

Ex 2.1

a) $\epsilon = \sqrt{\frac{1}{2N} \ln \frac{2m}{\delta}}$

$\therefore \epsilon^2 = \frac{1}{2N} \ln \frac{2m}{\delta}$

$\therefore N = \frac{1}{2\epsilon^2} \ln \frac{2m}{\delta}$

$\Rightarrow \cancel{F_{\alpha}} \text{ & } \epsilon \leq 0.05$

$\Rightarrow N \geq \frac{1}{2\epsilon^2} \frac{2m}{\delta}$

$m=1, \delta=0.03 \text{ & } \epsilon \leq 0.05$

$\therefore N \geq \frac{1}{2 \times (0.05)^2} \frac{\ln 2 \times 1}{0.03}$

= ~~839.941~~

b) $m=100, \delta=0.03 \text{ & } \epsilon \leq 0.05$

$\therefore N \geq \frac{1}{2 \times (0.05)^2} \frac{\ln 2 \times 100}{0.03}$

$\Rightarrow N \geq 1760.97$

c) $m=1000, \delta=0.03 \text{ & } \epsilon \leq 0.05$

$\therefore N \geq \frac{1}{2 \times (0.05)^2} \frac{\ln 2 \times 1000}{0.03}$

$\therefore N \geq 2682.009$

Problem 2.3

a) for positive rays the growth function is

$$\Rightarrow N+1$$

The dichotomies for negative rays

$$\Rightarrow \cancel{N-1}$$

The dichotomies for positive rays

$$\Rightarrow N+1$$

Adding the two

Positive dichotomies + Negative dichotomies

$$\Rightarrow N+1 + N-1$$

$$\Rightarrow 2N$$

$$\therefore m_H(N) = 2N$$

The largest values of all shatter is

given by by 2^N

$$N=1.$$

$$\therefore m_H(1) = 2 \quad | \quad 2^{(1)} = 2.$$

$$m_H(2) = 2 \times 2 = 4 \quad | \quad 2^{(2)} = 4$$

$$m_H(3) = 2 \times 3 = 6 \quad | \quad 2^{(3)} = 8.$$

i. The largest value for which $m_H(N) = 2^N$

$$\therefore dvc = 2$$

b) Growth function for positive intervals is given by $\frac{N^2}{2} + \frac{N}{2} + 1$

~~For~~ For negative intervals we will add only $N - 2$ new dichotomies.

~~f₂~~ ~~N=3~~ we add (+1, -1, +1)

& ~~f₂~~ N=4 we add (+1, -1, +1, +1)
& (+1, +1, -1, +1)

\therefore Total dichotomies = $\frac{N^2}{2} + \frac{N}{2} + 1 + N - 2$.

$$= \frac{N^2}{2} + \frac{3N}{2} - 1 \quad \text{for } N \geq 1$$

\therefore Computing the largest value for which $m_N(N) = 2^N$

~~m₂(2)~~

We will not use ~~N=1~~ as for 1 we already generated two dichotomies.

$$\therefore N = 2.$$

$$m_2(2) = \frac{2^2}{2} + \frac{2 \times 3}{2} - 1 \quad | \quad 2^2 = 4 \\ = 4$$

$$m_3(3) = \frac{3^2}{2} + \frac{3 \times 3}{2} - 1 \quad | \quad 2^3 = 8$$

$$= 8$$

$$m_4(4) = \frac{4^2}{2} + \frac{4 \times 3}{2} - 1 = 13 \quad | \quad 2^4 = 16.$$

$$\therefore dvc = 3 \quad (\text{for } m_4(4) = 2^4)$$

c) To find $m_N(N)$ we map \mathbb{R}^d onto $[0, +\infty]$

$$\therefore \phi : (x_1, \dots, x_d) \mapsto r = \sqrt{x_1^2 + \dots + x_d^2}$$

\therefore this problem reduces to positive interval problem.

\therefore we can conclude that:

$$m_N(N) = \frac{N^2}{2} + \frac{N}{2} + 1$$

\therefore From b) we have $d_{NC} = 323$

Problem 2.8

when $m_N(N) = 2^N$ then $dvc = +\infty$

if $m_N \neq 2^N$ then dvc is a finite value

and if $m(N)$ is bounded by dvc using $\underline{N^{dvc+1}}$

$$\rightarrow \therefore m(N) = 1+N$$

$$\begin{array}{l} m(1) = 2 \\ m(2) = 3 \end{array} \quad \left| \begin{array}{l} 2^1 = 2 \\ 2^2 = 4 \end{array} \right.$$

$$\therefore dvc = 1$$

so it is bounded by $N^{(1)}+1 = N+1$ for all N .

$\therefore m(N) = N+1$ is a possible growth function.

$$\rightarrow m(N) = 1+N+N(N-1)/2$$

$$m(1) = 2 \quad \left| \begin{array}{l} 2^1 = 2 \end{array} \right.$$

$$m(2) = 4 \quad \left| \begin{array}{l} 2^2 = 4 \end{array} \right.$$

$$m(3) = 7 \quad \left| \begin{array}{l} 2^3 = 8 \end{array} \right.$$

$$\therefore dvc = 2.$$

$\therefore m(N)$ must be bounded by N^2+1 for all N .

$\therefore m(N) = 1+N+N(N-1)/2$ is a ~~possible~~ possible growth function.

$$\rightarrow m_N(N) = \sum_{k=1}^{\sqrt{N}}$$

$$m(1) = 2 \quad \left| \begin{array}{l} 2^1 = 2 \end{array} \right.$$

$$m(2) = 2^{1+1} \quad \left| \begin{array}{l} 2^2 = 4 \end{array} \right.$$

$$\approx 2$$

$$\therefore dvc = 1$$

i. It must be bounded by $N+1$

Computing for squares.

$$\left. \begin{array}{l} N=9 \\ 2^{\sqrt{9}} = 2^3 = 8 \end{array} \right\} \begin{array}{l} N+1 \\ 9+1 = 10 \end{array}$$
$$\left. \begin{array}{l} 2^{\sqrt{16}} = 2^4 = 16 \\ 16+1 = 17 \end{array} \right\} \begin{array}{l} 16+1 = 17 \\ 16+1 = 17 \end{array}$$

$$2^{\sqrt{25}} = 2^5 = 32 \quad | \quad 25+1 = 26$$

Hence $m_H(N) = 2^{\sqrt{N}}$ is not a possible growth function.

$$\rightarrow m_H(N) = 2^{N/2}$$

$$m_H(1) = 2^{1/2} \quad | \quad 2^1 = 2$$

$$\therefore dvc = 0.$$

\therefore It must be bounded by $N^0 + 1 = 1 + 1 = 2$.
for all N .

Computing for different ~~N~~ $N=1, 2, 4$

$$N=2 \quad | \quad \begin{array}{l} dvc \\ N+1 \end{array}$$

$$m_H(2) = 2^{4/2} = 2 \quad | \quad \begin{array}{l} 2 \\ 2 \end{array}$$

$$m_H(4) = 2^{4/2} = 4 \quad | \quad \begin{array}{l} 2 \\ 2 \end{array} \leftarrow \text{This} \quad \begin{array}{l} \text{violates bound} \\ \text{condition} \end{array}$$

$\therefore m_H(N) = 2^{N/2}$ is not a possible growth function.

Problem 2.12

Sample complexity is given by.

$$N \geq \frac{8}{\epsilon^2} \ln \left(\frac{4m_H(2N)}{\delta} \right)$$

$$N \geq \frac{8}{\epsilon^2} \ln \left(\frac{4(2N)^{d_{VC}} + 1}{\delta} \right)$$

$$d_{VC} = 10, \quad \epsilon = 0.05, \quad \delta = 0.05$$

$$\therefore N \geq \frac{8}{0.05^2} \ln \left(\frac{4[2 \times \cancel{N}]^{10} + 1}{0.05} \right)$$

Trying an initial guess of $N = 1000$

we get

$$N \geq \frac{8}{0.05^2} \ln \left(\frac{4[2 \times 1000]^{10} + 1}{0.05} \right)$$

$$\approx 2.57251 \times 10^5$$

Now we try $N = 2.57251 \times 10^5$ in above eqⁿ.

we get $N = 4.5957 \times 10^5$. which is rapidly converging.

Problem 2.22

$$\begin{aligned}
 E_p [E_{\text{out}}(g^{(D)})] &= E_D [E_{u,y} [g^{(D)}(u) \\
 &\quad - y(u))^2]] \\
 &= E_{u,y} [E_D [(g^{(D)}_u - y(u))^2]] \\
 &= E_{u,y} [E_D [g^{(D)}(u)^2] - 2\bar{g}(u)y(u) + y(u)^2]
 \end{aligned}$$

where using Fubini's theorem.

$$\begin{aligned}
 &E_D [g^{(D)}(u)^2] - 2\bar{g}(u)y(u) + y(u)^2 \\
 &= E_D [g^{(D)}(u)^2] - \frac{1}{2} 2\bar{g}(u)(f(u) + \epsilon) \\
 &\quad + (f(u) + \epsilon)^2 \\
 &= (E_D [g^{(D)}(u)^2] - \bar{g}(u)^2) + (\bar{g}(u)^2 - 2\bar{g}(u) \\
 &\quad + f(u)^2) + \epsilon^2 - 2(\bar{g}(u) - f(u))\epsilon \\
 &= (E_D [g^{(D)}(u)^2] - 2E_D [g^{(D)}(u)]\bar{g}(u) \\
 &\quad + \bar{g}(u)^2) + (\bar{g}(u) - f(u))^2 + \epsilon^2 \\
 &\quad - 2(\bar{g}(u) - f(u))\epsilon \\
 &= E_D [g^{(D)}(u)^2 - 2g^{(D)}(u)\bar{g}(u) + \bar{g}(u)^2] \\
 &\quad + (\bar{g}(u) - f(u))^2 + \epsilon^2 - 2(\bar{g}(u) - f(u))\epsilon \\
 &= \text{var}(u) + \text{bias}(u) + \epsilon^2 - 2(\bar{g}(u) - f(u))\epsilon
 \end{aligned}$$

Taking expectation relative to (x, y)

$$E_{\theta} [\text{Err}_{\text{out}}(g^{(D)})] = E_{x,y} [\text{var}(u)] + E_{x,y} [\text{bias}(u)] \\ + E_{x,y} [\epsilon^2] - 2 E_{x,y} [(\bar{g}(x) - f(x))\epsilon]$$

$$= E_x [\text{var}(u)] + E_u [\text{bias}(u)] + E_x [E_{\epsilon} [\epsilon^2 | x]]$$

$$- 2 E_{x,y} [(\bar{g}(x) - f(x))\epsilon]$$

$$= \text{var} + \text{bias} + E_x [E_{\epsilon} [\epsilon^2]] \quad \cancel{- 2 E_{x,y} [(\bar{g}(x) - f(x))\epsilon]}$$

$$- 2 E_x [E_{\bar{g}(x)} [(\bar{g}(x) - f(x))\epsilon | x]]$$

$$= \cancel{\text{var} + bias + var_{\epsilon}[\epsilon]} \quad \cancel{2 E_x [(\bar{g}(x) - f(x))\epsilon]} \\ \cancel{- 2 E_x [E_{\bar{g}(x)} [(\bar{g}(x) - f(x))\epsilon | x]]}$$

$$= \text{var} + \text{bias} \quad \text{var}_{\epsilon}[\epsilon] - 2 E_x [(\bar{g}(x) - f(x)) E_{\epsilon}[\epsilon]]$$

$$= \text{var} + \text{bias} + \epsilon^2 \quad (\text{since } E_{\epsilon}[\epsilon] = 0)$$

~~Quest 6~~ Prove that selecting the hypothesis h that maximizes the likelihood $\prod_{n=1}^N P(y_n|x_n)$ is equivalent to minimizing the cross-entropy error.

Proof - The term likelihood is that we will get output y from input x given the target distribution $P(y|x)$ was indeed captured by hypothesis $h(x)$

$$\therefore P(y|x) = \begin{cases} h(x) & \text{for } y = +1 \\ 1 - h(x) & \text{for } y = -1 \end{cases}$$

Using $h(x) = \sigma(w^T x)$, we have.

$$P(y|x) = \sigma(y w^T x) \rightarrow \textcircled{1}$$

Now the method of maximum likelihood selects the hypothesis h which maximizes probability given by.

$$\prod_{n=1}^N P(y_n|x_n).$$

We can equivalently minimize it by putting

$$-\frac{1}{N} \ln \left(\prod_{n=1}^N P(y_n|x_n) \right)$$

$$= -\frac{1}{N} \sum_{n=1}^N \ln \left(\frac{1}{P(y_n|x_n)} \right)$$

$\frac{1}{N} \ln(\cdot)$ is a monotonically decreasing function.
 N & substituting $P(y_n/x_n)$ from eqn 8
we get

$$\frac{1}{N} \sum_{n=1}^N \ln \left(\frac{1}{\delta(y_n w^T x_n)} \right)$$

Since we are ~~mini~~ trying to minimize the term.
with respect to w allows us to call it as an
error measure. Substituting

$$\frac{1}{\delta(y_n w^T x_n)} = \frac{1}{1 + e^{-y_n w^T x_n}}$$

we have

$$E_m = \frac{1}{N} \sum_{n=1}^N \ln (1 + e^{-y_n w^T x_n})$$

i.e. the pointwise error measure is

$$e(x_n, y_n) = \ln (1 + e^{-y_n w^T x_n})$$

it has a negative exponential hence it will be
small if $y_n w^T x_n$ is large and vice versa.

This implies that $\text{sign}(w^T x_n) = y_n$.

∴ Error indeed ~~varies~~ varies w to classify each x_n correctly,
which is nothing but maximizing the
likelihood of $\prod_{n=1}^N P(y_n/x_n)$.

Thus the hypothesis h that minimizes the likelihood.

$\prod_{n=1}^N P(y_n/x_n)$ is equivalent to minimizing the
cross entropy error.

Ques) Derive the gradient of the in-sample error $\nabla E_{in}(w(t))$ using in gradient descent algorithm.

$$\text{Solutn} - E_{in}(w) = \frac{1}{N} \sum_{n=1}^N \ln(1 + e^{-y_n w^T x_n})$$

Taking gradient we have.

$$\nabla E_{in}(w) = \frac{1}{N} \sum_{n=1}^N \frac{d}{dw} \ln(1 + e^{-y_n w^T x_n}).$$

$$= \frac{1}{N} \sum_{n=1}^N \frac{1}{1 + e^{-y_n w^T x_n}} \frac{d}{dw} (1 + e^{-y_n w^T x_n})$$

$$= \frac{1}{N} \sum_{n=1}^N \frac{-y_n w^T x_n}{1 + e^{-y_n w^T x_n}} e^{-y_n w^T x_n}$$

$$= \frac{1}{N} \sum_{n=1}^N \frac{-y_n x_n}{e^{y_n w^T x_n} + e^{-y_n w^T x_n}} e^{-y_n w^T x_n}$$

$$\boxed{\nabla E_{in}(w) = -\frac{1}{N} \sum_{n=1}^N \frac{y_n x_n}{1 + e^{-y_n w^T x_n}}}$$

Exercise 3.6

a) $E_{in}(w)$ is given by.

$$E_{in}(w) = \frac{1}{N} \sum_{n=1}^N \ln \left(\frac{1}{\delta(y_n w^T x_n)} \right)$$

$$\text{we have } h(x) = \delta(w^T x). \quad \rightarrow ①$$

$$\& \text{ we have } \delta(-s) = 1 - \delta(s) \quad \rightarrow ②$$

$$\therefore E_{in}(w) = \frac{1}{N} \sum_{n=1}^N \ln \frac{1}{\delta(+1 \cdot w^T x_n)} \quad [\because y_n = +1]$$

$$+ [\because y_n = -1] \ln \frac{1}{\delta(-1 \cdot w^T x_n)}.$$

using eqⁿ ②

~~$$E_{in}(w) = \frac{1}{N} \sum_{n=1}^N \ln \frac{1}{1 - \delta(w^T x_n)}$$~~

$$E_{in}(w) = \frac{1}{N} \sum_{n=1}^N [\delta(y_n = +1)] \ln \frac{1}{\delta(w^T x_n)} \\ + [\delta(y_n = -1)] \ln \frac{1}{1 - \delta(w^T x_n)}$$

Using eqⁿ ①

$$E_{in}(w) = \frac{1}{N} \sum_{n=1}^N [\delta(y_n = +1)] \ln \frac{1}{h(x)} \\ + [\delta(y_n = -1)] \ln \frac{1}{1 - h(x)}$$

3.6

b) $h(x) = \Theta(w^T x)$.

$$E_{\text{in}}(w) = \frac{1}{N} \sum_{n=1}^N \mathbb{E}[y_n = +1] \ln \left(\frac{1}{h(x_n)} \right)$$

$$+ [y_n = -1] \ln \left(\frac{1}{1 - h(x_n)} \right)$$

Pattung $h(x_n) = \Theta(w^T x_n)$

$$= \frac{1}{N} \sum_{n=1}^N [y_n = +1] \ln \left(\frac{1}{\Theta(w^T x_n)} \right)$$

$$+ [y_n = -1] \ln \left(\frac{1}{1 - \Theta(w^T x_n)} \right)$$

We know, $1 - \Theta(w^T x_n) = \Theta(-w^T x_n)$

$$\Rightarrow \frac{1}{N} \sum_{n=1}^N [y_n = +1] \ln \left(\frac{1}{\Theta(w^T x_n)} \right)$$

$$+ [y_n = -1] \ln \left(\frac{1}{\Theta(-w^T x_n)} \right).$$

$$\Rightarrow \frac{1}{N} \sum_{n=1}^N \ln \left(\frac{1}{\Theta(y_n w^T x_n)} \right) +$$

$$\ln \left(\frac{1}{\Theta(-y_n w^T x_n)} \right)$$

$$\Rightarrow \frac{1}{N} \sum_{n=1}^N \ln \left(\frac{e^{y_n w_{nn}}}{1 + e^{y_n w_{nn}}} \right)$$

$$\Rightarrow \frac{1}{N} \sum_{n=1}^N \ln \left(\frac{1 + e^{y_n w_{nn}}}{e^{-y_n w_{nn}}} \right)$$

$$\Rightarrow \frac{1}{N} \sum_{n=1}^N \ln \left(1 + e^{-y_n w_{nn}} \right).$$

Exercise 3.13

a) Second order polynomial is given by

$$f_2(x) = (1, x_1, x_2; x_1^2, x_1 x_2, x_2^2). \rightarrow ①$$

Solving the eqn of parabola.

$$(x_1 - 3)^2 + x_2 = 1$$

$$x_1^2 - 6x_1 + 9 + x_2 = 0$$

Equating with ① we have weight vector.

$$\vec{w} = [w_0, w_1, w_2, w_3, w_4, w_5]$$

$$b) \text{ The circle } (x_1 - 3)^2 + (x_2 - 4)^2 = 1$$

$$x_1^2 - 6x_1 + 9 + x_2^2 - 8x_2 + 16 = 1$$

$$24 + x_1^2 - 6x_1 - 8x_2 + x_2^2 = 0$$

Equating with eqn ① we have

$$\vec{w} = [w_0, w_1, w_2, w_3, w_4, w_5]$$

$$c) \text{ The ellipse } 2(x_1 - 3)^2 + (x_2 - 4)^2 = 1$$

$$2[x_1^2 - 6x_1 + 9] + [x_2^2 - 8x_2 + 16] = 1$$

$$2x_1^2 - 12x_1 + 18 + x_2^2 - 8x_2 + 16 = 1$$

$$33 + 2x_1^2 - 12x_1 - 8x_2 + x_2^2 = 0.$$

Equating with eqn ① we have

$$\vec{w} = [w_0, w_1, w_2, w_3, w_4, w_5]$$

Problem 3.16

a)

		true classification	
		+1 (correct person)	-1 (contender)
you say	+1	c_1	c_a
	-1	c_2	0

$$\text{cost(accept)} = c_a \cdot P(y = +1/x)$$

$$+ c_a P(y = -1/x)$$

~~c_a~~

Now $g(x) = P[y = +1/x]$.

\therefore for negative probability
 $1 - g(x) = 1 - g(x)$

$$\therefore \text{cost accept} = c_a g(x)$$

$$\text{cost(reject)} = c_2 P[y = +1/x] + 0 \cdot P[y = -1/x]$$

$$= c_2 g(x)$$

b) $g(x)$ for accepting the person is

$$\text{cost (accept)} = \text{cost (reject)}$$

$$\Rightarrow C_a(1-g(x)) = C_r g(x)$$

$$\therefore g(x) = \frac{C_a}{C_a + C_r}$$

\therefore ~~our~~ threshold will be

$$k = \frac{C_a}{C_a + C_r}$$

c) For supermarket

~~$C_a = 10, C_r = 10$~~

$$C_r = 1$$

$$\therefore k = \frac{1}{1+10} = \frac{1}{11}$$

makes sense since we

are penalizing more for false reject;
when $g(x) < k$ which is very close to 0 we are rejecting
as we want to avoid false rejects.

~~For CIA~~

$$C_r = 1$$

$$C_a = 1000$$

$$\therefore k = \frac{1000}{1000+1} = \frac{1000}{1001}$$

which makes sense since

are penalizing more for false accept.
when $g(x) > k$ which is very close to 1 we
are accepting since we want to avoid
false accept.