# Lecture 14:

## PAC learnability

## Uniform convergence

# Recall:

## Probably Approximately Correct (PAC)

$X$    instance space      set of all instances

     (input space / domain)

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

$C$    concept class      a collection of labeling functions

$c^*$    target concept      $c^* \in C$   and   label all instances
                                       correctly

$D$    target distribution    distribution over instances

Sample / training data set
$$
\begin{cases}
\langle x_1, c^*(x_1) \rangle \\
\langle x_2, c^*(x_2) \rangle \\
\quad \vdots \\
\langle x_n, c^*(x_n) \rangle
\end{cases}
$$

# PAC learning    (Probably Approximately Correct)

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

$\forall c \in C, \quad \forall D$ over $X, \quad \forall \varepsilon, \delta \in (0, 0.5]$

$A$ receives $\varepsilon, \delta$, and samples $\langle x_1, c(x_1) \rangle$ $\dots, \langle x_n, c(x_n) \rangle$ where $x_i$'s are iid samples from $D$.

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

proper $[\in C]$

$$\text{err}(\hat{c}) \leq \varepsilon.$$

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

other notation

true error:

$$\text{err}(c) = \Pr_{(x, y) \sim D}\left[c(x) \neq y\right]$$

training error:

$$\hat{\text{err}}(c) = \frac{\#\ \text{Samples in } T \ \text{s.t } c(x_i) \neq y_i}{|T|}$$

fraction of samples in the training set that $c$ is mis-labeled.

**Example 2**   Boolean conjunctions :

$$X = \{0, 1\}^n \qquad\qquad \text{literals} \begin{cases} x_i \\ \overline{x}_i \end{cases}$$

$$\text{Conjunction} = \begin{cases} \text{literal} \\ \\ \text{literal} \;\wedge\; \text{conjunction} \end{cases}$$

$\longrightarrow$ logical and

concept :   a conjunction

example : $h(x) = x_1 \wedge \overline{x}_2$     $h((1,0,1)) = 1$

$x_5 (x_1, \dots, x_n)$     $h((0,0,1)) = 0$

$\mathcal{H}$ :   the set of all conjunction function

**Goal:**   PAC learning of $\mathcal{H}$

Suppose   we have   samples   of the form $\langle x, h^*(x) \rangle$ from a distribution $D$

$\hookrightarrow$ realizable

## Algorithm:

- start with $h = x_1 \wedge \bar{x}_1 \wedge x_2 \wedge \bar{x}_2 \wedge \ldots \wedge x_n \wedge \bar{x}_n$

- Try $m = ?$ examples

  - ignore negative example.
  - for positive example, remove inconsistent terms.

- Output $h$


Deleting an inconsistent literal

$h = \boxed{x_2} \wedge x_3 \wedge \bar{x}_4$

sample: $\langle (1, \boxed{0}, 1, 0), 1 \rangle$ $\rightarrow$ $x_2$ is inconsistant

$\downarrow$

we delete $x_2$ from $h$

$\downarrow$

new $h = x_3 \wedge \bar{x}_4$

Our goal is to analyze the performance of the algorithm.

First, we start by the error of the output hypothesis $\hat{h}$.

Initially, $h$ contains all literals. we only remove inconsistent literals. So, we never removes literals in $h^*$ from $h$. That is, $\hat{h}$ contains all the literals in $h^*$. This fact implies if $h^*(x) = 0$, $\hat{h}(x)$ must be zero too

$\Rightarrow$ Hence $\hat{h}$ always labels $x$ correctly if $h^*(x) = 0$

Now consider the rest of the domain elements $\alpha$ such that $h^*(\alpha) = 1$
If $\hat{h}$ makes a mistake (i.e. $\hat{h}(\alpha) = 0$), there must be a literal in $\hat{h}$, $\mathcal{l}$ that is inconsistent:

$$
\text{true error of } \hat{h} = err(\hat{h})
$$

$$
= \Pr_{\alpha \sim D} [\hat{h}(\alpha) \neq h^*(\alpha)]
$$

$$
= \Pr_{\alpha \sim D} \left[ \exists \text{ a literal } \mathcal{l} \in \hat{h} \text{ such that } \mathcal{l} \atop x_{\mathcal{l}} = 0 \text{ but } h^*(\alpha) = 1 \right]
$$

$$
= \sum_{\mathcal{l} \in \hat{h}} \Pr_{\alpha \sim D} [x_{\mathcal{l}} = 0 \text{ but } h^*(\alpha) = 1]
$$

call this $p(\mathcal{l})$

$$
= \sum_{\mathcal{l} \in \hat{h}} p(\mathcal{l}) \qquad *
$$

We call a literal bad iff $p(z)$

is at most $\dfrac{\varepsilon}{2n}$

bad $z$ $\longleftrightarrow$ $p(z) > \dfrac{\varepsilon}{2n}$

Using $*$ it is easy to see if

no bad literal survives in $\hat{h}$ then

$$\text{err}(\hat{h}) \leq \sum_{z \in \hat{h}} p(z) \leq 2n \cdot \dfrac{\varepsilon}{2n} \leq \varepsilon$$

$\Rightarrow$ hence the error of $\hat{h}$ is good

Now, let's focus on the probability
of $\text{err}(\hat{h}) > \varepsilon$

$$\Pr_T \left[ \text{outputting an inaccurate } \hat{h} \right]$$

training → $T$
set $= \Pr_T \left[ \text{err}(\hat{h}) > \varepsilon \right]$

$$\leq \Pr \left[ \exists \text{ a bad literal } z \text{ in } \hat{h} \right]$$

$$\leq 2n \cdot \Pr \left[ \begin{array}{l} \text{a bad literal survives} \\ \text{all the } m \text{ samples (not} \\ \text{been deleted )} \end{array} \right]$$

$$\leq 2n \cdot \left( 1 - p(z) \right)^m \quad \leftarrow \text{It is not hard}$$

to see that with

$$\leq 2n \left( 1 - \frac{\varepsilon}{2n} \right)^m$$

probability $p(z)$

$$\leq 2n \, e^{-\frac{\varepsilon m}{2n}} \quad \overset{?}{\leq} \delta$$

we delete $z$
at every
round.

by setting $m = \frac{2n}{\varepsilon} \log\left(\frac{2n}{\delta}\right)$

Hence our algorithm with prob.
$1-\delta$ output $\hat{h}$ that has
low error. $(\text{err}(\hat{h}) \leq \varepsilon)$

$\Rightarrow$ we PAC-learned $\mathcal{H}$ :)

# ERM

In both example we picked concepts $\hat{R}$ and $\hat{h}$ that were consistent with the samples in the training set

What we did is called :

ERM : Empirical Risk Minimization

comes from samples   error

ERM algorithm: it finds a concept $\hat{h}$ such that $\hat{err}(\hat{h}) = 0$
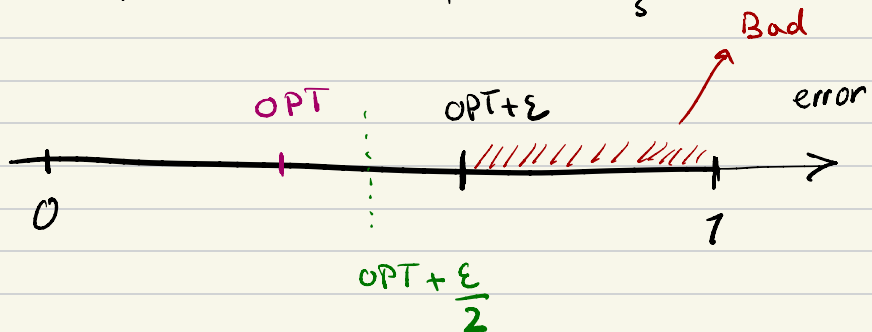
**+** Uniform convergence. (UC)

Class C has the uniform convergence property if $\forall \varepsilon, \delta \in (0,1)$, dist D $\exists$ m (as a function of $\varepsilon, \delta, \mathcal{H}$, but not D since we don't know D). s.t. for a training set of size m:

$$\Pr_{T \sim D^m}\left[\forall c \in C : \left|\hat{err}_T(c) - err(c)\right| \leq \varepsilon\right] \geq 1-\delta$$

Uniform convergence implies agnostic PAC learnability via EMR.

UC $\Rightarrow$ $\forall c \in C_B$ $\hat{err}_S(c) > OPT + \varepsilon/2$

UC $\Rightarrow$ $c^* =$ the best option) $\hat{err}_S(c^*) \leq OPT + \varepsilon$

Bad



OPT     OPT+ε     error

0                                          1

OPT + $\frac{\varepsilon}{2}$

There are two types of error in the agnostic setting:

$$\text{err}(\hat{c}) < \min_{c \in C} \text{err}(c) + \varepsilon$$

$\underbrace{\qquad\qquad}_{\varepsilon_{app} = \text{approximation error}}$  $\underbrace{\phantom{xx}}_{\varepsilon_{est} = \text{estimation error}}$

$\downarrow$

depends only to the choice of the class $C$

_ Is $C$ rich enough to capture how data is labeled?

$C$
larger    $\varepsilon_{app}$    $\varepsilon_{est}$
more complex

$\uparrow$    $\downarrow$    $\uparrow$