# THEORY OF MACHINE LEARNING

# LECTURE 3

PAC MODEL, GENERALIZATION

Homework 1:

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

**Naïve:** $h(x) = \ell(x) \;\; \forall x$

- **Formal definition of learning**

  - Formalizing generalization via "distributional assumption"

  - X: space of (all possible) inputs

    $- \; D:$ prob. distribution over $X$

  - Y: set of labels / outputs

    (Classification)

  - "Ground truth label" (**concept**).    $\ell: X \mapsto Y$ : function mapping inputs to outputs

- Goal of learning

  - "Learn" a hypothesis $h$ such that $h(x) = \ell(x)$ for all "inputs of interest"

  - Unknown probability distribution D over X; achieve small "risk" or "generalization error"

  - (Definition of risk):   $\Pr_{\{x \sim D\}} [\, h(x) \neq \ell(x) \,]$

Most common formal model to reason about learning

input sample $S = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$.

(unknown to learner)

- **Learnability (from examples).** [Suppose D is fixed.] We say that a concept class 'ℓ' 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

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

learning algo. $A : (X \times Y)^m \longrightarrow \mathcal{H} \rightsquigarrow$ possible output hypotheses.

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

- **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

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

# TODAY'S PLAN

$H: \{$ collection of functions $\}$.
$\qquad h: X \to Y$.

$\ell$
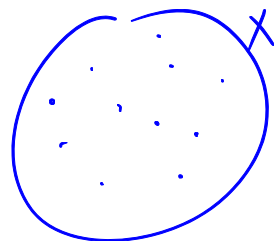
known hyp. class.

- Concept class (or class of hypothesis)

  - Assume that ground-truth label is (at least close to) a function in $H$

- "No free lunch theorem" (informal).  There is no "universal" (concept class agnostic) learning algorithm

- (Agnostic) PAC learning

- Finite classes are PAC learnable

# COMMON ML ASSUMPTIONS

*simple structure on class of concepts/hypothesis* *as long as you can choose the right set of features, ~~you can~~ you can predict label via linear separators.*

- (90s) Data is (approx.) linearly separable

- (these days) There exists 100-layer NN with width < ... that achieves low error on task

$$\begin{bmatrix} \ \end{bmatrix}$$

*( more complex ~~s~~ but still known — structure on hypothesis class.)*

- "Inductive bias" – assuming specific structure on concept

- What class of models do we use? *( for a given task.)* .

- Maybe.. we don't need to start with knowing a concept class
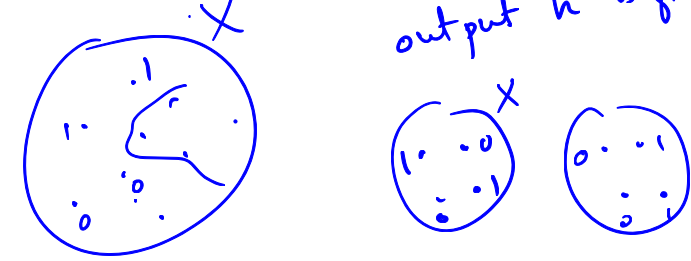
# NO FREE LUNCH THEOREM

- <u>Informal:</u> there is no "universal" learner, <u>even if it's allowed to be inefficient</u> (even for binary classification)

- **Theorem.** Let D be the uniform distribution on some input space X. Consider any (possibly randomized) algorithm A that uses < $|X|/2$ i.i.d. examples and produces $h : X \to \{0,1\}$. There exists a hypothesis $h$ for which A incurs risk > 1/10, with probability > 1/10.

  *0.9|x|*

  *|x|/16*

  *0.05*

- *(Recall def of "learnable" – fails with $\epsilon = 1/10$ and failure prob. 0.1 )*

  *$\epsilon \sim 1/20$*

  *X*

  *(∴ h produced by the theorem is not learnable--).*
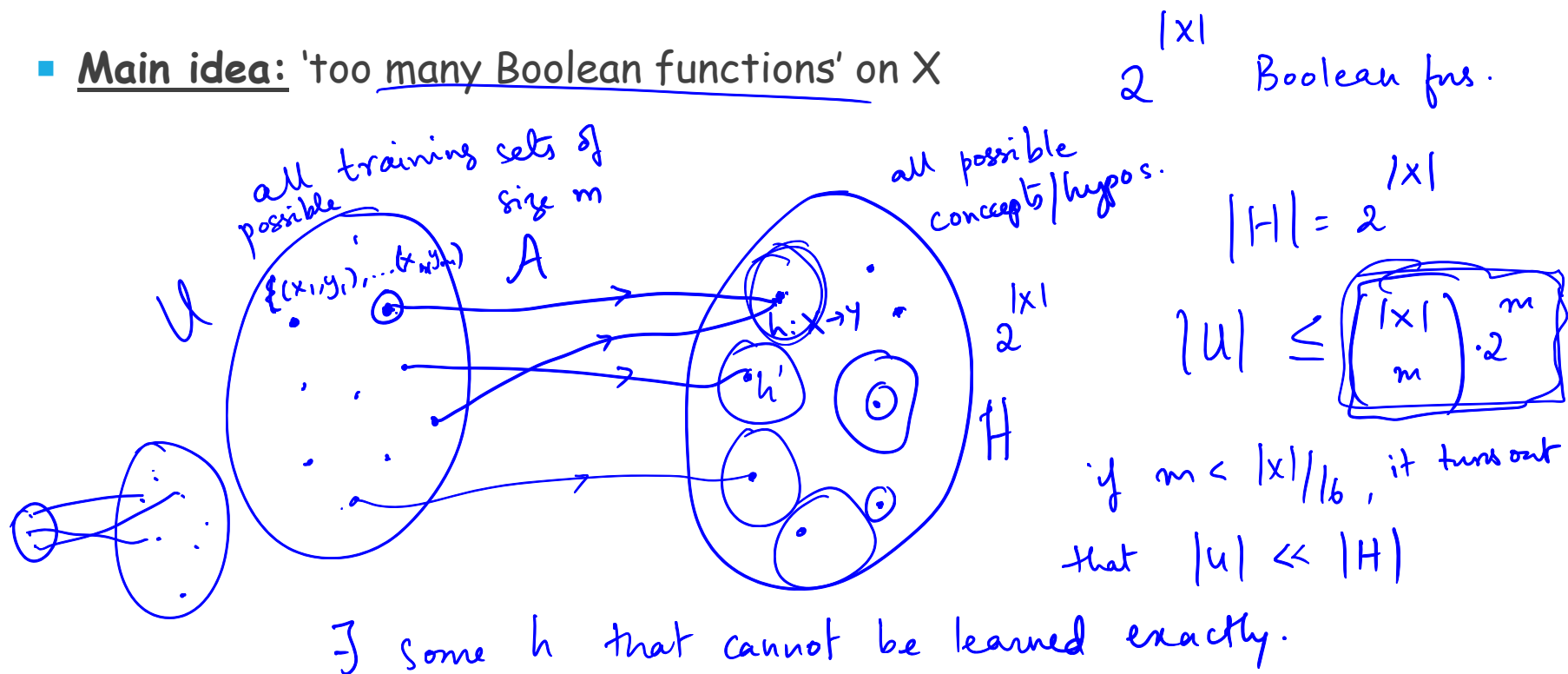
  *|x|/16*

# INFORMAL PROOF

y input to A { ie, $\{(x_1,y_1), (x_2,y_2), \ldots, (x_m,y_m)\}$ is given, output h is fixed.

- <u>Extra assumption:</u>  suppose A is deterministic; will show theorem with weaker constants

- **Main idea**: 'too <u>many Boolean functions</u>' on X

$2^{|X|}$ Boolean fns.

all training sets of size m

all possible training sets $\{(x_1,y_1), \ldots, (x_m,y_m)\}$   A

all possible concepts/hypos.

$h: X \to Y$

$2^{|X|}$

$H$

$|H| = 2^{|X|}$

$|U| \leq \binom{|X|}{m} \cdot 2^m$

if $m < |X|/16$, it turns out that $|U| \ll |H|$

$\exists$ some h that cannot be learned exactly.

# PROOF

$\rightarrow \left[ \text{Gilbert - Vershamov. bound} \right)$.

– But ... the goal is to obtain an $h'$ that agrees with $h$ on 90% of inputs.. ($\Rightarrow$ risk is $< 0.1$)

$\cdot h'$

$\cdot h$

$h$ is a $|X|$-bit string.

(exercise) how many $|X|$-bit strings differ from $h$ in $\leq |X|/10$ places?

$$\binom{|X|}{1} + \binom{|X|}{2} + \dots \rightarrow \binom{|X|}{|X|/10} \sim \binom{|X|}{|X|/10} \ll 2^{|X|}$$

Can verify: $\binom{|X|}{m} \cdot 2^m \cdot \binom{|X|}{|X|/10} < 2^{|X|}$

$\Rightarrow \exists h$ that differs in $\geq \frac{|X|}{10}$ places from ALL possible outputs of algo..

$\Rightarrow \exists h$ s.t. risk $> 0.1$ (with prob 1.)

## PAC LEARNING

- <u>Moral:</u> must suppose H is a *known* class of hypotheses (concept class)

- **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, takes $m(\epsilon, \delta)$ samples and produces $h$ with risk at most $\epsilon$ with probability at least $1 - \delta$.

- (The sample size must not depend on D)

# EVERY FINITE CLASS IS PAC LEARNABLE