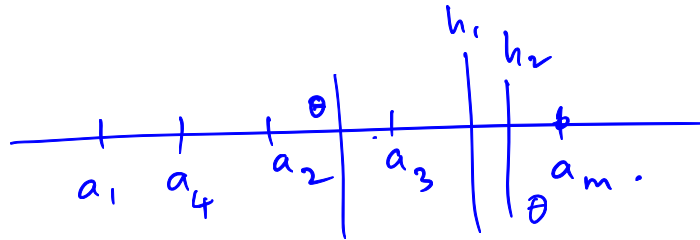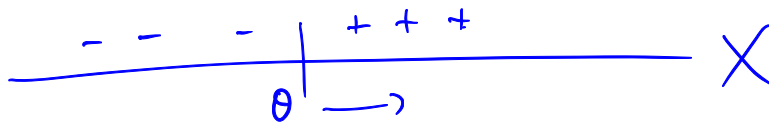# THEORY OF MACHINE LEARNING

# LECTURE 7
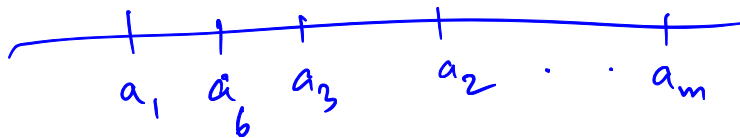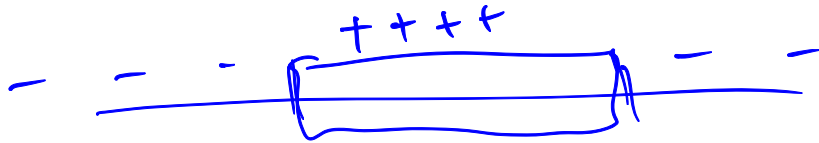
VC DIMENSION, FUNDAMENTAL THEOREM (CONTD.)

# LAST WEEK

PAC: Probably Approximately Correct.

- <u>Representative sample:</u> for a hyp class H and distribution D over X, S is called "representative" if

  for all $h \in H$, | (avg error on S)(h) – risk_D (h) | $\leq \epsilon$

- **Observation.** If training data happened to be a representative sample, ERM is an agnostic learning algorithm

  If S is representative, h automatically generalizes.

- **Observation 2.** For a finite hypothesis class, a *random* sample of size ~ log |H| is representative ( for any distribution D over a potentially infinite X.)

  w.h.p. ($\geq 1-\delta$)

  $$\frac{\log(|H|/\delta)}{\epsilon^2}$$

- Chernoff + Union bound

# WHAT ABOUT INFINITE CLASSES?

- <u>Note</u>:  if sample is representative, we are good!

  *~ How do you prove that a "small" sample is $\varepsilon$-rep. when $H$ is infinite?*

- "Growth function" of a class: total number of distinct ways in with H can label a set of m points (# distinct "sign patterns")

- H is the class of l.t.f. in 1D:  (m+1)

- H is the class of intervals in 1D: O(m^2)

  *$|S| = m$*

- H is the class of axis-parallel rectangles? messy, but O(m^4)

- H is the class of "convex polygons" in 2D:   2^m (exponential)

$$- - - - | + + +$$

$$\theta \quad \longrightarrow$$

X

$T_H(m)$: max # of distinct classifications for a __set S of__ size $m$

$(m+1)$

$$h_1 \quad h_2$$

$$a_1 \quad a_4 \quad \theta \quad a_2 \quad a_3 \quad | \quad | \quad a_m \quad . \quad \theta$$

$$- - - - \boxed{+ + + +} \quad - - -$$

$\sim O(m^2)$ distinct classifications.

$$a_1 \quad a_6 \quad a_3 \quad a_2 \quad \cdot \quad \cdot \quad a_m$$

Convex polygons in 2D:

$$T_H(m) = 2^m.$$

$$-$$
$$- \quad + \quad$$
$$- \quad -$$

$$- \quad b \quad b \quad + \quad + \quad + \quad c$$

$$+ \quad \cdot \quad \cdot$$
$$- \quad \cdot \quad \cdot \quad + \quad$$
$$+ \quad + \quad$$
$$\cdot \quad -$$

given any potential
classification, we can
obtain it via a ~~convex~~
hypothesis in ~~our~~ the class.
of convex polygons.

$$\Gamma_H(m) = 2^m.$$

Very rough statement: if $H$ has the property that
any $h \in H$ is determined by $\sim d$ "parameters",
then the growth fn $\sim n^d$.

# TODAY

- "Small" growth function => hypothesis class is learnable!

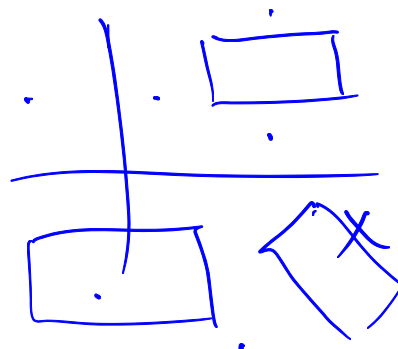  ( a "small" sample is $\varepsilon$-representative whp.).

- How to bound growth function? VC dimension & Sauer's lemma
  (attributed to Sauer and Shelah, VC)

  if VC dim $= d$, growth function $\leq O(n^d)$

  — defn of $\varepsilon$-rep was "deterministic".

  $\mathcal{H}$: axis parallel rectangles

# LEARNABILITY IN TERMS OF THE GROWTH FUNCTION

- Theorem:  Suppose $\tau_H(m)$ be the growth function of a hypothesis class H. Then for any X, D, if we take a sample S of size $m$, with prob. 1-$\delta$,

$$\forall h \in H:$$

$$\text{risk}_D(h)$$

$$\sup_{h \in H} |err(h, S) - err(h, D)| \leq \frac{4 + \sqrt{\log \tau_H(2m)}}{\delta \sqrt{2m}}$$

- In other words, if 'm' is chosen so that RHS < $\epsilon$, theorem implies that random sample of size m is $\epsilon$-representative

Recall: defn of $\epsilon$-rep: had an $\epsilon$ on the RHS. So if we choose

$$m \quad \text{s.t.} \quad \frac{4 + \sqrt{\log \tau_H (2m)}}{\delta \sqrt{2m}} \leq \epsilon$$

then Theorem implies that a random

Sample of size m is $\epsilon$-rep; with prob. $\geq 1-\delta$.

# EXAMPLES

Want:

$$\frac{4 + \sqrt{\log(2m+1)}}{8 \cdot \sqrt{2m}} \le \varepsilon$$

$$\frac{\log \frac{1}{\delta} + 1}{\varepsilon^2}$$

Suffices to set $m = \sqrt{\frac{4}{(\varepsilon\delta)^2} \cdot \log\left(\frac{1}{\varepsilon\delta}\right)}$

- H is the class of l.t.f. in 1D: $(m+1)$

$$m \sim \frac{O(1)}{\varepsilon^2} \cdot \log\left(\frac{1}{\varepsilon}\right)$$

- H is the class of intervals in 1D: $O(m^2)$

$$\frac{1}{3} \qquad \frac{2}{3}$$

- H is the class of axis-parallel rectangles? messy, but $O(m^4)$ → same bound, bigger constant.

- H is the class of "convex polygons" in 2D: $2^m$ (exponential)

$$\sqrt{\log(2m)^4} = \sqrt{4\log(2m)}$$

If growth fn $\le m^{\textcircled{d}}$,

**Obs 1:** If $\tau_H(m) \leq O(m^d)$, then to make

$$\frac{4 + \sqrt{\log(\tau_H(2m))}}{\delta \cdot \sqrt{2m}} \leq \varepsilon,$$ we just need to set

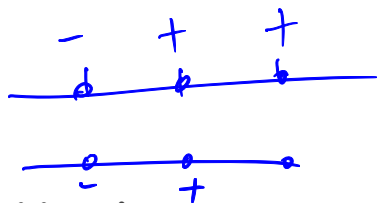$$m \geq \left(\frac{82 \cdot d}{(\varepsilon \delta)^2}\right) \log\left(\frac{d}{\varepsilon \delta}\right)$$

**Obs 2:** If $\tau_H(m) = (1.5)^m$, then the condition $2^{\sqrt{m}}$

$$\frac{4 + \sqrt{\log(\tau_H(2m))}}{\delta \sqrt{2m}} \leq \varepsilon \implies \frac{\boxed{4} + \sqrt{(2m)\log(1.5)}}{\delta \cdot \sqrt{2m}} \leq \varepsilon.$$

$$LHS \simeq \frac{1}{\delta}.$$

If growth fn is exp., then a sample cannot be guaranteed to be representative whp.
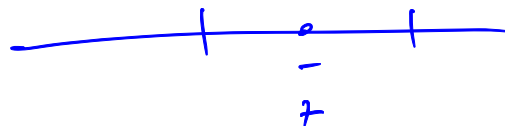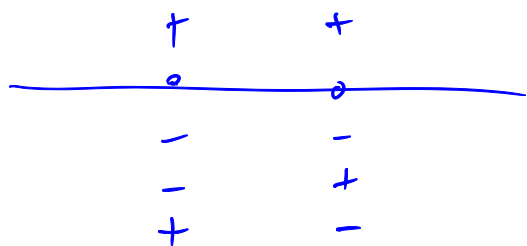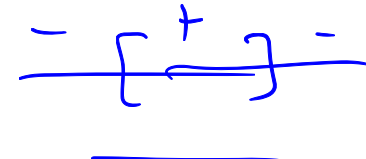
- **Shattering & VC dimension.**

Shattering: given a hyp. class $\mathcal{H}$ and a set $S$ of points, we say that $\mathcal{H}$ shatters $S$ if every possible classification on $S$ can be achieved by some $h \in \mathcal{H}$.

VC-dimension: VC-dimension of a class $\mathcal{H}$ of functions over $X$ is defined as the size of the largest set $S \subseteq X$ that is shattered by $\mathcal{H}$.

LTF:

What about $\mathcal{H} := \left\{ \begin{array}{l} \text{intervals on the line} \\ \text{(where label is + inside} \end{array} \right\}$. $\underline{\phantom{--}}\,\underline{\;[\;\;+\;\;]\;}\,\underline{\phantom{--}}$

in order to show VC-dim($\mathcal{H}$)=d  & - outside.

✓ +) ~~vc-dim~~ ∃ a set of size d ~~that~~ can be shattered → d=2

✓ 2)    No set of size (d+1) can be shattered! →

+,        ·   · -                    $\underline{\phantom{----}\,\overset{+}{\phantom{-}}\,\underline{[\;+\;]}\phantom{----}}$
                                              -        +

:          :                        +     -           +
-        ·                           |     |           |
         :

No set of 3 points can be shattered...

## SAUER-SHELAH LEMMA (VAPNIK-CHERVONENKIS)

- **Lemma.** Let H be a hypothesis class of finite VC dimension d. Then for every *m*, we have:

$$\tau_H(m) \leq \binom{m}{0} + \binom{m}{1} + \cdots + \binom{m}{d}$$

- Much better than exponential, for m large

- Proof by a clever inductive argument

# FUNDAMENTAL THEOREM OF (STAT) LEARNING THEORY

- <u>Theorem:</u>  The following statements are *equivalent:*

  - Class H is PAC learnable

  - Class H is *agnostically* PAC learnable

  - Class H has finite VC dimension


- Implies that if H has infinite VC dimension, it is **<u>not</u>** PAC learnable!  (same proof as no-free-lunch theorem)

## LEARNABILITY IN TERMS OF THE GROWTH FUNCTION

- <u>Theorem:</u> Suppose $\tau_H(m)$ be the growth function of a hypothesis class H. Then for any X, D, if we take a sample S of size $m$, with prob. 1-$\delta$,

$$\sup_{h \in H} |err(h, S) - err(h, D)| \leq \frac{4 + \sqrt{\log \tau_H(2m)}}{\delta \sqrt{2m}}$$

- In other words, if 'm' is chosen so that RHS < $\epsilon$, theorem implies that random sample of size m is $\epsilon$-representative

# OUTLINE -- THE TWO SAMPLE TRICK

- Want to show that a random sample is $\epsilon$-representative

- Take sample S, define event:

  A = Pr [ sample is not representative ]

- Way to "test" if S is not representative?

  - "Cross validation"

- Define new event S, S'

- "Swapping"