

F₀SML: Lecture 3

①

ex: Consider $X = \{0,1\}^n$, the set of Boolean vectors of length n . Let \mathcal{V}_n be the concept class formed as the power set of X , i.e. the set of all subsets of X . \mathcal{V}_n is quite large: $|X| = 2^n$
 $\Rightarrow |\mathcal{V}_n| = 2^{2^n}$. But, still finite.

How does our Theorem play with this concept class \mathcal{V}_n ? Well, to guarantee a consistent hypothesis, the concept class must be contained in the hypothesis class, i.e.

$$\mathcal{V}_n \subset H \Rightarrow |\mathcal{V}_n| \leq |H|.$$

↑
things we are trying to learn

↑
possible hypotheses

This puts a lower bound on the size of ~~the~~ any consistent hypothesis class:
 $2^{2^n} \leq |H|$. The sample bound in our Theorem is then quite bad:

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

$$\Rightarrow \frac{1}{\epsilon} \left(\log(2^{2^n}) + \log \left(\frac{1}{\delta} \right) \right)$$

$$= \frac{1}{\epsilon} \left(\log(2) \cdot 2^n + \log \left(\frac{1}{\delta} \right) \right).$$

This does not imply PAC learning (it doesn't imply not PAC learning either, though in

fact V_n is not PAC learnable.) Indeed, n is the cost of encoding an element of X , so the sample complexity guaranteed by our theorem is exponential in the size of an element of $X \Rightarrow$ PAC learning is not assured.

- Part of the problem here is requiring consistency. Indeed, this forces H to be quite large in the above example (size $\sim 2^{2^n}$). This suggests relaxing to inconsistent algorithms.

- To do so, recall Hoeffding's inequality:

Theorem (Hoeffding's Inequality): Let $\{X_i\}_{i=1}^m$ be independent random variables, with X_i taking values in $[a_i, b_i] \subset \mathbb{R}$. Then $\forall \epsilon > 0$, the sum $S_m = \sum_{i=1}^m X_i$ satisfies:

$$(1) P(|S_m - E(S_m)| \geq \epsilon) \leq \exp\left(-2\epsilon^2 / \sum_{i=1}^m (b_i - a_i)^2\right)$$

$$(2) P((S_m - E(S_m)) \leq -\epsilon) \leq \exp\left(-2\epsilon^2 / \sum_{i=1}^m (b_i - a_i)^2\right).$$

Remark: Hoeffding's Inequality is an example of a concentration inequality, which describes how a random variable (in this case S_m) concentrates around its mean (in this case $E(S_m)$). We will find use for such results over and over again. They are major tools in statistical and machine learning.

- From Hoeffding, we immediately get:

(3)
Corollary: Fix $\epsilon > 0$. For any hypothesis $h: X \rightarrow \{0,1\}$, the following inequalities hold:

$$(1) \mathbb{P}_{S \sim D} \left(\hat{R}_S(h) - R(h) \geq \epsilon \right) \leq \exp(-2\epsilon^2 m)$$

$$(2) \mathbb{P} \left(\hat{R}_S(h) - R(h) \leq -\epsilon \right) \leq \exp(-2\epsilon^2 m)$$

Remark: Taking a union yields

$$\mathbb{P} \left(|\hat{R}_S(h) - R(h)| \geq \epsilon \right) \leq 2 \exp(-2\epsilon^2 m)$$

In other words, the empirical error of a hypothesis h concentrates around its (true) generalization error. Clearly we expect the discrepancy to decay with the sample size m . Hoeffding establishes this decay is exponential.

To rephrase:

Theorem (Generalization bound, $|H|=1$): Fix a hypothesis $h: X \rightarrow \{0,1\}$. Then for

$$\text{any } \delta > 0, \mathbb{P} \left(R(h) \leq \hat{R}_S(h) + \sqrt{\frac{\log(\frac{2}{\delta})}{2m}} \right) \geq 1 - \delta.$$

The question, then, is how to generalize this to account for $|H| \geq 2$. As in the case where A is consistent, we shall require a uniform bound over $h \in H$.

Theorem (Learning bound, $|H|$ finite, inconsistent): Let $|H| < \infty$ be a finite hypothesis class. Then $\forall \delta > 0$, (4)

$$\mathbb{P}\left(R(h) \leq \hat{R}_S(h) + \sqrt{\frac{\log|H| + \log\frac{2}{\delta}}{2m}}, \forall h \in H\right) \geq 1 - \delta.$$

Proof: Let $\{h_i\}_{i=1}^{|H|}$ be the (finitely many) elements of H . Then:

$$\begin{aligned} & \mathbb{P}(\exists h \in H \mid |\hat{R}_S(h) - R(h)| > \epsilon) \\ &= \mathbb{P}\left(\bigcup_{i=1}^{|H|} |R_S(h_i) - R(h_i)| > \epsilon\right) \end{aligned}$$

$$\leq \sum_{i=1}^{|H|} \mathbb{P}(|R_S(h_i) - R(h_i)| > \epsilon)$$

$$\leq 2 \cdot |H| \cdot \exp(-2n\epsilon^2) \quad \text{by our corollary to Hoeffding.}$$

Setting $2 \cdot |H| \cdot \exp(-2n\epsilon^2)$ to δ and solving for ϵ gives the desired result. \blacksquare