# Lecture 15

Uniform convergence

overfitting

PAC learnability of finite classes.

No free lunch theorem

+ some proof witness

Last lecture:

Recall:

+ 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.

ERM could go very wrong if we
overfit.

$$\hat{R}(x) = \begin{cases} y_i & x=x_i \in T \\ 0 & x=x_i \in T \end{cases}$$

0 empirical error $\Big\}$ error 1 on any dist
with a continuous domain

ERM has really bad error!

*
ERM works for a finite class $C$ if we have enough samples.

- Problem setup:

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

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

Realizable case

Assume $\exists \; c^* \in C$ s.t. $err(c^*) = 0$

- Goal

find $\hat{c} \in C$ s.t. with probability $1 - \delta$, $err(\hat{c}) \leq \varepsilon$.

- Proof

Bad hypotheses $C_B := \{ c \in C \mid err(c) > \varepsilon \}$

$$\hat{err}_T(c) := \frac{|\{(x,y) \in T \mid c(x) \neq y\}|}{|T|}$$

Misleading training samples

$$M := \{T \mid \exists c \in C_B \text{ s.t. } \hat{err}_T(c) > 0\}$$

upon observing $T$, we may pick $c$ that
is a bad choice, but it "looked"
good from ERM perspective, since
$\hat{err}_T(c) = 0$.

Our goal is to show observing a
dataset $T \in M$ happens only with
probability $\delta$.
This is sufficient to prove ✱.

fix $\quad c \in C_B$

what is the probability of

$\hat{err}_T (c) = 0$

$$\Pr_{T \sim D^m} \left[ \hat{err}_T (c) = 0 \right]$$

$$= \Pr_{T \sim D^m} \left[ \forall (x,y) \in T . \; c(x) = y \right]$$

iid samples $\qquad \Longrightarrow = \left( \Pr_{(x,y) \sim D} \left[ c(x) = y \right] \right)^m$

err $(c) > \varepsilon \qquad \longrightarrow < (1-\varepsilon)^m \qquad \leq e^{-\varepsilon m}$

Now, we are ready to bound

$$\Pr_{T \sim D^m} \left[ T \in \mathcal{M} \right]$$

$$= \Pr_{T \sim D^m} \left[ \exists c \in C_B \text{ s.t. } \hat{err}_T(c) = 0 \right]$$

$$= \sum_{c \in C_B} \Pr_{T \sim D^m} \left[ \hat{err}_T(c) = 0 \right]$$

$$\leq |C_B| \cdot e^{-\varepsilon m} \leq |C| \cdot e^{-\varepsilon m}$$

set $\quad m = \dfrac{\log(|C|/\delta)}{\varepsilon}$

$$\Rightarrow \Pr \left[ \text{outputting a misleading } c \right]$$

$$\leq \delta$$

# The agnostic case :

what if there is no perfect $c \in C$?

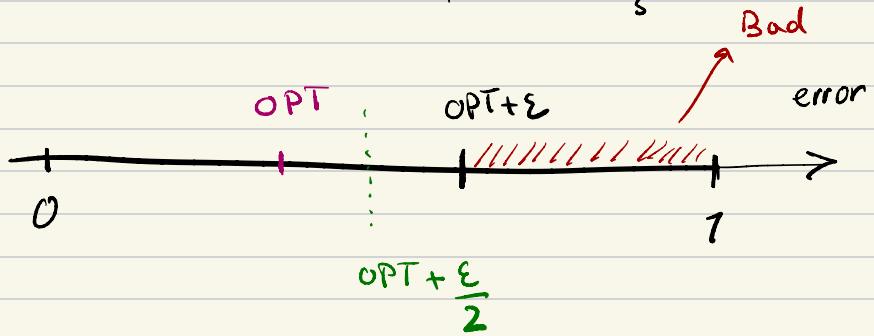$$\forall \, c \in C \qquad err(c) \quad > \quad 0$$

Goal

Find $\hat{c} \in C$ s.t.

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

$\underbrace{\min_{c \in C} err(c)}_{= OPT}$

the best possible option

## Uniform convergence implies agnostic PAC learnability via EMR.

$$UC \implies \forall c \in C_B \qquad \hat{err}_s(c) > OPT + \varepsilon/2$$

$$UC \implies c^* = \text{the best option}) \quad \hat{err}_s(c^*) \leq OPT + \varepsilon$$

Bad



Exercise!

Suppose we have a finite class $C$, and $m = O\left(\dfrac{(\log |C|/\delta)}{\varepsilon^2}\right)$. then w.p. at least $1-\delta$, for all $c \in C$, we have:

$$\left| \hat{err}_s(c) - err(c) \right| < \varepsilon/2$$

No free lunch theorem says if
there is no universal learner $_0^0$
for a complex $C$ even when
$\varepsilon_{app}$ is $0$, $\varepsilon_{est} \gg$ constant
with some constant probability

[ unless we have $\Omega (|X|)$ samples]

Suppose we have a set of 2m points

There are $2^{2m}$ possible labelings of these 2m points.

Suppose C is the class of $2^{2m}$ func. that assigns these labelings to these points.

Fix a labeling of the points ↗

Now assume D is the uniform distribution on the 2m points with their label.

$T \leftarrow$ Draw m samples from D

(WLOG assume they are unique)

How many function in C label T correctly? $\quad 2^m$

$$P := \left\{ c \in C \mid \hat{err}_T(c) = 0 \right\}$$

↳ promising hypothese. $\quad |P| = 2^{m/2}$

How many of them has error $< \varepsilon$ ?

$c$ is misleading if $\begin{cases} \mathrm{err}(c) > \varepsilon \\ \text{and } \hat{\mathrm{err}}_T(c) = 0 \end{cases}$

$$\mathcal{M} := \left\{ c \in C \,\middle|\, \mathrm{err}(c) > \varepsilon \ \& \ \hat{\mathrm{err}}_T(c) = 0 \right\}$$

$$|\mathcal{M}| = \frac{|\mathcal{M}|}{|P|} \cdot |P|$$

$$= 2^m \cdot \Pr_{c \sim_u P}\left[ c \in \mathcal{M} \right] \qquad \textcolor{green}{c \text{ makes}}$$

<span style="color:magenta">a random concept<br>in $P$</span>

$$\qquad\qquad\qquad\qquad\qquad \textcolor{green}{\geq m \cdot \varepsilon}$$
$$\textcolor{green}{\text{mistakes}}$$
$$\textcolor{green}{\text{in expectation}}$$

$$= 2^m \cdot \Pr\left[ \frac{\#\,\text{mistake}}{m} < \varepsilon \right]$$

$$= 2^m \left( 1 - \Pr\left[ \frac{\#\,\text{mistakes}}{m} < \tfrac{1}{2} - (\tfrac{1}{2} - \varepsilon) \right] \right)$$

$$\geq 2^m \left( 1 - e^{\left( -2m (\tfrac{1}{2} - \varepsilon)^2 \right)} \right)$$

$$\nearrow$$
Hoeffding bound

$$\geq 2^{m/2} \cdot 0.99$$
$$\nearrow$$
$$\varepsilon \leq \tfrac{1}{4} \qquad\qquad m \geq 40$$

⟹ 0.99 % of the promising concept
are bad!