



CMU 15-251

LEARNING THEORY

TEACHERS:

VICTOR ADAMCHIK

ARIEL PROCACCIA (THIS TIME)

MIDTERM 2

- Material: from graphs to online algorithms
- Format: similar to midterm 1
 - Short questions
 - One HW question
 - Two long questions
- We'll post a practice exam tonight or tomorrow which will be solved in recitation



THE PAC MODEL

- **PAC** = probably approximately correct
- Introduced by Valiant [1984]
- Learner can do well on training set but badly on new samples
- Establish guarantees on accuracy of learner when **generalizing** from examples



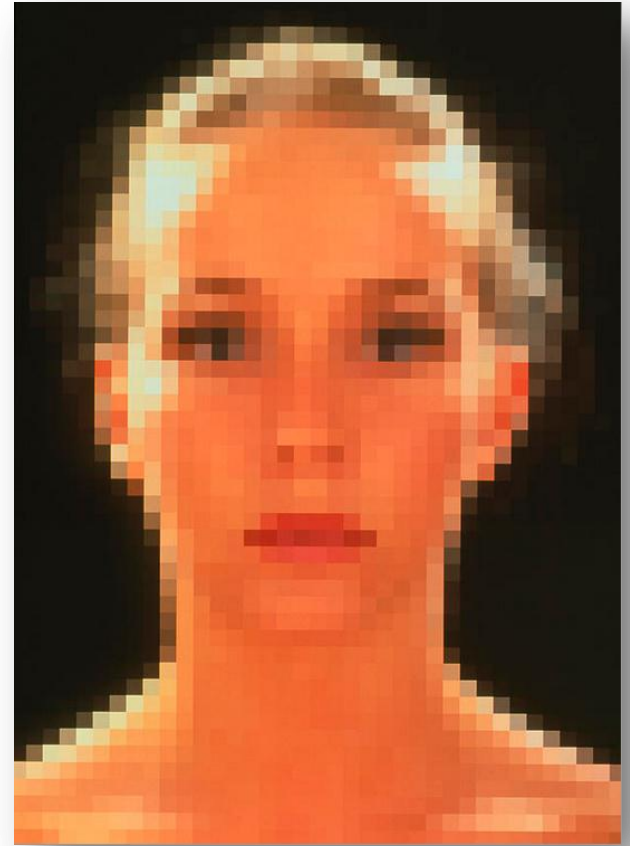
THE PAC MODEL

- Input space X
- D distribution over X : unknown but fixed
- Learner receives a set S of m instances x_1, \dots, x_m , sampled according to D
- Concept class C of functions $h: X \rightarrow \{+, -\}$
- Assume target function $c_t \in C$
- Training examples $Z = \{(x_i, c_t(x_i))\}$



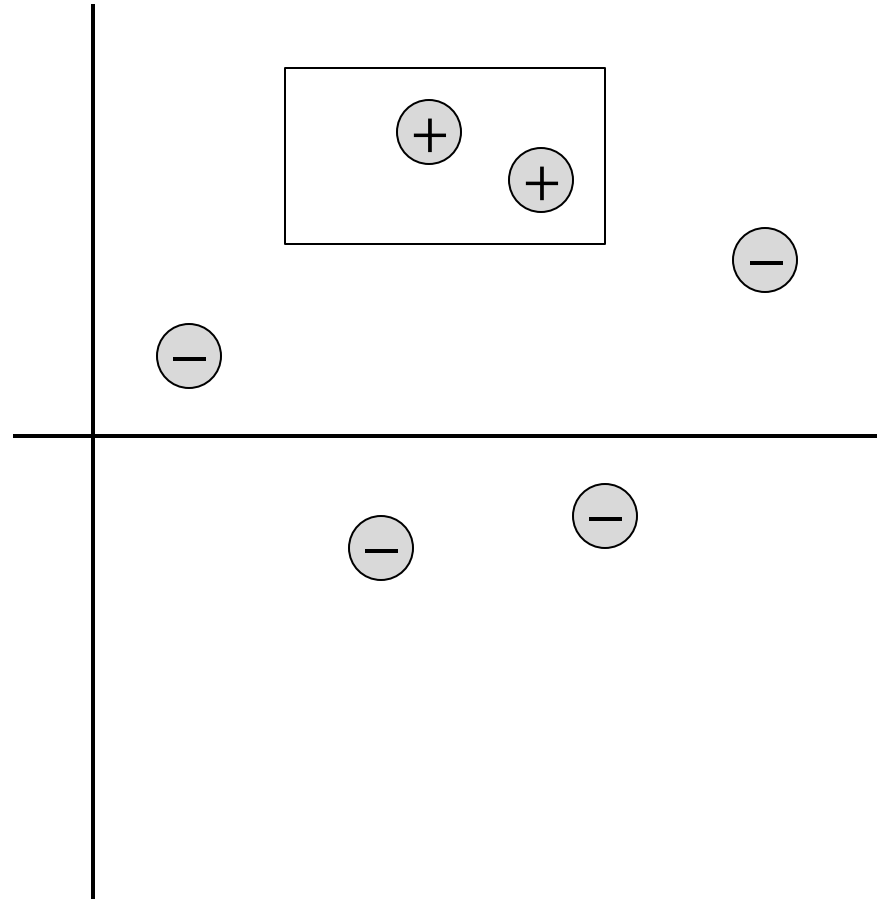
EXAMPLE: FACES

- $X = \mathbb{R}^n$
- Each $x \in X$ is a vector of colors, one per pixel
- $c_t(x) = +$ iff x is a picture of a face
- Training examples: Each is a picture labeled “face” or “not face”



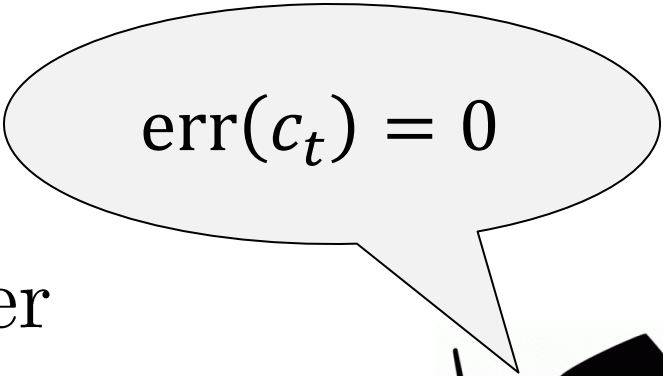
EXAMPLE: RECTANGLE LEARNING

- $X = \mathbb{R}^2$
- $\mathcal{C} =$ axes-aligned rectangles
- $h(x) = +$ iff x is contained in h



THE PAC MODEL

- The **error** of concept h is
$$\text{err}(h) = \Pr_{x \sim D} [x: c_t(x) \neq h(x)]$$
- Given **accuracy** parameter $\epsilon > 0$, would like to find concept h with $\text{err}(h) \leq \epsilon$
- Given **confidence** parameter $\delta > 0$, would like to achieve $\Pr[\text{err}(h) \leq \epsilon] \geq 1 - \delta$


$$\text{err}(c_t) = 0$$



THE PAC MODEL

- A **learning algorithm** L is a function from training examples to \mathcal{C} such that: for every $\epsilon, \delta > 0$ there exists $m_0(\epsilon, \delta)$ such that for every $m \geq m_0$ and every D , if m examples Z are drawn from D and $L(Z) = h$ then

$$\Pr[\text{err}(h) \geq \epsilon] \leq 1 - \delta$$

- \mathcal{C} is **learnable** if there is a learning algorithm for \mathcal{C}

$m_0(\epsilon, \delta)$ is independent
of D !



RECTANGLES ARE LEARNABLE

- Learning algorithm: given training set, return tightest fit for positive examples
- **Theorem:** axes-aligned rectangles are learnable with $m_0(\epsilon, \delta) \geq \frac{4}{\epsilon} \ln \frac{4}{\delta}$



RECTANGLES ARE LEARNABLE*

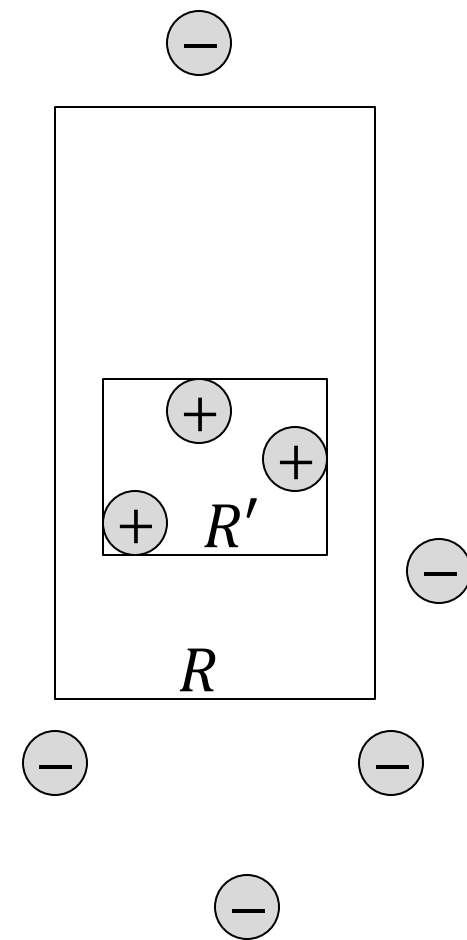
- **Proof:**

- Target rectangle R
- Recall: our learning algorithm returns the **tightest-fitting** R' around the positive examples

- For region E , let

$$w(E) = \Pr_{x \sim D} [x \in E]$$

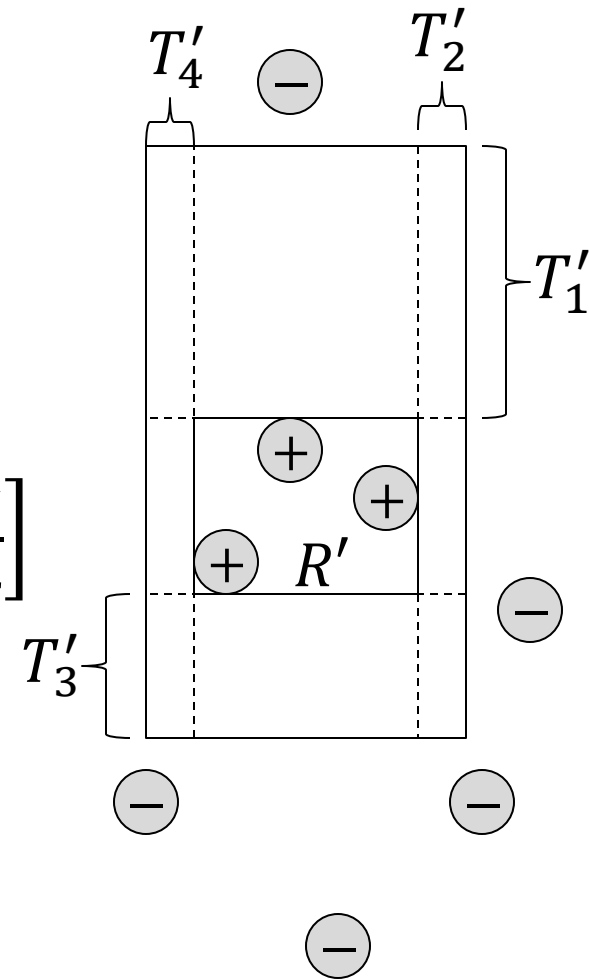
- $\text{err}(R') = w(R \setminus R')$ (**why?**)



*Just for fun

RECTANGLES ARE LEARNABLE*

- **Proof (cont.):**
 - Divide $R \setminus R'$ into four strips T'_1, T'_2, T'_3, T'_4
 - $\text{err}(R') \leq \sum_{i=1}^4 w(T'_i)$
 - We will estimate $\Pr \left[w(T'_i) \geq \frac{\epsilon}{4} \right]$

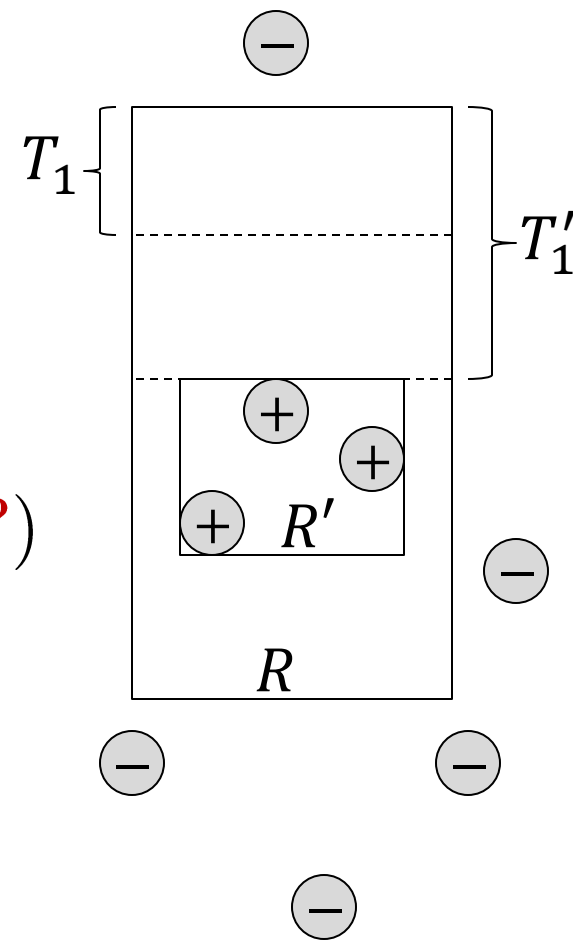


*Just for fun

RECTANGLES ARE LEARNABLE*

- Proof (cont.):

- Focusing wlog on T'_1 , define a strip T_1 such that $w(T_1) = \frac{\epsilon}{4}$
- $w(T'_1) \geq \frac{\epsilon}{4} \Leftrightarrow T_1 \subseteq T'_1$
- $T_1 \subseteq T'_1 \Leftrightarrow x_1, \dots, x_m \notin T_1$ (why?)
- $w(T'_1) \geq \frac{\epsilon}{4} \Leftrightarrow x_1, \dots, x_m \notin T_1$
- $\Pr[x_1, \dots, x_m \notin T_1] = \left(1 - \frac{\epsilon}{4}\right)^m$



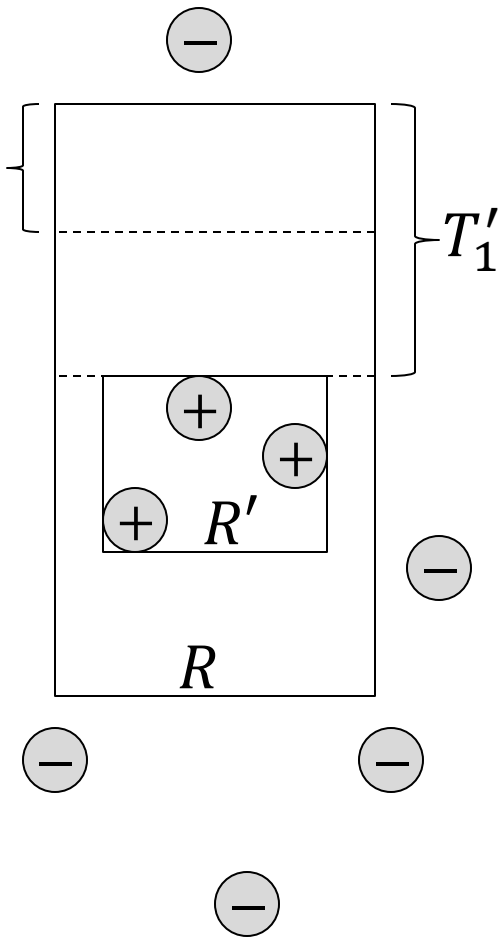
*Just for fun

RECTANGLES ARE LEARNABLE*

- Proof (cont.):

- $\Pr[w(R \setminus R') \geq \epsilon] \leq 4 \left(1 - \frac{\epsilon}{4}\right)^m T_1$

- So we want $4 \left(1 - \frac{\epsilon}{4}\right)^m \leq \delta$, and with a bit of algebra we get the desired bound ■



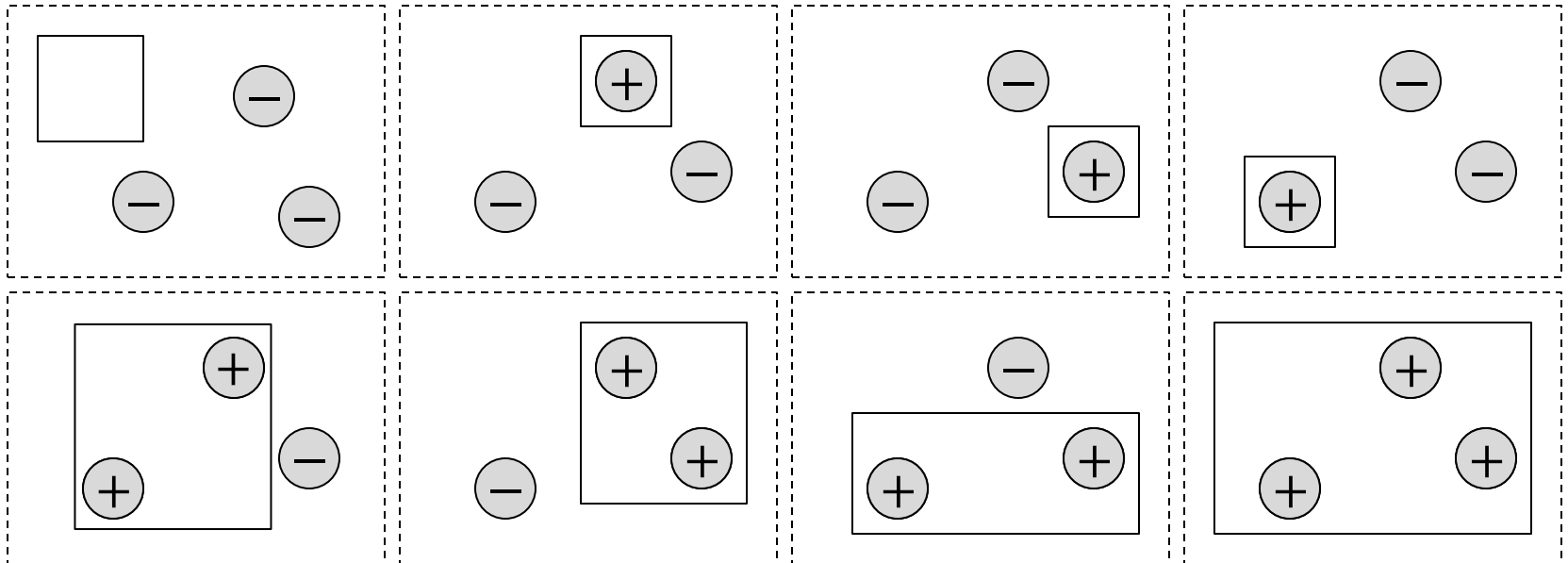
*Just for fun

VC DIMENSION

- We would like to obtain a more general result
- Let $S = \{x_1, \dots, x_m\}$
- $\Pi_C(S) = \{(h(x_1), \dots, h(x_m)) \mid h \in C\}$



VC DIMENSION




$$\Pi_C(S) = \{(-, -, -), (-, -, +), (-, +, -), (-, +, +), \\ (+, -, -), (+, -, +), (+, +, -), (+, +, +)\}$$



VC DIMENSION

- $X =$ real line
- $\mathcal{C} =$ intervals; points inside interval are labeled by $+$, outside by $-$

• **Note:** what is $|\Pi_{\mathcal{C}}(S)|$ for $S =$ 

1. 1

2. 2

3. 3

4. 4



VC DIMENSION

• **Note:** what is $|\Pi_C(S)|$ for $S =$ 

1. 5

2. 6

3. 7

4. 8



VC DIMENSION

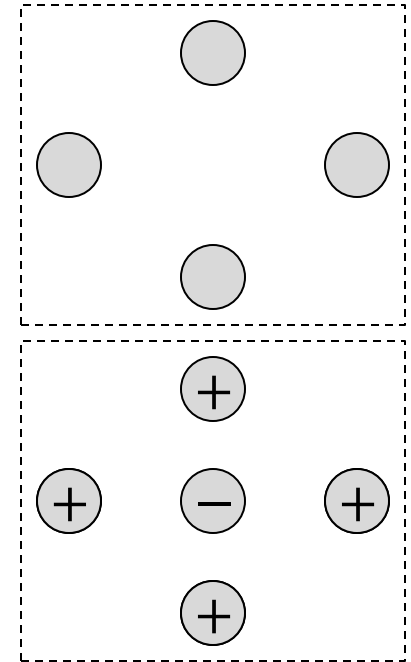
- S is **shattered** by \mathcal{C} if $|\Pi_{\mathcal{C}}(S)| = 2^m$
- The **VC dimension** of \mathcal{C} is the cardinality of the largest set that is shattered by \mathcal{C}

How do we prove
upper and lower
bounds?



EXAMPLE: RECTANGLES

- There is an example of four points that can be shattered
- For any choice of five points, one is “internal”
- A rectangle cannot label outer points by 1 and inner point by 0
- VC dimension is 4



VC DIMENSION

- **Note:** $X =$ real line, $\mathcal{C} =$ intervals, what is $\text{VC-dim}(\mathcal{C})$?

1. 1 3. 3

② 2 4. ∞

- **Note:** $X =$ real line, $\mathcal{C} =$ unions of intervals, what is $\text{VC-dim}(\mathcal{C})$?

1. 2 3. 4

2. 3 ④ ∞



SAMPLE COMPLEXITY

- **Theorem:** a concept class \mathcal{C} with $\text{VC-dim}(\mathcal{C}) = \infty$ is not PAC learnable
- **Theorem:** Let \mathcal{C} with $\text{VC-dim}(\mathcal{C}) = d$. Let L be an algorithm that produces an $h \in \mathcal{C}$ that is **consistent** with the given samples S . Then L is a learning algorithm for \mathcal{C} with $m_0 = c_0 \left(\frac{1}{\varepsilon} \log \frac{1}{\delta} + \frac{d}{\varepsilon} \log \frac{1}{\varepsilon} \right)$



WHAT WE HAVE LEARNED

- Definitions
 - PAC model
 - Error, accuracy, confidence
 - Learning algorithm
 - $\Pi_C(S)$, shattering
 - VC-dimension

