

Machine Learning HW2

Bill901174 傅啟恩

1. By one-sided Hoeffding's inequality,

$$P(\mu > \nu + \epsilon) \leq \exp(-2\epsilon^2 N), \text{ where in this case,}$$

$$\nu = \frac{c_m}{N_m} \text{ (sampled prob.), real probability } \mu = \mu_m, \text{ choose } \epsilon = \sqrt{\frac{\ln t - \frac{1}{2} \ln \delta}{N_m}}$$

$$\begin{aligned} \exp(-2\epsilon^2 N) &= \exp\left(-2 \cdot \frac{\ln \frac{t}{\sqrt{\delta}}}{N_m} N_m\right) \\ &= \exp\left(\ln \frac{\delta}{t^2}\right) = \frac{\delta}{t^2} = \delta t^{-2} \end{aligned}$$

$$\Leftrightarrow P\left(\mu_m > \frac{c_m}{N_m} + \sqrt{\frac{\ln t - \frac{1}{2} \ln \delta}{N_m}}\right) \leq \delta t^{-2} \quad \#$$

2.

$$\sum_{t=1}^{\infty} \sum_{m=1}^M \exp\left(-2 \cdot \frac{\ln \frac{tM}{\sqrt{\delta}}}{N_m} N_m\right)$$

$$= \sum_{t=1}^{\infty} \frac{\delta t^{-2}}{M^2} \cdot M \leq \delta M^{-2} \frac{\pi^2}{6} \leq \delta$$

for $m = 1, 2, \dots, M$, $t = M+1, M+2, \dots, \infty$

$$P\left(\mu_m \leq \frac{c_m}{N_m} + \sqrt{\frac{\ln t - \frac{1}{2} \ln \delta + \ln M}{N_m}}\right) \geq 1 - \delta \quad \#$$

3.

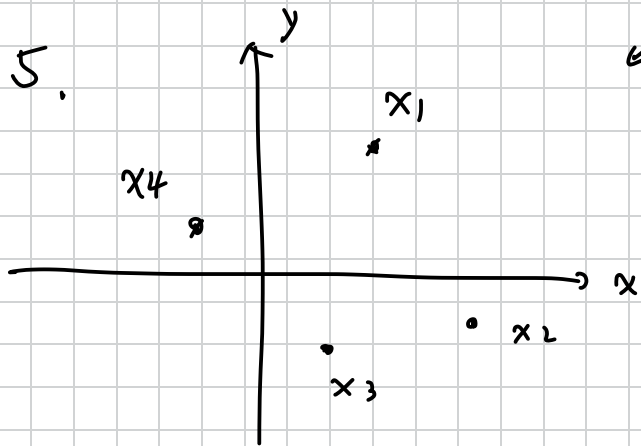
if all tickets
For each ticket, A/C or A/D or B/C or B/D (4 kinds)
would cause there's some "green"

\therefore the probability would be $\left(\frac{1}{2}\right)^5 \times 4 \text{ kinds} - \left(\frac{1}{4}\right)^5 \times 4$

$$= \frac{31}{256} \quad \#$$

4. $P(\text{five tickets contain 5 green 2's})$ would be
all tickets choose B or D

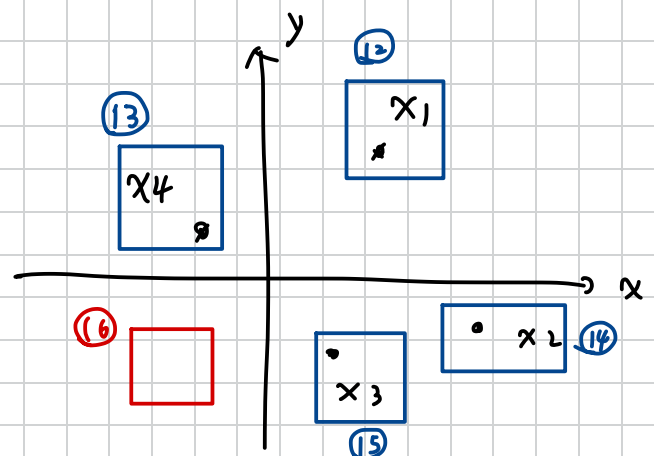
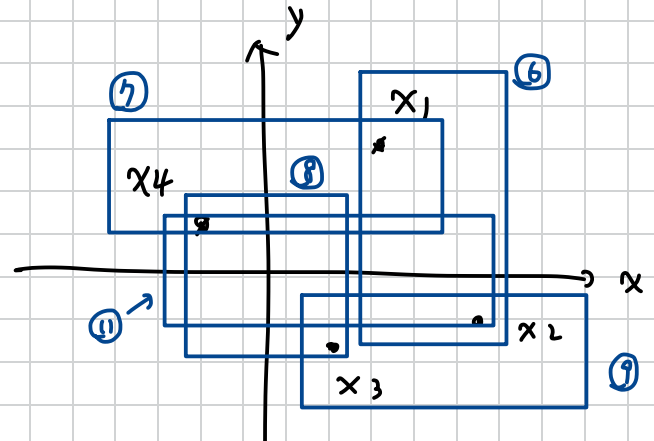
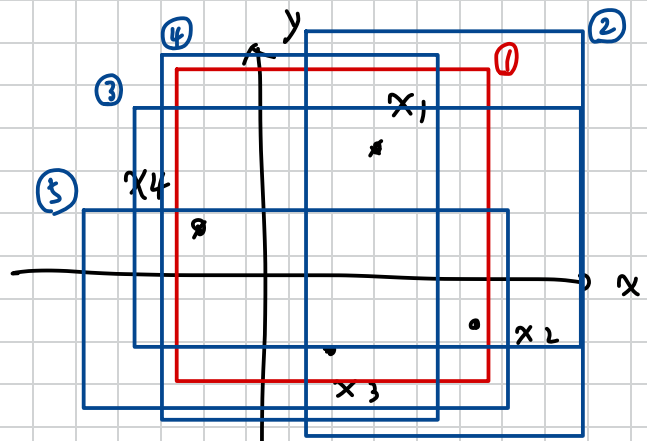
$$= \left(\frac{1}{2}\right)^5 = \frac{1}{32} *$$



← If x_1, x_2, x_3, x_4 are mapped to \mathbb{R}^2 like this,

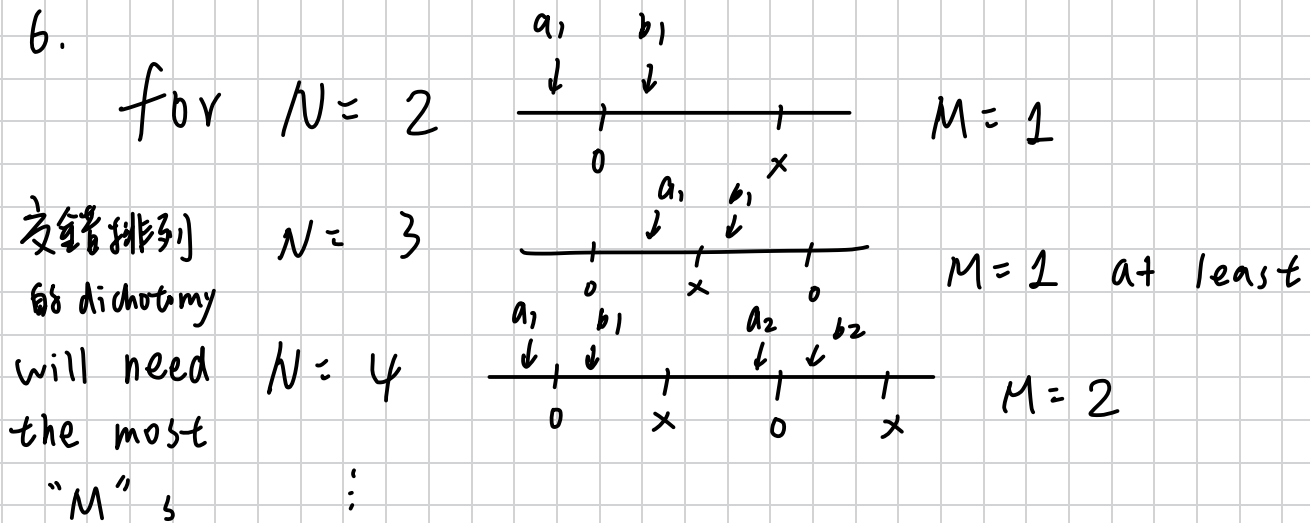
(choose "0" as $-1 \in$ in the rect
"x" as $+1 \in$ outside "

	x_1	x_2	x_3	x_4
①	0	0	0	0
②	0	0	0	x
③	0	0	x	0
④	0	x	0	0
⑤	x	0	0	0
⑥	0	0	x	x
⑦	0	x	x	0
⑧	x	x	0	0
⑨	x	0	0	x
⑩	0	x	0	x
⑪	x	0	x	0
⑫	0	x	x	x
⑬	x	x	x	0
⑭	x	0	x	x
⑮	x	x	0	x
⑯	x	x	x	x



\Rightarrow for 4-point vectors x_1, x_2, x_3, x_4 , there exists some x_1, x_2, x_3, x_4 that can be shattered by the hypothesis. ($N=4$ is NOT the break point)

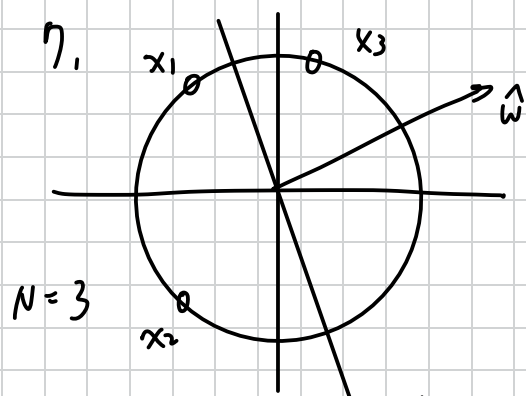
\Leftrightarrow The VC dimension is no less than 4!
(= break point - 1)



for N inputs, there exists at least $\lfloor N/2 \rfloor$ non-adjacent intervals that marked as "0" or "x" as shown above, so we need $a_1, b_1, a_2, b_2 \dots a_M, b_M$, where $M \geq \lfloor \frac{N}{2} \rfloor$ to shatter the hypothesis. Therefore, the break point of the hypothesis is $2^*M + 2$, the VC Dimension =

$$2^*M + 2 - 1 = \underline{2^*M + 1}^*$$

(for $M=1$, hypothesis cannot shatter $2^*M + 2$
= 4 inputs data because it would contain
 $0 \ x \ 0 \ x$
case.



$$H_0: \{h: h(x) = \text{sign}(w_1 x_1 + w_2 x_2)\}$$

$$\text{Think as } \hat{w} = (w_1, w_2)$$

$$\hat{x} = (x_1, x_2)$$

$$h(x) = \text{sign}(\hat{w}^T \hat{x})$$

We can know that w^\perp separate the \mathbb{R}^2 into 2 sides, one that causes $\text{sign}(\hat{w}^T \hat{x}) > 0$ with \hat{x} falls on the side, the other causes $\text{sign}(\hat{w}^T \hat{x}) \leq 0$

x_1 x_2 x_3 Initial state

x x 0

0 x 0 touch x_1

0 0 0 touch x_2

0 0 x touch x_3

x 0 x touch x_1

x x x touch x_2

~~x x x touch x_3~~

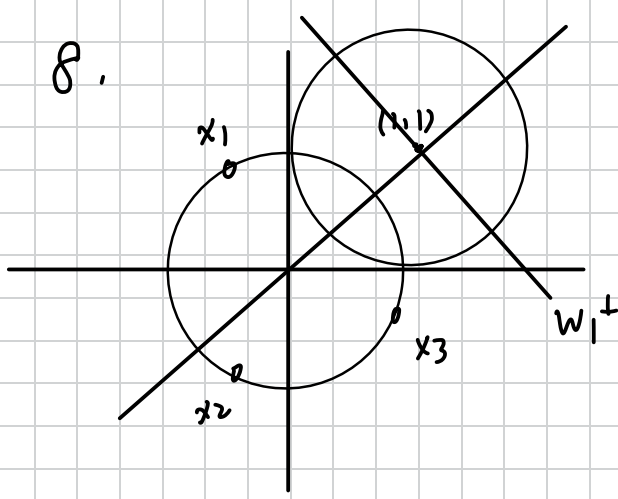
table 1

\Rightarrow In this figure shown above, we suppose at Initial state, $x_1 = "x"$, $x_2 = "x"$, and $x_3 = "0"$, as we rotate the w^\perp by counterclockwise, we will then create different "dichotomies" that satisfy the hypothesis (for some \hat{w}) After rotating 360° , w^\perp will touch $2^* N$ points, causes $2^* N + \underbrace{1}_{\text{repeat}}$

$$= 2^* N \text{ dichotomies}$$

$$\Rightarrow \underline{M_H(N) = 2^* N}$$

8.



$$H_0: \{ h: h_0(x) = \text{sign}(w_1 x_1 + w_2 x_2) \}$$

$$H_1: \{ h: h_1(x) = \text{sign}(w_1(x_1 - 1) + w_2(x_2 - 1)) \}$$

From problem 7, $N=3$ inputs has

$$VC \text{ Dimension} = 2, \quad m_H(3) = 2^2 = 4$$

for origin-passing perceptrons.

If we consider $H = H_0 \cup H_1$, we can create lines (w_1^\perp)

such that append 2 dichotomies $(x_1 = x_2 = x_3 = "0" / "x")$,

\Rightarrow there exists some x_1, x_2, x_3 such that hypothesis H can shatter the case $(2^2 + 2 = 6 = 2^{(N=3)})$

\Rightarrow VC Dimension of H is at least 3.

Now we consider $N=4$ points, because in PLA

we know that $N=4$ cannot be shattered, and the

set of PLA ^(perceptions) lines $(w^\perp) \supset \{ \text{origin-passing} \cup (1,1)\text{-passing lines } (w^\perp \text{ in } H) \}$

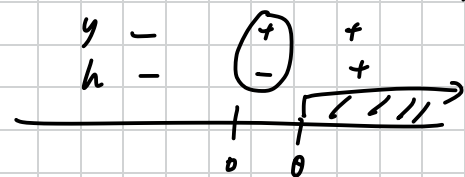
\Rightarrow Thus, if PLA cannot shatter $N=4$, hypothesis

H must also cannot shatter $N=4$

\Rightarrow VC Dimension of $H = 4 - 1 = 3$ #

9. $h_{S,\theta}(x) = S \cdot \text{sign}(x - \theta)$, $y = \text{sign}(x) + 10\% \text{ flip}$

Case 1: $S = +1$



$$E_{\text{out}}(h_{S,\theta}) = P(h_{S,\theta}(x) \neq y)$$

$$= 0.9 P(\text{sign}(x - \theta) \neq \text{sign}(x)) + 0.1 \times P(\text{sign}(x - \theta) = \text{sign}(x))$$

I. Consider $\theta > 0$

$$P(\text{sign}(x - \theta) \neq \text{sign}(x)) = P(0 \leq x \leq \theta) = \frac{\theta}{2}$$

$$P(\text{sign}(x - \theta) = \text{sign}(x)) = 0.5 + P(x \geq \theta)$$

$$= 0.5 + 0.5(1 - \theta)$$

$$= 1 - 0.5\theta$$

II. Consider $\theta < 0$

$$P(\text{sign}(x - \theta) \neq \text{sign}(x)) = P(\theta \leq x \leq 0) = \frac{-\theta}{2}$$

$$P(\text{sign}(x - \theta) = \text{sign}(x)) = 0.5 + P(x \leq \theta)$$

$$= 0.5 + \frac{\theta + 1}{2}$$

$$= 1 + 0.5\theta$$

Conclusion

$$\Rightarrow P(\text{sign}(x - \theta) \neq \text{sign}(x)) = \frac{|\theta|}{2}$$

$$P(\text{sign}(x - \theta) = \text{sign}(x)) = 1 - 0.5|\theta|$$

$$E_{\text{out}}(h_{S,\theta}) = 0.9 \times \frac{|\theta|}{2} + 0.1(1 - 0.5|\theta|)$$

$$= 0.45|\theta| + 0.1 - 0.05|\theta|$$

$$= 0.1 + 0.4|\theta|$$

case 2: $s = -1$

$$\begin{aligned} E_{\text{out}}(h_{s, \theta}) &= 0.9 \cdot (1 - 0.5|\theta|) + 0.1 \times \frac{|\theta|}{2} \\ &= 0.9 - 0.45|\theta| + 0.05|\theta| \\ &= 0.9 - 0.4|\theta| \end{aligned}$$

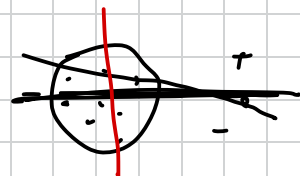
Combine case 1+2

$$\Rightarrow E_{\text{out}}(h_{s, \theta}) = 0.5 - 0.4s + 0.4s \cdot |\theta| \quad \#$$

13.

Suppose $C(N, d)$ means the number of dichotomies for N points on \mathbb{R}^d .

By the Cover's Theorem in the reference,



$$C(N+1, d) = C(N, d) + C(N, d-1), \text{ which means}$$

if we anchor one points

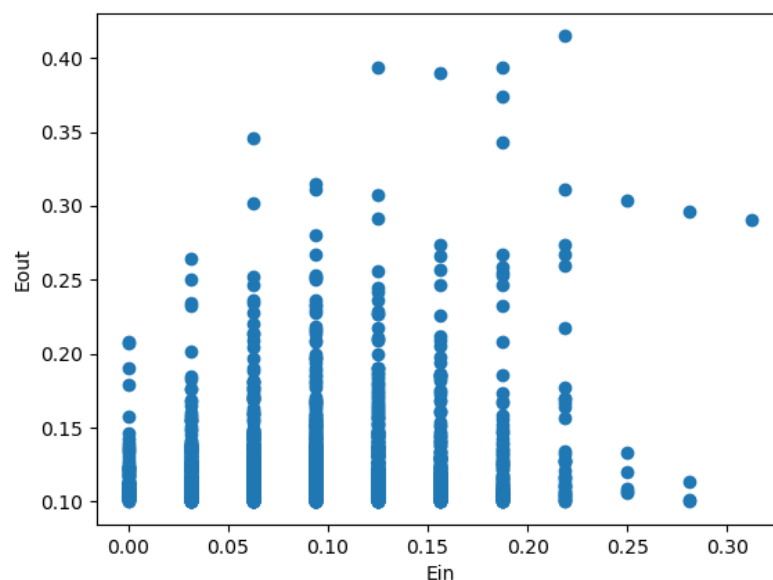
$$\begin{aligned} m_H(N) &= 2 \sum_{i=0}^{d-1} C_{i}^{N-1} \quad \text{for } k=1, \hat{a}_1 = \hat{o} \\ &= 2 \sum_{i=0}^{d-1} C_{i}^{N-1} = 2 \times (C_0^{N-1} + C_1^{N-1}) \\ &= 2 \times (1 + N-1) = 2N \end{aligned}$$

\Rightarrow if we anchor k points ($0 \leq k < d$)

$$m_{\tilde{H}}(N) = 2 \sum_{i=0}^{d-k} C_{i}^{N-1}$$

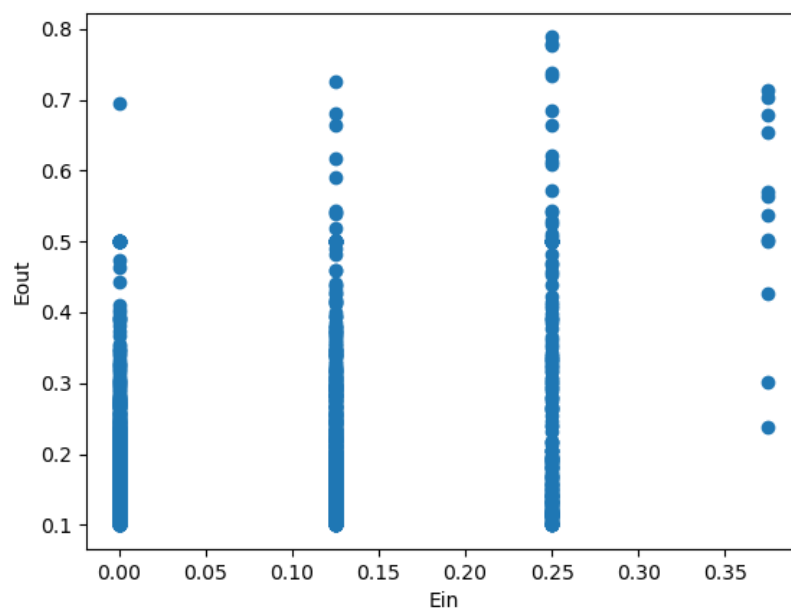
#

10. Median of $E_{out}(g) - E_{in}(g)$: 0.03747820081411739



11. Median of $E_{out}(g) - E_{in}(g)$: 0.12087622283412079

We saw that by reducing the size of x from 32 to 8, the number of different E_{in} become lower, and the range of E_{out} becomes larger from approximately (0, 0.45) to (0, 0.8), and also the median of $E_{out} - E_{in}$ becomes 4 times larger than problem 10, which means E_{out} is more different from E_{in} in smaller dataset in decision stump.



12. Median of $E_{out}(g) - E_{in}(g)$: 0.004120047286207934

By randomly chosen h_s, θ as g , with s uniformly sampled from $\{-1, +1\}$ and θ uniformly sampled from $[-1, 1]$, we saw that E_{in} has more possibilities, which range from $(0, 1.0)$, and so does E_{out} , and the median of $E_{out} - E_{in}$ becomes very small compared to problem 10, 11, which means E_{out} is very closed to E_{in} in each case. We can also see that E_{out} is higher as E_{in} goes higher, this means E_{out} and E_{in} share the same growing tendency.

