

---

---

---

---

---



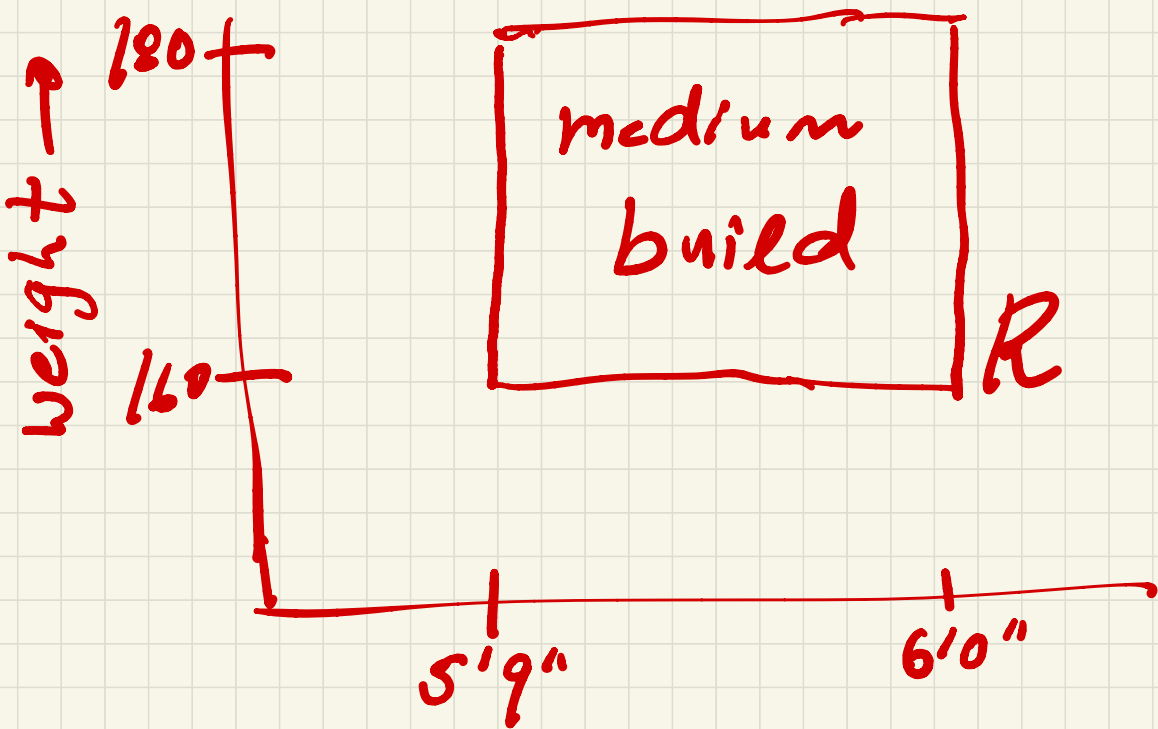
# A "Toy" ML Problem:

## Rectangles in $\mathbb{R}^2$

- Aliens arrive from outer space
- You'd like to teach them the concept of "medium build" for adult males (assume binary)
- You can label but not describe

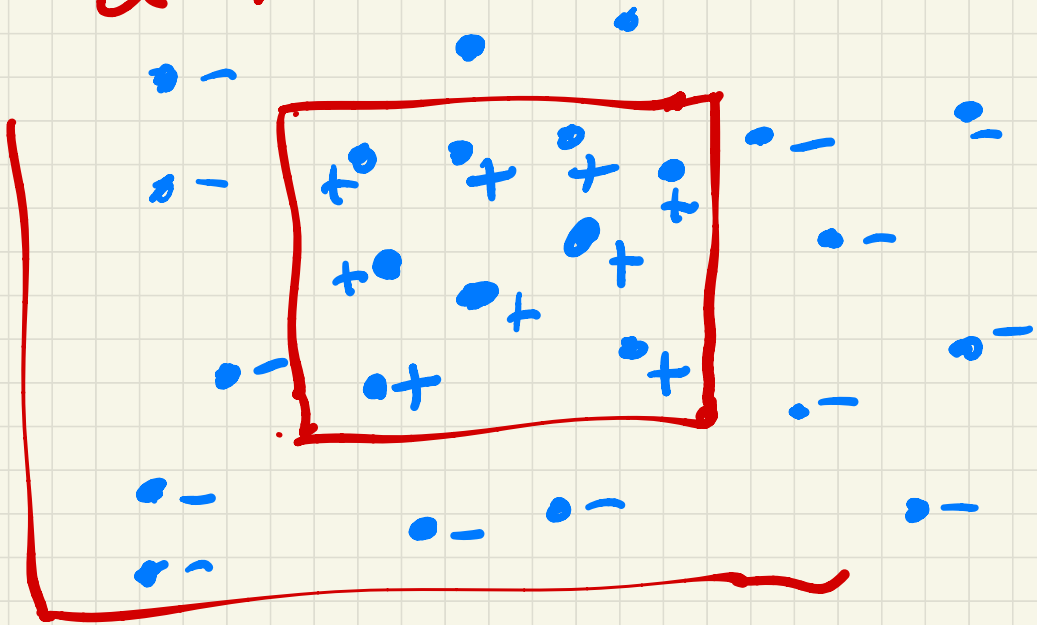
Let's formalize this:

- You (teacher): rectangle  $R$  in  $x$ - $y$  plane:



You generate

data for aliens:



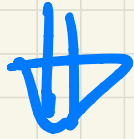
Assume: points drawn  
i.i.d from  $\mathcal{P}$  over  $\mathbb{R}^2$

Note: strong assumptions  
on  $\mathcal{R}$ , no assumptions on  $\mathcal{P}$

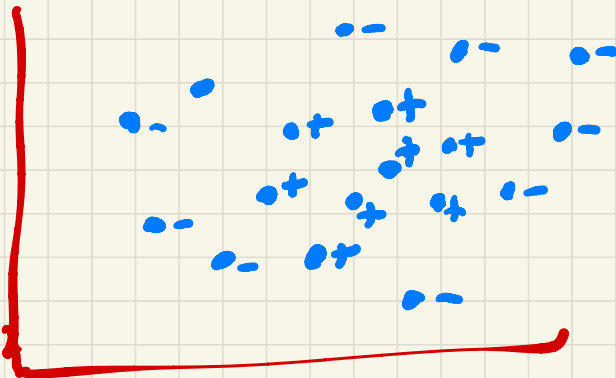


# Players:

- Input domain:  $\mathbb{R}^2$
- Model class: rectangles (binary functions)
- "Target" rectangle  $R$
- Input distribution  $P$



Data:

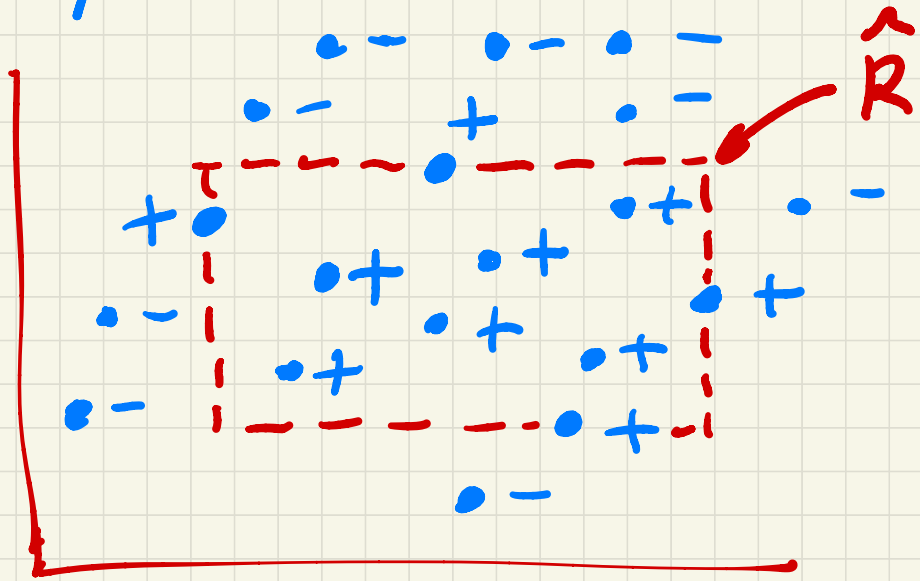


Alien (learner) goal:  
from data, learn  
a "good" hypothesis  
rectangle  $\hat{R}$ .

- What should "good" mean?
- What algorithm should  
alien use?

# A Proposed Algo:

$\hat{R}$  = "tightest fit" to positive (+) examples.



- Note: exploiting our assumption that "ground truth" ( $R$ ) is a rectangle!

What can we say about  $\hat{R}$ ?

Claim: Viewed as sets,  
 $\hat{R} \subseteq R$ .

What about something  
stronger/more interesting?

Remember points are  
drawn i.i.d. from  $P$ .

Let's define the error  
of  $\hat{R}$  w.r.t.  $R \notin P$ :

$$\varepsilon(\hat{R}) \triangleq P_{x \sim P} [\hat{R}(x) \neq R(x)]$$

(as functions)

$$= P[\hat{R} \Delta R]$$

(as sets)

Claim: With "high probability",  
 $\varepsilon(\hat{R})$  is "small" as long  
as sample is "large enough!"

# Analysis

Two inputs/parameters:

- small  $\delta > 0$ :  
"with high prob" =  
with prob  $\geq 1 - \delta$  w.r.t.  
draw of sufficiently  
large sample  $S$

- small  $\epsilon > 0$ :  
" $\epsilon(\hat{R})$  small" =  
 $\epsilon(\hat{R}) \leq \epsilon$

Goal: Show that if  $|S| = m$   
is large enough, then  
w.p.  $\geq 1 - \delta$ ,  $\epsilon(\hat{R}) \leq \epsilon$ .

Remark: Note that

$$E_S[\epsilon(\hat{R})] \leq (1-\delta)\epsilon + \delta \cdot 1$$

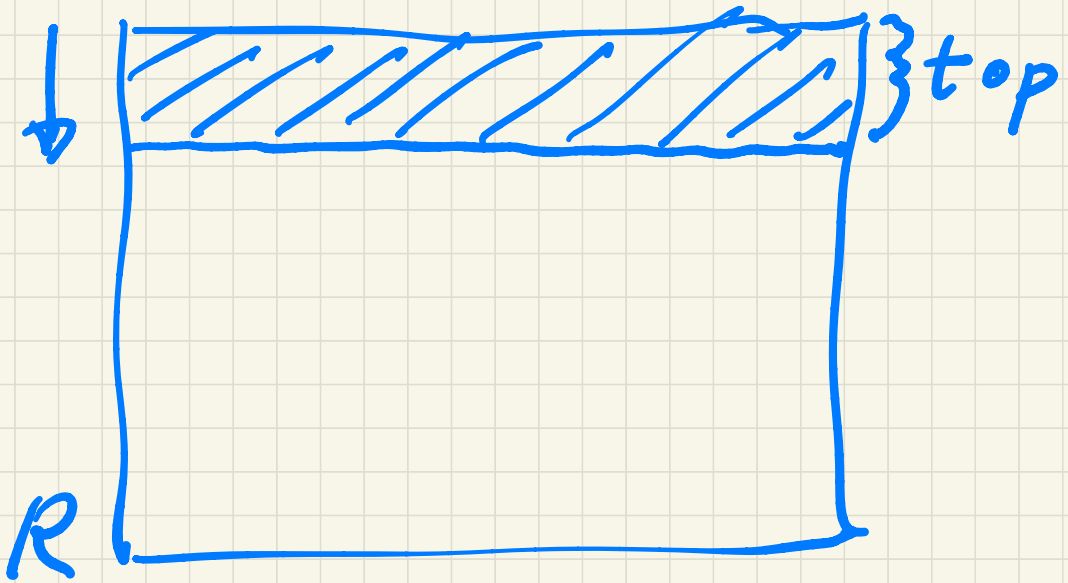
So why have both  $\epsilon$  &  $\delta$ ?

$\delta$ : Bounds prob. of a wildly unrepresentative sample  $S$

$\epsilon$ : Bounds error on representative samples

$\hat{R}$  is "probably ( $\geq 1-\delta$ )  
Approximately ( $\leq \epsilon$ )"  
Correct

Let's define 4 subsets  
of  $R$  (w.r.t.  $P$ ):

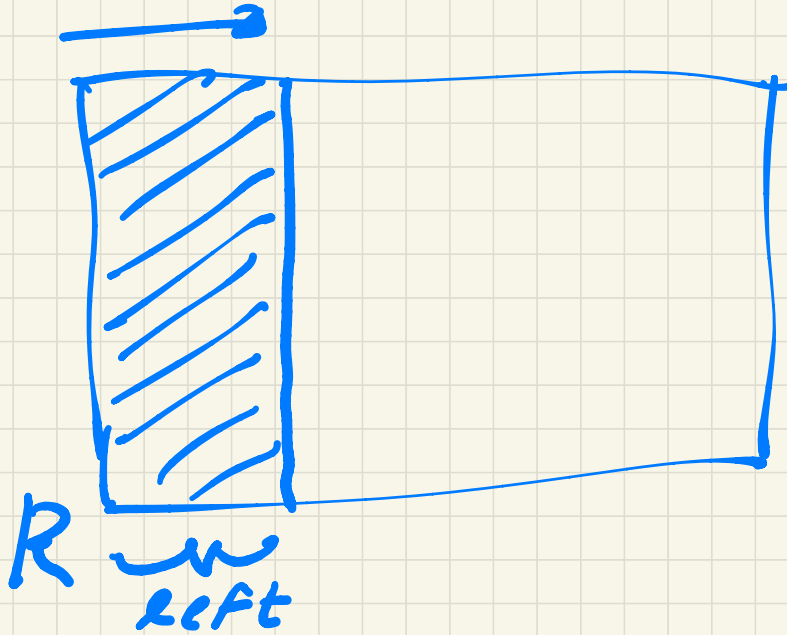


$$P_{x \sim P}[x \in \text{top}] = \epsilon/4$$

(Q: What if  $P[R] < \epsilon/4$ ?  
Assume not for now.)



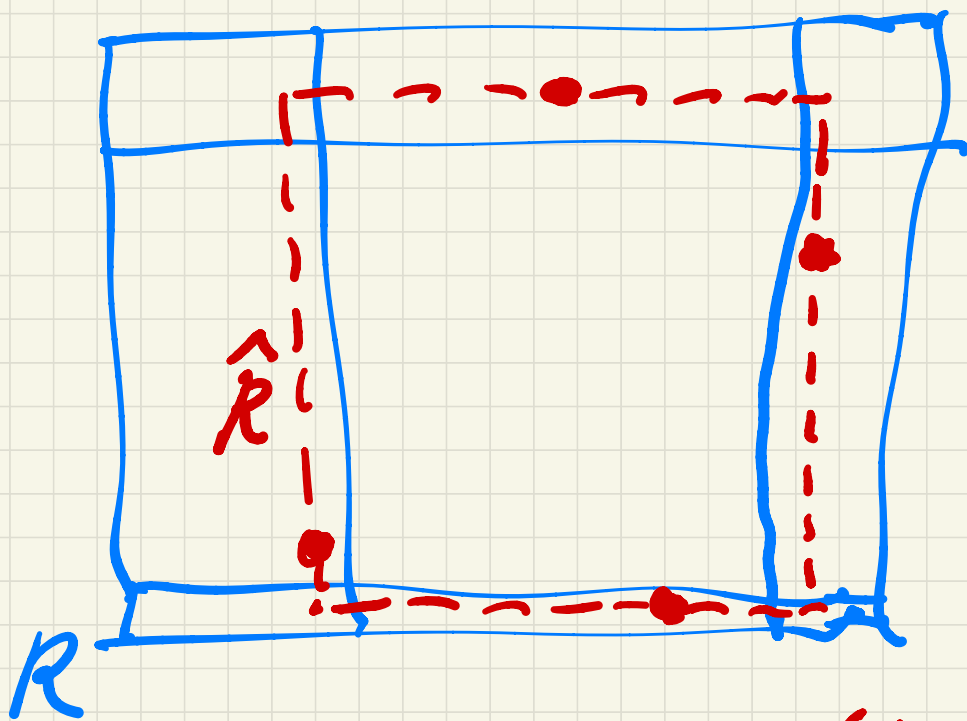
Similarly:



$$P_{r_{xup}}[x \in \text{left}] = \epsilon/4$$

Similarly for  
right, bottom.

If sample hits all  
of top, bottom, left, right:



$$\begin{aligned} \text{Then } \mathbb{E}(\hat{R}) &\leq \frac{\epsilon}{4} + \frac{\epsilon}{4} + \frac{\epsilon}{4} + \frac{\epsilon}{4} \\ &= \epsilon. \end{aligned}$$

So let's define a **bad**  
**sample**  $S$  as one s.t.  
 $S$  **misses** any of  $l, r, t, b$ .  
Goal: bound  $\Pr[S \text{ is bad}]$   
by  $\delta$ .

- Let  $n = |S| = \text{sample size}$
- Remember  $S$  i.i.d.  
wrt  $P$
- $\Pr[S \text{ misses top}]$   
 $= (1 - \epsilon/4)^n$  (indep.)
- Same for  $b, l, r$

$$\therefore \Pr[S \text{ misses } \underline{\text{any}} \text{ of } t, b, l, r]$$

$$\leq \Pr[S \text{ misses top}] + \Pr[S \text{ misses bottom}] + \Pr[S \text{ misses left}] + \Pr[S \text{ misses right}]$$

(union bound:

$$\Pr[A \text{ or } B] \leq \Pr[A] + \Pr[B])$$

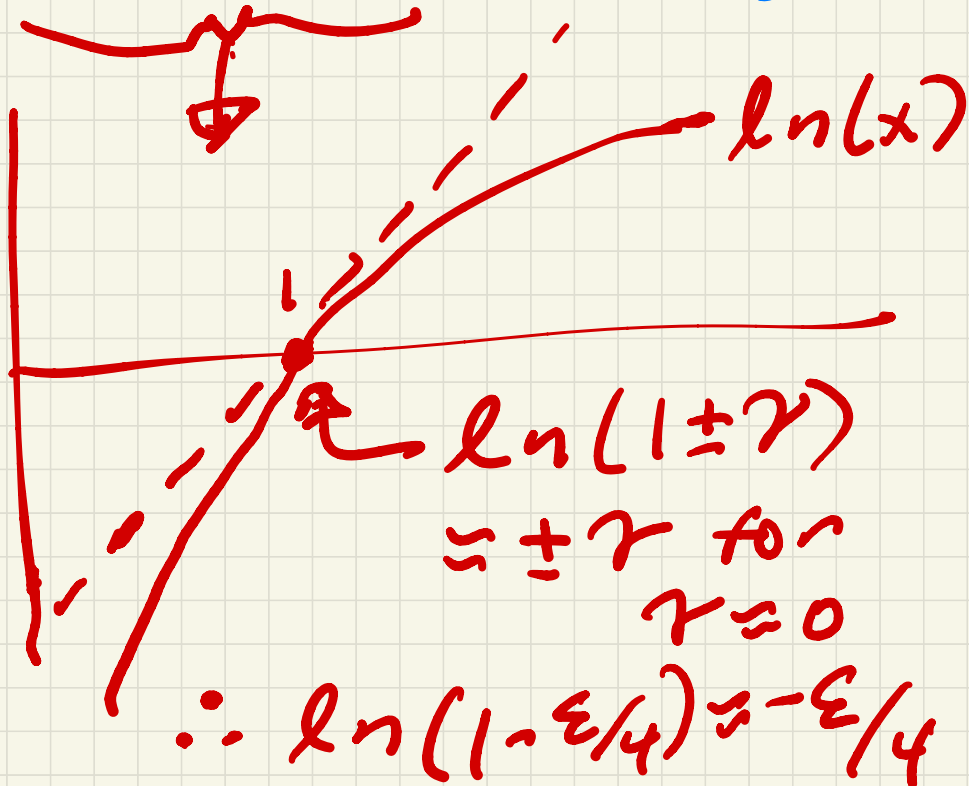
$$\leq 4 \cdot (1 - \epsilon/4)^m$$

• So  $\Pr[S \text{ bad}] \leq$

$4(1 - \epsilon/4)^m$ , set  $\leq \delta$ :

$$4(1 - \epsilon/4)^m \leq \delta$$

$$m \ln(1 - \epsilon/4) \leq \ln(\delta/4)$$



$$-m\varepsilon/4 \leq \ln(\delta/4)$$

$$m\varepsilon/4 \geq \ln(4/\delta)$$

$$m \geq \frac{4}{\varepsilon} \ln(4/\delta)$$

As long as  $\delta$  is this large, w.p.  $\geq 1-\delta$ ,

$$\varepsilon(\hat{R}) \leq \varepsilon.$$

Oh wait... what  
if e.g.  $P[R] < \epsilon/4$ ?

So have a **fast** algo  
with **small** sample  
complexity and a  
**rigorous** analysis.

# Proof overview:

- specify algo
- define "bad" events for algo
- bound prob. of each bad event
- take union bound
- set less than  $\delta$ , do algebra



# Extensions?

- Rectangles in  $\mathbb{R}^d$ ?
- Parallelograms in  $\mathbb{R}^2$ ?
- Circles? Triangles?
- Union of 2 rectangles?
- 

Next Up:

A General Model.