

Machine Learning Exercises (chapter 3)

2. Part 1: We define S as a training set of $X \times \{0, 1\}$. If the training set contains a positive instance (x^+), the algorithm returns h_{x^+} , otherwise, it returns h^- which is an ERM (we assume the fact that there can be at most a point with label 1) $\rightarrow L_S(h_S) = 0$

2. Part 2: (PAC learnability) - Assume that $x^+ \notin S$, if $D(\{x^+\}) \leq \epsilon$ then

$L_{(D, f)}(h) \leq \epsilon$ for all h in H . Suppose that $D(\{x^+\}) > \epsilon$,

then: $\forall x' \in X$ st. $x' \neq x^+ \rightarrow D(\{x'\}) \leq 1 - \epsilon$. then we have:

$$\begin{aligned} \{S|_X : L_{(D, f)}(h_S) > \epsilon\} &= \{S|_X : x^+ \notin S|_X, D(\{x^+\}) > \epsilon\} = \\ &= \{S|_X : \forall x' \in S|_X D(\{x'\}) \leq 1 - \epsilon\} \end{aligned}$$

$$\begin{aligned} \Rightarrow D^m(\{S|_X : L_{(D, f)}(h_S) > \epsilon\}) &= D^m(\{S|_X : \forall x' \in S|_X D(\{x'\}) \leq 1 - \epsilon\}) \leq \\ &\leq (1 - \epsilon)^m \leq e^{-\epsilon m} \end{aligned}$$

$$\text{Let } \delta \in (0, 1) \text{ st. } e^{-\epsilon m} \leq \delta \Rightarrow m \geq \frac{\log(\frac{1}{\delta})}{\epsilon}$$

$$\text{Hence } \mathcal{H} \text{ is PAC learnable with } m_H \leq \frac{\log(\frac{1}{\delta})}{\epsilon}$$

3. Assume that ERM algorithm A takes S as training set and it returns the tightest circle that contains all positive instances.

Define h_S as output function and its radius by r_S , we having

Realizability assumption (existing $h^* \in H$ that its radius is r^*)

Now Let $\epsilon, \delta \in (0, 1)$, we suppose that there is a r which $r \leq r^*$

$$\text{then: } D_X(\{x: r \leq \|x\| \leq r^*\}) = \epsilon / E = \{x \in \mathbb{R}^2, r \leq \|x\| \leq r^*\}$$

we suppose E then we have:

$$\begin{aligned} P(L_D(h_S) > \epsilon) &= P(x_i \in S \text{ s.t. } x_i \notin E) = \prod_{i=1}^m (1 - P(x_i \in E))^m = \\ &= (1 - \epsilon)^m \leq e^{-m\epsilon} \leq \delta \end{aligned}$$

Now Let $\delta \in (0, 1)$, then we have:

$$\begin{aligned} &e^{-m\epsilon} \leq \delta \\ \text{or} \\ &m \geq \frac{\log(1/\delta)}{\epsilon} \end{aligned}$$

Hence H is PAC-learnable and $m_H(\epsilon, \delta) \leq \frac{\log(1/\delta)}{\epsilon}$

5. H_B : set of bad hypotheses /

we suppose the $h \in H_B$ s.t. $L_{(\bar{D}_m, f)}(h) \geq \epsilon$ then we have:

$$P[L_{(S, f)}(h) = 0] = \prod_{i=1}^m P_{x \sim D_i}[f(x) = h(x)] =$$

$$= \prod_{i=1}^m P_{x \sim D_i}[f(x) = h(x)] \leq \left(\frac{1}{m} \sum_{i=1}^m P_{x \sim D_i}[f(x) = h(x)] \right)^m =$$

(by the geometric-arithmetic)

$$= (P_{x \sim \bar{D}_m}[f(x) = h(x)])^m \leq (1 - \epsilon)^m \leq e^{-\epsilon m}$$

Therefore we have,

$$P[\exists h \in H: L_{(\bar{D}_m, f)}(h) \geq \epsilon, L_{(S, f)}(h) = 0] \leq |H_B| e^{-\epsilon m} \leq |H| e^{-\epsilon m}$$

6. We define $h_S := A(S)$, D : distribution over $X \times \{0, 1\}$ then we

$$\text{have: } D^m(\{S \mid L_D(h_S) > \min_{h \in H} L_D(h) + \epsilon\}) < \delta$$

Suppose $f \in H$: f is a true labeling function.

\downarrow
(determine probability of x over y .)

If $m > m_H(\epsilon, \delta)$, $A(S)$ will perform and output h . then:

$$D_X^m(\{S \mid L_{(D_H, f)}(h) > \epsilon\}) < \delta$$

H is PAC-learnable

with realizability: $\min L_D(h) = 0$, hence A is a PAC learner of H

7. Suppose $x \in X$ and p : conditional probability of a positive label given x , then we have:

$$\begin{aligned} P[f_0(x) \neq Y | X=x] &= P[Y=0 | X=x] \mathbb{1}_{[p > \frac{1}{2}]} + P[Y=1 | X=x] \mathbb{1}_{[p < \frac{1}{2}]} \\ &= (1-p) \mathbb{1}_{[p > \frac{1}{2}]} + p \mathbb{1}_{[p < \frac{1}{2}]} \\ &= \min\{p, 1-p\} \end{aligned}$$

Suppose that g is classifier $g: X \rightarrow \{0, 1\}$:

$$\begin{aligned} P[g(x) \neq Y | X=x] &= P[g(x)=0 | X=x] P[Y=1 | X=x] + P[g(x)=1 | X=x] P[Y=0 | X=x] \\ &= P[X=x] p + P[g(x)=1 | X=x] (1-p) \geq \\ &\geq P[g(x)=0 | X=x] \min\{p, 1-p\} + P[g(x)=1 | X=x] \min\{p, 1-p\} \\ &= \min\{p, 1-p\} \end{aligned}$$

Therefore:

$$\begin{aligned} L_D(g) &= E_{(x,Y) \sim D} [\mathbb{1}_{g(x) \neq Y}] = E_{x \sim D_x} [E_{Y \sim D_{Y|x}} [\mathbb{1}_{g(x) \neq Y} | X=x]] \geq \\ &\geq E_{x \sim D_x} [E[\mathbb{1}_{f_0(x) \neq Y} | X=x]] = E[\mathbb{1}_{f_0(x) \neq Y}] = L_D(f_0) \end{aligned}$$

Hence:

$$L_D(f_0) \leq L_D(g) \quad \checkmark$$

Machine Learning Exercises (chapter 5)

5.2 - a) pros and cons: (Algorithm A)

Pros: It has small E_{est} , because of less complexity (less prone to overfitting)

This Algorithm can easily interpret models in a plot (dim: 2d)

Cons: Inductive bias might be too large (high E_{app}) and we can't use A, P, I in our model

(Algorithm B)

Pros: Smaller ~~than~~ inductive bias, reducing risk of underfitting.

B has

It has small E_{app} .

Cons: It has larger E_{est} . because of that our ^{complex} model may lead to overfitting

2.b) Increasing size of S (training data set) leads to $L_S(h_S)$ is a better estimate of $L_P(h_S)$. It means lower E_{est} , because of that B is better than A.

E_{app} can be reduced by choosing Algorithm B (Complexity)

E_{est} can be decreased with the size of S .