

Deep Learning: Homework #1

Due on March 7, 2021

Professor Zhen Li

Peng Deng

Problem 1

Cross entropy is often used as the objective function when training neural networks in classification problems. Suppose the training set includes N training pairs $D = \left\{ \left(\mathbf{x}_i^{(\text{train})}, y_i^{(\text{train})} \right) \right\}_{i=1}^N$, where $\mathbf{x}_i^{(\text{train})}$ is a training sample and $y_i^{(\text{train})} \in \{1, \dots, c\}$ is its corresponding class label. \mathbf{z}_i is the output of the network given input $\mathbf{x}_i^{(\text{train})}$ and the nonlinearity of the output layer is softmax. \mathbf{z}_i is a c dimensional vector, $z_{i,k} \in [0, 1]$ and $\sum_{k=1}^c z_{i,k} = 1$. The questions are as follows.

- (1) Write the objective function of cross-entropy with softmax activation function, and calculate the gradient of hidden-to-output and input-to-hidden weights as we introduced in class.
- (2) Verify it is equivalent to the negative log-likelihood on the training set, assuming the training samples are independent.

Solution

Subproblem (1)

We use one-hot encoding to encode the target value for the training sample $\mathbf{x}_i^{(\text{train})}$ as $\mathbf{t}_i (i = 1, 2, \dots, N)$, which is a c dimensional vector. Besides, it is defined as follows

$$t_{i,k} = \begin{cases} 1 & (k = y_i^{(\text{train})}) \\ 0 & (k \neq y_i^{(\text{train})}) \end{cases} \quad (1)$$

- The objective function can be written as follow. We use \mathbf{W} to denote all the weights in the network.

$$J(\mathbf{W}) = -\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^c t_{i,k} \log(z_{i,k}) \quad (2)$$

- The gradient of hidden-to-output weights is calculated as follow. We use net to denote the net value in the output layer. We set δ to denote the sensitivity in the output layer. We use w to denote the weights of hidden-to-output.

$$\begin{aligned} \frac{\partial J}{\partial w_{kj}} &= \sum_{i=1}^N \frac{\partial J}{\partial net_{i,k}} \frac{\partial net_{i,k}}{\partial w_{kj}} \\ &= \sum_{i=1}^N -\delta_{i,k} \frac{\partial net_{i,k}}{\partial w_{kj}} \end{aligned} \quad (3)$$

Then, we can calculate $\delta_{i,k}$ and δ_i as follow

$$\begin{aligned} \delta_{i,k} &= -\frac{\partial J}{\partial net_{i,k}} = -\frac{\partial J}{\partial z_{i,k}} \frac{\partial z_{i,k}}{\partial net_{i,k}} \\ \delta_i &= -\frac{\partial J}{\partial net_i} = -\frac{\partial J}{\partial \mathbf{z}_i} \frac{\partial \mathbf{z}_i}{\partial net_i} \end{aligned} \quad (4)$$

Set $m = y_i^{(\text{train})}$, then we have $t_{i,m} = 1$, and $\mathbf{t}_i = (0, \dots, 1, \dots, 0)$. Then we can get $\frac{\partial J}{\partial \mathbf{z}_i}$ as follow

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{z}_i} &= \left(\frac{\partial J}{\partial z_{i,1}} \quad \frac{\partial J}{\partial z_{i,2}} \quad \dots \quad \frac{\partial J}{\partial z_{i,c}} \right) \\ &= \left(0 \quad 0 \quad \dots \quad -\frac{1}{N} \frac{1}{z_{i,m}} \quad \dots \quad 0 \right) \end{aligned} \quad (5)$$

Then we can calculate $\frac{\partial \mathbf{z}_i}{\partial \mathbf{net}_i}$ as follow

$$\begin{aligned} \frac{\partial \mathbf{z}_i}{\partial \mathbf{net}_i} &= \begin{pmatrix} \frac{\partial z_{i,1}}{\partial \mathbf{net}_{i,1}} & \frac{\partial z_{i,1}}{\partial \mathbf{net}_{i,2}} & \cdots & \frac{\partial z_{i,1}}{\partial \mathbf{net}_{i,m}} & \cdots & \frac{\partial z_{i,1}}{\partial \mathbf{net}_{i,c}} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \frac{\partial z_{i,m}}{\partial \mathbf{net}_{i,1}} & \frac{\partial z_{i,m}}{\partial \mathbf{net}_{i,2}} & \cdots & \frac{\partial z_{i,m}}{\partial \mathbf{net}_{i,m}} & \cdots & \frac{\partial z_{i,m}}{\partial \mathbf{net}_{i,c}} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \frac{\partial z_{i,c}}{\partial \mathbf{net}_{i,1}} & \frac{\partial z_{i,c}}{\partial \mathbf{net}_{i,2}} & \cdots & \frac{\partial z_{i,c}}{\partial \mathbf{net}_{i,m}} & \cdots & \frac{\partial z_{i,c}}{\partial \mathbf{net}_{i,c}} \end{pmatrix} \\ &= \begin{pmatrix} \frac{\partial z_{i,1}}{\partial \mathbf{net}_{i,1}} & \frac{\partial z_{i,1}}{\partial \mathbf{net}_{i,2}} & \cdots & \frac{\partial z_{i,1}}{\partial \mathbf{net}_{i,m}} & \cdots & \frac{\partial z_{i,1}}{\partial \mathbf{net}_{i,c}} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ -z_{i,m} \cdot z_{i,1} & -z_{i,m} \cdot z_{i,2} & \cdots & z_{i,m} - z_{i,m}^2 & \cdots & -z_{i,m} \cdot z_{i,c} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \frac{\partial z_{i,c}}{\partial \mathbf{net}_{i,1}} & \frac{\partial z_{i,c}}{\partial \mathbf{net}_{i,2}} & \cdots & \frac{\partial z_{i,c}}{\partial \mathbf{net}_{i,m}} & \cdots & \frac{\partial z_{i,c}}{\partial \mathbf{net}_{i,c}} \end{pmatrix} \end{aligned} \quad (6)$$

By combination of equation 5 and equation 6, we can recalculate $\delta_{i,k}$ and δ_i as follow

$$\begin{aligned} \delta_i &= -\frac{\partial J}{\partial \mathbf{z}_i} \frac{\partial \mathbf{z}_i}{\partial \mathbf{net}_i} \\ &= -\begin{pmatrix} 0 & 0 & \cdots & -\frac{1}{N} \frac{1}{z_{i,m}} & \cdots & 0 \end{pmatrix} \begin{pmatrix} \frac{\partial z_{i,1}}{\partial \mathbf{net}_{i,1}} & \frac{\partial z_{i,1}}{\partial \mathbf{net}_{i,2}} & \cdots & \frac{\partial z_{i,1}}{\partial \mathbf{net}_{i,m}} & \cdots & \frac{\partial z_{i,1}}{\partial \mathbf{net}_{i,c}} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ -z_{i,m} \cdot z_{i,1} & -z_{i,m} \cdot z_{i,2} & \cdots & z_{i,m} - z_{i,m}^2 & \cdots & -z_{i,m} \cdot z_{i,c} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \frac{\partial z_{i,c}}{\partial \mathbf{net}_{i,1}} & \frac{\partial z_{i,c}}{\partial \mathbf{net}_{i,2}} & \cdots & \frac{\partial z_{i,c}}{\partial \mathbf{net}_{i,m}} & \cdots & \frac{\partial z_{i,c}}{\partial \mathbf{net}_{i,c}} \end{pmatrix} \\ &= \begin{pmatrix} -\frac{1}{N} z_{i,1} & -\frac{1}{N} z_{i,2} & \cdots & \frac{1}{N} (1 - z_{i,m}) & \cdots & -\frac{1}{N} z_{i,c} \end{pmatrix} \\ &= \frac{1}{N} \begin{pmatrix} 0 & 0 & \cdots & 1 & \cdots & 0 \end{pmatrix} - \begin{pmatrix} z_{i,1} & z_{i,2} & \cdots & z_{i,m} & \cdots & z_{i,c} \end{pmatrix} \\ &= \frac{1}{N} (\mathbf{t}_i - \mathbf{z}_i) \\ \delta_{i,k} &= \frac{1}{N} (t_{i,k} - z_{i,k}) \end{aligned} \quad (7)$$

Back to equation 3, then we can recalculate the gradient of hidden-to-output weights as follow

$$\begin{aligned} \frac{\partial J}{\partial w_{kj}} &= \sum_{i=1}^N -\delta_{i,k} \frac{\partial \mathbf{net}_{i,k}}{\partial w_{kj}} \\ &= \sum_{i=1}^N -\frac{1}{N} (t_{i,k} - z_{i,k}) y_{i,j} \\ &= -\frac{1}{N} \sum_{i=1}^N (t_{i,k} - z_{i,k}) y_{i,j} \end{aligned} \quad (8)$$

Thus, we can write down the matrix format of the gradient of hidden-to-output weights as follow. We suppose there are n_H unit in the hidden layer.

$$\frac{\partial J}{\partial \mathbf{w}} = -\frac{1}{N} \sum_{i=1}^N (\mathbf{t}_i - \mathbf{z}_i)^\top \mathbf{y}_i \quad (9)$$

Where $(\mathbf{t}_i - \mathbf{z}_i)^\top$ is a $c \times 1$ column vector and \mathbf{y}_i is a $1 \times n_H$ row vector.

- We suppose $\mathbf{x}_i^{(\text{train})}$ is a d dimensional vector, and we use q to denote the q^{th} unit. The gradient of input-to-hidden weights is calculated as follow. We use net' to denote the net value in the hidden layer. We set δ' to denote the sensitivity in the output layer. We use w' to denote the weights of hidden-to-output.

$$\begin{aligned} \frac{\partial J}{\partial w'_{jq}} &= \sum_{i=1}^N \frac{\partial J}{\partial \text{net}'_{i,j}} \frac{\partial \text{net}'_{i,j}}{\partial w'_{jq}} \\ &= \sum_{i=1}^N -\delta'_{i,j} \frac{\partial \text{net}'_{i,j}}{\partial w'_{jq}} \end{aligned} \quad (10)$$

Then, we can calculate $\delta'_{i,j}$ and δ'_i as follow

$$\begin{aligned} \delta'_{i,j} &= -\frac{\partial J}{\partial \text{net}'_{i,j}} = -\frac{\partial J}{\partial y_{i,j}} \frac{\partial y_{i,j}}{\partial \text{net}'_{i,j}} \\ \delta'_i &= -\frac{\partial J}{\partial \text{net}_i} = -\frac{\partial J}{\partial \mathbf{y}_i} \frac{\partial \mathbf{y}_i}{\partial \text{net}_i} \end{aligned} \quad (11)$$

Then we can get $\frac{\partial J}{\partial \mathbf{y}_i}$ as follows

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{y}_i} &= \frac{\partial J}{\partial \mathbf{z}_i} \frac{\partial \mathbf{z}_i}{\partial \mathbf{y}_i} = \frac{\partial J}{\partial \mathbf{z}_i} \frac{\partial \mathbf{z}_i}{\partial \text{net}_i} \frac{\partial \text{net}_i}{\partial \mathbf{y}_i} \\ &= -\delta_i \frac{\partial \text{net}_i}{\partial \mathbf{y}_i} = -\delta_i \mathbf{w} \end{aligned} \quad (12)$$

Where δ_i is a $1 \times c$ row vector and \mathbf{w} is a $c \times n_H$ matrix. We can get $\frac{\partial \mathbf{y}_i}{\partial \text{net}_i}$ as follow

$$\begin{aligned} \frac{\partial \mathbf{y}_i}{\partial \text{net}_i} &= \begin{pmatrix} \frac{\partial y_{i,1}}{\partial \text{net}'_{i,1}} & \frac{\partial y_{i,1}}{\partial \text{net}'_{i,2}} & \dots & \frac{\partial y_{i,1}}{\partial \text{net}'_{i,n_H}} \\ \frac{\partial y_{i,2}}{\partial \text{net}'_{i,1}} & \frac{\partial y_{i,2}}{\partial \text{net}'_{i,2}} & \dots & \frac{\partial y_{i,2}}{\partial \text{net}'_{i,n_H}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_{i,n_H}}{\partial \text{net}'_{i,1}} & \frac{\partial y_{i,n_H}}{\partial \text{net}'_{i,2}} & \dots & \frac{\partial y_{i,n_H}}{\partial \text{net}'_{i,n_H}} \end{pmatrix} \\ &= \begin{pmatrix} y_{i,1} - y_{i,1}^2 & -y_{i,1} \cdot y_{i,2} & \dots & -y_{i,1} \cdot y_{i,n_H} \\ -y_{i,2} \cdot y_{i,1} & y_{i,2} - y_{i,2}^2 & \dots & -y_{i,2} \cdot y_{i,n_H} \\ \vdots & \vdots & \ddots & \vdots \\ -y_{i,n_H} \cdot y_{i,1} & -y_{i,n_H} \cdot y_{i,2} & \dots & y_{i,n_H} - y_{i,n_H}^2 \end{pmatrix} \\ &= \text{diag}(\mathbf{y}_i) - \mathbf{y}_i^\top \mathbf{y}_i \end{aligned} \quad (13)$$

By combination of equation 12 and equation 13, we can recalculate $\delta'_{i,k}$ and δ'_i as follow

$$\begin{aligned}\delta'_i &= -\frac{\partial J}{\partial \mathbf{y}_i} \frac{\partial \mathbf{y}_i}{\partial \mathbf{net}_i'} \\ &= \delta_i \mathbf{w} \frac{\partial \mathbf{y}_i}{\partial \mathbf{net}_i'} \\ &= \frac{1}{N} (\mathbf{t}_i - \mathbf{z}_i) \mathbf{w} (\text{diag}(\mathbf{y}_i) - \mathbf{y}_i^\top \mathbf{y}_i)\end{aligned}\tag{14}$$

Back to equation 10, then we can recalculate the gradient of input-to-hidden weights as follow

$$\begin{aligned}\frac{\partial J}{\partial \mathbf{w}'} &= -\sum_{i=1}^N \delta_i'^\top \mathbf{x}_i = -\sum_{i=1}^N \left(\delta_i \mathbf{w} \frac{\partial \mathbf{y}_i}{\partial \mathbf{net}_i'} \right)^\top \mathbf{x}_i \\ &= -\frac{1}{N} \sum_{i=1}^N (\text{diag}(\mathbf{y}_i) - \mathbf{y}_i^\top \mathbf{y}_i)^\top \mathbf{w}^\top (\mathbf{t}_i - \mathbf{z}_i)^\top \mathbf{x}_i \\ &= -\frac{1}{N} \sum_{i=1}^N (\text{diag}(\mathbf{y}_i) - \mathbf{y}_i^\top \mathbf{y}_i) \mathbf{w}^\top (\mathbf{t}_i - \mathbf{z}_i)^\top \mathbf{x}_i\end{aligned}\tag{15}$$

Subproblem (2)

We suppose $z_{i,k} \in [0, 1]$ and $\sum_{k=1}^c z_{i,k} = 1$, which is the probability of $\mathbf{x}_i^{(\text{train})}$ belongs to class k . Thus, we can get the likelihood of $\mathbf{x}_i^{(\text{train})}$ as follow

$$p_i = \prod_{k=1}^c (z_{i,k})^{t_{i,k}}\tag{16}$$

Then we can get the likelihood function as follow

$$L = \prod_{i=1}^N p_i = \prod_{i=1}^N \prod_{k=1}^c (z_{i,k})^{t_{i,k}}\tag{17}$$

Then we can get the log-likelihood function as follow

$$\begin{aligned}\log L &= \sum_{i=1}^N \log \left(\prod_{k=1}^c (z_{i,k})^{t_{i,k}} \right) \\ &= \sum_{i=1}^N \sum_{k=1}^c t_{i,k} \log z_{i,k}\end{aligned}\tag{18}$$

According to equation 2 and equation 18, we can find that

$$\begin{aligned}J(\mathbf{W}) &= -\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^c t_{i,k} \log(z_{i,k}) \\ \log L &= \sum_{i=1}^N \sum_{k=1}^c t_{i,k} \log z_{i,k}\end{aligned}\tag{19}$$

Because we know that $J(\mathbf{W})$ can be expressed as follow as well

$$J(\mathbf{W}) = -\sum_{i=1}^N \sum_{k=1}^c t_{i,k} \log(z_{i,k})\tag{20}$$

Thus, we can see that $J(\mathbf{W})$ is equivalent to the negative log-likelihood on the training set.

Problem 2

x_1 and x_2 are two input variables, and y is the target variable to be predicted. The network structure is shown in Figure 1(a). $h_{11} = f_{11}(x_1)$, $h_{12} = f_{12}(x_2)$, and $y = g(h_{11}, h_{12})$. The questions are as follows.

- (1) Assuming $x_1 \in [0, 1]$ and $x_2 \in [0, 1]$, in order to obtain the decision regions in Figure 1(b), decide functions f_{11} , f_{12} , and g .
- (2) Now we extend the range of x_1 and x_2 to $[0, 2]$. Please add one more layer to Figure 1(a) in order to obtain the decision regions in Figure 1(c).
- (3) Although the decision boundaries in Figure 1(c) look complicated, there exist regularity and global structure. Please explain such regularity and global structure. Based on your observation, draw the decision boundaries when the range of x_1 and x_2 to $[0, 4]$.
- (4) Following the question above and assuming the range of x_1 and x_2 is extended to $[0, 2^n]$, draw the network structure and the transform function in each layer, in order to obtain the decision regions with the same regularity and global structure in Figure 1(b) and (c). The complexity of computation units should be $O(n)$.
- (5) Assuming the range of x_1 and x_2 is $[0, 2^n]$ and only one hidden layer is allowed, specify the network structure and transform functions.

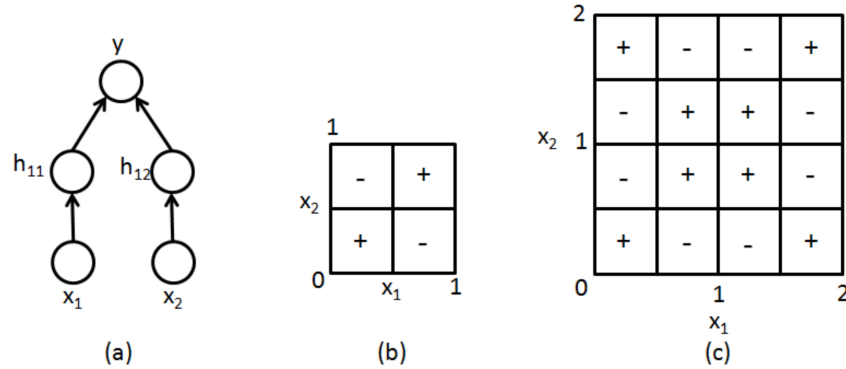


Figure 1 Figure for problem 2

Solution

Subproblem (1)

We can decide the functions as follows

$$\begin{aligned}
 f_{11}(x) &= x - \frac{1}{2} \implies h_{11} = f_{11}(x_1) = x_1 - \frac{1}{2} \implies h_{11} \in \left[-\frac{1}{2}, \frac{1}{2}\right] \\
 f_{12}(x) &= x - \frac{1}{2} \implies h_{12} = f_{12}(x_2) = x_2 - \frac{1}{2} \implies h_{12} \in \left[-\frac{1}{2}, \frac{1}{2}\right] \\
 g(u, v) &= u \cdot v \implies y = g(h_{11}, h_{12}) = h_{11} \cdot h_{12}
 \end{aligned} \tag{21}$$

Thus, when $y > 0$, we classify the point to positive class, and when $y < 0$, we classify the point to negative class.

Subproblem (2)

When we extend the range of x_1 and x_2 to $[0, 2]$, we can add one more layer to Figure 1(a) and then we can obtain the decision regions in Figure 1(c). The new network structure is showed as Figure 2.

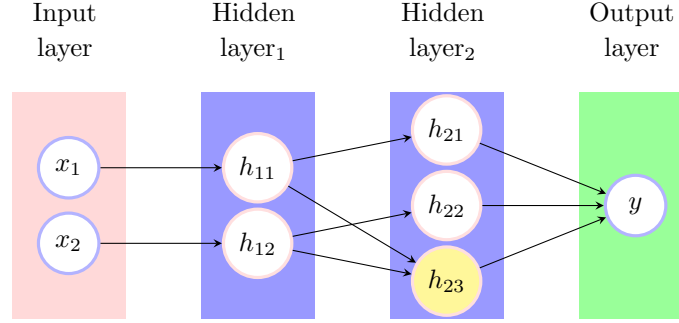


Figure 2 The new network structure

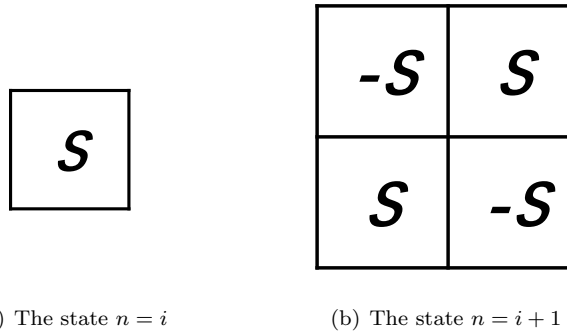
Besides, the relationships between units in different layers are decided as follows. The yellow unit is used to transmit symbols $(+, -)$.

$$\begin{aligned}
 h_{11} &= x_1 - 1 \implies h_{11} \in [-1, 1] \\
 h_{12} &= x_2 - 1 \implies h_{12} \in [-1, 1] \\
 h_{21} &= h_{11} - \text{sgn}(h_{11}) \cdot \frac{1}{2} \implies h_{21} \in \left[-\frac{1}{2}, \frac{1}{2}\right] \\
 h_{22} &= h_{12} - \text{sgn}(h_{12}) \cdot \frac{1}{2} \implies h_{22} \in \left[-\frac{1}{2}, \frac{1}{2}\right] \\
 h_{23} &= h_{11} \cdot h_{12} \\
 y &= h_{21} \cdot h_{22} \cdot h_{23}
 \end{aligned} \tag{22}$$

Thus, when $y > 0$, we classify the point to positive class, and when $y < 0$, we classify the point to negative class.

Subproblem (3)

As we can see, there is a recurrence process in this problem and we can get the recurrence formula between two states as follow. We define the decision boundaries when $n = i$ as Figure 3(a), and then we can get decision boundaries when $n = i + 1$ as Figure 3(b). Thus, we get the recurrence formula of the decision boundaries between state $n = i$ and state $n = i + 1$ as Figure 3, which is the regularity and global structure. Then we can draw the decision boundaries when the range of x_1 and x_2 to $[0, 4]$ as Figure 4, and we can see in this situation, $n = 2$.

Figure 3 The recurrence formula of the decision boundaries between state $n = i$ and state $n = i + 1$

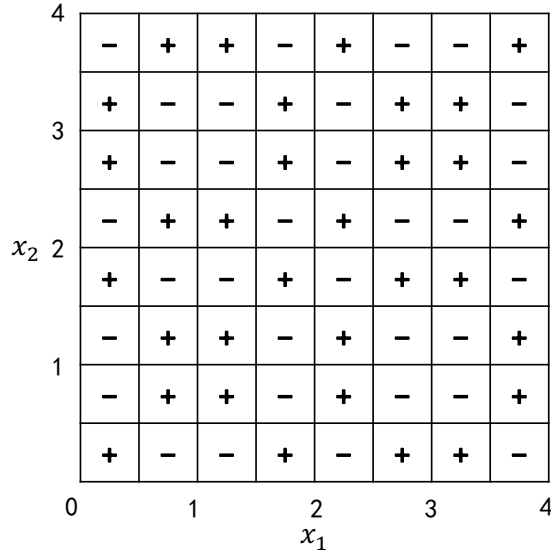


Figure 4 The decision boundaries when the range of x_1 and x_2 to $[0, 4]$

Subproblem (4)

When the range of x_1 and x_2 is extended to $[0, 2^n]$, the network structure is showed as Figure 5.

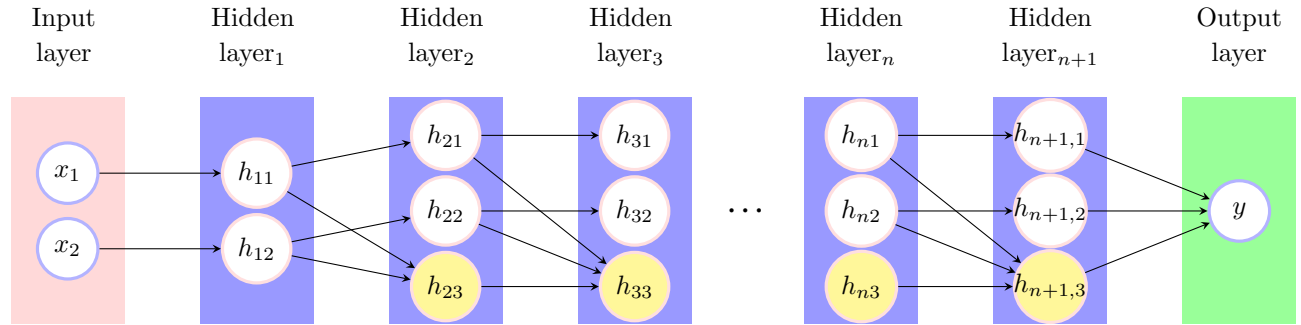


Figure 5 The network structure

The relationships between units in different layers are decided as follows. The yellow unit is used to transmit symbols (+, -).

◦ For input to hidden layer

$$\begin{aligned} h_{11} &= x_1 - 2^{n-1} \implies h_{11} \in [-2^{n-1}, 2^{n-1}] \\ h_{12} &= x_2 - 2^{n-1} \implies h_{12} \in [-2^{n-1}, 2^{n-1}] \end{aligned} \quad (23)$$

◦ For hidden layer 1 to hidden layer 2

$$\begin{aligned} h_{21} &= h_{11} - \text{sng}(h_{11}) \cdot 2^{n-2} \implies h_{21} \in [-2^{n-2}, 2^{n-2}] \\ h_{22} &= h_{12} - \text{sng}(h_{12}) \cdot 2^{n-2} \implies h_{22} \in [-2^{n-2}, 2^{n-2}] \\ h_{23} &= h_{11} \cdot h_{12} \end{aligned} \quad (24)$$

- For hidden layer i to hidden layer $i + 1$, $i = 2, 3, \dots, n$

$$\begin{aligned} h_{i+1,1} &= h_{i1} - \text{sng}(h_{i1}) \cdot 2^{n-i-1} \implies h_{i+1,1} \in [-2^{n-i-1}, 2^{n-i-1}] \\ h_{i+1,2} &= h_{i2} - \text{sng}(h_{i2}) \cdot 2^{n-i-1} \implies h_{i+1,2} \in [-2^{n-i-1}, 2^{n-i-1}] \\ h_{i+1,3} &= h_{i1} \cdot h_{i2} \cdot h_{i3} \end{aligned} \quad (25)$$

- For the last hidden layer (hidden layer $n+1$) to the output layer

$$y = h_{n+1,1} \cdot h_{n+1,2} \cdot h_{n+1,3} \quad (26)$$

Thus, when $y > 0$, we classify the point to positive class, and when $y < 0$, we classify the point to negative class.

Then, we can get the complexity of computation units as follow. We define the number of units as Q .

$$\begin{aligned} Q &= Q_{\text{input}} + Q_{\text{hidden}} + Q_{\text{output}} \\ &= 2 + (n + 1) \cdot 3 - 1 + 1 \\ &= 3n + 5 \end{aligned} \quad (27)$$

Thus, we can see the the complexity of computation units is $O(n)$. With only one hidden layer, the network structure is specified as Figure 6.

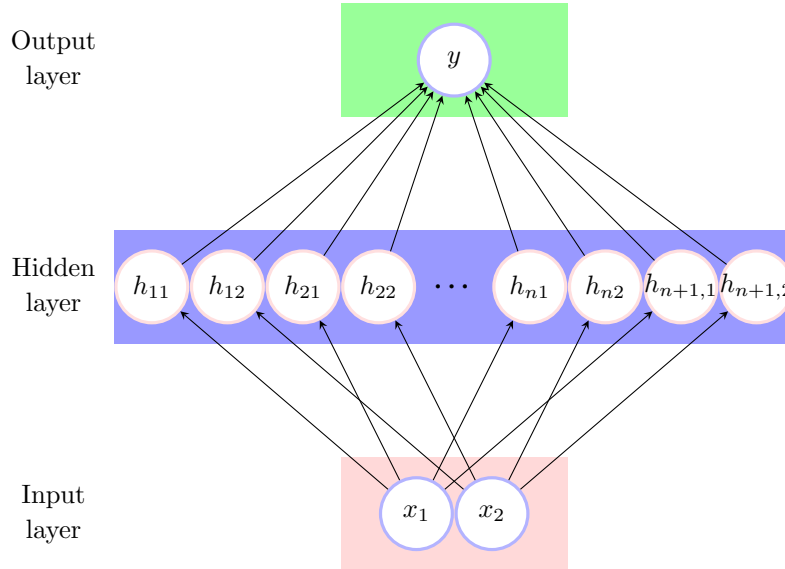


Figure 6 The network structure with only one hidden layer

The relationships between units in different layers are decided as follows.

- For input layer to h_{11}, h_{12}

$$\begin{aligned} h_{11} &= x_1 - 2^{n-1} \implies h_{11} \in [-2^{n-1}, 2^{n-1}] \\ h_{12} &= x_2 - 2^{n-1} \implies h_{12} \in [-2^{n-1}, 2^{n-1}] \end{aligned} \quad (28)$$

- For input layer to $h_{i+1,1}, h_{i+1,2}$, $i = 1, 3, \dots, n$.

$$\begin{aligned} h_{i+1,1} &= h_{i1} - \text{sng}(h_{i1}) \cdot 2^{n-i-1} \implies \text{we can compute it iteratively and } h_{i+1,1} \in [-2^{n-i-1}, 2^{n-i-1}] \\ h_{i+1,2} &= h_{i2} - \text{sng}(h_{i2}) \cdot 2^{n-i-1} \implies \text{we can compute it iteratively and } h_{i+1,2} \in [-2^{n-i-1}, 2^{n-i-1}] \end{aligned} \quad (29)$$

◦ For the hidden layer to the output layer

$$y = \prod_{i=1}^{n+1} h_{i1} \cdot h_{i2} \quad (30)$$

Thus, when $y > 0$, we classify the point to positive class, and when $y < 0$, we classify the point to negative class.

Problem 3

Siamese neural nets have an interesting architecture – the same parameters and functions are used to evaluate 2 inputs. As one might expect, Siamese nets are useful to train similarity metrics, evaluations of how “close” inputs are. These nets have been applied to facial recognition tasks with a good deal of success, but in this example, we’ll see how to train Siamese nets to learn a distance metric for two inputs, e.g., word vectors. One might imagine training a net to map word vectors across languages, discover synonyms or antonyms, etc.

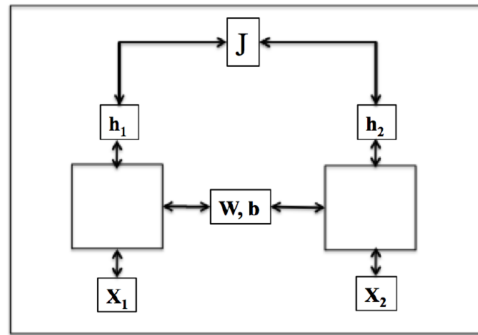


Figure 7 The Siamese neural nets

Here is one such model to evaluate how similar two inputs are using Euclidean distance. There are two inputs $x_1, x_2 \in R^n$, shared parameters $W \in R^{m \times n}$ and $b \in R^m$, and a single hidden layer associated with each input:

$$\begin{aligned} h_1 &= \sigma(Wx_1 + b) \\ h_2 &= \sigma(Wx_2 + b) \end{aligned}$$

We evaluate the distance between the two activations h_1, h_2 using Euclidean distance as our similarity metric. The model objective J is

$$J = \frac{1}{2} \|h_1 - h_2\|_F^2 + \frac{\lambda}{2} \|W\|_F^2$$

where λ is a given regularization parameter. (The Frobenius norm $\|\cdot\|_F$ is a matrix norm defined by $\|A\|_F = \sqrt{\sum_{i,j} |A_{ij}|^2}$). The questions are as follows.

- (1) Calculate the gradient $\nabla_W J$ and $\nabla_b J$.
- (2) Write out the (vanilla) gradient descent update rules for the model parameters for a single training example (with arbitrary step size α).
- (3) If $W \in R^{10 \times 5}$ and $b \in R^{10 \times 1}$, how many parameters does the model have?

Now imagine you wanted to see how ReLU/sigmoid nonlinearities might affect training on single input. But instead of training two separate nets, you want to train a pseudo-Siamese net like Figure 8.

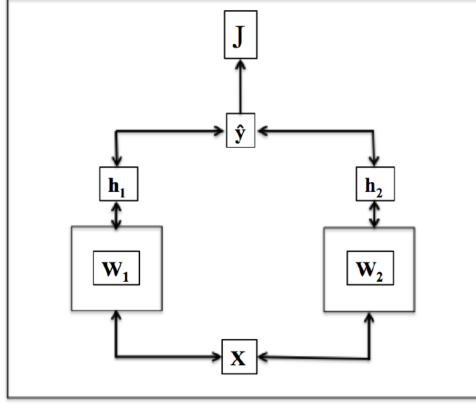


Figure 8 The pseudo-Siamese neural nets

The whole model would be changed to

$$\begin{aligned} \mathbf{h}_1 &= \sigma(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) \\ \mathbf{h}_2 &= \text{Relu}(\mathbf{W}_2 \mathbf{x} + \mathbf{b}_2) \\ \hat{\mathbf{y}} &= \text{softmax}(\mathbf{W}_3(\mathbf{h}_1 + \mathbf{h}_2) + \mathbf{b}_3) \end{aligned}$$

where $\mathbf{x} \in R^n$, $\mathbf{W}_1, \mathbf{W}_2 \in R^{m \times n}$, $\mathbf{W}_3 \in R^{k \times m}$, $\mathbf{b}_1, \mathbf{b}_2 \in R^m$ and $\mathbf{b}_3 \in R^k$. We calculate this model for N examples and k classes with cross-entropy loss

$$J = -\frac{1}{N} \sum_{j=1}^N \sum_{i=1}^k y_j^i \log(\hat{y}_j^i)$$

where \mathbf{y}_j is the one-hot vector for example j with all probability mass on the correct class and $\hat{\mathbf{y}}_j$ are the softmax scores for example j .

(4) Calculate $\nabla_{\mathbf{h}_1} J$, $\nabla_{\mathbf{h}_2} J$, and $\nabla_{\mathbf{x}} J$.

(5) For $\mathbf{W}_1, \mathbf{W}_2$, which one is likely to train faster, please explain it.

Solution

Subproblem (1)

◦ For $\nabla_{\mathbf{W}} J$

$$\begin{aligned} \frac{\partial J}{\partial w_{ji}} &= \frac{\frac{\partial}{2} \|\mathbf{h}_1 - \mathbf{h}_2\|_F^2}{\mathbf{h}_1 - \mathbf{h}_2} \cdot \frac{\partial(\mathbf{h}_1 - \mathbf{h}_2)}{\partial w_{ji}} + \frac{\partial \frac{\lambda}{2} \|\mathbf{W}\|_F^2}{\partial w_{ji}} \\ &= (\mathbf{h}_1 - \mathbf{h}_2)^\top \left(\frac{\partial \mathbf{h}_1}{\partial w_{ji}} - \frac{\partial \mathbf{h}_2}{\partial w_{ji}} \right) + \lambda w_{ji} \\ &= (\mathbf{h}_1 - \mathbf{h}_2)^\top \begin{pmatrix} \frac{\partial h_{11}}{\partial w_{ji}} - \frac{\partial h_{21}}{\partial w_{ji}} \\ \frac{\partial h_{12}}{\partial w_{ji}} - \frac{\partial h_{22}}{\partial w_{ji}} \\ \vdots \\ \frac{\partial h_{1m}}{\partial w_{ji}} - \frac{\partial h_{2m}}{\partial w_{ji}} \end{pmatrix} + \lambda w_{ji} = (\mathbf{h}_1 - \mathbf{h}_2)^\top \begin{pmatrix} 0 \\ 0 \\ \vdots \\ \sigma'(\mathbf{w}_j \mathbf{x}_1 + \mathbf{b}_j) \cdot x_{1i} - \sigma'(\mathbf{w}_j \mathbf{x}_2 + \mathbf{b}_j) \cdot x_{2i} \\ \vdots \\ 0 \end{pmatrix} + \lambda w_{ji} \\ &= (h_{1j} - h_{2j}) \cdot \left(\sigma'(\mathbf{w}_j \mathbf{x}_1 + \mathbf{b}_j) \cdot x_{1i} - \sigma'(\mathbf{w}_j \mathbf{x}_2 + \mathbf{b}_j) \cdot x_{2i} \right) + \lambda w_{ji} \end{aligned} \tag{31}$$

Thus, we can get the gradient $\nabla_{\mathbf{W}} J$ as follow

$$\nabla_{\mathbf{W}} J = \left[(\mathbf{h}_1 - \mathbf{h}_2) \circ \sigma'(\mathbf{W}\mathbf{x}_1 + \mathbf{b}) \right] \cdot \mathbf{x}_1^\top - \left[(\mathbf{h}_1 - \mathbf{h}_2) \circ \sigma'(\mathbf{W}\mathbf{x}_2 + \mathbf{b}) \right] \cdot \mathbf{x}_2^\top + \lambda \mathbf{W} \quad (32)$$

Where \mathbf{h}_1 , \mathbf{h}_2 , and \mathbf{b} are $m \times 1$ column vectors, and \mathbf{x}_1 , \mathbf{x}_2 are $n \times 1$ column vectors.

◦ For $\nabla_{\mathbf{b}} J$

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{b}} &= \frac{\frac{1}{2} \|\mathbf{h}_1 - \mathbf{h}_2\|_F^2}{\mathbf{h}_1 - \mathbf{h}_2} \cdot \frac{\partial (\mathbf{h}_1 - \mathbf{h}_2)}{\partial \mathbf{b}} \\ &= (\mathbf{h}_1 - \mathbf{h}_2)^\top \cdot \left(\frac{\partial \mathbf{h}_1}{\partial \mathbf{b}} - \frac{\partial \mathbf{h}_2}{\partial \mathbf{b}} \right) \\ &= (\mathbf{h}_1 - \mathbf{h}_2)^\top \cdot \begin{pmatrix} \sigma'(\mathbf{w}_1 \mathbf{x}_1 + b_1) - \sigma'(\mathbf{w}_1 \mathbf{x}_2 + b_1) & 0 & \cdots & 0 \\ 0 & \sigma'(\mathbf{w}_2 \mathbf{x}_1 + b_2) - \sigma'(\mathbf{w}_1 \mathbf{x}_2 + b_1) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma'(\mathbf{w}_m \mathbf{x}_1 + b_m) - \sigma'(\mathbf{w}_1 \mathbf{x}_2 + b_1) \end{pmatrix} \\ &= (\mathbf{h}_1 - \mathbf{h}_2)^\top \cdot \left[\text{diag}(\sigma'(\mathbf{W}\mathbf{x}_1 + \mathbf{b})) - \text{diag}(\sigma'(\mathbf{W}\mathbf{x}_2 + \mathbf{b})) \right] \end{aligned}$$

Thus, we can get $\nabla_{\mathbf{b}} J$ as follows

$$\begin{aligned} \nabla_{\mathbf{b}} J &= \left(\frac{\partial J}{\partial \mathbf{b}} \right)^\top \\ &= \left[\text{diag}(\sigma'(\mathbf{W}\mathbf{x}_1 + \mathbf{b})) - \text{diag}(\sigma'(\mathbf{W}\mathbf{x}_2 + \mathbf{b})) \right] \cdot (\mathbf{h}_1 - \mathbf{h}_2) \end{aligned} \quad (33)$$

Where \mathbf{h}_1 , \mathbf{h}_2 , and \mathbf{b} are $m \times 1$ column vectors, and \mathbf{x}_1 , \mathbf{x}_2 are $n \times 1$ column vectors.

Subproblem (2)

We can write out the gradient descent update rules for the model parameters for a single training example (\mathbf{x}_1 and \mathbf{x}_2) as follow

$$\begin{aligned} \mathbf{W} &= \mathbf{W} - \alpha \nabla_{\mathbf{W}} J \\ &= \mathbf{W} - \alpha \left(\left[(\mathbf{h}_1 - \mathbf{h}_2) \circ \sigma'(\mathbf{W}\mathbf{x}_1 + \mathbf{b}) \right] \cdot \mathbf{x}_1^\top - \left[(\mathbf{h}_1 - \mathbf{h}_2) \circ \sigma'(\mathbf{W}\mathbf{x}_2 + \mathbf{b}) \right] \cdot \mathbf{x}_2^\top + \lambda \mathbf{W} \right) \\ \mathbf{b} &= \mathbf{b} - \alpha \nabla_{\mathbf{b}} J \\ &= \mathbf{b} - \alpha \left(\left[\text{diag}(\sigma'(\mathbf{W}\mathbf{x}_1 + \mathbf{b})) - \text{diag}(\sigma'(\mathbf{W}\mathbf{x}_2 + \mathbf{b})) \right] \cdot (\mathbf{h}_1 - \mathbf{h}_2) \right) \end{aligned} \quad (34)$$

Subproblem (3)

Because $\mathbf{W} \in R^{10 \times 5}$ and $\mathbf{b} \in R^{10 \times 1}$, and there is a hyper-parameter λ , so the number of parameters in the model is

$$10 \times 5 + 10 \times 1 + 1 = 61 \quad (35)$$

Subproblem (4)

We have the loss function as follow

$$\begin{aligned} J &= -\frac{1}{N} \sum_{j=1}^N \sum_{i=1}^k y_j^i \log(\hat{y}_j^i) \\ &= -\frac{1}{N} \sum_{j=1}^N J_j \\ \implies J_j &= \sum_{i=1}^k y_j^i \log(\hat{y}_j^i) \end{aligned} \quad (36)$$

Which means that J_j is the loss for the j^{th} sample. Then we can compute the gradient for sample j . Besides, we define $\hat{\mathbf{y}}_j = \text{softmax}(\mathbf{net}_j)$. We suppose that the j^{th} sample belongs to the class m , which means $y_j^m = 1$. Besides, we suppose $\mathbf{h}_1, \mathbf{h}_2$ is the vector generated by the j^{th} sample.

◦ For $\nabla_{\mathbf{h}_1} J$.

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{h}_1} &= \frac{\partial \left(-\frac{1}{N} \sum_{j=1}^N J_j \right)}{\partial \mathbf{h}_1} \\ &= -\frac{1}{N} \cdot \frac{\partial J_j}{\partial \mathbf{h}_1} \\ &= -\frac{1}{N} \cdot \frac{\partial J_j}{\partial \hat{\mathbf{y}}_j} \cdot \frac{\partial \hat{\mathbf{y}}_j}{\partial \mathbf{net}_j} \cdot \frac{\partial \mathbf{net}_j}{\partial \mathbf{h}_1} \end{aligned} \quad (37)$$

We can calculate $\frac{\partial J_j}{\partial \hat{\mathbf{y}}_j}$ as follow

$$\begin{aligned} \frac{\partial J_j}{\partial \hat{\mathbf{y}}_j} &= \frac{\partial \log(\hat{y}_j^m)}{\partial \hat{\mathbf{y}}_j} \\ &= \begin{pmatrix} 0 & 0 & \cdots & \frac{1}{\hat{y}_j^m} & \cdots & 0 \end{pmatrix} \end{aligned} \quad (38)$$

We can calculate $\frac{\partial \hat{\mathbf{y}}_j}{\partial \mathbf{net}_j}$ as follow

$$\begin{aligned} \frac{\partial \hat{\mathbf{y}}_j}{\partial \mathbf{net}_j} &= \begin{pmatrix} \frac{\partial \hat{y}_j^1}{\partial \mathbf{net}_j^1} & \frac{\partial \hat{y}_j^2}{\partial \mathbf{net}_j^2} & \cdots & \frac{\partial \hat{y}_j^m}{\partial \mathbf{net}_j^m} & \cdots & \frac{\partial \hat{y}_j^k}{\partial \mathbf{net}_j^k} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \end{pmatrix} \\ &= \begin{pmatrix} -\hat{y}_j^1 \hat{y}_j^m & -\hat{y}_j^2 \hat{y}_j^m & \cdots & \hat{y}_j^m - (\hat{y}_j^m)^2 & \cdots & -\hat{y}_j^k \hat{y}_j^m \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \end{pmatrix} \end{aligned} \quad (39)$$

We can calculate $\frac{\partial \mathbf{net}_j}{\partial \mathbf{h}_1}$ as follow

$$\begin{aligned} \frac{\partial \mathbf{net}_j}{\partial \mathbf{h}_1} &= \frac{\partial (\mathbf{W}_3 (\mathbf{h}_1 + \mathbf{h}_2) + \mathbf{b}_3)}{\partial \mathbf{h}_1} \\ &= \frac{\partial (\mathbf{W}_3 \mathbf{h}_1)}{\partial \mathbf{h}_1} = \mathbf{W}_3 \end{aligned} \quad (40)$$

Then, by interesting equation 38, equation 39, and equation 40 into equation 37, we can compute $\frac{\partial J_j}{\partial \mathbf{h}_1}$ as follow

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{h}_1} &= -\frac{1}{N} \cdot \frac{\partial J_j}{\partial \hat{\mathbf{y}}_j} \cdot \frac{\partial \hat{\mathbf{y}}_j}{\partial \mathbf{net}_j} \cdot \frac{\partial \mathbf{net}_j}{\partial \mathbf{h}_1} \\ &= -\frac{1}{N} \cdot \begin{pmatrix} 0 & 0 & \cdots & \frac{1}{\hat{y}_j^m} & \cdots & 0 \end{pmatrix} \cdot \begin{pmatrix} -\hat{y}_j^1 \hat{y}_j^m & -\hat{y}_j^2 \hat{y}_j^m & \cdots & \hat{y}_j^m - (\hat{y}_j^m)^2 & \cdots & -\hat{y}_j^k \hat{y}_j^m \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \end{pmatrix} \cdot \mathbf{W}_3 \\ &= -\frac{1}{N} \cdot \begin{pmatrix} -\hat{y}_j^1 & -\hat{y}_j^2 & \cdots & 1 - \hat{y}_j^m & \cdots & -\hat{y}_j^k \end{pmatrix} \cdot \mathbf{W}_3 \\ &= -\frac{1}{N} \cdot (\mathbf{y}_j - \hat{\mathbf{y}}_j)^\top \cdot \mathbf{W}_3 \end{aligned} \quad (41)$$

Finally, we can get $\nabla_{\mathbf{h}_1} J$ as follow

$$\begin{aligned}\nabla_{\mathbf{h}_1} J &= \left(\frac{\partial J}{\partial \mathbf{h}_1} \right)^\top \\ &= -\frac{1}{N} \mathbf{W}_3^\top (\mathbf{y}_j - \hat{\mathbf{y}}_j)\end{aligned}\quad (42)$$

Where \mathbf{y}_j and $\hat{\mathbf{y}}_j$ are $k \times 1$ column vector.

◦ For $\nabla_{\mathbf{h}_2} J$, it is similar to $\nabla_{\mathbf{h}_1} J$. Thus we can get it as follow

$$\begin{aligned}\nabla_{\mathbf{h}_2} J &= \left(\frac{\partial J}{\partial \mathbf{h}_2} \right)^\top \\ &= -\frac{1}{N} \mathbf{W}_3^\top (\mathbf{y}_j - \hat{\mathbf{y}}_j)\end{aligned}\quad (43)$$

◦ For $\nabla_{\mathbf{x}} J$, we consider \mathbf{x} is the j^{th} sample, as \mathbf{x}_j . Besides, in the follow equation, \mathbf{h}_1 and \mathbf{h}_2 are also generated by the j^{th} sample.

$$\begin{aligned}\frac{\partial J}{\partial \mathbf{x}_j} &= \frac{\partial \left(-\frac{1}{N} \sum_{j=1}^N J_j \right)}{\partial \mathbf{x}_j} \\ &= -\frac{1}{N} \cdot \frac{\partial J_j}{\partial \mathbf{x}_j} \\ &= -\frac{1}{N} \cdot \left(\frac{\partial J_j}{\partial \mathbf{h}_1} \cdot \frac{\partial \mathbf{h}_1}{\partial \mathbf{x}_j} + \frac{\partial J_j}{\partial \mathbf{h}_2} \cdot \frac{\partial \mathbf{h}_2}{\partial \mathbf{x}_j} \right)\end{aligned}\quad (44)$$

We can calculate $\frac{\partial \mathbf{h}_1}{\partial \mathbf{x}_j}$ as follows. We define $\hat{\mathbf{net}}_j = \mathbf{W}_1 \mathbf{x}_j + \mathbf{b}_1$.

$$\begin{aligned}\frac{\partial \mathbf{h}_1}{\partial \mathbf{x}_j} &= \frac{\partial \mathbf{h}_1}{\partial \hat{\mathbf{net}}_j} \cdot \frac{\partial \hat{\mathbf{net}}_j}{\partial \mathbf{x}_j} \\ &= \begin{pmatrix} \frac{\partial \sigma(\hat{\mathbf{net}}_{j,1})}{\partial \hat{\mathbf{net}}_{j,1}} & 0 & \cdots & 0 \\ 0 & \frac{\partial \sigma(\hat{\mathbf{net}}_{j,2})}{\partial \hat{\mathbf{net}}_{j,2}} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \frac{\partial \sigma(\hat{\mathbf{net}}_{j,m})}{\partial \hat{\mathbf{net}}_{j,m}} \end{pmatrix} \cdot \mathbf{W}_1 \\ &= \text{diag} \left(\sigma'(\hat{\mathbf{net}}_j) \right) \cdot \mathbf{W}_1 \\ &= \text{diag} \left(\sigma'(\mathbf{W}_1 \mathbf{x}_j + \mathbf{b}_1) \right) \cdot \mathbf{W}_1\end{aligned}\quad (45)$$

Similarly, we can get $\frac{\partial \mathbf{h}_2}{\partial \mathbf{x}_j}$ as follows.

$$\frac{\partial \mathbf{h}_2}{\partial \mathbf{x}_j} = \text{diag} \left(\text{ReLU}'(\mathbf{W}_2 \mathbf{x}_j + \mathbf{b}_1) \right) \cdot \mathbf{W}_2 \quad (46)$$

Back to equation 44, we can calculate $\frac{\partial J}{\partial \mathbf{x}_j}$ as follow

$$\begin{aligned}\frac{\partial J}{\partial \mathbf{x}_j} &= -\frac{1}{N} \cdot \left[(\mathbf{y}_j - \hat{\mathbf{y}}_j)^\top \cdot \mathbf{W}_3 \cdot \text{diag} \left(\sigma'(\mathbf{W}_1 \mathbf{x}_j + \mathbf{b}_1) \right) \cdot \mathbf{W}_1 + (\mathbf{y}_j - \hat{\mathbf{y}}_j)^\top \cdot \mathbf{W}_3 \cdot \text{diag} \left(\text{ReLU}'(\mathbf{W}_2 \mathbf{x}_j + \mathbf{b}_1) \right) \cdot \mathbf{W}_2 \right] \\ &= -\frac{1}{N} \cdot (\mathbf{y}_j - \hat{\mathbf{y}}_j)^\top \cdot \mathbf{W}_3 \cdot \left[\text{diag} \left(\sigma'(\mathbf{W}_1 \mathbf{x}_j + \mathbf{b}_1) \right) \cdot \mathbf{W}_1 + \text{diag} \left(\text{ReLU}'(\mathbf{W}_2 \mathbf{x}_j + \mathbf{b}_1) \right) \cdot \mathbf{W}_2 \right]\end{aligned}\quad (47)$$

Thus, we can get $\nabla_{\mathbf{x}_j} J$ as follow

$$\begin{aligned}\nabla_{\mathbf{x}_j} J &= \left(\frac{\partial J}{\partial \mathbf{x}_j} \right)^\top \\ &= -\frac{1}{N} \cdot \left[\mathbf{W}_1^\top \text{diag} \left(\sigma'(\mathbf{W}_1 \mathbf{x}_j + \mathbf{b}_1) \right) + \mathbf{W}_2^\top \text{diag} \left(\text{ReLU}'(\mathbf{W}_2 \mathbf{x}_j + \mathbf{b}_1) \right) \right] \cdot \mathbf{W}_3^\top \cdot (\mathbf{y}_j - \hat{\mathbf{y}}_j)\end{aligned}\tag{48}$$

Subproblem (5)

\mathbf{W}_2 is likely to train faster. For the sigmod function, the gradient in the positive and negative saturation regions is close to 0, which may cause the vanishing gradient problem and then the training process would be very slow. On the contrary, for the ReLU function, it doesn't have a tendency to saturate, and there will be no particularly small gradients. Thus, the training process would be faster.