

Problem 2.1 : In Equation (2.1), set $\delta = 0.03$, and let $E(M, N, \delta)$

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

a) For $M=1$, how many examples do we need to make $E \leq 0.05$?

$$\rightarrow E = \sqrt{\frac{1}{2N} \ln \frac{2M}{\delta}}$$

$$\Rightarrow N = \frac{1}{2E^2} \ln \frac{2M}{\delta} \Rightarrow N = \frac{1}{2 \times (0.05)^2} \ln \frac{2 \times 1}{0.03}$$

$$\Rightarrow N = \frac{1}{0.005} \ln (66.67) = 839.94$$

\therefore for $E \leq 0.05$, we have $N \geq 839.94$

b) For $M=100$, how many examples do we need to make $E \leq 0.05$

$$\rightarrow N = \frac{1}{2 \times (0.05)^2} \ln \frac{2 \times 100}{0.03} = \frac{1}{0.005} \ln (6666.67) = 1760.97$$

\therefore for $E \leq 0.05$, we have $N \geq 1760.97$

c) For $M=10,000$, how many examples do we need to make $E \leq 0.05$

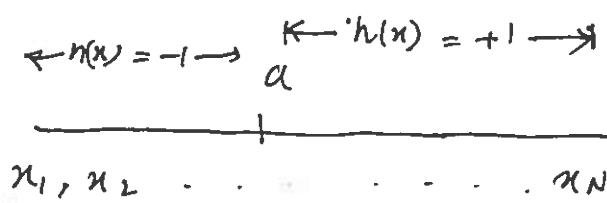
$$\rightarrow N = \frac{1}{2 \times (0.05)^2} \ln \frac{2 \times 10,000}{0.03} = \frac{1}{0.005} \times \ln (666666.67) = 2682.009$$

\therefore for $E \leq 0.05$, we have $N \geq 2682.009$

Problem 2.3 :- Compute the max^m no. of dichotomies, $m_H(N)$, for these learning models, and consequently compute d_{vc} , the vc dimension.

a) Positive or Negative Ray : \mathcal{H} contains the functions which are +1 on $[a, \infty)$ (for some a) together with those that are +1 on $(-\infty, a)$ (for some a).

→ Positive Rays



For these N points, growth funcⁿ $m_H(N) = \frac{N+1}{2}$ ($N-1$ sigma and for negative rays we have $N-1$. $+1$ on each end)
So, total dichotomies for 2^N .

$$m_H(N) = 2^N$$

Largest value of N for which $m_H(N) = 2^N$ is $\underline{\underline{2}}$

For positive rays $\underline{\underline{d_{rc} = 1}}$ & for negative rays $\underline{\underline{d_{rc} = 1}}$

as $m_H(N) = N+1$, break point = $\underline{\underline{2}}$

$$\begin{aligned} k+1 &< 2 \\ 2+1 &< 2^2 \text{ i.e.} \end{aligned}$$

(b) Positive or Negative Intervals : - For this $m_H(N) = \frac{N^2}{2} + \frac{N}{2} + 1$

Break point = 3 as

$$\begin{aligned} \frac{1}{2} + \frac{1}{2} + 1 &< 2 & \times \\ 2 + 1 + 1 &< 2 & \times \\ \frac{3}{2} + \frac{3}{2} + 1 &< 2^3 & \times \end{aligned}$$

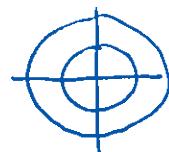
$$\therefore \underline{\underline{d_{rc} = 3}}$$

(c) Two concentric spheres in \mathbb{R}^d . $a \leq \sqrt{x_1^2 + \dots + x_d^2} \leq b$.

From the equation we can figure out that varying a and b , the circles will be centered around the origin. So the growth function will be like the positive and negative intervals.

$$\therefore m_H = \frac{N^2}{2} + \frac{N}{2} + 1$$

$$\Rightarrow \underline{\underline{d_{rc} = 3}}$$



Problem 2.8 : Which of the following are possible growth functions $m_H(N)$ for some hypothesis set

1. $1+N$:- For a finite VC dimension $m_H \leq 2^N$ and bounded by $N^{d_{VC}} + 1$

$$m_H(N) = 1+N \leq 2^N$$

This is true for $N=2$ (breakpoint), $d_{VC} = 1$

$$m_H(2) = 3 \leq 2^3 \text{ and is bounded by } N+1$$

\therefore This is a possible growth function. Ex - Positive Ray

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

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

This is true for $N=3$ (breakpoint), $d_{VC} = 2$

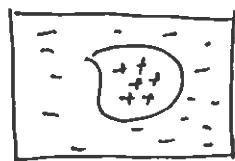
$$m_H(3) = 1+\frac{3}{2}+\frac{9}{2}=7 \leq 8 \text{ and is bounded by } N^2+1$$

\therefore This is a possible growth function. Ex - Positive Intervals

$$3. \underline{2^N} :- \text{For an infinite VC-dimension, } m_H = 2^N$$

For convex sets, there are infinite break points and $m_H = 2^N$

Ex :-



\therefore This is a possible growth function

$$4. \underline{\underline{2^{\sqrt{N}}}} :- m_H = 2^{\sqrt{N}} \leq 2^N$$

This is true for $N=2$, $d_{VC} = 1$

But, the bound $N+1$ fails for $N=25$. \therefore This is not a possible growth function.

$$5. \underline{2^{N/2}} :- m_H = 2^{N/2} \leq 2^N \cdot \text{This is true for } N=0 \text{ as } 1 \leq 2$$

But, the bound $N^0+1=2$, fails for $N=4$. \therefore This is not a possible growth function

Exercise

Problem 2.12 :- For an f with $dvc = 10$, what sample size do you need to have 95% confidence that your generalization error is at most .05?

→ We have found given the equation :-

$$Eout(f) \leq Ein(f) + \sqrt{\frac{8}{N} \ln \frac{4((2N)^{dvc} + 1)}{\delta}}$$

$$\text{Here, } E = \sqrt{\frac{8}{N} \ln \frac{4((2N)^{dvc} + 1)}{\delta}}$$

Confidence = 95%. means $\delta = 0.05$

$$E = 0.05$$

for error to be at most .05

$$N \geq \frac{8}{E^2} \ln \frac{4((2N)^{dvc} + 1)}{\delta}$$

$$\Rightarrow N \geq \frac{8}{(0.05)^2} \ln \frac{4((2 \times N)^{10} + 1)}{0.05}$$

The general rule of thumb is $N \geq 10 dvc$

We take a random guess of $N=1000$ initially in R.H.S

$$\Rightarrow N \geq \frac{8}{(0.05)^2} \ln \frac{4((2000)^{10} + 1)}{0.05} \approx 2.57 \times 10^5$$

Iteratively,

$$N \geq \frac{8}{(0.05)^2} \ln \left(\frac{4((2 \times 2.57 \times 10^5)^{10} + 1)}{0.05} \right) \approx 5.6 \times 10^5$$

$$\text{Iteratively } N \geq \frac{8}{(0.05)^2} \ln \left(\frac{4((2 \times 5.6 \times 10^5)^{10} + 1)}{0.05} \right) \approx 4.59 \times 10^5$$

$$\text{Iteratively } N \geq \frac{8}{(0.05)^2} \ln \left(\frac{4 \times ((2 \cdot 4.59 \times 10^5)^{10} + 1)}{0.05} \right) \approx 4.53 \times 10^5$$

$$\text{Iteratively } N \geq \frac{8}{(0.05)^2} \ln \left(\frac{4 \times ((2 \cdot 4.53 \times 10^5)^{10} + 1)}{0.05} \right) \approx 4.529 \times 10^5$$

Stratitely, $N \geq \frac{8}{(0.5)^2} \ln \left(\frac{4((2 \times 4.529 \times 10^5)^{10} + 1)}{0.05} \right) \approx 4.53 \times 10^6$

\therefore The estimate of $N \approx \underline{4.529 \times 10^5}$

Problem 2.22 : - When there is noise in the data, $E_{out}(g^{(D)}) = E_{x,y} [(g^{(D)}(x) - y(x))^2]$, where $y(x) = f(x) + \epsilon$. If ϵ is a zero mean noise random variable with variance δ^2 , show that the bias-variance decomposition becomes $R_D [E_{out}(g^{(D)})] = \delta^2 + bias + var$.

→ Expected value of the error over the entire space may be represented as -

$$\begin{aligned} R_D [E_{out}(g^{(D)})] &= E_D [E_{x,y} [(g^{(D)}(x) - y(x))^2]] \\ &= E_{x,y} [E_D [(g^{(D)}(x) - y(x))^2]] \end{aligned}$$

To evaluate this, we need the average hypothesis,

$$\bar{g}(x) = E_D [g^{(D)}(x)]$$

$$\Rightarrow E_D [E_{out}(g^{(D)})] = E_{x,y} [E_D [g^{(D)}(x)^2] - 2\bar{g}(x)y(x) + y(x)^2]$$

$$\begin{aligned} \text{Now, } E_D [g^{(D)}(x)^2] - 2\bar{g}(x)y(x) + y(x)^2 &= E_D [g^{(D)}(x)^2] - 2\bar{g}(x)(f(x) + \epsilon) + (f(x) + \epsilon)^2 \\ &= E_D [g^{(D)}(x)^2] - \bar{g}(x)^2 + \bar{g}(x)^2 - 2\bar{g}(x)f(x) - \\ &\quad 2\bar{g}(x)\epsilon + f(x)^2 + \epsilon^2 + 2f(x)\epsilon \end{aligned}$$

$$\begin{aligned} &= (E_D [g^{(D)}(x)^2] - \bar{g}(x)^2) + (\bar{g}(x)^2 - 2\bar{g}(x)f(x) + f(x)^2) \\ &\quad + \epsilon^2 - 2\bar{g}(x) - f(x)\epsilon \end{aligned}$$

$$= E_D \left[\left(g^{(D)}(x) - \hat{g}(x) \right)^2 + \epsilon^2 + 2(\hat{g}(x) - f(x))\epsilon \right]$$

as $E_D(g^{(D)}(x)) = \hat{g}(x)$

$$= E_D \left[g^{(D)}(x) - 2g^{(D)}(x)\hat{g}(x) + \hat{g}(x)^2 \right] + (\hat{g}(x) - f(x))^2 + \epsilon^2 - 2(\hat{g}(x) - f(x))\epsilon$$

$$= E_D \underbrace{\left[(g^{(D)}(x) - \hat{g}(x))^2 \right]}_{\text{var}(x)} + \underbrace{(\hat{g}(x) - f(x))^2}_{\text{bias}(x)} + \epsilon^2 - 2(\hat{g}(x) - f(x))\epsilon$$

Now,

$$\begin{aligned} E_D[E_{\text{out}}(g)] &= E_{n,y} [\text{var}(x) + \text{bias}(x) + \epsilon^2 - 2\hat{g}(x)f(x)\epsilon] \\ &= E_{n,y}[\text{var}(x)] + E_{n,y}[\text{bias}(x)] + E_{n,y}[\epsilon^2] - 2E_{n,y}[\hat{g}(x)f(x)\epsilon] \end{aligned}$$

$$\rightarrow \boxed{E_D[E_{\text{out}}(g)] = \text{var} + \text{bias} + \sigma^2}$$

Q6: 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 $E_{\text{in}}(w) = \frac{1}{N} \sum_{n=1}^N \ln(1 + e^{-y_n w^T x_n})$

→ likelihood is the probability of getting y_1, y_2, \dots, y_N in D from the corresponding x_1, x_2, \dots, x_N .

$$P((y_1, y_2, \dots, y_N) | (x_1, x_2, \dots, x_N)) = \prod_{n=1}^N P(y_n|x_n)$$

likelihood for one point,

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

Substitute $h(x) = \Theta(w^T x)$

From the equation of logistic function $\Theta(s) = \frac{e^s}{1+e^s}$
 $\therefore \Theta(-s) = 1 - \Theta(s)$

where $s = w^T x$

$$\Rightarrow P(y|x) = \Theta(y w^T x) \quad \text{where } y = \{-1, +1\}$$

Multiplying the likelihood for all the points,

$$\rightarrow P(y_1, y_2, \dots, y_N | x_1, x_2, \dots, x_N) = \prod_{n=1}^N \Theta(y_n w^T x_n)$$

Maximizing the likelihood,

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

↑
 \because the algorithm is increasing, so the log of this is also maximized
multiplying by $\frac{1}{N}$ is in line with the algorithm growing proportionally.

\therefore The above equation is the maximized likelihood.

To minimize, we have error measure,

$$E_{\text{in}}(w) = -\frac{1}{N} \ln \left(\prod_{n=1}^N \Theta(y_n w^T x_n) \right)$$

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

$$\therefore \Theta(s) = \frac{e^s}{1+e^s}$$

$$\Rightarrow \Theta(s) = \frac{1}{1+e^{-s}}$$

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

In sample error of
Logistic Regression or Cross Entropy Error.

Q7. Derive the gradient of the in-sample error $\nabla E_{in}(w(t))$ used in the gradient descent algorithm.

→ The in-sample error measure for logistic regression,

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

The gradient of this can be calculated as follows

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

as we are differentiating
w.r.t $w(t)$

$$\Rightarrow \nabla E_{in}(w(t)) = \frac{1}{N} \sum_{n=1}^N \left(\frac{e^{-y_n w^T x_n}}{1 + e^{-y_n w^T x_n}} \right) (-y_n x_n) \cdot e^{-y_n w^T x_n}$$

$$= \frac{1}{N} \sum_{n=1}^N \Theta(-y_n w^T x_n) (-y_n x_n)$$

$$= -\frac{1}{N} \sum_{n=1}^N y_n x_n \Theta(-s)$$

$$\therefore \phi(-x) = \frac{1}{1+e^{-x}}$$

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

Exercise 3.13 (a) $z = \phi_2(x)$. How can we use a hyperplane \tilde{w} in z to represent the following boundaries in x ?

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

$$\Rightarrow x_1^2 + 9 - 6x_1 + x_2 = 1$$

$$\Rightarrow -6x_1 + x_2 + x_1^2 + 8 = 0$$

Now $\phi_2(x) = (1, x_1, x_2, x_1^2, x_1 x_2, x_2^2)$

Hence the weights are $(-6, 1, 1, 0, 8)$ and intercept is 8. This satisfies the plane in the feature space (x_1, x_2, x_1^2) . This shows that the boundary is linear in the feature space $\phi_2(x)$

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

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

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

weights are $(-6, -8, 1, 0, 1)$ corresponding to $\phi_2(x)$

and the intercept is 24. This shows that a circular boundary is linear in this feature space (x_1, x_2, x_1^2, x_2^2)

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

$$\Rightarrow 2(x_1^2 + 9 - 6x_1) + x_2^2 + 16 - 8x_2 = 1$$

$$\Rightarrow -12x_1 - 8x_2 + 2x_1^2 + x_2^2 + 23 = 0$$

Here we see that the weights are $(-12, -8, 2, 0, 1)$ and the intercept is 23 for the feature space (x_1, x_2, x_1^2, x_2^2) . This shows that the ellipse is linearly separable in $\phi_2(x)$.

Problem 3.16 :- $g(x) = P[y = +1 | x]$

Cost Matrix

		True Classification	
		+1	-1
You say	+1	0	c_a
	-1	c_r	0

$$\begin{aligned}
 (a) \text{ cost (accept)} &= 0 \cdot P[y = +1 | x] + c_a \cdot P[y = -1 | x] \\
 &= c_a P[y = -1 | x] \\
 &= c_a (1 - g(x))
 \end{aligned}$$

$$\begin{aligned}
 \text{cost (reject)} &= c_r \cdot P[y = +1 | x] + 0 \cdot P[y = -1 | x] \\
 &= c_r g(x)
 \end{aligned}$$

$$\begin{aligned}
 (b) \text{ cost (accept)} &= \text{cost (reject)} \quad (\text{As initially we consider that both events are equally likely}) \\
 \Rightarrow c_a (1 - g(x)) &= c_r g(x)
 \end{aligned}$$

$$\Rightarrow c_r g(x) = c_a - c_a g(x)$$

$$\Rightarrow g(x) = \frac{c_a}{c_a + c_r} = k, \text{ where } k = \text{threshold}$$

(c) Supermarket example

		t	
		+1	-1
h	+1	0	1
	-1	10	0

CIA Example

		t	
		+1	-1
h	+1	0	1000
	-1	2	0

In the Supermarket example, we don't want any False rejects.

so $k = 1/11$. Hypothesis is rejected when $g(x) < k \approx 0$

In the CIA example, we don't want any false accept.

so $k = 1000/1001$ so Hypothesis is rejected when $g(x) > k \approx 1$.

Exercise-3.6 :- (a) Generally, if we are learning from ± 1 data to predict

a noisy target $p(y|x)$ with candidate hypothesis h , show that the maximum likelihood method reduces to the task of finding h

$$\text{that minimizes } E_{\text{ml}}(w) = \sum_{n=1}^N [[y_n = +1]] \ln \frac{1}{h(x_n)} + [[y_n = -1]] \ln \frac{1}{1-h(x_n)}$$

→ The method of maximum likelihood selects the hypothesis h which maximises the probability $p(y|x)$

$$\min E_{\text{ml}}(w) = \frac{1}{N} \sum_{n=1}^N \ln \frac{1}{p(y_n|x_n)} \quad \text{as} \quad \prod_{n=1}^N p(y_n|x_n) \\ = -\frac{1}{N} \ln \prod_{n=1}^N p(y_n|x_n)$$

For Pointwise error

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

for ± 1 data we have

$$E_{\text{ml}}(w) = \sum_{n=1}^N [[y_n = +1]] \ln \frac{1}{p(y_n|x_n)} + \sum_{n=1}^N [[y_n = -1]] \ln \frac{1}{p(y_n|x_n)}$$

$$\Rightarrow E_{\text{ml}}(w) = \sum_{n=1}^N [[y_n = +1]] \ln \frac{1}{h(x_n)} + [[y_n = -1]] \ln \frac{1}{1-h(x_n)}$$

(k) For the case $h(x) = \Theta(w^T x)$, argue that minimizing the in sample error in part (a) is equivalent to minimizing the one in 3.9

$$\rightarrow h(x) = \Theta(w^T x)$$

In sample error,

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

$$= \sum_{n=1}^N [y = \pm 1] \ln \frac{1}{\Theta(w^T x)} + [y = -1] \ln \frac{1}{1-\Theta(w^T x)}$$

from the question, we see that for the two probability distributions $\{p, 1-p\}$ and $\{q, 1-q\}$, the cross entropy error is $p = [y_n = +1]$ and $q = h(x)$

$$\Rightarrow \text{Cross Entropy Error} = p \log \frac{1}{q} + (1-p) \log \frac{1}{1-q} \quad \text{--- (1)}$$

From Equation 3.9,

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

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

$$= [y = \pm 1] \ln \frac{1}{\Theta(w^T x_n)} + [y = -1] \ln \frac{1}{1-\Theta(w^T x)}$$

$$= [y = \pm 1] \ln \frac{1}{h(x)} + [y = -1] \ln \frac{1}{1-h(x)}$$

$$\Leftrightarrow p \log \frac{1}{q} + (1-p) \log \frac{1}{1-q} \quad \therefore \text{Minimizing the 2 equations is equivalent.}$$