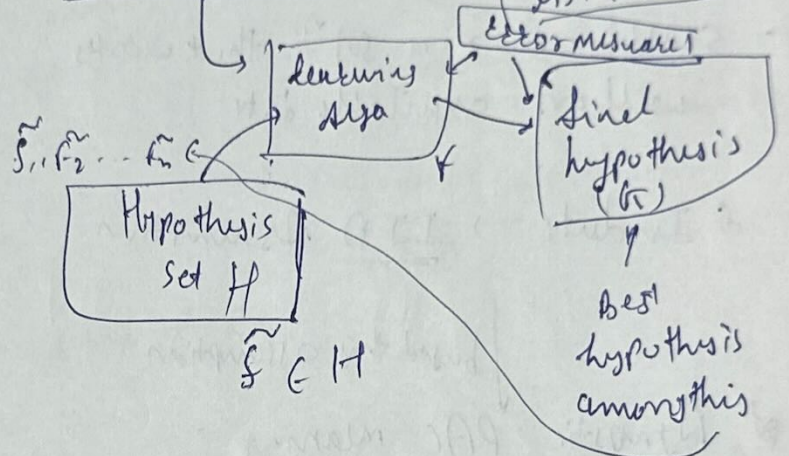


* Theory of generalization

Unknown Target Function.
 $f: \mathcal{X} \rightarrow \mathcal{Y}$

Training example

Unknown
I/P
Distribution



* A statistical learning framework:-

→ Domain - \mathcal{X} , Label - \mathcal{Y}

- training data

- data Generation Model:-

Prob. dist.ⁿ over \mathcal{X} is D .

label funⁿ $f: \mathcal{X} \rightarrow \mathcal{Y}$

- measures ~~fun~~ success.

$$L_{D, P}(h) = \sum_{x \sim D} P[h(x) \neq f(x)]$$

\nwarrow \nearrow
 D, P h
 \nearrow \nwarrow
 dist^n hypothesis

\downarrow
 $P(h) \neq f(x)$
 \downarrow
 $\text{H} \neq \text{means } 'y'$

this is about miss classification rate

*ERM Framework (Empirical Risk Minimization)

training error

$$L_S(h) = \frac{|\{i \in [m] : h(x_i) \neq y_i\}|}{m}$$

$$\frac{L_{\text{total with miss classify points}}}{\text{total instances}}$$

- search for a solⁿ that works well on available data

Inductivity \rightarrow IID Assumption

based this assumption

* Agnostic PAC meaning:

probably approximately correct

\rightarrow H is PAC learnable if

two conditions should be followed

- $m_H : (0,1)^2 \rightarrow N$
- $\epsilon, \delta \in (0,1)$, D over X, Y

error failure prob.

$$E_{in}(h) = \frac{1}{N} \sum_{n=1}^N [h(x_n) \neq f(x_n)]$$

training error expectation

$$E_{out}(h) = P[h(x) \neq f(x)] \text{ (out of sample error)}$$

kind of testing error expectation (unknown to us)

for fixed hypo. h

$$P[|E_{in}(h) - E_{out}(h)| > \epsilon] \leq \delta$$

prob. of failure

⑧ Let Hoeffding's Inequality:

$$P\left(\left|\frac{1}{m} \sum_{i=1}^m \theta_i - \mu\right| > \epsilon\right) \leq 2 \exp\left(\frac{-2m\epsilon^2}{(a-b)^2}\right)$$

Sample mean

mean of IID R.V.
original mean

(m = no. of samples)

this quantity we want to this quantity should be small.

$\theta_i \rightarrow$ samples (R.V.)

$m \rightarrow$ no. of samples

$a, b \rightarrow$ bounds for prob of R.V.

$$\text{so } \delta \geq 2 \exp\left(\frac{-2m\epsilon^2}{(a-b)^2}\right)$$

$$\therefore \frac{1}{2} \log \delta \geq \frac{-2m\epsilon^2}{(a-b)^2}$$

$$m\epsilon^2 \geq \frac{1}{2} \log \frac{1}{\delta} \cdot (a-b)^2$$

$$\therefore m \geq \frac{1}{2\epsilon^2} \log \frac{1}{\delta} \cdot (a-b)^2$$

we want δ should be small. So: $\therefore m \geq \frac{1}{2\epsilon^2} \log \frac{1}{\delta} \cdot (a-b)^2$

$$\therefore m \geq O\left(\frac{1}{\epsilon^2} \log \frac{1}{\delta}\right)$$

means ϵ, δ are const. so.

we need increase successful prob. then m should be increase.

hook \Rightarrow learning from data

by abou Ismail.

$$A \quad E_{in}(h) = \frac{1}{N} \sum_{i=1}^N 1(h(x_i) \neq f(x_i))$$

↓
Training error (ERM)

$$E_{out}(h) = P[h(x) \neq f(x)]$$

↓ Marble experiments:

Bagging with ~~random~~ ~~data~~ marble

$$V = m$$

$$U = mu$$

$$P(|E_{in}(g) - E_{out}(g)| > \epsilon) \leq 2Me^{-2\epsilon^2 N}$$

$$E_{in}(g) - E_{out}(g) > \epsilon$$

$$E_{out}(g) \leq E_{in}(g) + \sqrt{\frac{1}{2N} \ln \frac{2M}{\delta}}$$

generalization bound

$$\delta \leq 2Me^{-\frac{2m\epsilon^2}{c}}$$

$$\frac{\delta}{2M} \leq \exp\left(-\frac{2m\epsilon^2}{c}\right)$$

($\because m = \text{no. of hypothesis}$)

($\because N = \text{no. of samples}$)

$$\ln \frac{\delta}{2M} \geq -\frac{2m\epsilon^2}{c}$$

$$\frac{1}{2m} \ln \frac{\delta}{2M} \leq \epsilon^2$$

$$\sqrt{\frac{1}{2m} \ln \frac{\delta}{2M}} \leq \epsilon$$

$$\ln \frac{2M}{\delta} \leq 2m\epsilon^2$$

$$\epsilon \geq \sqrt{\frac{1}{2m} \ln \frac{2M}{\delta}}$$

what if $M = \infty$?? almost always the case!!!!

→ So we need to represent the infinite No. to finite no.

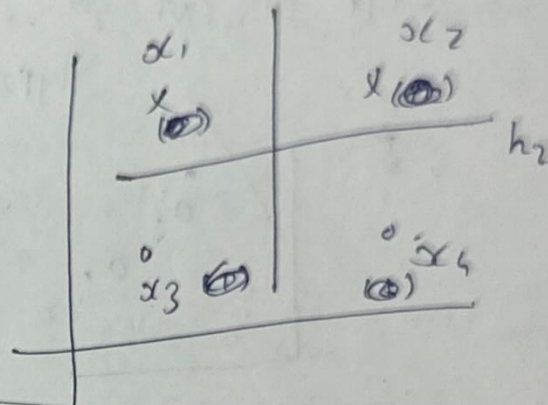
* Def A: for $x_1, x_2, \dots, x_n \in N$ the dichotomies generated by H on these points are defined by

$$H(x_1, \dots, x_n) = \{ h(x_1), \dots, h(x_n) \mid h \in H \}$$

* ~~Def B~~

$$H = \{h_1, h_2\}$$

$$\begin{aligned} h_1(x_1, x_2, x_3, x_4) &= (0, 0, 1, 1) \\ &= (1, 0, 1, 0) \\ h_2 &= (0, 0, 1, 1) \end{aligned}$$



$$H = \{(1, 0, 1, 0), (0, 0, 1, 1)\}$$

this is finite set because the dichotomies m_H^n we can create dichotomies is 2^n in this case only

* def B: the growth function for a hypothesis set H is

(this funⁿ is real no.)
 p^+

$$m_H(n) = \max_{x_1, \dots, x_n} |H(x_1, x_2, \dots, x_n)|$$

$$m_H(n) \leq 2^n$$

\uparrow
Hypothesis size bound

① if $m_H(n) = 2^n$, H shatters the points x_1, \dots, x_n

Ex

3 points x_1, x_2, x_3

$$\begin{aligned} H_1 &= \{h_1, h_2, h_3\} \\ H_2 &= \{h_1, h_2, h_3, h_4\} \\ H_3 &= \{h_1, h_2, h_3, h_4, h_5\} \end{aligned}$$

\uparrow
 h_1, h_2, h_3

it shatters the

but the set

x_1	x_2	x_3
0	0	0
0	0	1
0	1	0
0	1	1
1	0	0
1	0	1
1	1	0
1	1	1

8

cannot shatter

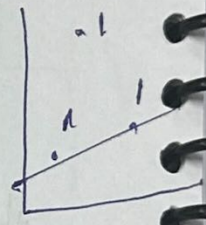
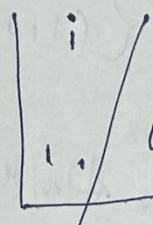
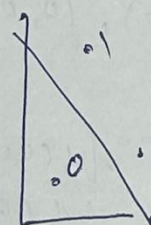
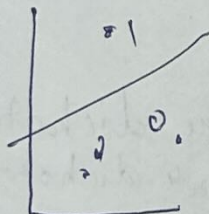
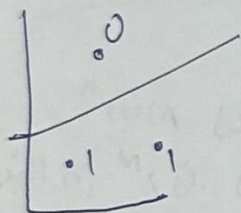
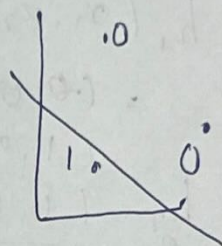
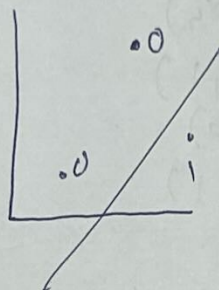
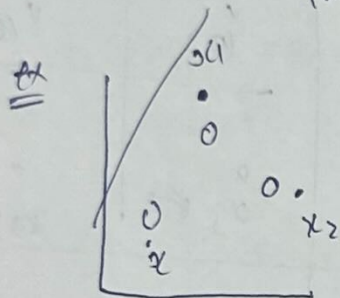
size it cannot all classify the

can classify the all class points.

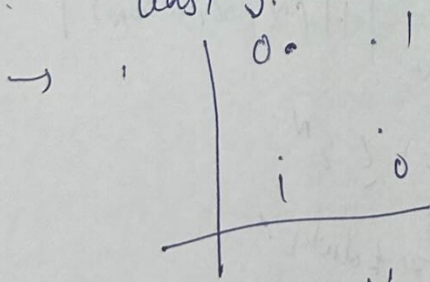
(*) VC dimension: VC-dimension of a hypo. set H denoted by $d_{VC}(H)$ or d_{VC} is largest value of N for which growth funⁿ $m_H(N) = 2^N$.

⊖ if $m_H(N) = 2^N$ for $\forall N$, $\rightarrow d_{VC}(H) = \infty$

\rightarrow if $m_H(N) = 2^N \rightarrow$ then N is VC dimⁿ of that Hypothesis.



\hookrightarrow we can classify all point so. $\forall d_{VC} = \infty$ at least 3.



$\nexists d_{VC} = 4$ not possible

\rightarrow for 2D point VC dimension is 3.

(*) If VC-dimension $(H) = n$ then H shatters every set of size n & there exists a set of size $(n+1)$ (Radon's theorem) which can not shatter.

\hookrightarrow point are d d dim. then we can shatters $d+1$ dimension.

\rightarrow sine waves $\rightarrow VC \dim(H) = \infty$

$\rightarrow NN \rightarrow \text{comb } VC \dim(H) = 0$

* Important Result (Sauer-Shelah lemma theorem)
 then \exists : If $m_H(k) < 2^{\frac{k-1}{2}}$ for some value k then
 $m_H(k) \leq \sum_{i=1}^{k-1} \binom{N}{i}$ \rightarrow can't shutter b'coz of strictly less than

ex $k=4, N=3$
 $m_H(3) \leq \binom{3}{0} + \binom{3}{1} + \binom{3}{2} + \binom{3}{3}$

* by definition VC-dimension
 $m_H(N) \leq \sum_{i=0}^{d_{VC}} \binom{N}{i}$

ex 5 points - 2 dimen. $d_{VC} = 2+1 = 3$

$m_H(5) \leq \binom{5}{0} + \binom{5}{1} + \binom{5}{2} + \binom{5}{3}$

\rightarrow If $d_{VC}(H) = n$ \exists a set of size n which H can shutter and H cannot shutter any set of size $(n+1)$

* VC-generalization: VC-generalization bound for any tolerance $\delta > 0$

$E_{out}(g) \leq E_{in}(g) + \sqrt{\frac{\delta}{N} \ln \frac{4m_H(2^N)}{\delta}} \rightarrow 2^{d_{VC}}$

w.p. $1-\delta$
 with probabily.

\downarrow Now this is limit from we can apply Union bound.

⊙ if d_{VC} is increased then more complex the model, more complex hypothesis space than our model is overfit. test error is very high.