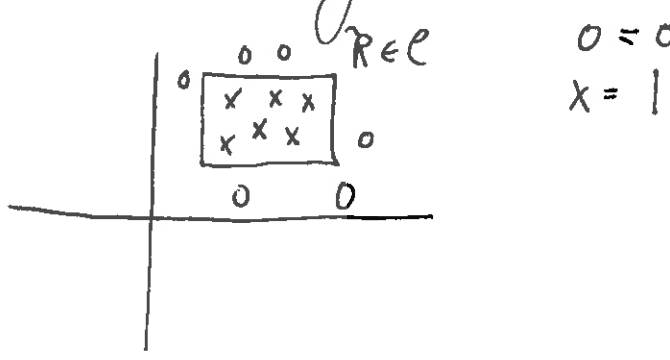


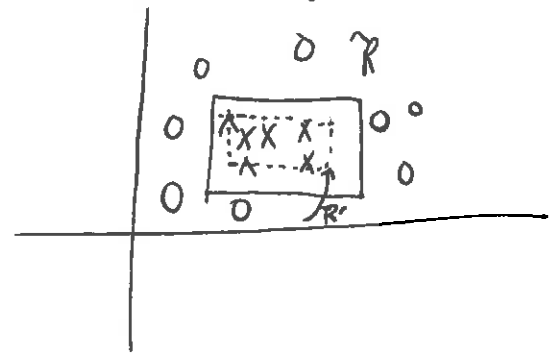
FoSML: Lecture 2

Remark: The concept class C is known to A , but $c \in C$ is not. The algorithm must work well (quantified by PAC definition) for all elements of C .

ex (Axis aligned rectangles). Let C be the concept class of all axis-aligned rectangles in \mathbb{R}^2 . Note that C is uncountably infinite, but is in some sense quite simple.



In order to show C is PAC-learnable, we must consider an algorithm that estimates these rectangles "well" (ϵ -close...) given a sample. Consider the algorithm that simply fits the tightest possible (area-minimizing) rectangle to the data sample (labeled S):



Intuitively $R \rightarrow \hat{R}$ as the number of samples increases. We want to quantify " $R \rightarrow \hat{R}$ " in terms of error of function estimation, and show that the convergence rate is sufficiently fast as a function of the number of samples n .

Let $R_S = \hat{R}$, i.e. the error region, which is a function of the sample S . To prove PAC-learnability, we must analyze $P_{S \sim D^n} [R(R_S) > \epsilon]$.

$$\Leftrightarrow \exp(-m\epsilon/4) \leq \delta/4$$

$$\Leftrightarrow -\frac{m\epsilon}{4} \leq \log(\delta/4)$$

$$\Leftrightarrow \frac{m\epsilon}{4} \geq \log(4/\delta)$$

$$\Leftrightarrow m \geq \frac{4}{\epsilon} \log\left(\frac{4}{\delta}\right).$$

So, we have shown that $m \geq \frac{4}{\epsilon} \log\left(\frac{4}{\delta}\right)$

$$\rightarrow \mathbb{P}(R(\mathcal{R}_S) \leq \epsilon)$$

$$= 1 - \mathbb{P}(R(\mathcal{R}_S) \geq \epsilon)$$

$$\geq 1 - \delta.$$

Note that $\frac{4}{\epsilon} \log\left(\frac{4}{\delta}\right) \Rightarrow$ poly logarithmic \Rightarrow PAC learnable.

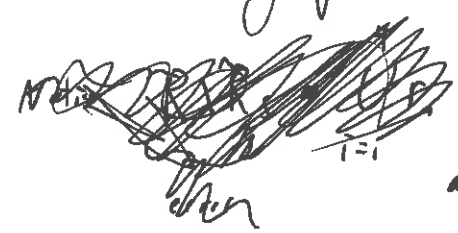
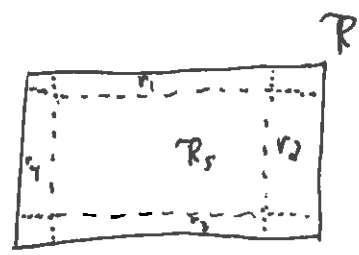
$\underbrace{\frac{4}{\epsilon}}_{\substack{\text{linear} \\ \text{in } \frac{1}{\epsilon}}} \underbrace{\log\left(\frac{4}{\delta}\right)}_{\substack{\text{logarithmic} \\ \text{in } \frac{1}{\delta}}}$

In the homework, you will see that other concept classes aren't quite so simple...

Note that our scheme for axis-aligned rectangles was consistent, meaning in this case that there was no error on the training set.

We now consider a generalization error estimate that holds for any finite concept

In this particular example, we can analyze this probability geometrically.



The r_i are chosen to have probability area exactly $\frac{\epsilon}{4}$

(prob = area under a uniform distribution)

In particular, due to the rectangle geometry, we may argue that

$$P_{S \sim D}(R(R_s) > \epsilon) \leq P\left[\bigcup_{i=1}^4 \{R_s \cap r_i = \emptyset\}\right]$$

that intersecting with R_s

Indeed, $R(R_s) > \epsilon \Rightarrow$ ~~at least one of these~~ ~~at least one of these~~ exterior rectangles ~~is not intersecting with R_s~~

Recalling $P(\bigcup_{i=1}^m A_i) \leq \sum_{i=1}^m P(A_i)$ (union bound), we get

$$P_{S \sim D}(R(R_s) > \epsilon) \leq \sum_{i=1}^4 P_{S \sim D}(\{R_s \cap r_i = \emptyset\})$$

$$\leq 4 \max_{i=1 \rightarrow 4} P_{S \sim D}(R_s \cap r_i = \emptyset)$$

$$\leq 4 \left(1 - \frac{\epsilon}{4}\right)^m$$

$$\leq 4 \exp(-m \epsilon/4), \quad \text{where we note } (1 - \epsilon) \leq \exp(-\epsilon)$$

So, to ensure $4 \exp(-m \epsilon/4) \leq \delta$, it suffices to solve for m :

$$4 \exp(-m \epsilon/4) \leq \delta$$

class and consistent algorithm.

Theorem (Learning bound, finite H , consistent algorithm): Let $|H| < \infty$, $h \in H$ s.t. $h: X \rightarrow Y$.

Let \mathcal{A} be a consistent algorithm, i.e. for any target ~~concept~~ concept $c \in H$, and for any iid. sample S , the hypothesis h_S produced by \mathcal{A} satisfies $\hat{R}_S(h_S) = 0$.

Then, $\forall \epsilon, \delta > 0$, $m \geq \frac{1}{\epsilon} (\log |H| + \log(\frac{1}{\delta}))$

$$\Rightarrow \mathbb{P}_{S \sim D} (R(h_S) \leq \epsilon) \geq 1 - \delta.$$

Remark: It is equivalent to say $R(h_S) \leq \frac{1}{m} (\log |H| + \log(\frac{1}{\delta}))$ with probability at least $1 - \delta$.

Proof: Let $H_\epsilon = \{h \in H \mid R(h) > \epsilon\}$, for any $\epsilon > 0$. Now, if we have m training samples and h is correct on all of them (i.e. h is consistent), this happens with probability

$$\mathbb{P}(\hat{R}_S(h) = 0) \leq (1 - \epsilon)^m.$$

$\hat{R}_S(h) = 0$ consistent, i.e. 0 empirical error
 $\mathbb{P}(\text{wrong on one point}) = \epsilon$.

By a union bound, we see $\mathbb{P}(\exists h \in H_\epsilon \mid \hat{R}_S(h) = 0)$
 $= \mathbb{P}(\{\hat{R}_S(h_1) = 0\} \cup \{\hat{R}_S(h_2) = 0\} \cup \dots \cup \{\hat{R}_S(h_{|H_\epsilon|}) = 0\})$
 $\leq \sum_{h \in H_\epsilon} \mathbb{P}(\hat{R}_S(h) = 0)$

$$\leq \sum_{h \in \mathcal{H}_\varepsilon} (1-\varepsilon)^m$$

$$= |\mathcal{H}_\varepsilon| (1-\varepsilon)^m$$

$$\leq |\mathcal{H}| (1-\varepsilon)^m$$

$$\leq |\mathcal{H}| e^{-m\varepsilon}$$

Therefore, $\mathbb{P}_{S \sim D} (R(h_S) \leq \varepsilon)$

$$= 1 - \mathbb{P}(R(h_S) \geq \varepsilon)$$

$$\geq 1 - \delta$$

provided $|\mathcal{H}| e^{-m\varepsilon} \leq \delta$

$$\Leftrightarrow e^{-m\varepsilon} \leq \frac{\delta}{|\mathcal{H}|}$$

$$\Leftrightarrow -m\varepsilon \leq \log(\delta) - \log(|\mathcal{H}|)$$

$$\Leftrightarrow m\varepsilon \geq \log\left(\frac{1}{\delta}\right) + \log(|\mathcal{H}|)$$

$$\Leftrightarrow m \geq \frac{1}{\varepsilon} \left(\log\left(\frac{1}{\delta}\right) + \log(|\mathcal{H}|) \right)$$