

HW 4

↳ Hardik Prajapati (2678294168)

$$N = 4601$$

$$N_{T_r} = 3000$$

$$N_{Test} = 1601$$

$$\delta (\text{tolerance}) = 0.1$$

Given:- If H is a linear perceptron classifier in D Dimensions
 $d_{vc}(H) = D+1$

Here, $D = 57$ features

a) $d_{vc}(H) = D+1 = 58$

b) $E_{out}(h_g) \leq E_{in}(h_g) + \varepsilon_{vc}$

$E_{in}(h_g) \rightarrow$ Training Data Set

$$\varepsilon_{vc} = \sqrt{\frac{8}{N} \ln 4 \left[\frac{(2N)^{d_{vc}} + 1}{\delta} \right]}$$

$$= \sqrt{\frac{8}{3000} \ln 4 \left[\frac{(2(3000))^{58} + 1}{(0.1)} \right]}$$

$$= 1.1642$$

$$c) D = 10$$

$$N_{T_r} = 10,000$$

$$\Rightarrow d_{vc} = 11$$

$$\varepsilon_{vc} = \sqrt{\frac{8}{(10,000)} \ln 4 \left[\frac{(2(10,000))^{d_{vc}} + 1}{0.1} \right]}$$

$$= 0.3001$$

$$d) \quad \varepsilon_{vc} = 0.1$$

$$\delta = 0.1$$

$$D = 10$$

$$d_{vc} = 11$$

$$\varepsilon_{vc} = \sqrt{\frac{8}{N_{T_r}} \ln 4 \left[\frac{(2N_{T_r})^{d_{vc}} + 1}{\delta} \right]}$$

$$\Rightarrow 0.1 = \sqrt{\frac{8}{N_{T_r}} \ln 4 \left[\frac{(2N_{T_r})^{d_{vc}} + 1}{0.1} \right]}$$

$$\Rightarrow N_{T_r} = 110000$$

$$e) E_{in}(h_g) \rightarrow E_{test}(h_g)$$

$H = \text{Hypothesis set} = \{h_g\}$

Cardinality = 1

$$f) E_{\text{out}}(h_g) \leq E_{\text{test}}(h_g) + \varepsilon$$

$$D = 57$$

$$N_{\text{test}} = 1601$$

$$\underline{\varepsilon_n = \sqrt{\frac{1}{2N} \ln \frac{2M}{\delta}}}$$

$$\Rightarrow \underline{\varepsilon_n = \sqrt{\frac{1}{2(1601)} \ln \frac{2(1)}{0.1}}}$$

$$= 0.0305$$

$$2) X = \{x\} \times R^d$$

To show:- $d_{VC} \geq d+1$ & $d_{VC} \leq d+1 \Rightarrow d_{VC} = d+1$

$$a) x_0 = [0, 0, \dots, 0]_{1 \times d}^T$$

$$x_1 = [1, 0, \dots, 0]_{1 \times d}^T$$

$$x_2 = [0, 1, 0, \dots, 0]_{1 \times d}^T$$

$$x_d = [0, 0, \dots, 1]_{1 \times d}^T$$

} $d+1$ points

Linear Perception Model :- $\text{sgn} [\omega^T x + \omega_0] \quad \{+1, -1\}$

Augmented column of 1

Now,

$$\begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ 1 & 1 & 0 & \cdots & 0 \\ 1 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & 0 & \cdots & 1 \end{bmatrix} \quad (d+1) \times (d+1)$$

$$\text{Let } \omega' = [\omega_0, \omega_1, \omega_2, \dots, \omega_d]_{(d+1) \times 1}$$

$$h(\underline{x}) = \underline{x} \omega'$$

$$\begin{bmatrix} h(x_0) \\ h(x_1) \\ \vdots \\ h(x_d) \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ 1 & 1 & 0 & \cdots & 0 \\ 1 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & 0 & \cdots & 1 \end{bmatrix}_{(d+1) \times (d+1)} \begin{bmatrix} \omega_0 \\ \omega_1 \\ \vdots \\ \omega_d \end{bmatrix}_{(d+1) \times 1}$$

$$\Rightarrow h(\underline{x}) = \begin{bmatrix} \omega_0 \\ \omega_0 + \omega_1 \\ \omega_0 + \omega_1 + \omega_2 \\ \vdots \\ \omega_0 + \omega_d \end{bmatrix}$$

$$\text{Let, } \omega_0 = 0.5 h(x_0)$$

$$\Rightarrow \begin{bmatrix} h(x_0) \\ h(x_1) \\ \vdots \\ h(x_d) \end{bmatrix} = \text{sgn.} \begin{bmatrix} 0.5^* h(x_0) \\ 0.5^* h(x_1) + w_1 \\ \vdots \\ 0.5^* h(x_d) + w_d \end{bmatrix}$$

Now, if $w_i = h(x_i)$ $\forall i \in \{1, \dots, d\}$

Then, Linear perception can classify all these $d+1$ points correctly for their different combinations.

\Rightarrow The perception model can atleast shatter all the $d+1$ points.

$$\Rightarrow d_{VC} \geq d+1$$

b) Now, Let $x_{d+1} = [1, 0, 0, \dots, 1]_{1 \times d}^T$
 \uparrow
 Linear combination of x_1, x_2, \dots, x_d .

$$\Rightarrow X = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 1 & 1 & 0 & \dots & 0 \\ 1 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & 0 & \dots & 1 \\ 1 & 1 & 0 & \dots & 1 \end{bmatrix}_{(d+2) \times (d+1)}$$

Let $w_0 = 0.5 * h(x_0)$

$$\Rightarrow \begin{bmatrix} h(x_0) \\ h(x_1) \\ \vdots \\ h(x_d) \\ h(x_{d+1}) \end{bmatrix} = \text{Sgn} \begin{bmatrix} 0.5^* h(x_0) \\ 0.5^* h(x_1) + w_1 \\ \vdots \\ 0.5^* h(x_d) + w_d \\ 0.5^* h(x_0) + w_1 + w_d \end{bmatrix}$$

Again, $w_i = h(x_i) \quad \forall i \in \{1, \dots, d\}$

Now,

i) Say, $h(x_0) = 1, h(x_1) = 1, h(x_d) = 1$
 Assume $h(x_{d+1}) = 1$
 $\Rightarrow h(x_{d+1}) = \text{sgn}[(0.5)(1) + 1 + 1]$

$$= 1$$

\Rightarrow Correctly classified.

ii) Say, $h(x_0) = 1, h(x_1) = 1, h(x_d) = 1$

Assume $h(x_{d+1}) = -1$

$$\Rightarrow h(x_{d+1}) = \text{sgn}[(0.5)(1) + 1 + 1]$$

$$= 1$$

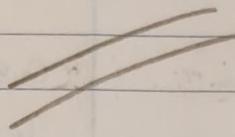
\Rightarrow Mis-classified.

\Rightarrow The linear perceptron cannot shatter $d+2$ points

$$\Rightarrow d_{VC} \leq d+1.$$

Ans:- From part (a) & part (b), it can be concluded that VC Dimension of Linear perception (with $d+1$ parameters, counting w_0) is exactly equal to $d+1$.

$$\text{Hence, } d_{VC} = d+1$$



3. AML Pb 2.24

$$x \in \overline{[-1, 1]}$$

$$D = \{\{x_1, x_1^2\}, \{x_2, x_2^2\}\}$$

$$\text{Hypothesis set: } \{h(x) = ax + b \mid a \in \mathbb{R}, b \in \mathbb{R}\}$$

a) $g^{(D)}(x) = \text{line fitting points } \{\{x_1, x_1^2\}, \{x_2, x_2^2\}\}$

$$\text{Line: } x_1^2 = ax_1 + b$$

$$x_2^2 = ax_2 + b$$

$$\underline{x_1^2 - x_2^2 = a(x_1 - x_2)}$$

$$\Rightarrow a = \frac{x_1^2 - x_2^2}{x_1 - x_2}$$

$$\Rightarrow a = x_1 + x_2, \text{ if } x_1 \neq x_2$$

$$\Rightarrow b = x_1^2 - \{x_1 + x_2\} x_1$$

$$\Rightarrow b = x_1^2 - x_1^2 - x_1 x_2$$

$$\Rightarrow b = -x_1 x_2$$

$$\Rightarrow g^{(D)}(x) = (x_1 + x_2)x - x_1 x_2$$

~~g(x)~~

a.2) $\bar{g}(x) = E_D[(x_1 + x_2)x - x_1 x_2]$

$$= E_D[x_1 x] + E_D[x_2 x] - E_D[x_1 x_2]$$

↳ Now, Here 'x' is not Random Variable.

↳ Also, x_1 & x_2 are 2 independent Random variables, both having uniform Distribution [-1, 1]

$$\text{Hence, } E_D[x_1] = E_D[x_2] = 0$$

$$\Rightarrow \bar{g}(x) = x E_D[x_1] + x E_D[x_2] - E_D[x_1] E_D[x_2]$$

$$\Rightarrow \bar{g}(x) = 0$$

~~g(x)~~

- b) \hookrightarrow Fix \underline{x}
- ↳ Sample (Random) 2 points from uniform distribution
[-1, 1]
- Repeat for 1000 times {
- \hookrightarrow Fit $y = (x_1 + x_2)x + x_1 \epsilon_1$ & compute $g^D(\underline{x})$
- \hookrightarrow Compute average of $g^D(\underline{x})$, which will equal to $\bar{g}(\underline{x})$

Now, for out of sample error

- Repeat for 5,000 times {
- \hookrightarrow Random Sample \underline{x} from Unif [-1, 1]
- \hookrightarrow Compute $g^D(\underline{x})$ as shown above.
- \hookrightarrow We will have array for each \underline{x} . Take average of it, $\bar{g}(\underline{x})$
- \hookrightarrow Compute $E_D [(\bar{g}(\underline{x}) - g^D(\underline{x}))^2]$
- \hookrightarrow Compute $E_D [(\bar{g}(\underline{x}) - f(\underline{x}))^2]$
- \hookrightarrow Compute $E_D [(g^D(\underline{x}) - f(\underline{x}))^2]$

\hookrightarrow Take average & compute

$$E_x [E_D [(\bar{g}(\underline{x}) - g^D(\underline{x}))^2]] \rightarrow \text{Variance}$$

$$E_x [(\bar{g}(\underline{x}) - f(\underline{x}))^2] \rightarrow \text{Bias}$$

$$E_x [E_D [(g^D(\underline{x}) - f(\underline{x}))^2]] \rightarrow \text{Out of Sample Error}$$

c) \rightarrow Simulation

d) From (a), $\bar{g}(\underline{x}) = 0$

$$\hookrightarrow \text{Var} = E_x \left[E_D \left[(\bar{g}(\underline{x}) - g^D(\underline{x}))^2 \right] \right]$$

$$= E_x \left[E_D \left[(0 - (x_1 + x_2) \underline{x} + x_1 x_2)^2 \right] \right]$$

$$= E_x \left[E_D \left[(x_1 + x_2)^2 \underline{x}^2 + x_1^2 x_2^2 - 2x_1 x_2 (x_1 + x_2) \underline{x} \right] \right]$$

$$= E_D \left[E_x \left[(x_1 + x_2)^2 \underline{x}^2 + x_1^2 x_2^2 - 2x_1 x_2 (x_1 + x_2) \underline{x} \right] \right]$$

$$= E_D \left[(x_1 + x_2)^2 E_x[\underline{x}^2] + x_1^2 x_2^2 E_x[1] - 2x_1 x_2 (x_1 + x_2) E_x[\underline{x}] \right]$$

Now,

$$E_x[\underline{x}] = 0, E_x[\underline{x}^2] = \frac{1}{3}$$

$$E_D[x_1] = E_D[x_2] = 0, E_D[x_1^2] = E_D[x_2^2] = \frac{1}{3}$$

$$\Rightarrow \text{Var} = E_D \left[(x_1 + x_2)^2 \left(\frac{1}{3} \right) + x_1^2 x_2^2 - 0 \right]$$

$$= \frac{1}{3} E_D[x_1^2 + x_2^2 + 2x_1 x_2] + E_D[x_1^2 x_2^2]$$

$$= \frac{1}{3} \{ E_D[x_1^2] + E_D[x_2^2] + 2E_D[x_1] E_D[x_2] \} + E_D[x_1^2] E_D[x_2^2]$$

$$\begin{aligned}
 \Rightarrow \text{Var} &= \frac{1}{3} \left\{ \frac{1}{3} + \frac{1}{3} + 0 \right\} + \left\{ \frac{1}{3} \right\} \left\{ \frac{1}{3} \right\} \\
 &= \frac{2}{9} + \frac{1}{9} \\
 &= \frac{3}{9} \\
 &= \underline{\underline{\frac{1}{3}}}
 \end{aligned}$$

$$\begin{aligned}
 \hookrightarrow \text{Bias} &= E_x \left[(\bar{g}(\underline{x}) - f(\underline{x}))^2 \right] \\
 &= E_x \left[(0 - \underline{x}^2)^2 \right] \\
 &= E_x \left[\underline{x}^4 \right] \\
 &= \underline{\underline{\frac{1}{5}}}
 \end{aligned}$$

$$\begin{aligned}
 \hookrightarrow E_D \{ E_{\text{out}} \} &= E_x \left[E_D \left[(g^{(0)}(\underline{x}) - f(\underline{x}))^2 \right] \right] \\
 &= E_x \left[E_D \left[((x_1 + x_2)\underline{x} - x_1 x_2 - \underline{x}^2)^2 \right] \right]
 \end{aligned}$$

$$\Rightarrow E_D\{E_{\text{out}}\} = E_x\left[E_D\left[\{(x_1+x_2)\underline{x} - x_1x_2\}^2 + \underline{x}^4 - 2\underline{x}\{(x_1+x_2)\underline{x} - x_1x_2\}\right]\right]$$

$$= E_x\left[\underbrace{E_D\left[\{(x_1+x_2)\underline{x} - x_1x_2\}^2\right]}_{\frac{1}{3} = (\text{Var}_m)} + E_D[\underline{x}^4] - 2\underline{x} E_D[(x_1+x_2)\underline{x} - x_1x_2]\right]$$

$$= E_x\left[\frac{1}{3}\right] + \underbrace{E_x[\underline{x}^4]}_{\frac{1}{5} = \text{Bias}} + 2 E_D[(x_1+x_2)\underline{x} - x_1x_2] \underbrace{E_x[\underline{x}]}_0$$

$$\Rightarrow E_D\{E_{\text{out}}\} = \frac{1}{3} + \frac{1}{5}$$

$$= \frac{8}{15}$$

//

$$\hookrightarrow \text{Var}(\underline{x}) = E_D\left[\left(g^{(0)}(\underline{x}) - \bar{g}(\underline{x})\right)^2\right]$$

$$= E_D\left[\left((x_1+x_2)\underline{x} + x_1x_2\right)^2\right]$$

$$= E_D\left[\left(x_1+x_2\right)\underline{x}^2 + x_1^2x_2^2 - 2x_1x_2(x_1+x_2)\underline{x}\right]$$

$$\begin{aligned}
 \Rightarrow \text{Var}(\underline{x}) &= \underline{x}^2 E_D[(x_1 + x_2)^2] + E_D[x_1^2]E_D[x_2^2] \\
 &\quad - 2 \underline{x} E_D[x_1^2 x_2 + x_2^2 x_1] \\
 &= \underline{x}^2 \left\{ E_D[x_1^2] + E_D[x_2^2] + 2 E_D[x_1] E_D[x_2] \right\} + \binom{4}{3} \binom{1}{3} \\
 &\quad - 2 \underline{x} \left\{ E_D[x_1^2] E_D[x_2] + E_D[x_2^2] E_D[x_1] \right\} \\
 &= \underline{x}^2 \left\{ \frac{1}{3} + \frac{1}{3} + 0 \right\} + \frac{1}{9} - 2 \underline{x} \left\{ 0 + 0 \right\} \\
 &= \underline{x}^2 \left\{ \frac{2}{3} \right\} + \frac{1}{9}
 \end{aligned}$$

$$\hookrightarrow \text{Bias}(\underline{x}) = (\bar{g}(\underline{x}) - f(\underline{x}))^2$$

$$= (0 - \underline{x}^2)^2$$

$$= \underline{x}^4$$

4 AML Pb 2.13 (a), (b)

(a) $H = \{h_1, h_2, \dots, h_M\}$

To prove $d_{vc}(H) \leq \log_2 M$

Now, $m_H(n) \leq n^{d_{vc}} + 1$

For 'd_{vc}' no. of points, growth function is

$$m_H(d_{vc}) = 2^{d_{vc}} \quad \text{---(i)}$$

Let us assume that $d_{vc}(H) > \log_2 M$

$$\Rightarrow 2^{d_{vc}(H)} > M \quad \text{---(ii)}$$

\Rightarrow From (i) & (ii),

$$m_H(d_{vc}) = 2^{d_{vc}} > M$$

\Rightarrow No. of Hypothesis that shatter 'd_{vc}' no. of points is larger than M.

\Rightarrow But, $H = \{h_1, h_2, \dots, h_M\}$ has only finite 'M' hypothesis.

\Rightarrow Contradiction to our Assumption.

$$\Rightarrow d_{vc}(H) \leq \log_2 M$$

Hence Proved

(b) Hypothesis Sets :- H_1, H_2, \dots, H_k

finite VC Dim :- $d_{VC}(H_k)$

Derive & Prove - tightest upper bound & lower bound on
 $d_{VC}(\bigcap_{k=1}^K H_k)$

Now, $d_{VC}(\bigcap_{k=1}^K H_k) \geq 0$ -(i)

As it can't be negative integer. Also, it can be zero as the intersection can be an empty \emptyset set.

By intuition,

Let us assume that

$$d_{VC}(\bigcap_{k=1}^K H_k) > \min_k d_{VC}(H_k)$$

\Rightarrow There are $a = d_{VC}(\bigcap_{k=1}^K H_k)$ # of points that intersection can shatter.

\Rightarrow Since, it's an intersection, the respective Hypothesis which was able to shatter must be present in all $H_k (k=1 \dots K)$

$\Rightarrow d_{VC}(\bigcap_{k=1}^K H_k) > \min_k d_{VC}(H_k)$ is a

contradiction.

\Rightarrow Hence, $d_{VC}(\bigcap_{k=1}^K H_k) \leq \min_k d_{VC}(H_k)$

Ans:-

$$0 \leq d_{VC}(\bigcap_{k=1}^K H_k) \leq \min_k d_{VC}(H_k) //$$

5. Ans Pb. n.h (a)-(c)

$$x = [-1, 1] \quad \text{Uniform} \quad p(x) = \frac{1}{2}$$

$$f(x) = \sum_{q=0}^{\infty} a_q L_q(x) \quad , \quad L_q(x) \sim \text{Legendre polynomials}$$

$$D = (x_1, y_1); \dots; (x_n, y_n) \quad , \quad y_n = f(x_n) + \sigma \varepsilon_n$$

y_n & ε_n ~ i.i.d Normal R.V

Legendre polynomials:-

$$L_0(x) = 1$$

$$L_1(x) = x$$

$$L_k(x) = \frac{2k-1}{k} x L_{k-1}(x) - \frac{k-1}{k} L_{k-2}(x)$$

Specified Values :- Q, N, σ

a) We normalize $f \Rightarrow E_{a,x} [f(x)] = 0$

$$\text{Var}[f(x)] = E_{a,x} [f^2(x)] = 1$$

$$\Rightarrow E_{a,x} [y_n] = E_{a,x} [f(x_n) + \sigma \varepsilon_n]$$

$$= \underbrace{E_{a,x} [f(x_n)]}_0 + \underbrace{\sigma E_{a,x} [\varepsilon_n]}_0$$

$$\Rightarrow E_{\alpha, \sigma} [y_n] = 0$$

Also,

$$\text{Var}[y_n] = \text{Var}[f(x) + \sigma \varepsilon_n]$$

$$= \text{Var}[f(x)] + \sigma^2 \text{Var}[\varepsilon_n]$$

$$+ 2\sigma \text{Cov}(f(x), \varepsilon_n)$$

$$= 1 + \sigma^2 + 2\sigma \text{Cov}(f(x), \varepsilon_n)$$

if $f(x)$ & ε_n are not correlated.

$$\Rightarrow \text{Var}[y_n] = 1 + \sigma^2 + 0$$

\Rightarrow for Different Order of Polynomial, comparison becomes easy bet. $f(x)$ & noise term.

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$$

$g_2 \rightarrow$ Best fit Hypothesis from H_2 [Second order polynomial]

$g_{10} \rightarrow$ Best fit Hypothesis from H_{10} [10^m Order polynomial]

\Rightarrow Comparing with Linear Regression model,

$$Y = \omega X \rightarrow \omega^* = (X^T X)^{-1} X^T Y$$

Here, $X = \Sigma$

$$f(\underline{x}) = \sum_{q=0}^{\infty} a_q L_q(x) = \Phi_q(\underline{x})$$

$$\Rightarrow \text{Input Matrix} = \begin{bmatrix} \phi_a(x_1) \\ \phi_a(x_2) \\ \vdots \\ \phi_a(x_n) \end{bmatrix} \quad \leftarrow \text{Let it be represented by } \phi_a \quad (N) \times (n+1)$$

$$\Rightarrow w^* = (\phi_a^\top \phi_a)^{-1} \phi_a^\top y \quad \leftarrow \text{By Least Squares Best fit}$$

$$\Rightarrow g_a(x) = \phi_a(x) (\phi_a^\top \phi_a)^{-1} \phi_a^\top y$$

\Leftrightarrow For a single experiment,

$$E_{\text{out}} = E_{x,y} [(g_{10}(x) - y(x))^2]$$

$$= E_{x,y} [g_{10}^2(x) + y^2(x) - 2 g_{10}(x) y(x)]$$

$$= E_x [g_{10}^2(x)] + E_{x,y} [y^2(x)] - 2 \underbrace{E_{x,y} [g(x)] E_{x,y} [y(x)]}_{\text{Var}(y)}$$

$$= E_x [g_{10}^2(x)] + (1 + \sigma^2) - 0$$

$$= E_x [g_{10}^2(x)] + (1 + \sigma^2)$$

$\rightarrow \begin{cases} g_1, f(x) & \text{are not} \\ \epsilon_n & \text{correlated} \end{cases}$

5 (i) Fig 4.3(a), $\sigma^2 = 0.5$, $N = 60 \sim 130$.

With fixed Noise level, $\sigma^2 = 0.5$, as we increase Data Points, $N = 60$ to 130 , overfitting decreases.

As the signal f is normalized, noise level is calibrated to signal & noise level σ^2 plays a major factor.

With smaller N , & more complex, larger model, noise is more susceptible to deviate signal from true function.

For a small Dataset, the model is less capable to distinguish noise from signal & hence tries to fit the noise, leading to overfit.

ii) Fig 4.3(a), $N=100$, $\sigma^2 = 0$ to $\sigma^2 = 2$.

With increase in Noise, Overfitting increases.

As f is normalized i.e. $\text{Var}(f) = 1$ noise level is calibrated to signal & has major influence.

More large, complex models (H_1 compared to H_2) are more susceptible to noise as, it has more ways to deviate from true function. Hence, increase in overfitting.

(iii) Fig 4.3(b), $N = 75$, $Q_f = 0$ to $Q_f = 100$

With higher target function complexity Q_f , it will

be more susceptible to noise as it has more ways to deviate from true signal even though there is non-linearity.

Hence with increase in $Q_f (> 15)$, overfitting increases.

With $N = 75$, with lower range of target complexity, overfitting decreases ($Q_f = 0$ to $Q_f = 15$)

With $N = 75$, & $Q_f > 15$, it overfits the data compared to H_2 .

With $Q_f = 0$ to $Q_f = 15$, it seems like the signal was able to distinguish between noise & true function due to non-linearity of $Q_f(x)$. Hence overfitting decreases from $Q_f = 0$ to $Q_f = 15$.