# THEORY OF MACHINE LEARNING

# LECTURE 3 4 ·

## PAC MODEL, VC DIMENSION

# RECAP – VALIANT'S THEORY OF (SUPERVISED) LEARNING

$$(\text{concept}/\text{hyp})$$
$$\text{labeling function}$$

- **Learnability (from examples).** [Suppose D is fixed.] We say that a concept class is "learnable" if there exists an [efficient] algorithm **A** with the property: for all $\epsilon > 0$, there exists $m$ (number of samples) such that when given $m$ _i.i.d._ samples from D along with their labels, **A** produces a hypothesis $h$ with risk less than $\epsilon$, with prob. >= 0.9

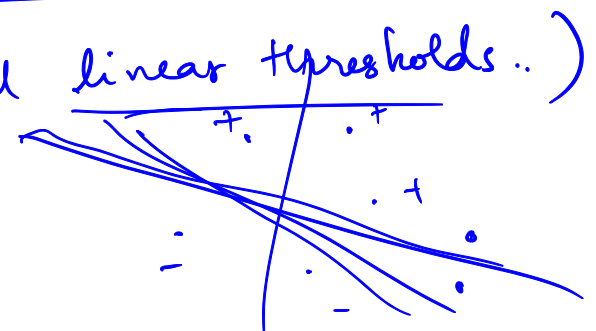$$R_{D, \ell}(h) = \Pr_{x \sim D} \left( h(x) \neq \ell(x) \right)$$

- (Recall, risk = expected error on sample from distribution)

- Beyond examples? (technically yes, e.g., teacher/student)

# RECAP: NO FREE LUNCH THEOREM

*inductive bias.*

- <u>Motivation:</u>  do we really need to restrict the hypothesis/concept class *before starting* learning? – yes!

- <u>No free lunch:</u> (informal) there is no "universal" learner, <u>even if it's allowed to be inefficient</u>  (even for binary classification under a uniform distribution, unless it "sees most of the labels")

- Proof via a counting argument – too many hypotheses

$\mathcal{H}$ : - class of hypotheses  (all *linear thresholds*..)

    - all 100 layer NNs with width 100k.

## TODAY'S PLAN

- Definition. (Agnostic) PAC learning

- Finite classes are PAC learnable

- Dealing with infinite classes: 'growth function' and VC dimension

# PAC LEARNING (REALIZABLE CASE)

(over some domain X) . $\mathcal{H}$: bunch of hypotheses over X.

- **Learnability of a concept class.** A concept class H is *PAC learnable* (over domain X) if there exists an algorithm **A** that for all $\epsilon, \delta > 0$ and distributions D, has the following property:

  $\exists$ a function $m(\epsilon, \delta)$

  $2^{1/\epsilon^2} \cdot 2^{2^{1/\delta}}$

  $N \quad \frac{\log \frac{d}{\delta}}{\epsilon^2}$ ?

  - given $m(\epsilon, \delta)$ samples $(x, f(x))$, where $x \sim$ D and f is a (unknown) function (in H) it outputs h with risk at most $\epsilon$ with probability at least $1 - \delta$.

    $0.05$ $\quad \frac{2}{3}$ (label function is in $\mathcal{H}$.) .

- (The sample size must not depend on D)

- As such h need not belong to H (improper learning)

  $f \in \mathcal{H}$ $\qquad (x, f(x))$ $\quad$ X

  - goal is to find true label function f
  - declare success, if we find h that has risk $\le \epsilon$ $\qquad h \rightsquigarrow f$ .

Concept ≡ hypothesis.

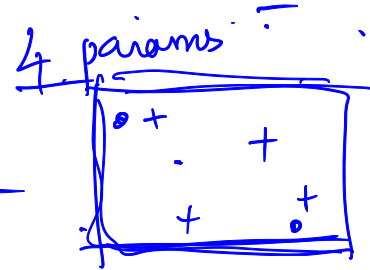Example of $\mathcal{H}$:    $X = \mathbb{R}^2$

3 params

$\text{Sign}\left(\frac{ax+by+c}{a}\right)$    $\mathcal{H}:$ all linear classifiers.

$ax+by+c$

$\text{Sign}\left(-(ax+by+c)\right)$

$X = \mathbb{R}^2$

$\mathcal{H}:$ all rectangles
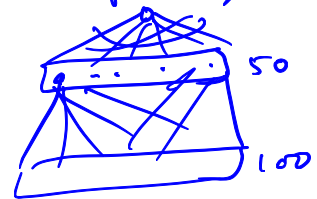
4 params

$X: \mathbb{R}$

1 param

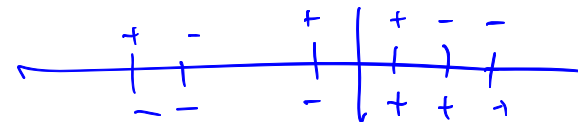$\mathcal{H}:$ all threshold functions

$X = \mathbb{R}^{100}$  ;  $\mathcal{H}:$ all depth 2 nns with 100 inputs,

every gate is ReLU.

50

100

# PAC LEARNING (NON-REALIZABLE CASE)

(agnostic)

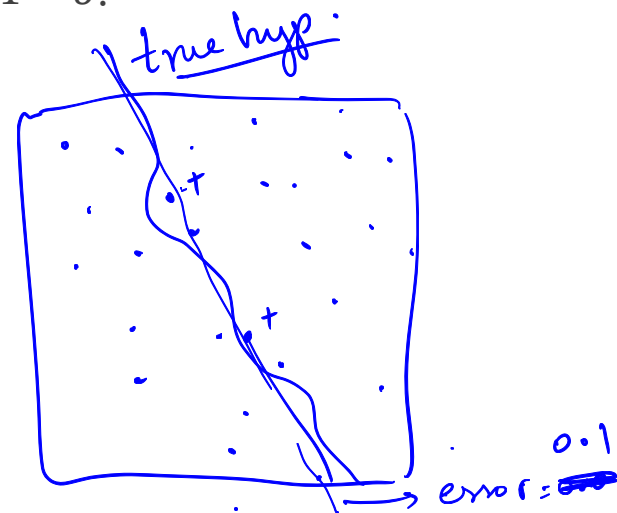- **Learnability of a concept class.** A concept class H is **agnostically** _PAC learnable_ (over domain X) if there exists an algorithm **A** that for all $\epsilon, \delta > 0$ and distributions D, has the following property:

  - given $m(\epsilon, \delta)$ samples $(x, f(x))$, where $x \sim D$ and $f$ is a (unknown) function <u>not</u> <u>necessarily in H</u>, it outputs $h$ with risk at most $\epsilon$ **more than the risk of the** $h'$ in H that is **"closest" to f,** with probability at least $1 - \delta$.

- (The sample size must not depend on D)

- Again, $h$ need not belong to H (improper learning)

true hyp.

→ can find $\underline{h}$ with risk $\leq 0.1 + \epsilon$.

"weakens" the inductive bias.

error := 0.1

0.1

# EVERY FINITE CLASS IS PAC LEARNABLE (EVEN AGNOSTIC)

$$\mathcal{H}: \{h_1, h_2, \ldots, h_N\} \; ; \; \text{input space } X \; (\text{potentially infinite}).$$

Theorem: $\mathcal{H}$ is PAC-learnable.

$$(x, f(x))$$

- Suppose H has only finitely many hypotheses (input space X may still be infinite)

example: $(x, f(x))$

Sample $S$: collection of these examples.

- **Generic algorithm:** underline{empirical risk minimization} (ERM)

  - get (m examples) (we'll fix m later)
  - find $h \in \mathcal{H}$ that is consistent with all these examples.
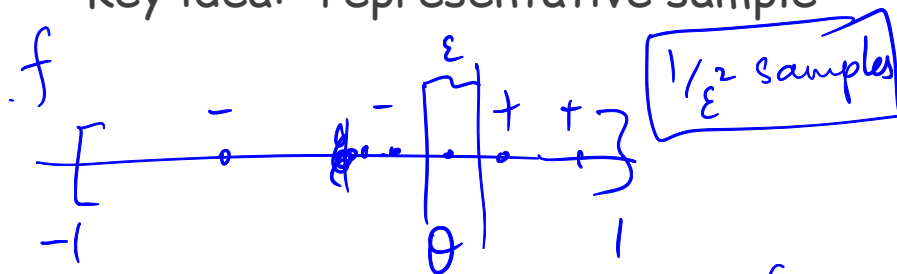  - output (one such) $h$.

  - get m examples
  - find $h \in \mathcal{H}$ that minimizes "empirical risk"
  - output $h$.

- Key idea: "representative sample"

$1/\varepsilon^2$ samples
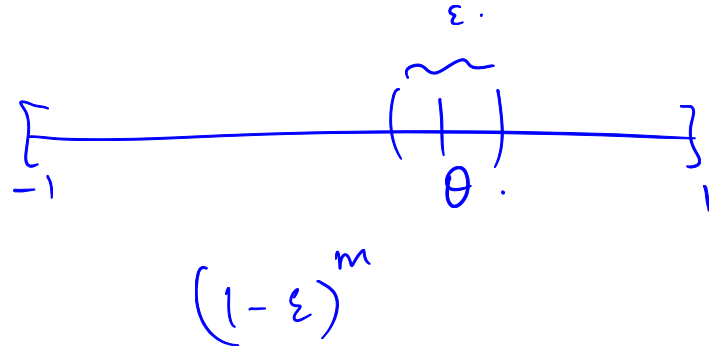
"empirical risk"
of $h$ & a sample $= \frac{1}{|S|} \cdot \sum_{x \in S} \mathbb{1}_{h(x) \neq f(x)}$
$S$
(minimizing training error)

# When is ERM bad?

- if there are too few examples ($ie$, $m$ is too small)

- just got unlucky with the examples?

$$-1 \qquad [\overbrace{(\ |\ )}^{\varepsilon} \underset{\theta}{} ]\qquad 1$$

$$(1-\varepsilon)^m$$

# REPRESENTATIVE SAMPLE

- Let H be a hypothesis class and X be an input space with a distribution D on it, and let f be a target function. Sample $S \subseteq X$ is said to be $\epsilon -$"representative" if **for all** h in H, we have:

$$\left| \frac{1}{|S|} \text{error}(S, h) - \text{risk}_D(h, f) \right| < \epsilon$$

$$\Pr_{x \sim D}\left[ h(x) \neq f(x) \right].$$

empirical
risk.
(risk on examples.

$$\frac{1}{|S|} \cdot \sum_{x \in S} \mathbb{1}[h(x) \neq f(x)]$$

true risk
w.r.t. $D$.

# RANDOM SAMPLE IS REPRESENTATIVE WHP!

- **Chernoff bound (Hoeffding).** Suppose $X_1, X_2, \ldots X_n$ are n iid samples from a distribution with mean $\mu$ and <u>support</u> [a, b]. Then we have

$$\Pr\left[\left|\frac{1}{n}(X_1 + \cdots + X_n) - \mu\right| > \epsilon\right] \leq 2\exp\left(-\frac{\epsilon^2 n}{(a-b)^2}\right)$$

# WHAT ABOUT INFINITE CLASSES?

- Note:  as long as sample is representative, we are good!

- What if we can divide hypotheses into finitely many "classes"?

- Example of threshold functions on a line

# GROWTH FUNCTION OF A CLASS

- For a class H and an input space X, we can define a notion of "growth function"