

1. 因為每一層 layer 都是 fully connected，其中每一層 layer 都有一個常數節點不會受 weight 影響，第一層數量為 10，最後一層為 1，對於層數 L 的 neural network，總 weight 數量為：

$$\sum_{l=0}^{L-2} d^{(l)} * (d^{(l+1)} - 1) + d^{L-1}$$

Constraint 為 $\sum_{l=1}^{L-1} d^{(l)} = 36$

這個不等式不是很容易求，只好使用窮舉法找出最小與最大 weight 數量，其程式碼如下：

```

1  import numpy as np
2
3  maxw = -1
4  minw = 10000
5  minstruc = []
6  maxstruc = []
7
8  def layer_structure( number, unitnum, hidden = []):
9      global maxw
10     global minw
11     global minstruc
12     global maxstruc
13     if np.sum(hidden) != None:
14         if np.sum(hidden) == unitnum:
15             wnum = 10*(hidden[0]-1) + hidden[len(hidden) - 1]
16             for i in range(len(hidden) - 1):
17                 wnum += hidden[i]*(hidden[i+1]-1)
18             if maxw < wnum:
19                 maxw = wnum
20                 maxstruc = hidden
21                 #print('max:', maxw, maxstruc)
22             if minw > wnum:
23                 minw = wnum
24                 minstruc = hidden
25                 #print('min:', minw, minstruc)
26             return
27         elif np.sum(hidden) > unitnum:
28             return
29         for i in number:
30             layer_structure(number, unitnum, hidden + [i])
31
32 maxw = -1
33 minw = 10000
34 number = np.arange(2, 36)
35 layer_structure(number, 36)
36
37 print('final min:', minw, minstruc)
38 print('final max:', maxw, maxstruc)
39
40

```

最後得出的最小 weight 數為 46，layer 排列為

{10, 2,2,2,2,2,2,2,2,2,2,2,2,2,2,2,2,1}

2. 最後得出的最大 weight 數為 510，layer 排列為 {10, 22, 14, 1}

3.

$$\begin{aligned}
 \nabla_w \text{err}_n(w) &= \nabla_w (\|x_n - w w^T x_n\|^2) \\
 &= \nabla_w (x_n^T x_n - 2 w w^T x_n x_n + (w w^T x_n)^2) \\
 &= -4 w^T x_n x_n + \nabla_w (x_n^T w w^T w w^T x_n) \quad (x_n^T w = w^T x_n \text{ 都是同一條線}) \\
 &= -4 w^T x_n x_n + \nabla_w ((x_n^T w)^2 (w^T w)) \\
 &= -4 w^T x_n x_n + 2 x_n^T w \cdot (w^T w) \cdot x_n + 2 w \cdot (x_n^T w)^2
 \end{aligned}$$

4.

$$\begin{aligned}
 E_n(w) &= \frac{1}{N} \sum_{n=1}^N \|x_n - w w^T (x_n + \varepsilon_n)\|^2 \\
 &= \frac{1}{N} \sum_{n=1}^N (x_n^T x_n - 2(x_n + \varepsilon_n)^T w w^T x_n + (x_n + \varepsilon_n)^T w w^T w w^T (x_n + \varepsilon_n)) \\
 &= \frac{1}{N} \sum_{n=1}^N (x_n^T x_n - 2x_n^T w w^T x_n + x_n^T w w^T w w^T x_n) + \frac{1}{N} \sum_{n=1}^N (-2\varepsilon_n^T w w^T x_n + \varepsilon_n^T w w^T w w^T \varepsilon_n + 2\varepsilon_n^T w w^T w w^T x_n) \\
 &= \frac{1}{N} \sum_{n=1}^N \|x_n - w w^T x_n\|^2 + \frac{1}{N} \sum_{n=1}^N ((-2w w^T x_n + 2w w^T w w^T x_n) \cdot \varepsilon_n + \varepsilon_n^T w w^T w w^T \varepsilon_n) \\
 &\quad \text{已知 } E(\varepsilon_n) = 0, E(\varepsilon_n \varepsilon_n^T) = I_n \\
 E(E_n(w)) &= \frac{1}{N} \sum_{n=1}^N \|x_n - w w^T x_n\|^2 + \frac{1}{N} \sum_{n=1}^N ((-2w w^T x_n + 2w w^T w w^T x_n) E(\varepsilon_n) + E(\frac{1}{N} \sum_{n=1}^N \varepsilon_n^T w w^T w w^T \varepsilon_n)) \\
 &= \frac{1}{N} \sum_{n=1}^N \|x_n - w w^T x_n\|^2 + \frac{1}{N} \sum_{n=1}^N E(\text{tr}(\varepsilon_n^T w w^T w w^T \varepsilon_n)) \\
 &= \frac{1}{N} \sum_{n=1}^N \|x_n - w w^T x_n\|^2 + \frac{1}{N} \sum_{n=1}^N \text{tr}(E(w w^T \varepsilon_n \cdot \varepsilon_n^T w w^T)) \\
 &= \frac{1}{N} \sum_{n=1}^N \|x_n - w w^T x_n\|^2 + \frac{1}{N} \sum_{n=1}^N \text{tr}(w w^T E(\varepsilon_n \varepsilon_n^T) w w^T) \\
 &= \frac{1}{N} \sum_{n=1}^N \|x_n - w w^T x_n\|^2 + \frac{1}{N} \sum_{n=1}^N N (w^T w)^2 \Rightarrow \mathcal{L}(w) = (w^T w)^2 \quad \#
 \end{aligned}$$

5.

$$\text{error function } E = \sum_{i=1}^d (g_i(x) - x_i)^2$$

$$\text{let } U = \begin{bmatrix} u_{11} & \dots & u_{1d} \\ \vdots & & \vdots \\ u_{d1} & \dots & u_{dd} \end{bmatrix} \Rightarrow g_i(x) = \sum_{k=1}^d u_{ik} \tanh\left(\sum_{k=1}^d u_{kj} x_k\right)$$

$$E = \sum_{i=1}^d \left(\sum_{j=1}^J u_{ij} \tanh\left(\sum_{k=1}^d u_{kj} x_k\right) - x_i \right)^2$$

6.

$$\text{for } w: E = \sum_{i=1}^d \left(\sum_{j=1}^J w_{ji}^{(2)} \tanh\left(\sum_{k=1}^d w_{kj}^{(1)} x_k\right) - x_i \right)^2$$

$$\frac{\partial E}{\partial w_{ij}^{(1)}} = \sum_{i=1}^d 2 \left(\sum_{j=1}^J w_{ji}^{(2)} \tanh\left(\sum_{k=1}^d w_{kj}^{(1)} x_k\right) - x_i \right) \cdot \frac{\partial \left(\sum_{j=1}^J w_{ji}^{(2)} \tanh\left(\sum_{k=1}^d w_{kj}^{(1)} x_k\right) - x_i \right)}{\partial w_{ij}^{(1)}} \approx C$$

$$= \sum_{i=1}^d 2 \left(\sum_{j=1}^J w_{ji}^{(2)} \tanh\left(\sum_{k=1}^d w_{kj}^{(1)} x_k\right) - x_i \right) w_{ji} \tanh'\left(\sum_{k=1}^d w_{kj}^{(1)} x_k\right) \cdot \frac{\partial \left(\sum_{k=1}^d w_{kj}^{(1)} x_k \right)}{\partial w_{ij}^{(1)}}$$

$$= \sum_{i=1}^d 2 \left(\sum_{j=1}^J w_{ji}^{(2)} \tanh\left(\sum_{k=1}^d w_{kj}^{(1)} x_k\right) - x_i \right) w_{ji} \tanh'\left(\sum_{k=1}^d w_{kj}^{(1)} x_k\right) \cdot x_i$$

$$\frac{\partial E}{\partial w_{ji}^{(2)}} = 2 \left(w_{ji}^{(2)} \tanh\left(\sum_{k=1}^d w_{kj}^{(1)} x_k\right) - x_i \right) \cdot \frac{\partial \left(w_{ji}^{(2)} \tanh\left(\sum_{k=1}^d w_{kj}^{(1)} x_k\right) - x_i \right)}{\partial w_{ji}^{(2)}} \approx C$$

$$= 2 \left(w_{ji}^{(2)} \tanh\left(\sum_{k=1}^d w_{kj}^{(1)} x_k\right) - x_i \right) \cdot \tanh\left(\sum_{k=1}^d w_{kj}^{(1)} x_k\right)$$

$$\frac{\partial E}{\partial u_{ij}} = \frac{\partial \left(\sum_{j=1}^J u_{ij} \tanh\left(\sum_{k=1}^d u_{kj} x_k\right) - x_i \right)}{\partial u_{ij}} + \frac{\partial \left(\sum_{i=1}^d \left(\sum_{j=1}^J u_{ij} \tanh\left(\sum_{k=1}^d u_{kj} x_k\right) - x_i \right)^2 \right)}{\partial u_{ij}}$$

$$= 2 \left(u_{ij} \tanh\left(\sum_{k=1}^d u_{kj} x_k\right) - x_i \right) \cdot \tanh\left(\sum_{k=1}^d u_{kj} x_k\right) + u_{ij} \tanh'\left(\sum_{k=1}^d u_{kj} x_k\right) \cdot x_i$$

$$+ \sum_{i=1}^d 2 \left(u_{ij} \tanh\left(\sum_{k=1}^d u_{kj} x_k\right) - x_i \right) \cdot u_{ij} \tanh'\left(\sum_{k=1}^d u_{kj} x_k\right) \cdot x_i$$

$$= \sum_{i=1}^d 2 \left(u_{ij} \tanh\left(\sum_{k=1}^d u_{kj} x_k\right) - x_i \right) \cdot u_{ij} \tanh'\left(\sum_{k=1}^d u_{kj} x_k\right) \cdot x_i + 2 \left(u_{ij} \tanh\left(\sum_{k=1}^d u_{kj} x_k\right) - x_i \right) \tanh\left(\sum_{k=1}^d u_{kj} x_k\right)$$

take $w_{ij}^{(1)}, w_{ji}^{(2)}$ as u_{ij}

$$\frac{\partial E}{\partial w_{ij}^{(1)}} + \frac{\partial E}{\partial w_{ji}^{(2)}} = \sum_{i=1}^d 2 \left(u_{ij} \tanh\left(\sum_{k=1}^d u_{kj} x_k\right) - x_i \right) \cdot u_{ij} \tanh'\left(\sum_{k=1}^d u_{kj} x_k\right) \cdot x_i + 2 \left(u_{ij} \tanh\left(\sum_{k=1}^d u_{kj} x_k\right) - x_i \right) \tanh\left(\sum_{k=1}^d u_{kj} x_k\right)$$

$$= \frac{\partial E}{\partial u_{ij}} \quad \#$$

7

兩個點的 1-Nearest Neighbor 即為求任意點離哪一個點比較近,

因此 hypothesis 為 2 點的中垂面.

對於平面上的點 x

$$\begin{aligned} (x^+ - x^-, z) \cdot (x - \frac{x^+ + x^-}{2}, 0) &= 0 \\ \Rightarrow (x^+ - x^-) \cdot x - \frac{\|x^+\|^2 - \|x^-\|^2}{2} &= 0 \\ g_{LW}(x) = \text{sign}(w^T x + b) &= \text{sign}((x^+ - x^-) \cdot x - \frac{\|x^+\|^2 - \|x^-\|^2}{2}) \end{aligned}$$

8.

已知 $\text{sign}(x) = \text{sign}(ax)$ when $a > 0$, 且 $\exp(x) > 0$

$$\begin{aligned} g_{\text{GRAFNET}}(x) &= \text{sign}(\beta_+ \exp(-\|x - M_+\|^2) + \beta_- \exp(-\|x - M_-\|^2)) \\ &= \text{sign}(\exp(\|x - M_-\|^2) \cdot (\beta_+ \exp(-\|x - M_+\|^2) + \beta_- \exp(-\|x - M_-\|^2))) \\ &= \text{sign}(\beta_+ \exp(\|x - M_-\|^2 - \|x - M_+\|^2) + \beta_-) \quad (\because \beta_+ > 0) \\ &= \text{sign}(\exp(\|x - M_-\|^2 - \|x - M_+\|^2) + \frac{\beta_-}{\beta_+}) \quad (\because \beta_+ > 0) \end{aligned}$$

若 $\exp(\|x - M_-\|^2 - \|x - M_+\|^2) + \frac{\beta_-}{\beta_+} > 0$:

$$\exp(\|x - M_-\|^2 - \|x - M_+\|^2) > -\frac{\beta_-}{\beta_+}$$

$$\Leftrightarrow \|x - M_-\|^2 - \|x - M_+\|^2 > \ln(-\frac{\beta_-}{\beta_+})$$

$$\Leftrightarrow \|x - M_-\|^2 - \|x - M_+\|^2 - \ln(-\frac{\beta_-}{\beta_+}) > 0$$

$$\text{同理, 若 } \exp(\|x - M_-\|^2 - \|x - M_+\|^2) + \frac{\beta_-}{\beta_+} \leq 0 \Leftrightarrow \|x - M_-\|^2 - \|x - M_+\|^2 - \ln(-\frac{\beta_-}{\beta_+}) \leq 0$$

$$\Rightarrow \text{sign}(\exp(\|x - M_-\|^2 - \|x - M_+\|^2) + \frac{\beta_-}{\beta_+}) = \text{sign}(\|x - M_-\|^2 - \|x - M_+\|^2 - \ln(-\frac{\beta_-}{\beta_+}))$$

$$= \text{sign}((x^T x - 2M_+^T x + M_+^T M_+) - (x^T x - 2M_-^T x + M_-^T M_-) - \ln(-\frac{\beta_-}{\beta_+}))$$

$$= \text{sign}((2(M_- - M_+)^T x + (M_-^T M_- - M_+^T M_+ - \ln(-\frac{\beta_-}{\beta_+}))) \neq$$

9.

現在要 fix V 求 $E_m(W_m) = \sum_{(x_n, r_{n,m}) \in D_m} (r_{n,m} - W_m^T V_n)^2$ 的最小值.

$$\nabla_{W_m} E_m(W_m) = \sum_{(x_n, r_{n,m}) \in D_m} -2(r_{n,m} - W_m^T V_n) V_n$$

已知 $V_n = 1$, 令上式為 0:

$$\nabla_{W_m} E_m(W_m) = \sum_{(x_n, r_{n,m}) \in D_m} -2(r_{n,m} - W_m) = 0 \Rightarrow \sum_{(x_n, r_{n,m}) \in D_m} r_{n,m} = \sum_{(x_n, r_{n,m}) \in D_m} W_m$$

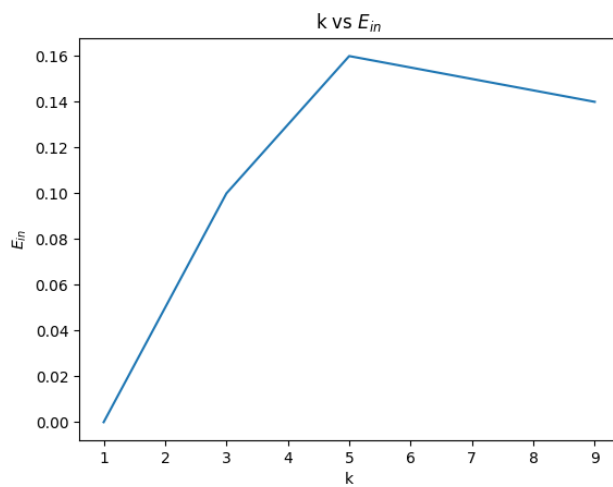
$$W_m = \frac{\sum_{(x_n, r_{n,m}) \in D_m} r_{n,m}}{\sum_{(x_n, r_{n,m}) \in D_m} 1} = m\text{-th 電影的平均分數.} \#$$

10.

$$V_{N+1}^T W_m = \frac{1}{N} \sum_{n=1}^N V_n^T W_m = \frac{1}{N} \sum_{n=1}^N r_{n,m}$$

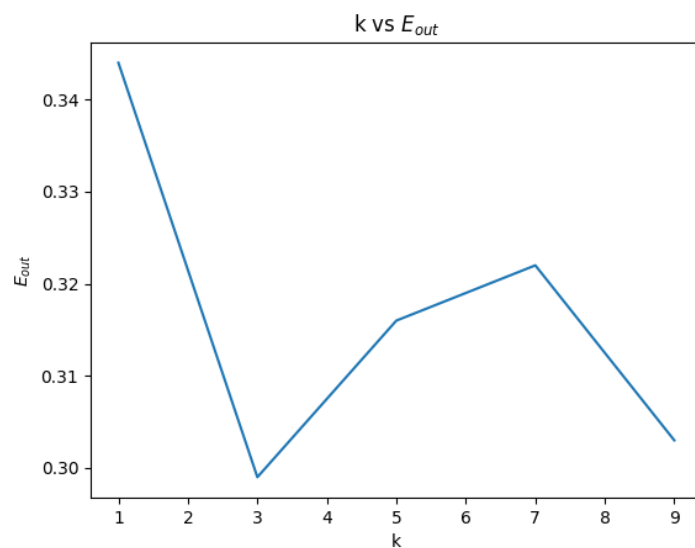
$\frac{1}{N} \sum_{n=1}^N r_{n,m}$ 最大的電影即為平均分數最高的電影. #

11.



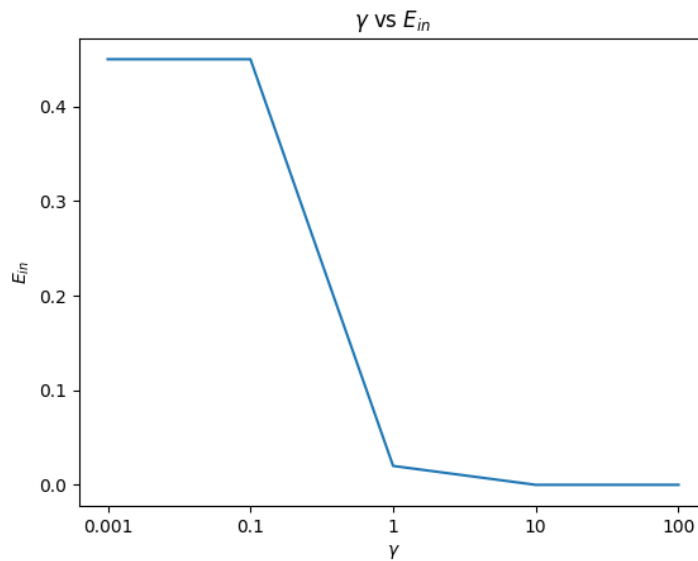
對於任意 sample 點 x ，當 $k = 1$ 時距離最近的點就是自己，所以 $E_{in} = 0$ 。若 $k \neq 1$ 則也會受周遭點的影響，因為不能保證 `hw4_train.dat` 數據 1 與 -1 與附近分佈具相關性，因此 E_{in} 跟 k 並沒有明顯的正負相關性。

12.



跟 E_{in} 不同，當 $k = 1$ 時 E_{out} 反而最大，因為最近的點不一定會跟自己的 y 值一樣，可以看出 k Nearest Neighbor 的表現並不好，而且做 `predict` 也很吃運算資源。

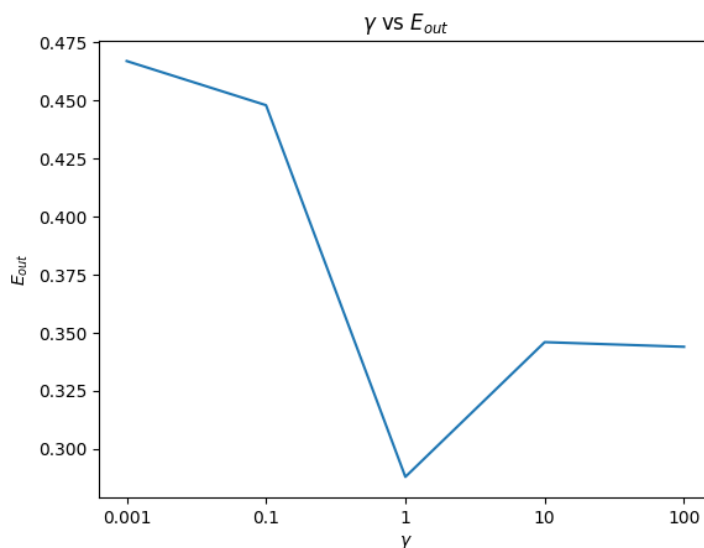
13.



$$g_{\text{uniform}} = \text{sign}(\sum_{m=1}^N y_m \exp(-\gamma \|x - x_m\|^2))$$

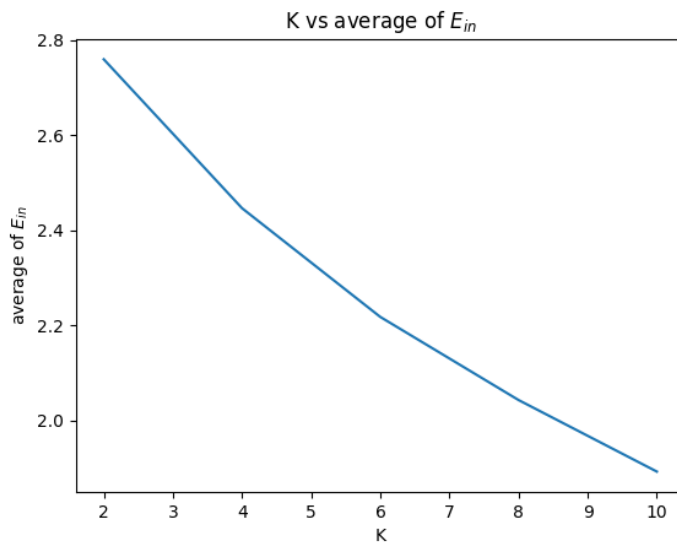
可以看出當 γ 越大 E_{in} 越小。因為 γ 越大代表與 x 越接近的點所擁有的權重越大，而最接近 x 的 sample 即為 x 自己，所以當 γ 夠大時 x 會 dominant 整個決策，此時 $E_{in} = 0$ 。

14.



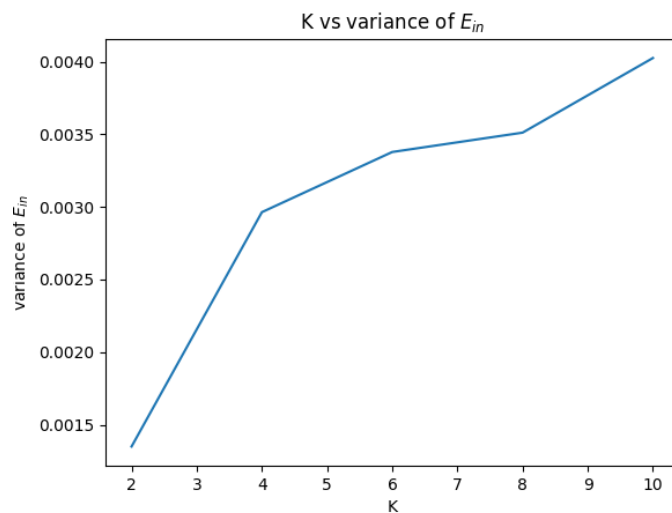
γ 跟 E_{out} 的相關性相對來說就沒有上題那麼明顯了，只知道 γ 越大 g_{uniform} 的值也會跟 x 相近的點越接近。理論上來說會呈負相關，但沒有一定的嚴格遞減性。

15.



可以看出 K 越大平均 E_{in} 越小，因為 K 越大能夠分越多群， μ_m 跟 x_m 之間的平均距離當然也會越小。

16.



K 越大 E_{in} 的 variance 就越大。因為當 K 很小時，分的群體多樣性也低，因此會收斂到差不多的位置， E_{in} 也會差不多。當 K 變大時，群體多樣性增加，最後收斂的 μ_m 位置受到起始 random 位置的影響也會變大，因此 E_{in} 也較容易差很多。

17.

$$\text{for } \Delta \geq 2, N \geq 3\Delta \log_2 \Delta \Rightarrow N+1 < 2^N$$

先證 $N^0 < 2^N$:

$$\text{當 } N = 3\Delta \log_2 \Delta:$$

$$N^0 < 2^N \Leftrightarrow N^0 < (2^{\log_2 \Delta})^{3\Delta}$$

$$\Leftrightarrow N^0 < \Delta^{3\Delta} \Leftrightarrow N < \Delta^3$$

$$\Leftrightarrow 3\Delta \log_2 \Delta < \Delta^3$$

$$\text{當 } \Delta = 2, 3 < 2^3 \text{ 成立}$$

$$\text{取 } f(\Delta) = \Delta^3 - 3\log_2 \Delta$$

$$f'(\Delta) = 3\Delta - \frac{3}{\Delta \ln 2} > 0 \text{ for } \Delta \geq 2$$

$$\Rightarrow f(\Delta) \text{ 是嚴格遞增 for } \Delta \geq 2 \Rightarrow \Delta^3 - 3\log_2 \Delta > 0 \Leftrightarrow 3\log_2 \Delta < \Delta^3 \Leftrightarrow N^0 < 2^N \text{ for } \Delta \geq 2 - (1)$$

$$N^0 < 2^N \Leftrightarrow \Delta \ln N < N \ln 2$$

$$\text{取 } f(N) = \Delta \ln N - N \ln 2$$

$$\therefore N^0 < 2^N \text{ for } N = 3\Delta \log_2 \Delta$$

$$\therefore f(3\Delta \log_2 \Delta) < 0 \text{ for } N = 3\Delta \log_2 \Delta$$

$$f'(N) = \frac{\Delta}{N} - \ln 2 = 0 \Rightarrow N = \frac{\Delta}{\ln 2}$$

$$\Rightarrow \text{當 } N > \frac{\Delta}{\ln 2}, f'(N) = \frac{\Delta}{N} - \ln 2 < \frac{\Delta}{\frac{\Delta}{\ln 2}} - \ln 2 = 0$$

$$\Rightarrow f(N) \text{ 是嚴格遞減 for } N > \frac{\Delta}{\ln 2}$$

$$\text{又 } 3\Delta \log_2 \Delta = \frac{3\Delta \ln \Delta}{\ln 2} > \frac{\Delta}{\ln 2}$$

$$\Rightarrow f(N) \leq f(3\Delta \log_2 \Delta) < 0 \text{ for } N = 3\Delta \log_2 \Delta$$

$$\Rightarrow N^0 < 2^N$$

根據下題, $\Delta = 3cd+1$, N 為 sample 數, 兩者皆為整數, 因此 $N^0+1 < 2^N$ 也成立

18.

We know that bounding function $BCN(k)$:

maximum possible $m_H(N)$ when $\text{rank point} = k$

且对 N 组数据, d 维 input, d 维 perceptrons 能表示组合数 $\leq BCN(d+1)$

而在题中第 0 层到第 1 层 (4 节点即为 d 维 perceptrons,

3 节点能表示组合数 $\leq BCN(d+1)^3$

已知 $BCN(k) \leq \sum_{i=0}^k C_i^N \leq N^{k-1} + 1$ (ML foundation hw 2)

\Rightarrow 3 节点能表示组合数 $\leq BCN(d+1)^3 \leq (N^{d+1} + 1)^3$

要证 H_{3A} 的 VC dimension $\leq 3(3cd+1+1) \log_2(3cd+1+1)$:

let $\Delta = 3cd+1+1$ for $d \geq 0 \Rightarrow \Delta \geq 4 \geq 2$

let $N \geq 3\Delta \log_2 \Delta$, 则根据上题:

$$N^{3cd+1+1} + 1 < 2^N$$

3 节点能表示的组合数 $\leq BCN(d+1)^3 \leq (N^{d+1} + 1)^3$

$$\leq (N^{d+1} + 1)^3 = N^{3cd+1} + 3N^{2cd+1} + 3N^{d+1} + 1$$

对于: $N \geq 3\Delta \log_2 \Delta \geq 3 \cdot 4 \log_2 4 = 24$ ($\because \Delta \geq 4$)

$$\begin{aligned} \Rightarrow (N^{d+1} + 1)^3 &= N^{3cd+1} + 3N^{2cd+1} + 3N^{d+1} + 1 < N^{3cd+1} + N^{3cd+1} + N^{3cd+1} + 1 \\ &= 3N^{3cd+1} + 1 < N^{3cd+1+1} + 1 < 2^N \end{aligned}$$

$\Rightarrow H_{3A}$ 无法 shatter $3\Delta \log_2 \Delta$ 个类

$$VC \text{ dimension} < 3\Delta \log_2 \Delta = 3(3cd+1+1) \log_2(3cd+1+1) \#$$