

## Lecture 4

What can we learn privately?

Recall:

Probably Approximately Correct (PAC)

$X$  instance space      set of all instances  
(input space / domain)

$h: X \rightarrow \{+1, -1\}$  concept      a function to label elements

$\mathcal{H}$  concept class      a collection of labeling functions

$h^*$  target concept       $h^* \in \mathcal{H}$  and label all instances correctly

$D$  target distribution      distribution over instances

sample / training data set

}	$(x_1, h^*(x_1))$
	$(x_2, h^*(x_2))$
	$\vdots$
	$(x_m, h^*(x_m))$

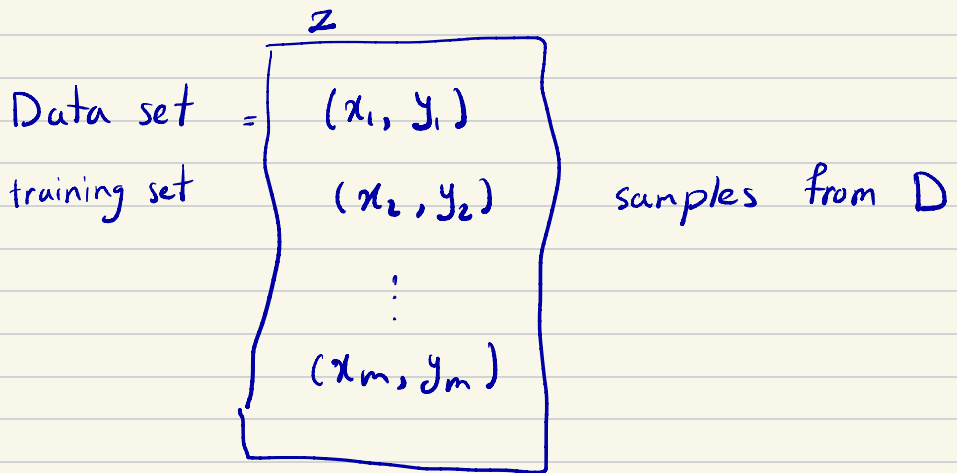
Agnostic case  $h^*$  DOES NOT exist

training set =

}	$(x_1, y_1)$
	$\vdots$
	$(x_m, y_m)$

# Recall

$$\text{true error : } \text{err}(h) := \Pr_{(x,y) \sim D} [ h(x) \neq y ]$$



$$\hat{\text{err}}(h) = \frac{|\{ i \mid i \in [m] : h(x_i) \neq y_i \}|}{m}$$

ERM : empirical risk minimizer

$$\text{output} := \arg \min_h \hat{\text{err}}(h)$$

↳ training error

Recall

## PAC learning (Probably Approximately Correct)

Suppose that we have a concept class  $\mathcal{H}$  over  $X$ . We say that  $\mathcal{H}$  is **PAC learnable** if there exists an algorithm  $A$  s.t.:

$$\forall h \in \mathcal{H}, \forall D \text{ over } X, \forall \alpha, \beta \in (0, 0.5]$$

$A$  receives  $\alpha, \beta$ , and samples  $\langle x_1, y_1 \rangle, \dots, \langle x_m, y_m \rangle$  where  $x_i$ 's are iid samples from  $D$ .

Then, w. p.  $\geq 1 - \beta$ ,  $A$  outputs  $\hat{h}$  s.t.

$$\text{err}(\hat{h}) \leq \min_{h \in \mathcal{H}} \text{err}(h) + \alpha$$

The probability is taken over the randomness in the samples and any internal coin flips of  $A$ .



\*

ERM works for a finite class  $\mathcal{H}$  if we have enough samples.

- Problem setup:

samples  $(x_1, y_1), \dots, (x_m, y_m) \sim D$

$$h \in \mathcal{H} : \text{err}(h) := \Pr_{(x,y) \sim D} [h(x) \neq y]$$

- Goal

find  $\hat{h} \in \mathcal{H}$  s.t. with probability  $1 - \beta$ ,  $\text{err}(\hat{h}) \leq \min_{h \in \mathcal{H}} \text{err}(h) + \alpha$

proof via uniform convergence.

$$\text{set } \alpha' = \frac{\alpha}{3}, \beta' = \frac{\beta}{2}$$

we show if  $m \geq O\left(\frac{\log(|\mathcal{H}|) + \log(1/\beta')} then$

$$\Pr[\forall h \in \mathcal{H} : |\hat{\text{err}}(h) - \text{err}(h)| < \alpha'] \geq 1 - \beta'$$

For all  $h \in \mathcal{H}$ , we define:

$X_h := \#$  samples in  $z$  where  $h(x_i) \neq y_i$

$$\hat{\text{err}}(h) = \frac{X_h}{m}$$

Each sample is mislabeled w.p.

$$\text{err}(h) = \Pr_{(x,y) \sim D} [h(x) \neq y].$$

$\Rightarrow X_h$  is a binomial random variable.

$$X_h \sim \text{Bin}(m, \text{err}(h))$$

Via Hoeffding bound:

$$\Pr [ |\hat{\text{err}}(h) - \text{err}(h)| > \alpha' ] =$$

$$\Pr [ \left| \frac{X_h}{m} - \text{err}(h) \right| > \alpha' ] \leq 2e^{-\frac{m\alpha'^2}{3}}$$

$$\Pr [ \exists h \in \mathcal{H} : | \hat{\text{err}}(h) - \text{err}(h) | > \alpha' ]$$

$$\leq \sum_{h \in \mathcal{H}} \Pr [ | \hat{\text{err}}(h) - \text{err}(h) | > \alpha' ]$$

$$\leq |\mathcal{H}| \cdot 2 \cdot e^{-\frac{m \alpha'^2}{3}} \leq \beta'$$

This holds for all  $\leftarrow$

$$m \geq \frac{3 (\ln(2|\mathcal{H}|) + \ln(1/\beta'))}{\alpha'^2}$$

$$\Rightarrow \Pr [ \forall h \in \mathcal{H} : | \hat{\text{err}}(h) - \text{err}(h) | \leq \alpha' ]$$

$$\geq 1 - \beta'$$

The guarantee of the  
uniform convergence

## Recall

Uniform convergence (UC) implies agnostic PAC learnability via ERM.

Outcome of ERM minimizes the  $\hat{\text{err}}(h)$ :

$$* \hat{h}_{\text{ERM}} := \arg \min_{h \in \mathcal{H}} \hat{\text{err}}(h)$$

$$h^* := \arg \min_{h \in \mathcal{H}} \text{err}(h)$$

$$* \Rightarrow \hat{\text{err}}(\hat{h}_{\text{ERM}}) \leq \hat{\text{err}}(h^*)$$

$$\text{UC} \stackrel{\text{w.p. } 1-\beta}{\Rightarrow} \left\{ \begin{array}{l} \hat{\text{err}}(\hat{h}_{\text{ERM}}) \leq \text{err}(h_{\text{ERM}}) + \alpha' \\ \hat{\text{err}}(h^*) \leq \text{err}(h^*) + \alpha' \end{array} \right.$$

$$\Rightarrow \hat{\text{err}}(\hat{h}_{\text{ERM}}) \leq \text{err}(h^*) + 2\alpha'$$

$$\leq \min_{h \in \mathcal{H}} \text{err}(h) + 2\alpha'$$

as desired in PAC learning.

## Exponential mechanism:

Help us to pick an element privately

→ an element that maximize a score  
or utility

Algorithm:

Input: dataset  $Z$ , set of options  
score function  $u$ , privacy  
parameter  $\epsilon$ ,

Output: pick  $h \in \mathcal{H}$  with prob.

$$\propto e^{\frac{\epsilon \cdot u(h, Z)}{2\Delta}}$$

$\Delta$  is sensitivity of  $h$ :

$$\Delta u := \max_{\substack{h, Z, Z'}} |u(h, Z) - u(h, Z')|$$

↙  
neighboring

## Proof of privacy:

for all  $h$ :

$$\frac{e^{\frac{\epsilon \cdot u(h, z)}{2\Delta}}}{e^{\frac{\epsilon \cdot u(h, z')}{2\Delta}}} = e^{\frac{\epsilon |h(z) - h(z')|}{\Delta}} \leq e^{\epsilon/2}$$

$$\Pr[\text{output} = h | z] = \frac{e^{\frac{\epsilon \cdot u(h, z)}{2\Delta}}}{\sum_{h'} e^{\frac{\epsilon \cdot u(h, z)}{2\Delta}}}$$

$$\leq \frac{e^{\frac{\epsilon \cdot u(h, z')}{2\Delta}} \cdot e}{\sum_{h'} e^{\frac{\epsilon \cdot u(h, z')}{2\Delta}} \cdot e^{-\epsilon/2}}$$

$$= e \cdot \Pr[\text{output} = h | z']$$

## Score performance of output $\hat{h}$

Thm  $\Pr [ u(\hat{h}) \leq \text{OPT} - \frac{2\Delta}{\epsilon} (\ln(H) + t) ]$   
 $\leq e^{-t}$

$\rightarrow \text{OPT} = \max_h u(h, z)$

proof:

$c := \text{OPT} - \frac{2\Delta}{\epsilon} (\ln(H) + t)$

$$\Pr [ u(\hat{h}, z) \leq c ] = \sum_{h: u(h, z) \leq c} \Pr [\text{pick } h | z]$$

$$\leq \frac{\sum_{h: u(h, z) \leq c} e^{\frac{\epsilon u(h, z)}{2\Delta}}}{\sum_{h \in \mathcal{H}} e^{\frac{\epsilon u(h, z)}{2\Delta}}}$$

$$\leq \frac{|\mathcal{H}| e^{\frac{\epsilon c}{2\Delta}}}{e^{\frac{\epsilon \text{OPT}}{2\Delta}}}$$

$$= \exp \left( \ln(H) + \frac{\epsilon}{2\Delta} (c - \text{OPT}) \right) \leq e^{-t}$$

# Private PAC learning

Exponential mechanism chooses an item with high utility. We can use this mechanism to mimic ERM. We privately select a concept that labels a large number of samples correctly.

For all  $h \in \mathcal{H}$ , we define  $u(h, z)$  to be the number of correctly labeled samples:

$$u(h, z) = \left| \left\{ i \mid h(x_i) = y_i \right\} \right|$$

↓  
training set

Note that  $\hat{\text{err}}(h) = 1 - \frac{u(h, z)}{m}$

ERM minimizes  $\hat{\text{err}}(h)$ .

Exponential mechanism privately maximizes  $u(h, z)$



Note  $\hat{h}_{ERM}$  maximizes  $u(h, z)$ :

$$\text{OPT} = u(\hat{h}_{ERM}, z)$$

Let  $\hat{h}_m$  be the outcome of the exponential mechanism.

— we have shown w.p.  $1 - \beta'$ :  $\Delta$  of  $u$  is 1

$$u(\hat{h}_m, z) > u(\hat{h}_{ERM}, z) - \frac{2\Delta}{\epsilon} (\ln(1/\beta') + \ln(1/\beta'))$$

$\underbrace{\hspace{10em}}_{\text{OPT}}$

$$= m - m \cdot \text{err}(\hat{h}_m)$$

$$\Rightarrow m - m \cdot \text{err}(\hat{h}_m) > m - m \cdot \text{err}(h_m) - \frac{2}{\epsilon} (\ln(1/\beta') + \ln(1/\beta'))$$

$$\left(\times \frac{-1}{m}\right) \Rightarrow \text{err}(\hat{h}_m) < \text{err}(\hat{h}_{ERM}) + \frac{2}{m\epsilon} (\ln(1/\beta') + \ln(1/\beta'))$$

Thus, for all  $m \geq \frac{2(\ln(1/\beta') + \ln(1/\beta'))}{\epsilon \alpha'}$ , we get:

$$\text{err}(\hat{h}_m) \leq \text{err}(\hat{h}_{ERM}) + \alpha' \quad \textcircled{1}$$

if  $m \geq \Theta\left(\frac{\log(|\mathcal{H}|) + \log(1/\beta')}{\alpha'^2}\right)$ , with probability of  $1 - \beta'$  uniform convergence holds:

$$\forall h \quad |\hat{\text{err}}(h) - \text{err}(h)| < \alpha' \quad (2)$$

Finally lets analyze the error of  $\hat{h}_m$ :

Using (1) and (2), with probability  $1 - 2\beta'$ :

$$\text{err}(\hat{h}_m) \leq \hat{\text{err}}(\hat{h}_m) + \alpha' \quad \text{via (2)}$$

$$\leq \hat{\text{err}}(\hat{h}_{\text{ERM}}) + 2\alpha' \quad \text{via (1)}$$

$$\leq \hat{\text{err}}(h^*) + 2\alpha' \quad \text{using ERM}_h \text{ minimizes err.}$$

$$\leq \hat{\text{err}}(h^*) + \underbrace{3\alpha'}_{\alpha'} \quad \text{via (2)}$$

$$\Rightarrow \text{If } m \geq \Theta\left(\log(|\mathcal{H}|) + \log(1/\beta)\right) \cdot \left(\frac{1}{\alpha^2} + \frac{1}{\alpha\epsilon}\right)$$

$$\text{then } \Pr\left[\text{err}(\hat{h}_m) \leq \min_{h \in \mathcal{H}} \text{err}(h) + \alpha\right] \geq 1 - \beta$$

Note that for  $\epsilon > \alpha$  privacy is "free": That is the sample complexity increases by a constant factor.

Otherwise, it's linear in terms of  $\alpha$  and  $\epsilon$ .