

10.1

Let $\epsilon, \delta \in (0, 1)$. pick k "chunks" of size $m_H(\epsilon/2)$

Apply A on every k chunks, to obtain $\hat{h}_1, \dots, \hat{h}_k$.

Apply ERM over $\hat{H} = \{\hat{h}_1, \dots, \hat{h}_k\}$

training data \leadsto last chunk of size $\left\lceil \frac{2 \log(4k/\delta)}{\epsilon^2} \right\rceil$

with probability at least $1 - \delta/2$, $L_D(\hat{h}) \leq \min_{i \in [k]} L_D(h_i) + \epsilon/2$

Apply the union bound \leadsto probability at least $1 - \delta$

$$L_D(\hat{h}) \leq \min_{i \in [k]} L_D(h_i) + \epsilon/2$$

$$\leq \min_{h \in H} L_D(h) + \epsilon$$

10.4

(a) Let X be finite set of size n .

Let B be the class of all func $X \rightarrow \{0, 1\}$

$L(B, T) = B$. for any $T \geq 1$,

$$\text{VC dim}(B) = \text{VC dim}(L(B, T)) = \log 2^n = n$$

(b) B class of decision stump in \mathbb{R}^d

$$B = \{h_{j,b,\theta} : j \in [d], b \in \{-1, 1\}, \theta \in \mathbb{R}\}$$

where $h_{j,b,\theta}(x) = b \cdot \text{sign}(\theta - x_j)$

①

for each $j \in [d]$ let $B_j = \{h_{b,\theta} : b \in \{-1, 1\}, \theta \in \mathbb{R}\}$
 $h_{b,\theta}(x) = b \cdot \text{sign}(\theta - x_j)$, ($\text{VCdim}(B_j) = 2$)

$$B = \bigcup_{j=1}^d B_j \Rightarrow \text{VCdim}(B) \leq 16 + 2 \log d$$

$d = 2^k$ for some $k \in \mathbb{N}$. $A \in \mathbb{R}^{k \times d} \Rightarrow$ matrix whose
columns range over set $\{0, 1\}^k$ for each $i \in [k]$

let $x_i = A_{i, \rightarrow}$. ~~we claim~~ $C = \{x_1, \dots, x_k\}$ is shattered

let $I \subseteq [k]$. the instances $[k] \setminus I$ are labeled
negatively. exist an index j such that $A_{i,j} = x_{i,j}$

iff $i \in I$. Then, $h_{j, -1, 1/2}(x_i) = 1$ if $i \in I$. $= 1$

(c) for each $i \in [T_{k/2}]$, let $x_i = [i/k] A_{i, \rightarrow}$.

set $C = \{x_i : i \in [T_{k/2}]\}$ is shattered by $L(B_d, T)$

Let $I \subseteq [T_{k/2}]$. $I = I_1 \cup \dots \cup I_{T_2}$. I_1 is subset of

$\{(t-1)k+1, \dots, tk\}$. let j_t be the corresponding column

of A . $h(x) = \text{sign}((h_{j_1, -1, 1/2} + h_{j_2, 1, 3/2} + \dots + h_{j_{T_2}, -1, T_2 - 1/2})(x))$

Then $h(x_i) = 1$ iff $i \in I \Rightarrow$ observe that $h \in L(B_d, T)$

(2)

11.1

let h be the output of the described learning algorithm. $L_D(h) = 1/2$.

calculating Estimate $L_V(h)$. parity of S is 0

1. Fix some fold $\{(x_i, y_i)\} \subseteq S$.

- parity of $S \setminus \{x_i\}$ is 1. follows that $y_i = 0$.

trained using $S \setminus \{x_i\}$, outputs the constant predictor $h(x_i) = 1$. Hence the leave-one-out estimate using this fold is 1.

- parity of $S \setminus \{x_j\}$ is 0. it follows that $y_j = 1$.

when being trained using $S \setminus \{x_j\}$. the algorithm outputs the constant predictor $h(x_j) = 0$.

Averaging over folds, the estimate of the error of h is 1. consequently, the difference between the estimate and true error is $1/2$. parity of S is 0 is analyzed analogously.

11.2 H_k in the agnostic-Pac model provides the following bound for ERM hypothesis h :

$$L_D(h) \leq \min_{h \in H_k} L_D(h) + \sqrt{\frac{2(k+1+\log(1/\delta))}{m}}.$$

Assume that j is the minimal index which contains a hypothesis $h^* \in \arg\min_{h \in H_k} L_D(h)$. Fix some $r \in [k]$.

By Hoeffding's inequality, with probability at least $1 - \frac{\delta}{2k}$

$$|L_D(\hat{h}_r) - L_V(\hat{h}_r)| \leq \sqrt{\frac{1}{2\alpha m} \log \frac{4}{\delta}}.$$

Applying the union bound, we obtain that with probability at least $1 - \frac{\delta}{2}$, the following inequality holds for $\forall r \in [k]$:

$$\begin{aligned} L_D(\hat{h}) &\leq L_V(\hat{h}) + \sqrt{\frac{1}{2\alpha m} \log \frac{4k}{\delta}} \\ &\leq L_D(\hat{h}_r) + \sqrt{\frac{2}{\alpha m} \log \frac{4k}{\delta}}. \end{aligned}$$

① ~~with~~ with probability $1 - \frac{\delta}{2}$ we have:

$$L_D(\hat{h}) \leq L_D(\hat{h}_j) + \sqrt{\frac{2}{\alpha m} \log \frac{4k}{\delta}}.$$

$$\begin{aligned} \rightarrow L_D(\hat{h}_j) &\leq L_D(h^*) + \sqrt{\frac{2}{(1-\alpha)m} \log \frac{4|H_j|}{\delta}} \\ &= L_D(h^*) + \sqrt{\frac{2}{(1-\alpha)m} \log \frac{4|H_0|}{\delta}} \end{aligned}$$

combining the two last inequalities:

$$L_D(\hat{h}) \leq L_D(h^*) + \sqrt{\frac{2}{\alpha m} \log \frac{4k}{\delta}} + \sqrt{\frac{2}{(1-\alpha)m} \log \frac{4|H_j|}{\delta}}$$

we ~~can~~ conclude that

$$L_D(\hat{h}) \leq L_D(h^*) + \sqrt{\frac{2}{\alpha m} \log \frac{4k}{\delta}} + \sqrt{\frac{2}{(1-\alpha)m} (j + \log \frac{4}{\delta})}$$

comparing the two bounds.

if j is logarithmic in k , we achieve a logarithmic improvement

18.2

(a) information ~~for~~ gain for feature 1 is:

$$H(1/2) - \left(\frac{3}{4} H\left(\frac{2}{3}\right) + \frac{1}{4} H(0) \right) \approx 0.22$$

the information gain for feature 2 as well as feature 3

$$H(1/2) - \left(\frac{1}{2} H(1/2) + \frac{1}{2} H(1/2) \right) = 0$$

we won't be able to classify all three examples perfectly. since we have 4 examples in the training set, it follows that the training error is at least $1/4$.

(b)

