

Hand-written Part

学号: t11902210 姓名: 张凡 (交换生)

code question report 在后面

Problem 1

Swish is an activation function that can be viewed as a "modification" of the logistic function. Swish has been shown to be better than ReLU in some deep learning models. Recall that the logistic function is defined as

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

Now Swish is defined as

$$\varphi(s) = s \cdot \theta(s)$$

What is $\varphi'(s)$?

解:

$$\varphi(s) = s \cdot \theta(s)$$

$$\begin{aligned} \varphi'(s) &= \theta(s) + s \cdot \frac{d(\theta(s))}{ds} \quad \Leftarrow \frac{d(\theta(s))}{ds} = \theta'(s) = \left(\frac{1}{1 + e^{-s}} \right)' \\ &= \frac{1}{1 + e^{-s}} + \frac{s \cdot e^{-s}}{(1 + e^{-s})^2} &= -\frac{e^{-s} \cdot (-1)}{(1 + e^{-s})^2} = \frac{e^{-s}}{(1 + e^{-s})^2} \\ &= \frac{1 + e^{-s} + s \cdot e^{-s}}{(1 + e^{-s})^2} \\ &= \frac{1 + (s+1)e^{-s}}{(1 + e^{-s})^2} & \varphi'(s) = \frac{1 + (s+1)e^{-s}}{(1 + e^{-s})^2} \end{aligned}$$

Problem 2

In the class we briefly talked about the famous PageRank algorithm, which assigns a rank to each web page. Given a set of pages: $P = \{P_1, P_2, \dots, P_n\}$ and their incoming links: $\text{in}(P_j | P_j \text{ has a link to } P_i, i \neq j)$. The basic idea is that a web page should be ranked higher if it has 1) more incoming hyperlinks and 2) incoming links from high-ranked pages. The PageRank can be computed through the following iterative algorithm.

1. Initialize $\text{PageRank}_0(P_i) = \frac{1}{n}$ for each page.
2. Compute the updated PageRank for each page: $\text{PageRank}_k(P_i) = \frac{\sum_{P_j \in \text{in}(P_i)} \text{PageRank}_k(P_j)}{\text{out}(P_j)}$
where $\text{out}(P_j)$ denotes the out-degree of P_j .
3. Repeat step 2 until the ranks converge.

Now consider the following example with 3 pages:

$$P = \{P_1, P_2, P_3\} \quad \text{in}(P_1) = \{P_2, P_3\} \quad \text{in}(P_2) = \{P_3\} \quad \text{in}(P_3) = \{P_1\}$$

We define the following variables:

$$V_t = (\text{PageRank}_t(P_1), \text{PageRank}_t(P_2), \text{PageRank}_t(P_3))^T$$

$$P = \begin{bmatrix} 0 & 1 & 0.5 \\ 0 & 0 & 0.5 \\ 1 & 0 & 0 \end{bmatrix}.$$

where V_t is the ranks after the t -th iteration, and P is the transition matrix between pages. We can rewrite the update rule above as $V_t = PV_{t-1}$. Answer the following questions:

- (A) Based on the iterative algorithm above, please compute the values of v_1, v_2, v_3, v_4 and v_5 .

~~Step 1~~ according to the algorithm.

Step 1: $P = \{P_1, P_2, P_3\} \Rightarrow n=3 \Rightarrow V_0 = \left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3}\right)^T$

initialize $\text{PageRank}_0(P_i) = \frac{1}{3} \quad (i=1, 2, 3) \Rightarrow V_0 = \left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3}\right)$

Step 2-1: $V_1 = P \cdot V_0 = P V_0 = \left(\begin{array}{ccc|c} 0 & 1 & \frac{1}{2} & \frac{1}{3} \\ 0 & 0 & \frac{1}{2} & \frac{1}{3} \\ 1 & 0 & 0 & \frac{1}{3} \end{array}\right) \left(\begin{array}{c} \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \end{array}\right) = \left(\begin{array}{c} \frac{1}{2} \\ \frac{1}{6} \\ \frac{1}{3} \end{array}\right)$

Step 2-2: $V_2 = P V_1 = \left(\begin{array}{ccc|c} 0 & 1 & \frac{1}{2} & \frac{1}{2} \\ 0 & 0 & \frac{1}{2} & \frac{1}{6} \\ 1 & 0 & 0 & \frac{1}{3} \end{array}\right) \left(\begin{array}{c} \frac{1}{2} \\ \frac{1}{6} \\ \frac{1}{3} \end{array}\right) = \left(\begin{array}{c} \frac{1}{3} \\ \frac{1}{6} \\ \frac{1}{2} \end{array}\right)$

Step 2-3: $V_3 = P V_2 = \left(\begin{array}{ccc|c} 0 & 1 & \frac{1}{2} & \frac{1}{3} \\ 0 & 0 & \frac{1}{2} & \frac{1}{6} \\ 1 & 0 & 0 & \frac{1}{2} \end{array}\right) \left(\begin{array}{c} \frac{1}{3} \\ \frac{1}{6} \\ \frac{1}{2} \end{array}\right) = \left(\begin{array}{c} \frac{5}{12} \\ \frac{1}{4} \\ \frac{1}{3} \end{array}\right)$

Step 2-4: $V_4 = P V_3 = \left(\begin{array}{ccc|c} 0 & 1 & \frac{1}{2} & \frac{5}{12} \\ 0 & 0 & \frac{1}{2} & \frac{1}{4} \\ 1 & 0 & 0 & \frac{1}{3} \end{array}\right) \left(\begin{array}{c} \frac{5}{12} \\ \frac{1}{4} \\ \frac{1}{3} \end{array}\right) = \left(\begin{array}{c} \frac{5}{12} \\ \frac{1}{6} \\ \frac{5}{12} \end{array}\right)$

Step 2-5 $V_5 = P V_4 = \left(\begin{array}{ccc|c} 0 & 1 & \frac{1}{2} & \frac{5}{12} \\ 0 & 0 & \frac{1}{2} & \frac{1}{6} \\ 1 & 0 & 0 & \frac{5}{12} \end{array}\right) \left(\begin{array}{c} \frac{5}{12} \\ \frac{1}{6} \\ \frac{5}{12} \end{array}\right) = \left(\begin{array}{c} \frac{3}{8} \\ \frac{5}{24} \\ \frac{5}{12} \end{array}\right)$

In conclusion $\left\{ \begin{array}{l} V_0 = \left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3}\right)^T \\ V_1 = \left(\frac{1}{2}, \frac{1}{6}, \frac{1}{3}\right)^T \\ V_2 = \left(\frac{1}{3}, \frac{1}{6}, \frac{1}{2}\right)^T \\ V_3 = \left(\frac{5}{12}, \frac{1}{4}, \frac{1}{3}\right)^T \\ V_4 = \left(\frac{5}{12}, \frac{1}{6}, \frac{5}{12}\right)^T \\ V_5 = \left(\frac{3}{8}, \frac{5}{24}, \frac{5}{12}\right)^T \end{array} \right.$

- (B) Recall that in the class, we mentioned that the algorithm converges when $v^* = Pv^*$ solve this equation to derive v^* . Note that you should provide a normalized v^* as answer.

设 $V = \alpha \begin{pmatrix} v_1 \\ v_2 \\ v_3 \end{pmatrix}$ (待定) $PV^* = V^* \Rightarrow (P - E)V^* = 0$ (E 为单位矩阵)

$$P = \begin{pmatrix} 0 & 1 & \frac{1}{2} \\ 0 & 0 & \frac{1}{2} \\ 1 & 0 & 0 \end{pmatrix} \quad (P - E) = \begin{pmatrix} -1 & 0 & \frac{1}{2} & \frac{1}{2} \\ 0 & -1 & \frac{1}{2} \\ 1 & 0 & -1 \end{pmatrix}$$

$$(P - E)V^* = \alpha \begin{pmatrix} -v_1 + v_2 + \frac{1}{2}v_3 \\ -v_2 + \frac{1}{2}v_3 \\ v_1 - v_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$$

$$\begin{aligned} \det(P - E) &= \begin{vmatrix} 0 & 1 & -\frac{1}{2} \\ 0 & -1 & \frac{1}{2} \\ 1 & 0 & -1 \end{vmatrix} \\ &= -\begin{vmatrix} 0 & 1 & -\frac{1}{2} \\ 0 & -1 & \frac{1}{2} \\ 1 & 0 & 0 \end{vmatrix} \\ &= -\begin{vmatrix} 1 & \frac{1}{2} \\ -1 & \frac{1}{2} \end{vmatrix} = 1 \neq 0 \end{aligned}$$

$$\Rightarrow \begin{cases} -v_1 + v_2 + \frac{1}{2}v_3 = 0 \\ -v_2 + \frac{1}{2}v_3 = 0 \\ v_1 - v_3 = 0 \end{cases} \Rightarrow \begin{cases} 2v_2 = v_3 \\ v_1 = v_3 \end{cases}$$

$$V^* = \alpha \begin{pmatrix} 2v_2 \\ v_2 \\ 2v_2 \end{pmatrix} = \alpha \cdot v_2 \begin{pmatrix} 2 \\ 1 \\ 2 \end{pmatrix}$$

$$V^* \Rightarrow \text{normalized } \sqrt{2^2 + 1^2 + 2^2} = 3$$

解得 $V^* = \begin{pmatrix} \frac{2}{3} \\ \frac{1}{3} \\ \frac{2}{3} \end{pmatrix}$

Consider $L \in \{1, 2, 3\}$
and

Problem 3

For a $d^{(0)} - d^{(1)} - d^{(2)} - \dots - d^{(L)}$ fully-connected neural network, the total number of neurons is

$$D = \sum_{l=1}^L d^{(l)}$$
 and the total number of weight is $\sum_{l=1}^L d^{(l-1)}d^{(l)}$

which ignore the so-called 'bias' terms in the network. Let the number of inputs $d^{(0)} = 10$ and the total number of neurons $D = 100$. Over all possible choices of L , $d^{(1)}, d^{(2)}, \dots, d^{(L)}$. What is the maximum total number of weights (and what neural network architecture does it correspond to)? What is the minimum total number of weights (and what neural network architecture does it correspond to)?

If: maximum: 选用较大的 $d^{(l)}$ 和较小的 L . 使 ~~每层 $d^{(l-1)} - d^{(l)} - d^{(l+1)}$~~ $= 100$

$$\sum_{l=1}^L d^{(l-1)}d^{(l)} = \cancel{10 \cdot 10 \cdots 100 \cdot 1}$$

$$d^{(0)}d^{(1)} + d^{(1)}d^{(2)} + \dots + d^{(L-1)}d^{(L)}$$

数、-1式子 $a+b=100$ ($a>0, b>0$)

$$10000 = (a+b)^2 \geq a^2 + 2ab + b^2 \geq a^2 + b^2.$$

因此在总 D 等的情况下，分层越多 $\sum_{l=1}^L d^{(l-1)}d^{(l)}$ 越大。

直觉想法

最大情况下，设计方法 $\cancel{d^{(l-1)} - d^{(l)} - d^{(l+1)}}$

但见好 $d^{(0)}$ 很好

我们通常选择 $(L=1, 2, 3)$

$$\textcircled{1} L=1 \quad \sum_{l=1}^L d^{(l-1)}d^{(l)} = 10 \times 100 = 1000$$

$$\textcircled{2} L=2 \text{ 时 } \sum_{i=1}^2 d(i-1)d(i) = d(0)d(1) + d(1)d(2)$$

$$\begin{cases} d(1)=x \\ d(2)=y \end{cases} \quad W_2 = 10x + xy \leq -x^2 + 110x = -(x-55)^2 + 3025 \leq 3025$$

$$\textcircled{3} L=3 \text{ 时 } W_3 = 10d(1) + d(1)d(2) + d(2)d(3)$$

$$\begin{cases} x, y, z = d_1, d_2, d_3 \\ x+y+z=100 \end{cases}$$

$$d(1) > 0 \quad W = 10x + xy + yz = 10x + xy + (100-x-y)y$$

$$\frac{\partial W}{\partial x} = 10 \rightarrow 0 \quad \frac{\partial W}{\partial y} = 2x + 100 - x = x + 100 \quad = 10x + 100y - y^2 = 10 \cdot (100-y-z)$$

$$\text{由 } Z \text{ 的最小值 } z \geq 0. \quad y = 45, z = 1 \quad = -(y-50)^2 + 2500 + 10x + 100y - y^2$$

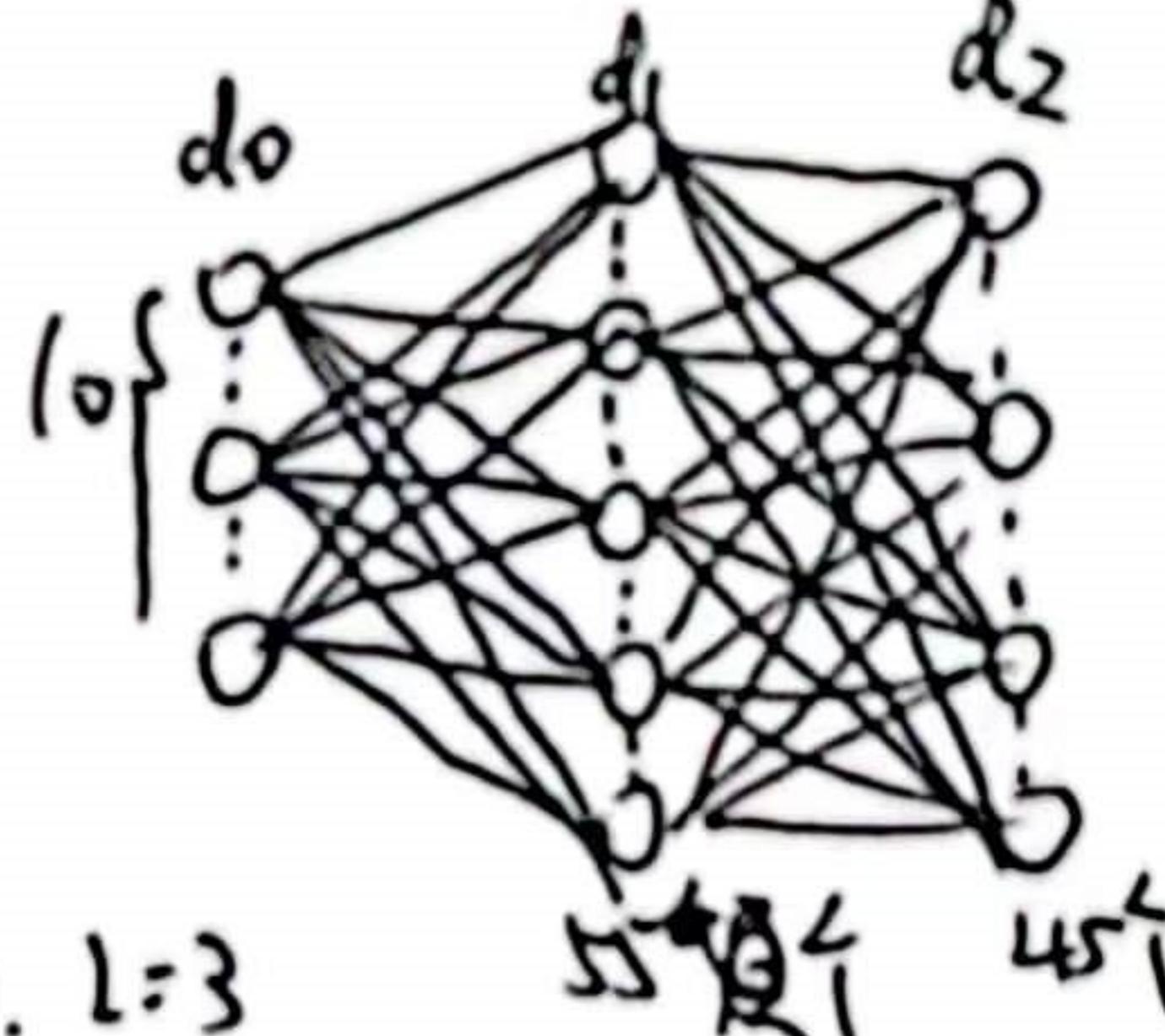
$$x = 54 \quad = -\frac{(x-50)^2 + 1000}{100} - y^2 - 10z + 90y$$

$$1000 < \frac{3025}{3025} < 3025 \quad = \frac{3000}{3025} - (y-45)^2 - 10z$$

$$3015 < \quad = 3025 - 0 - 10 = 3015$$

因此，在 $L=2$ 时，令 $d(1)=55, d(2)=45$ 取得最大值

神经网络架构：



$$\sum_{i=1}^L d(i-1)d(i) = 3025 \quad (\text{最大值})$$

minimize ($\frac{\text{最小值}}{\text{最大值}}$)：也就是说 $L=1, L=2, L=3$

$$\textcircled{1} L=1 \text{ 时, } \text{与求最大值相同} \quad \sum_{i=1}^1 d^{(1-i)}d^{(1)} = 10 \times 100 = 1000$$

$\textcircled{2} L=2 \text{ 时, 同样设 } d(1)=x, d(2)=y$

$$W_2 = -(x-55)^2 + 3025 \quad \leq |x-55| \text{ 最大}$$

$$x > 0 \quad x=1 \quad W_2 = -(1-55)^2 + 3025 = -1 + 110 \times 1 = -1 + 110 = 109$$

$\textcircled{3} L=3 \text{ 时, 同样设 } x, y, z = d_1, d_2, d_3$

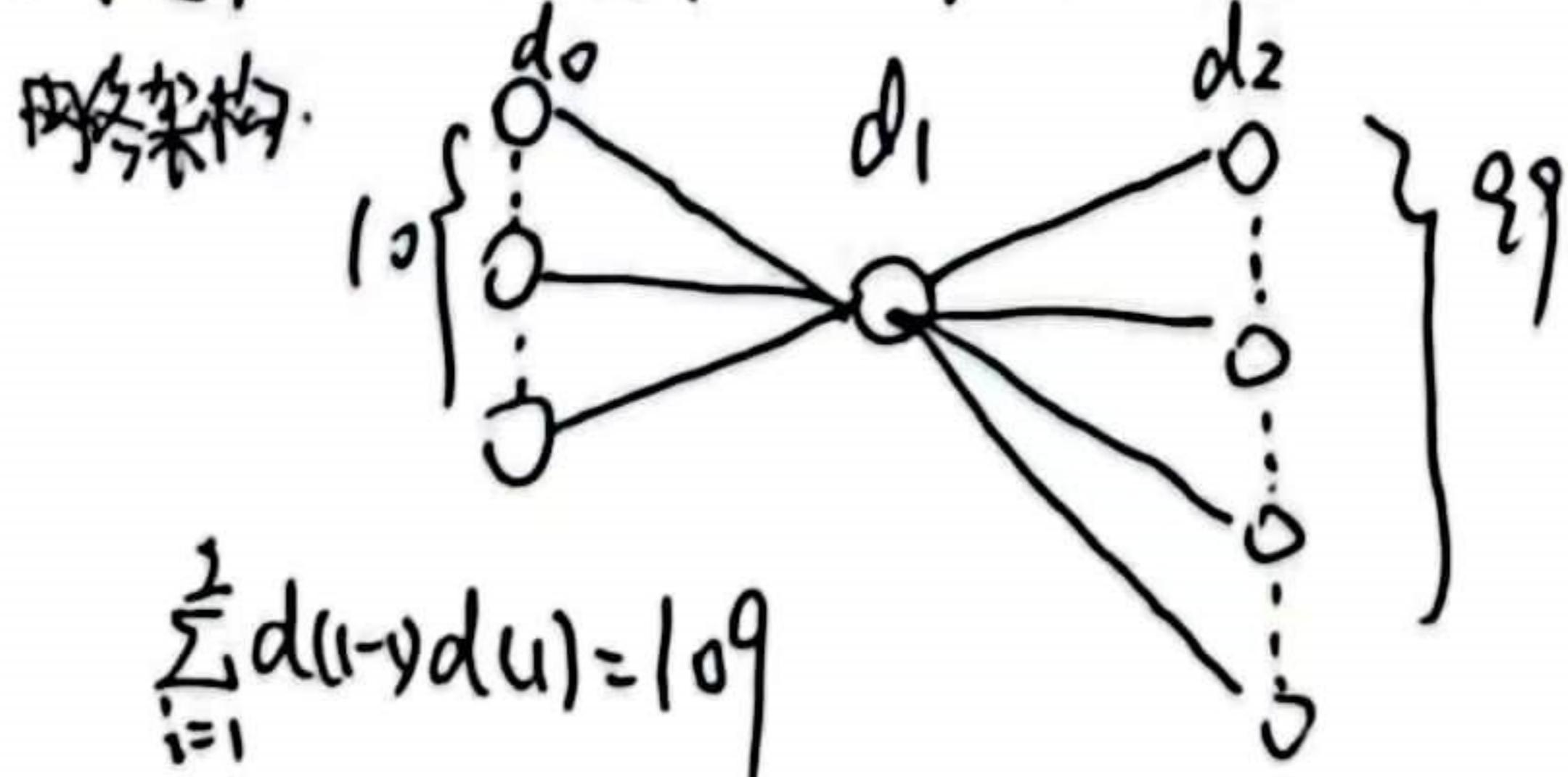
$$W = - (y-45)^2 - 10z + 3025 \quad x, y, z > 0$$

$z \text{ 的约束 } -10 < 0, |y-45| \text{ 最大值, } Z \text{ 的最大值 } x=1, y=1$

$$W = 10x + 10y + 10z = 109 \quad z = 100 - 2 = 98$$

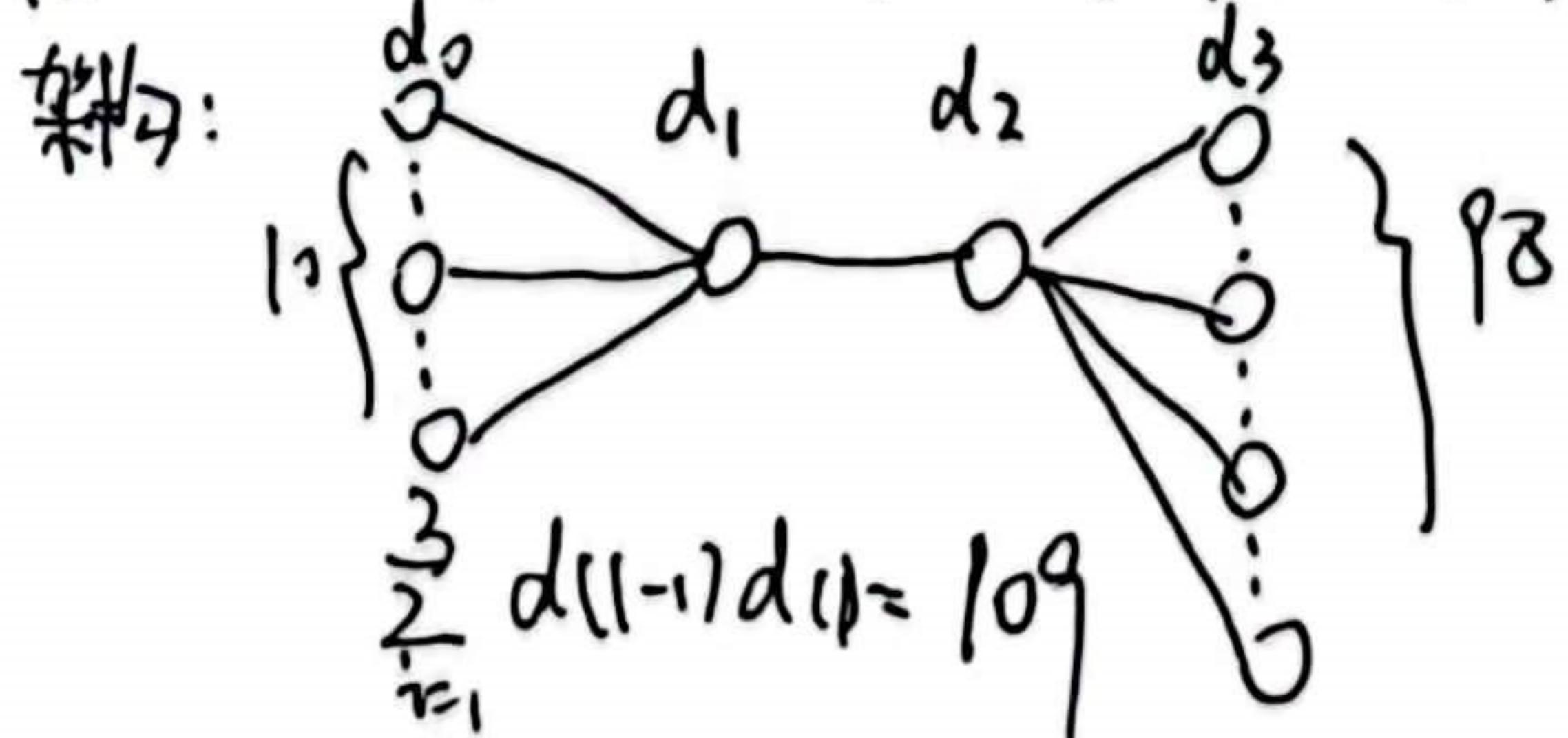
$1000 > 109 \Rightarrow 109 = 109 \quad \text{最小值有多种情况}$

情况1： $L=2, d(1)=1, d(2)=99$



$$\sum_{i=1}^L d(i-1)d(i) = 109$$

情况2： $L=3, d(1)=d(2)=1, d(3)=98$



$$\sum_{i=1}^L d(i-1)d(i) = 109$$