

1. by 投影法 p38, dvc  $\leq d+1$

故 3D perceptron 的 dvc  $\leq 3+1 = 4$

$$[a] \begin{pmatrix} X_1 & X_2 & X_3 \\ \begin{matrix} 1 & 1 & 1 \\ 2 & 4 & 3 \\ 3 & 3 & 3 \\ 4 & 2 & 3 \end{matrix} \end{pmatrix}$$

$$\rightarrow \begin{pmatrix} 1 & 1 & 1 \\ 0 & 2 & 1 \\ 0 & 0 & 0 \\ 0 & -2 & -1 \end{pmatrix} \rightarrow \begin{pmatrix} 1 & 1 & 1 \\ 0 & 2 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}$$

$$[b] \begin{pmatrix} X_1 & X_2 & X_3 & X_4 \\ \begin{matrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 4 & 4 \\ 1 & 3 & 3 & 2 \\ 1 & 4 & 2 & 3 \end{matrix} \end{pmatrix} \rightarrow L \cdot D$$

$$\rightarrow \begin{pmatrix} 1 & 1 & 1 & 1 \\ 0 & 1 & 3 & 3 \\ 0 & 0 & -4 & -5 \\ 0 & 0 & -8 & -7 \end{pmatrix} \rightarrow \begin{pmatrix} 1 & 1 & 1 & 1 \\ 0 & 1 & 3 & 3 \\ 0 & 0 & -4 & -5 \\ 0 & 0 & 0 & 3 \end{pmatrix} \rightarrow L \cdot I$$

(c)  $X_1 \ X_2 \ X_3 \ X_4$

$$\left( \begin{array}{cccc} 1 & 1 & 1 & 1 \\ 1 & 2 & 4 & 2 \\ 1 & 3 & 3 & 2 \\ 1 & 4 & 2 & 2 \end{array} \right) \rightarrow \left( \begin{array}{cccc} 1 & 1 & 1 & 1 \\ 0 & 1 & 3 & 1 \\ 0 & 2 & 2 & 1 \\ 0 & 3 & 1 & 1 \end{array} \right)$$

$$\rightarrow \left( \begin{array}{cccc} 1 & 1 & 1 & 1 \\ 0 & 1 & 3 & 1 \\ 0 & 0 & -4 & -1 \\ 0 & 0 & -8 & -2 \end{array} \right) \rightarrow \left( \begin{array}{cccc} 1 & 1 & 1 & 1 \\ 0 & 1 & 3 & 1 \\ 0 & 0 & -4 & -1 \\ 0 & 0 & 0 & 0 \end{array} \right)$$

$\rightarrow L \cdot D$

(D)  $d_{VC} \leq d+1, 5 > 4$  故排除.

(a)(c) 因線性相依使得某變量  $x_n$  只能由其它  $n-1$  個  $x$  決定，  
故無法產生  $2^n$  個 dichotomy.

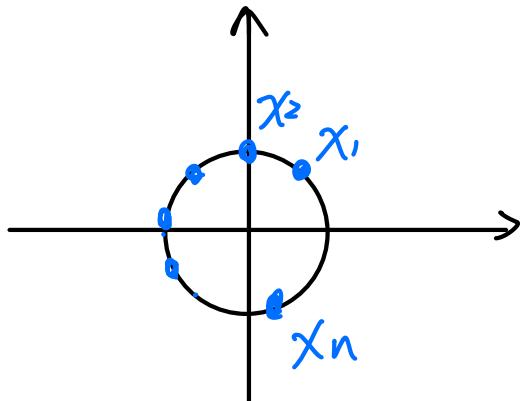
(a)  $W^T X_3 = W_0^T X_1 - W_x^T X_2 > 0 \quad d_{VC} = 2$

(c)  $W^T X_4 = W_0^T X_1 + W_0^T X_2 - W_x^T X_3 > 0 \quad d_{VC} = 3$

(d) L.I.,  $W^T X_4 \neq W^T X_1 + W^T X_2 + W^T X_3, d_{VC} = 4$

2. 由於  $w_0 = 0$ , 故 Data 可轉換至極座標表達,  
且長度不影響 data 關係. 故可以單位圓表示,

如下圖



由於 data 只由 "角度" 來分類, 故可將問題等價成

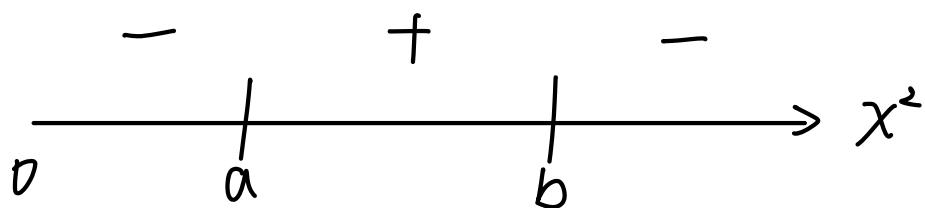


positive negative ray.

故  $m_M(N) = 2N$ .

$\Rightarrow \text{ans} : [\text{C}]$

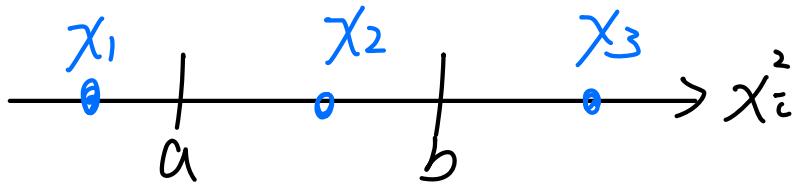
3. data 為圓心的距離  $\chi^2$ , 決定 "+", "-"  
因此問題可等價成一維 positive interval, 如下圖



$$\text{故 } m_H(N) = \binom{N+1}{2} + 1$$

$\Rightarrow \text{ans: } [a]$

4. 由於 positive interval 的  $dvc(M_H) = 2$  (by 第3題)



若  $\{x_1, x_2, x_3\} = \{+, -, +\}$  則無法 shatter  
故  $dvc(M_H) = 2$ .

$\Rightarrow \text{ans} : [d]$

5. [a] degree of freedom = 4, dvc = 4

1. 單一 positive intervals 決定了 "O" 区間後，若能從剩餘 X 中找出不相鄰 "O" 的 "O"，則無法 shatter
2. union positive intervals 決定了 2 個不相鄰 "O" 区間後，若能從剩餘 data 找出不相鄰 data 的組合，則無法 shatter

$n=4$  時  $\{ OXOX \}$  為所有 O 不相鄰組合，  
 $\{ OXXX \}$   
 $\{ OXXO \}$   
 $\{ XOZO \}$

此時無論將任意 X 變改為 O，皆會使原來的 O 相鄰  
故找不出第 3 個 O 区間，因此 union positive intervals  
能 shatter  $n=4$

$n=5$  時  
1. OXOX  
2. OXXOX  
3. OXXXO  
4. XOXOX  
5. XOXZO  
6. XXCXO

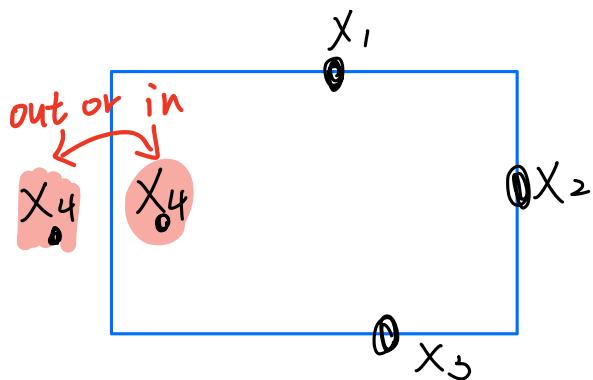
為所有 (不相鄰 O) 且 (O 区間長度皆為 1) 的組合

此時 1 和 3 號組合可以從 X 区間再選出一連續 "O"  
 $\Rightarrow OXOXO$ ，因此  $n=5$  為 breakpoint -  $dvc=4$

[b] degree of freedom = 4 - (長方形的 top-left  
 top-right  
 bottom-left  
 bottom-right)

(以上根據 lecture 4, P40 簡略估計，以下較詳細)

- 在  $n=4$  時，將長方形的任3邊以任3莫決定，如下



此時  $x_4$  可在矩形內 or 外，故在 3 莫決定後，  
 第 4 個仍可 shatter，所以  $n=4$  可 shatter.

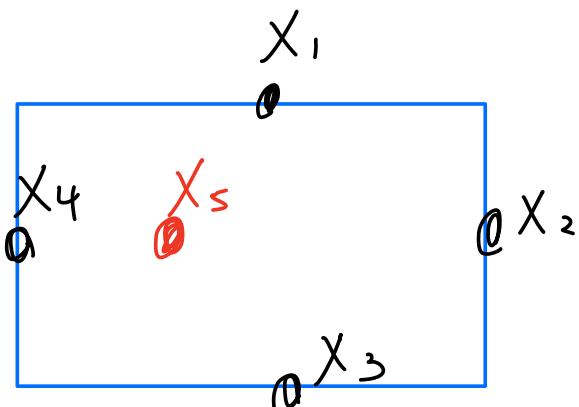
- 在  $n=5$  時，將長方形任 4 邊以  $x_1 \sim x_4$

決定，其中  $x_1$  決定  $\max y$  -  $x_2$  決定  $\max X$

$x_3$  決定  $\min y$  -  $x_4$  決定  $\min X$

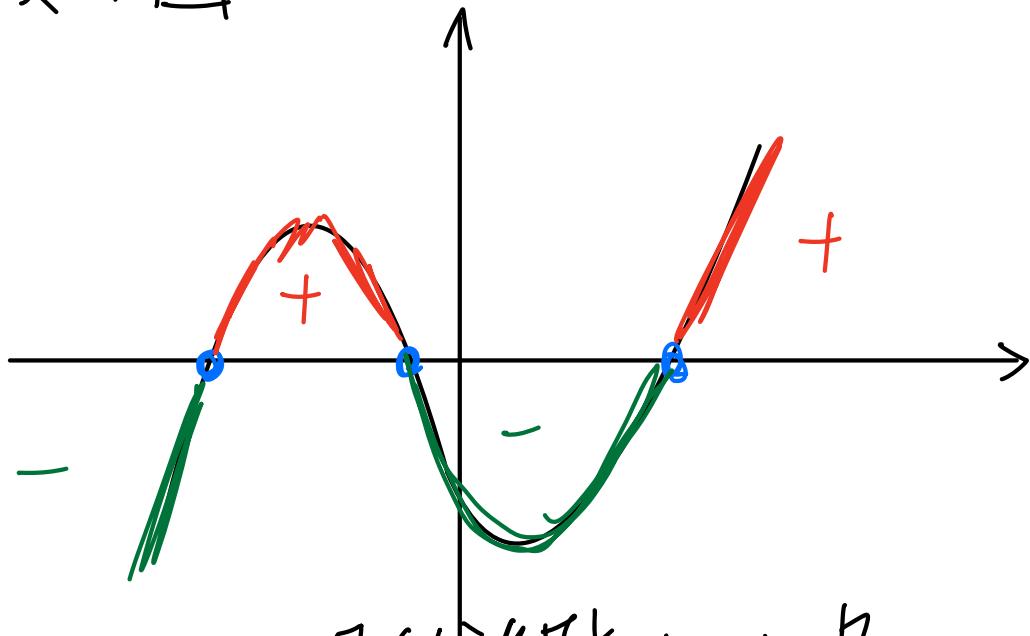
此時  $x_5$  必在矩形內，被  $x_1 \sim x_4$  決定 (如下圖)

故  $x_5$  無法被 shatter，因此  $dvc = 4$



[d] 3次多項式最多會有3個根

如下圖：



因  $x_1 \sim x_n$  最多變號  $n-1$  次

故  $x_1 \sim x_4$  最多變號 3 次

Ex:  $h(x_1), h(x_2), h(x_3), h(x_4) \Rightarrow -, +, -, +$

故  $n=4$  時可 shatter (變號 3 次)

而  $n=5$  時，最多會變號 4 次，

故  $n=5$  無法 shatter,  $d_{VC} = 4$ .

$$(c)(ii) X = \begin{bmatrix} -x_1 & - \\ -x_2 & - \\ \vdots & \\ -x_4 & - \end{bmatrix} \text{ when } n=3.$$

$$\text{special } X = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 0 \end{bmatrix} \Rightarrow L \cdot I \text{-invertible}$$

$$\text{for any } y = \begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{pmatrix}, \text{ find } w, \text{ s.t.}$$

$$\text{sign}(Xw) = y \Leftarrow (Xw) = y \Leftrightarrow w = X^{-1}y$$

其中  $w_0 > 0$ , 故:

$$X^{-1} = \begin{bmatrix} x'_{11} & x'_{12} & \cdots & x'_{14} \\ & \vdots & & \\ x'_{41} & x'_{42} & \cdots & x'_{44} \end{bmatrix}_{4 \times 4}$$

$$\Rightarrow w_0 = y_1 x'_{11} + y_2 x'_{12} + y_3 x'_{13} + y_4 x'_{14}$$

when  $y_i \cdot x_{1i} > 0$ , then  $w_0 > 0$

由於  $x'_{11} \sim x'_{14}$  可為 + or -, 故任意  $y$   
均可在  $w_0 > 0$  的前提下得到

一組  $x'_{11} \sim x'_{14}$  與  $y$  的內積為正,  
ex:  $y = \begin{bmatrix} 0 \\ x \\ 0 \end{bmatrix}$  then  $x'_{11} > 0, x'_{12} < 0, x'_{13} > 0, x'_{14} < 0$   
故  $w_0 > 0$  不改變 VC 維度,  
 $\Rightarrow dvc \geq 4$

$$(1) X = \begin{pmatrix} -x_1 \\ -x_2 \\ -x_3 \\ -x_4 \\ -x_5 \end{pmatrix}, L.D \text{ 不可逆.}$$

$$\text{special } X = \left| \begin{array}{ccccc} 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 1 \end{array} \right| \quad 5 \times 4$$

$\therefore$  (1) 已說明  $w_0 > 0$  不改變  $y$  的 dichotomy 次數

故存在某  $w$ , s.t.

$$w^T x_5 = \underbrace{w^T x_1}_0 - \underbrace{w^T x_2}_X - \underbrace{w^T x_3}_X - \underbrace{w^T x_4}_X > 0$$

(由於 L.D 故  $x_5$  為  $x_1 \sim x_4$  的線性組合)

因此  $n=5$  時  $x_5$  不 shatter

$$\Rightarrow dvc \leq 4, \text{ by (1), (2)} \Rightarrow dvc = 4$$

b. 每個 dichotomy 會對應 1 個 hypothesis

而  $n$  個 input 最多有  $2^n$  個 dichotomy.

因本題的 hypothesis = 1126

$$\Rightarrow 2^n < 1126 \Rightarrow n = 10.13$$

$$dvc(H) \leq \lg |H|$$

$$dvc(H) \leq \lg(1126)$$

$$\leq 10.13$$

- then dvc(H) 最大為 10.

7.

$$2M \cdot \exp(-2\epsilon^2 N) = \delta$$

$$\Rightarrow -2\epsilon^2 N = \ln\left(\frac{\delta}{2M}\right)$$

$$\Rightarrow \epsilon^2 N = \frac{1}{2} \ln\left(\frac{2M}{\delta}\right)$$

$$\Rightarrow \epsilon = \sqrt{\frac{1}{2N} \ln\left(\frac{2M}{\delta}\right)}$$

$$\because P[|E_{out}(g) - E_{out}(g^*)| < \epsilon] \geq 1 - \delta$$

1 - δ 的 upper bound

⇒ ans is [C]

8. positive ray 的  $m_1(N) = 2N$

$$\varepsilon = 0.1, \delta = 0.1$$

$$\Rightarrow f = 8 \cdot N \cdot \exp\left(-\frac{1}{8}(0.1)^2 \cdot N\right) = 0.1 = f$$

在 [a],  $f = 0.298$

~~[b]~~,  $f = 0.09$

[c],  $f = 0.029$

[d],  $f = 0.009$

[e] -  $f = 0.0028$

所以  $N = 11000$  是滿足  $f < 0.1$  的最少的  $N$  [b]

$$9 \cdot \textcircled{1} E(\omega) \simeq E(u) + b_E(u)^T (\omega - u) + \frac{1}{2} (\omega - u)^T A_E(u) (\omega - u)$$

$$\textcircled{2} \quad \omega = u + v, \quad V = \omega - u.$$

$$E(\omega) \simeq E(u) + b_E(u)^T \cdot V + \frac{1}{2} V^T A_E(u) \cdot V$$

$$\frac{\partial E(\omega)}{\partial \omega} = \frac{\partial}{\partial \omega} \left( E(u) + b_E(u)^T \cdot V + \frac{1}{2} V^T A_E(u) \cdot V \right)$$

$$\frac{\partial E(\omega)}{\partial \omega} = b_E(u) + V A_E(u)$$

$$\text{令 } \frac{\partial E(\omega)}{\partial \omega} = 0, \quad V = -b_E(u) A_E^{-1}(u)$$

$$\Rightarrow \text{ans : } [b]$$

$$10. \frac{\partial E_{in}(w)}{\partial w} = \frac{\partial}{\partial w} \frac{1}{N} \sum_{n=1}^N \ln(1 + \exp(-y_n w^T x_n))$$

$$\stackrel{?}{=} 1 + \exp(-y_n w^T x_n) = \square$$

$$-y_n w^T x_n = 0$$

$$= \frac{1}{N} \sum_{n=1}^N \left( \frac{1}{\square} \right) \frac{\partial}{\partial w} \square$$

$$= \frac{1}{N} \sum_{n=1}^N \frac{1}{1 + \exp(0)} \cdot (y_n x_n) \cdot \exp(0)$$

$$= \frac{1}{N} \sum_{n=1}^N \frac{\exp(0)}{1 + \exp(0)} \cdot (-y_n x_n)$$

$$= \frac{1}{N} \sum_{n=1}^N \frac{1}{1 + \exp(-0)} \cdot (-y_n x_n)$$

$$= \frac{1}{N} \sum_{n=1}^N h_t(y_n x_n) \cdot (-y_n x_n) = E_{in}(w)$$

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

$$= \frac{1}{N} \sum_{n=1}^N \frac{(-y_n x_n) \exp(y_n w^T x_n)}{(1 + \exp(y_n w^T x_n))^2} \cdot (-y_n x_n)$$

$$= \frac{1}{N} \sum_{n=1}^N \frac{\exp(y_n w^T x_n)}{1 + \exp(y_n w^T x_n)} \cdot \frac{1}{1 + \exp(y_n w^T x_n)} \cdot (x_n x_n^T)$$

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

$$\cdot (x_n x_n^T)$$

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

$\Rightarrow$  Ans: [d]

$$/\!/. \text{ by def } XX^+X = X - \textcircled{1}$$

$$X^+XX^+ = X^+ - \textcircled{2}$$

$$(XX^+)^T = XX^+ - \textcircled{3}$$

$$(X^+X)^T = X^+X - \textcircled{4}.$$

$$\text{(a) by } \textcircled{3}. \text{ 同乘 } X^T \Rightarrow X^T(XX^+)^T = X^TXX^+$$

$$\Rightarrow (XX^+X)^T = (X^TX)X^+$$

$$\text{by } \textcircled{1} \Rightarrow X = (X^TX)X^+$$

$$X^TX \text{ 可逆} \Rightarrow (X^TX)^{-1}X = X^+$$

$$(b) (XX^+)^2 = XX^+XX^+$$

by def  $XX^+X = X$

$$\text{then } XXXX^+ = XX^+$$

$$\begin{aligned}(XX^+)^3 &= (XX^+)^2(XX^+) \\ &= (XX^+)^2 = XX^+\end{aligned}$$

$$\text{so, } (XX^+)^k = XX^+$$

$$(c) U\Sigma V^T V\Sigma^+ U^T$$

$$= U\Sigma I\Sigma^+ U^T$$

$$= U\Sigma\Sigma^+ U^T = U I U^T = I$$

由於  $X$  為  $N$  by  $d+1$  matrix

then  $XX^+$  為  $N \times N$  matrix.

(d) 假設  $X_{6 \times 5}$  則  $X_{5 \times 6}^+$

by (c) 的結論.

故  $XX^+$  為  $I_{6 \times 6}$

而  $X_{6 \times 5}$  的 rank 最大為 5 (因第 6 列必為  
前 5 列的線性組合)

$$\Rightarrow \text{trace}(I_{6 \times 6}) = 6 \neq \text{rank}(X) = 5$$

$\Rightarrow$  (d) false.

$$12. D = \{(X_1, y_1), (X_2, y_2), \dots, (X_n, y_n)\}$$

$x_i \in \mathbb{R}^d, y_i \in \mathbb{R}$

$$\text{find } W^* = \operatorname{argmax} P(D|W^*)$$

$$P(D|W^*) = P(y_1, y_2, \dots, y_n | X_1, X_2, \dots, X_n, W^*)$$

$$= \prod_{i=1}^n P(y_i | X_i, W^*)$$

$$= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(y_i - w^T x_i)^2\right)$$

$$= \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^n \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - w^T x_i)^2\right)$$

改寫  $\sum_{i=1}^n (y_i - w^T x_i) \Rightarrow \begin{pmatrix} y_1 - w^T x_1 \\ y_2 - w^T x_2 \\ \vdots \\ y_n - w^T x_n \end{pmatrix} = y - w^T X$

$= y - X^T w$

$$= \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^n \exp\left(-\frac{1}{2\sigma^2} (y - X^T w)^T (y - X^T w)\right)$$

$$P(D|W*) = \left( \frac{1}{\sqrt{2\pi}\sigma^2} \right)^n \exp \left( -\frac{1}{2\sigma^2} (y - X^T w)^T (y - X^T w) \right)$$

What  $w^*$  can maximize  $P(D|w^*)$ ?

我們需要找出  $\underset{w}{\operatorname{argmin}} (y - X^T w)^T (y - X^T w)$

假設  $G(w) = (y - X^T w)^T (y - X^T w)$

$$\frac{\partial G(w)}{\partial w} = (-X^T(y - X^T w)) \times 2$$

$$= X^T(X^T w - y) \times 2$$

$$= 2(X X^T w - X^T y)$$

$$\text{令 } \frac{\partial G(w)}{\partial w} = 0 = X X^T w - X^T y$$

$$\Rightarrow X^T y = X X^T w$$

$$\Rightarrow w = (X X^T)^{-1} X^T y$$

所以當  $w^* = (X X^T)^{-1} X^T y$  時  $P(D|w^*)$  有  
最大值。



```
def generateData(tranum, valnum, addoutliers = 0):
    # 產生200筆 training data traVal
    # 產生5000筆 testing data testVal
    traLabel = np.random.normal(0, 1, tranum).reshape(tranum, 1)
    valLabel = np.random.normal(0, 1, valnum).reshape(valnum, 1)

    # traData = np.zeros((200,2))
    # valData = np.zeros((5000,2))
    traData = []
    valData = []
    for _ in traLabel:
        if _[0] >= 0:
            # 產生 var = 0.6 ,mean = [2 , 3]
            x1 = np.random.normal(2, 0.6 ** 0.5, 1)[0]
            x2 = np.random.normal(3, 0.6 ** 0.5, 1)[0]
        else:
            # 產生var = 0.4 ,mean = [0 , 4]
            x1 = np.random.normal(0, 0.4 ** 0.5, 1)[0]
            x2 = np.random.normal(4, 0.4 ** 0.5, 1)[0]
        traData.append([1, x1, x2])
    for _ in valLabel:
        if _[0] >= 0:
            # 產生 var = 0.6 ,mean = [2 , 3]
            x1 = np.random.normal(2, 0.6 ** 0.5, 1)[0]
            x2 = np.random.normal(3, 0.6 ** 0.5, 1)[0]
        else:
            # 產生var = 0.4 ,mean = [0 , 4]
            x1 = np.random.normal(0, 0.4 ** 0.5, 1)[0]
            x2 = np.random.normal(4, 0.4 ** 0.5, 1)[0]
        valData.append([1, x1, x2])
    traData = np.asarray(traData, dtype=np.float64)
    valData = np.asarray(valData, dtype=np.float64)

    if addoutliers :
        Noise = []
        NoiseLabel = [1 for _ in range(20)]
        for _ in range(20):
            x1 = np.random.normal(6, 0.3 ** 0.5, 1)[0]
            x2 = np.random.normal(0, 0.1 ** 0.5, 1)[0]
            Noise.append([1, x1, x2])

        Noise = np.asarray(Noise,dtype=np.float64)
        NoiseLabel = np.asarray(NoiseLabel,dtype=np.float64).reshape(20,1)

        traData = np.concatenate((traData,Noise),axis=0)
        traLabel = np.concatenate((traLabel,NoiseLabel),axis=0)

    traLabel[traLabel >= 0] = 1
    traLabel[traLabel < 0] = -1
    valLabel[valLabel >= 0] = 1
    valLabel[valLabel < 0] = -1

    return [[traData,traLabel],[valData,valLabel]]


def meanerr(iterNUM,noiseTF = 0):
    lin_errinBag = []
    lin_errInOutBag = []
    lin_log_errout01Bag = []
    lin_log_errout01Bag_noise = []

    for _ in range(iterNUM):

        DataSET = generateData(200, 5000, 0)
        M_plus = np.linalg.pinv(DataSET[0][0])
        W_lin = np.dot(M_plus, DataSET[0][1])
        errin = np.dot(DataSET[0][0], W_lin) - DataSET[0][1]
        sqrErrMean = np.sum(errin ** 2) / len(errin)
        lin_errinBag.append(sqrErrMean)

        #for Q14
        lin_errin01 = linValidate(W_lin, DataSET[0][0], DataSET[0][1])
        lin_errout01 = linValidate(W_lin, DataSET[1][0], DataSET[1][1])
        lin_errInOutBag.append(abs(lin_errout01-lin_errin01))

        #for Q15
        log_errout01= logisticGrad(500,DataSET[0][0],DataSET[0][1],DataSET[1][0],DataSET[1][1])
        lin_log_errout01Bag.append([lin_errout01,log_errout01])
        #for Q16
        if noiseTF == 1:

            DataSET = generateData(200, 5000, 1)
            M_plus = np.linalg.pinv(DataSET[0][0])
            W_lin = np.dot(M_plus, DataSET[0][1])

            lin_errout01 = linValidate(W_lin, DataSET[1][0], DataSET[1][1])
            log_errout01= logisticGrad(500,DataSET[0][0],DataSET[0][1],DataSET[1][0],DataSET[1][1])
            lin_log_errout01Bag_noise.append([lin_errout01,log_errout01])

        lin_errIn01 = sum(lin_errinBag) / iterNUM
        lin_errInOutDiff01 = sum(lin_errInOutBag)[0] / iterNUM
        lin_log_errout01Bag = np.array(lin_log_errout01Bag)
        lin_log_err = np.mean(lin_log_errout01Bag, axis=0).reshape(-1)

        lin_log_errout01Bag_noise = np.array(lin_log_errout01Bag_noise)
        lin_log_errout01Bag_noise = np.mean(lin_log_errout01Bag_noise, axis=0).reshape(-1)

    return lin_errIn01,lin_errInOutDiff01,lin_log_err,lin_log_errout01Bag_noise


def linValidate(W_lin, valDATA, valLabel):
    pred = np.dot(valDATA,W_lin)
    pred[pred >= 0] = 1
    pred[pred < 0] = -1
    err01 = sum(pred != valLabel) / len(valLabel)
    return err01


def logisticGrad(iterNUM, traData, traLabel, valData, valLabel):
    # traData 200,3
    # traLabel 200,1
    # weight 3,1
    W = np.zeros((3,1),dtype=np.float64)
    theta = lambda x: (1 / (1 + np.exp(-x)))

    for _ in range(iterNUM):

        # s = y_n * W^T * x_n
        s = np.dot(traData, W) * traLabel*(-1)
        grad = ( theta(s)* (-1)*traLabel * traData)
        grad = np.mean(grad, axis=0).reshape(3, 1)
        W = W - (0.1* grad)

    s = np.dot(valData, W)
    pred = theta(s)
    pred[pred >= 0.5] = 1
    pred[pred < 0.5] = -1
    res = sum(pred != valLabel) / len(valLabel)
    return res
```