis of AdaBoost (alg. for)

Today:

- analysis of AdaBoost: bound on error rate of its final hyp. over its input sample
- start unit on PAC learning in presence of noise
 - General framework, specific noise models:
 - malicious noise
 - random classification noise

Note: PS 4 due Tues Nov. 11

Questions?

Look at AdaBoost...

(with σ
+ list!)

Note: AdaBoost update rule: like WM \nearrow :

WM	AdaBoost
<ul style="list-style-type: none">• expert i• t^{th} trial $\left. \begin{array}{l} \downarrow \\ \downarrow \end{array} \right\}$• pred. of a_t \downarrow• expert i's wt a_t \downarrow• mult. update	<ul style="list-style-type: none">• i^{th} example in S• t^{th} run of weak learner• $h_t(x^i)$• $\mathcal{D}_t(x^i)$• mult. update

Twist: increase γ_t on mistakes
decrease

Interpret. of alg:

$$\begin{aligned} \text{Suppose } h_t(x^i) = y_i &\Rightarrow h_t(x^i) \cdot y_i = 1. \\ \Rightarrow \exp(-\alpha_t y_i h_t(x^i)) &= \exp(-\alpha_t) \\ &= \exp\left(-\frac{1}{2} \ln \frac{1-\epsilon_t}{\epsilon_t}\right) \\ &= \sqrt{\frac{\epsilon_t}{1-\epsilon_t}}. \end{aligned}$$

So in stage t ,
any x^i s.t. $h_t(x^i) = y_i$ has its wt (ignoring Z_t normaliz.)

• by $\sqrt{\frac{\epsilon_t}{1-\epsilon_t}} < 1$: wt goes down

OTOH, Suppose $h_t(x^i) \neq y_i \Rightarrow h_t(x^i) \cdot y_i = -1.$

$$\begin{aligned} \Rightarrow \exp(-\alpha_t y_i h_t(x^i)) &= \exp(\alpha_t) \\ &= \exp\left(\frac{1}{2} \ln \frac{1-\epsilon_t}{\epsilon_t}\right) \\ &= \sqrt{\frac{1-\epsilon_t}{\epsilon_t}} > 1 ; \text{ wt goes up.} \end{aligned}$$

Thm: Suppose run Ada. for T stages. Then its final hyp. H makes errors on \leq

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

frac. of the m points in S , where $\frac{1}{2} - \epsilon_t = \gamma_t$

is advantage of t^{th} weak hyp. h_t under \mathcal{D}_t ($1-x \leq e^{-x}$)

easy: $(1-4\gamma_t^2)^{\frac{1}{2}} \leq \exp((4\gamma_t^2) \cdot \frac{1}{2}) = \exp(-2\gamma_t^2)$

Cor: If each $\gamma_t \geq \gamma$, can run Ada. for $T \geq \frac{1}{2\gamma^2} \cdot \ln(1/\epsilon)$ stages, & achieve H s.t. wrong on $\leq \epsilon$ frac. of pts in S .
(solve $\epsilon \geq \exp(-2T\gamma^2)$).

(optimal!)

T follows from 3 claims:

$$H(x) = \text{sign}(f(x)),$$
$$f(x) = \sum_{t=1}^T \alpha_t h_t(x).$$

real value

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

Claim 2: $\frac{1}{m} \sum_{i=1}^m \exp(-y_i f(x_i)) = \prod_{t=1}^T z_t$.

Claim 3: For each $t \in T$, $z_t = \sqrt{1-4\gamma_t^2}$.

Proofs:

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

Pf: We'll show: for any i s.t. $H(x_i) \neq y_i$, have $\exp(-y_i f(x_i)) \geq 1$. (This does it b/c other i 's contrib. 0 on LHS, > 0 on RHS.)

Fix an i s.t. $H(x^i) \neq y_i$, i.e. $\text{sign}(f(x^i)) \neq y_i$.
 Suppose $y_i = +1$, so $f(x^i) < 0$; this means
 $-y_i f(x^i) > 0$, so $\exp(-y_i f(x^i)) > 1$. □

Claim 2: $\frac{1}{m} \sum_{i=1}^m \exp(-y_i f(x_i)) = \prod_{t=1}^T z_t$.

Pf: Consider $\mathcal{D}_{T+1}(x^i)$:

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

(Red bracket under $\mathcal{D}_1(x^i)$ with $\frac{1}{m}$ above it)

Sum over $i=1, \dots, m$: know $\sum_{i=1}^m \mathcal{D}_{T+1}(x^i) = 1$ \ddot{c} , so

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

$$\begin{aligned} \prod_{t=1}^T z_t &= \sum_{i=1}^m \frac{1}{m} \cdot \exp\left(-\sum_{t=1}^T \alpha_t y_i h_t(x^i)\right) \\ &= \sum_{i=1}^m \frac{1}{m} \exp\left(-y_i \sum_{t=1}^T \alpha_t h_t(x^i)\right) = \end{aligned}$$

$$= \frac{1}{m} \cdot \sum_{i=1}^m \exp(-y_i f(x^i)) \quad \blacksquare$$

(Now, in Claim 3, use def of α_t)

Claim 3: For each $t \in T$, $z_t = \sqrt{1 - 4r_t^2}$.

Pf: Fix a $t \in T$.

$a, b, c:$
 $\frac{a}{a+b+c}, \frac{b}{a+b+c}, \frac{c}{a+b+c}$
 sum to 1

$$\begin{aligned} z_t &= \sum_{i=1}^m \mathcal{D}_t(x^i) \exp(-\alpha_t y_i h_t(x^i)) \\ &= \sum_{i: h_t(x^i) \neq y_i} \mathcal{D}_t(x^i) \exp(-\alpha_t y_i h_t(x^i)) = A \\ &\quad + \sum_{i: h_t(x^i) = y_i} \mathcal{D}_t(x^i) \exp(-\alpha_t y_i h_t(x^i)) = B \end{aligned}$$

We saw: for an i in A , have

$$\exp(-\alpha_t y_i h_t(x^i)) = \sqrt{\frac{1 - \epsilon_t}{\epsilon_t}}. \quad \text{So } A =$$

$$\begin{aligned} \sqrt{\frac{1 - \epsilon_t}{\epsilon_t}} \cdot \sum_{i: h_t(x^i) \neq y_i} \mathcal{D}_t(x^i) &= \sqrt{\frac{1 - \epsilon_t}{\epsilon_t}} \cdot \epsilon_t \\ &= \sqrt{\epsilon_t (1 - \epsilon_t)}. \end{aligned}$$

We ^{also} saw: for an i in B ^{the} sum, have

$$\exp(-\alpha_t y_i h_t(x^i)) = \sqrt{\frac{\epsilon_t}{1 - \epsilon_t}}. \quad \text{So } B =$$

$$\sqrt{\frac{\epsilon_t}{1 - \epsilon_t}} \cdot \sum_{i: h_t(x^i) = y_i} \mathcal{D}_t(x^i) = \sqrt{\frac{\epsilon_t}{1 - \epsilon_t}} \cdot (1 - \epsilon_t)$$

$$= \sqrt{\epsilon_t(1-\epsilon_t)}.$$

$$\text{So } Z_t = A+B = 2\sqrt{\epsilon_t(1-\epsilon_t)}$$
$$= \sqrt{(2\epsilon_t)(2(1-\epsilon_t))}$$

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

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

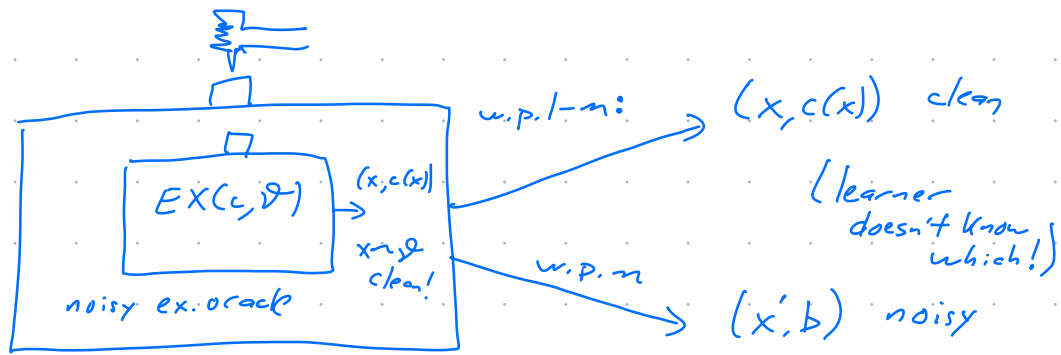
$$= \sqrt{(1-2\gamma_t)(1+2\gamma_t)} = \sqrt{1-4\gamma_t^2}.$$

Learning in Presence of Noise

Motivation: obvious: want more robust models.

Our noise framework: still is a target $c \in \mathcal{C}$,
an unknown dist. \mathcal{D} .

We assume learner has access to a $0 \leq \eta < \frac{1}{2}$
noisy example oracle: noise rate η



Diff. ^{specific} noise models: diff. assump. about (x', b) .

Goal: unchanged: still want to, w.p. $1-\delta$,
output h s.t. $err_{\mathcal{D}}[h, c] \leq \epsilon$.