

## Lecture 11

\* HW1 due today

\* I'll send a survey about HW1 & lectures, do file!

### RECAP

Last time:

\* Cryptographic assumptions  $\Rightarrow$  hardness of improper learning

\* Agnostically learning halfspaces is hard over some complexity assumptions (hardness of "refuting random  $k$ - $\text{XOR}$ ")

\* Agnostic learning is quite a challenging model.

Thm There is no efficient proper learning algorithm for agnostically learning rectangles in  $\mathbb{R}^d$  & conjunctions, unless  $\text{RP} = \text{P}$ .

Agnostic learning allows "worst-case noise" which makes learning hard.

Today: Learning under random noise.

## Learning with Random Classification Noise (RCN).

- \* Only the labels are noisy
- \* All labels are flipped with probability  $\eta$  (white-noise model)

ORACLE :  $EX^\eta(c, D)$

- Draws  $x \sim D$  from  $\mathcal{X}$
  - With prob.  $(1-\eta)$ , return  $(x, c(x))$
  - With prob.  $\eta$ , return  $(x, 1-c(x))$
- It costs one unit time to make call to  $EX^\eta(c, D)$

Criterion for success

with  $1-\delta$ , find hypothesis  $h$  s.t.

$$\text{error}(h; c, D) = \Pr_{x \sim D} [c(x) \neq h(x)] \leq \epsilon.$$

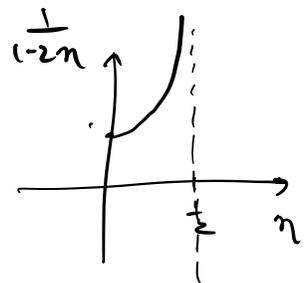
(comparing with true labels)

Restrict  $\eta$  to lie in  $[0, \frac{1}{2})$

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

As  $\eta \rightarrow \frac{1}{2}$ , learning becomes harder.

Allow factors of  $\left(\frac{1}{1-2\eta}\right)$  in complexity



bounds.

### Definition (PAC learning with RCM)

Let  $\mathcal{C}$  be a concept class and  $\mathcal{H}$  be a hypothesis class over  $\mathcal{X}^d$ . We say that  $\mathcal{C}$  is efficiently PAC learnable with RCM if:  $\exists$  a learning algo.  $A$  s.t.  $\forall c \in \mathcal{C}$ ,  $\forall$  dist.  $D$  over  $\mathcal{X}^d$ ,  $\forall \epsilon, \delta \in (0, \frac{1}{2})$ ,  $\forall \eta \in (0, \frac{1}{2})$  if  $A$  is given access to  $Ex^\eta(c, D)$  and inputs  $\epsilon, \delta$  &  $\eta_0$  s.t.  $\eta_0 \in (\eta, \frac{1}{2})$ , with prob  $\geq \delta$ ,  $A$  outputs  $h$  s.t.  $\text{error}(h; c, D) \leq \epsilon$  (prob. includes randomness in noisy labels).

We say that  $\mathcal{C}$  is efficiently PAC learnable if  $A$  runs in time  $\text{poly}(d, \frac{1}{\epsilon}, \frac{1}{\delta}, \frac{1}{1-2\eta_0})$ .

### Alg. for learning conjunctions

Previous alg. fails miserably on random noise!

Consider a robust algo. that relies on aggregate statistics instead of a single example.

• Suppose  $l$  is a literal ( $\bar{x}_i$  or  $x_i$ )

$$\text{Define } p_0(l) = \underset{a=0}{\text{Pr}} [l \text{ is } 0 \text{ in } a]$$

$$p_{\text{bad}}(l) = \underset{a=0}{\text{Pr}} [l \text{ is } 0 \text{ in } a \wedge c(a)=1]$$

Note that if  $l$  is in  $c(x)$ , then  $p_{\text{bad}}(l) = 0$ .

→  $l$  is significant if  $p_0(l) \geq \epsilon/8d$ .

→  $l$  is harmful if  $p_{\text{bad}}(l) \geq \epsilon/8d$ .

Claim: If  $h$  is conjunction of all significant literals that are not harmful, then  $\text{error}(h; c, D) \leq \epsilon$ .

Proof

$$\text{error}(h; c, D) = \underset{a=0}{\text{Pr}} [h(a)=1 \wedge c(a)=0] + \underset{a=1}{\text{Pr}} [h(a)=0 \wedge c(a)=1]$$

(caused if  $h(x)$  does not have some literal which  $c(x)$  has (i.e. that literal is set to 0 in example  $a$ ).

∴ caused due to insignificant literals

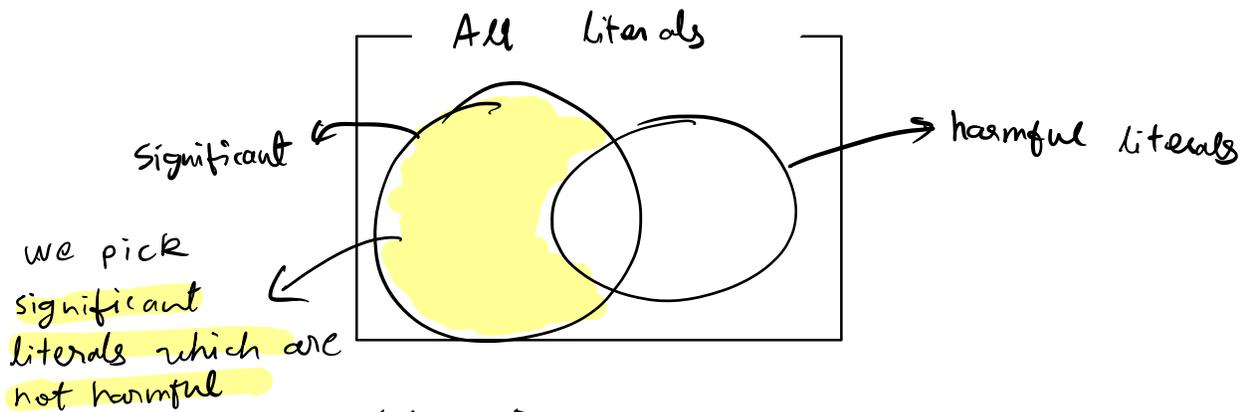
$$\leq 2d \cdot \left( \frac{\epsilon}{8d} \right) = \frac{\epsilon}{4}$$

(caused if  $h(x)$  has a literal which  $c(x)$  does not have (i.e. that literal is set to 0 in  $a$ )

$h(x)$  does not have any harmful literals

∴ caused due to significant literals which are not harmful

$$\leq 2d \cdot \left( \frac{\epsilon}{8d} \right) = \frac{\epsilon}{4}$$



$$\text{error}(h; c, D) \leq \epsilon/2$$

Still in noiseless case:

for any literal  $l$ ,  $P_+(l)$  &  $P_{\text{bad}}(l)$  can be estimated based on samples from  $E_{\pm}(c, D)$ . Can take a large set to get concentration with Chernoff.

This implies we can PAC learn conjunctions using this algorithm.

This algorithm can also be run with  $E_{\pm}^m(c, D)$ !

## Statistical Query Learning (Kearns ~90)

Do not have access to example oracle.

Statistical query oracle:  $STAT(C, D)$

query:  $(\phi, \tau)$

$$\phi: \mathcal{X} \times \{0, 1\} \rightarrow \{0, 1\}$$

$$\tau \in (0, 1)$$

$$\text{Let } P_\phi = P_{x \sim D} [\phi(x, c(x)) = 1]$$

$STAT(C, D)$  returns  $\hat{P}_\phi$  satisfying  $\hat{P}_\phi \in [P_\phi - \tau, P_\phi + \tau]$

### Examples

Consider a spam classification task.

"What fraction of emails labelled as spam (label is 1), have the words "Urgent" & "Free", but not "USC", and the number of words is  $\leq 20$ ?"

for our conjunction algorithm:

$$P_0(l) = P_{x \sim D} (l \text{ is } 0 \text{ in } x)$$

$$\phi: \mathcal{X} \times \{0, 1\} \rightarrow \{0, 1\}$$

$$\phi(x, y) = \begin{cases} 1 & \text{if } l \text{ is } 0 \text{ in } x \\ 0 & \text{otherwise} \end{cases}$$

$$P_{\phi_l} = \mathbb{E}_{x \sim D} [\phi_l(x, c(x))] = p_0(l)$$

$$P_{\text{bad}}(l) = P_{x \sim D} (l \text{ is } 0 \text{ in } a \text{ \& } c(a)=1)$$

$$\phi'_l(x, y) = \begin{cases} 1 & \text{if } l \text{ is } 0 \text{ in } a \text{ \& } y=1 \\ 0 & \text{otw} \end{cases}$$

$$P_{\phi'_l} = \mathbb{E}_{x \sim D} [\phi'_l(x, c(x))] = P_{\text{bad}}(l).$$

### Definition (SQ learning)

Let  $\mathcal{C}$  be a concept class and  $\mathcal{H}$  be a hypothesis class. We say  $\mathcal{C}$  is efficiently learnable from statistical queries (SQ) using  $\mathcal{H}$ , if  $\exists$  alg.  $A$ , and polynomials  $p(\cdot, \cdot)$ ,  $q(\cdot, \cdot)$  &  $r(\cdot, \cdot)$  s.t.  $\forall c \in \mathcal{C}$ ,  $\forall$  dist.  $D$  over  $\mathcal{X}^d$ ,  $\forall \epsilon \in (0, \frac{1}{2})$ , if  $A$  is given access to  $\text{STAT}(c, D)$  and input  $\epsilon$  then,

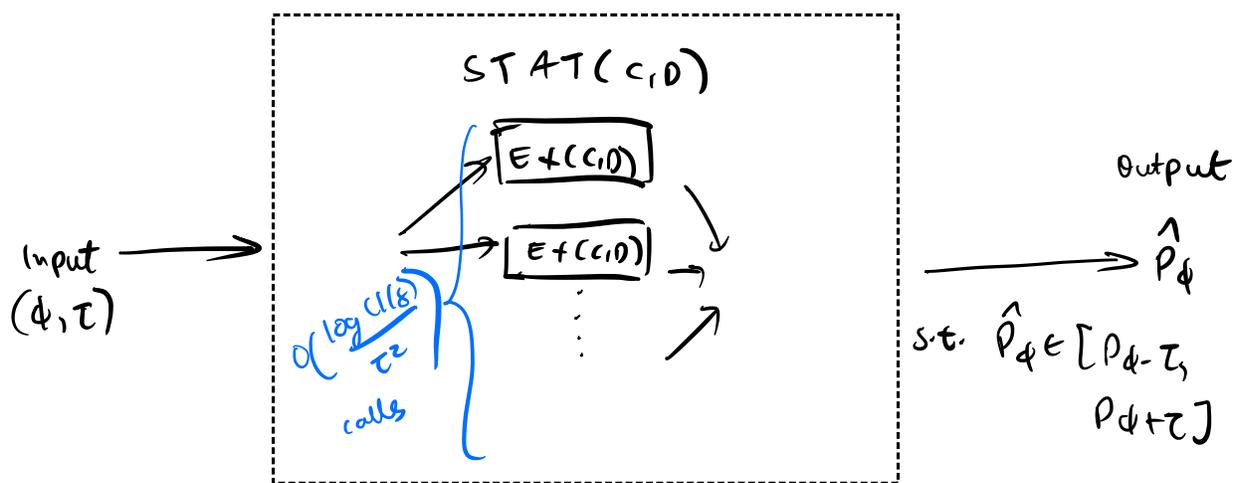
- If query  $(\phi, \tau)$  made by  $A$ , the predicate  $\phi$  can be evaluated in time  $q(\frac{1}{\epsilon}, d)$  and  $\frac{1}{\epsilon}$  is bounded by  $r(\frac{1}{\epsilon}, d)$ .
- $A$  runs in time  $p(\frac{1}{\epsilon}, d)$
- $A$  outputs  $h \in \mathcal{H}$  s.t.  $\text{error}(h; c, D) \leq \epsilon$ .

Note: no confidence parameter  $\delta$ . Before,  $\delta$  took care of the probability of seeing a set of "bad examples". But the  $\text{EX}(c, D)$  has been replaced by  $\text{STAT}(c, D)$ , which is deterministic.

**Thm** If  $\mathcal{L}$  is efficiently SQ learnable using  $\mathcal{H}$ , then  $\mathcal{L}$  is efficiently PAC learnable using  $\mathcal{H}$ .

Proof Exercise. More details below.

We can implement  $\text{STAT}(c, \delta)$  using  $E_{\tau}(c, \delta)$



By Chernoff, if we take  $O(\frac{\log(1/\delta)}{\tau^2})$  samples from  $E_{\tau}(c, \delta)$  then  $\hat{P}_{\phi}$  satisfies

$$\hat{P}_{\phi} \in [P_{\phi} - \tau, P_{\phi} + \tau] \text{ w.p. } (1 - \delta)$$

Aside: **Differential privacy**. If  $\mathcal{L}$  is efficiently SQ learnable, then  $\mathcal{L}$  is efficiently PAC learnable with differential privacy.

Thm: If  $\mathcal{L}$  is efficiently  $\mathcal{S}_2$  learnable, then  $\mathcal{L}$  is PAC learnable in the presence of RCN.

Thm The class of conjunctions is efficiently learnable in the  $\mathcal{S}_2$  model. Therefore, it is efficiently learnable with RCN.

Proof Our conjunction learning algo. can be implemented in the  $\mathcal{S}_2$  model. ■