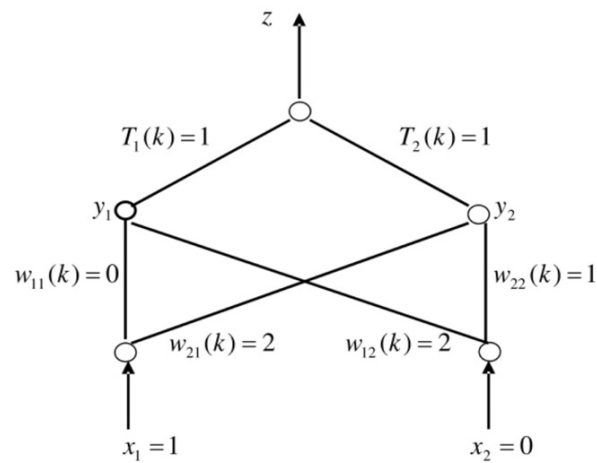


对如下的 BP 神经网络，学习系数  $\eta=1$ ，各点的阈值  $\theta=0$ 。作用函数为：

$$f(x) = \begin{cases} x & x \geq 1 \\ 1 & x < 1 \end{cases}。$$

输入样本  $x_1=1, x_2=0$ ，输出节点  $z$  的期望输出为 1，对于第  $k$  次学习得到的权值分别为  $w_{11}(k)=0, w_{12}(k)=2, w_{21}(k)=2, w_{22}(k)=1, T_1(k)=1, T_2(k)=1$ ，求第  $k$  次和  $k+1$  次学习得到的输出节点值  $z(k)$  和  $z(k+1)$ （写出计算公式和计算过程）。



计算如下：

1. 第  $k$  次训练的正向过程如下：

$$y_i = f(\sum_j w_{ij} x_j - \theta_i) = f(net_i)$$

$$y_1 = f(\sum_{j=1}^2 w_{1j} x_j - \theta) = f(net_1) = f(0 \times 1 + 2 \times 0) = f(0) = 1$$

$$y_2 = f\left(\sum_{j=1}^2 w_{2j} x_j\right) = f(\text{net}_2) = f(2 \times 1 + 1 \times 0) = f(2) = 2$$

$$O_l = f\left(\sum_i T_{li} y_i - \theta_l\right) = f(\text{net}_l)$$

$$z = f\left(\sum_{i=1}^2 T_i y_i\right) = f(\text{net}_l) = f(1 \times 1 + 1 \times 2) = f(3) = 3$$

$$E = \frac{1}{2} (1 - 3)^2 = 2$$

2. 第  $k$  次训练的反向过程如下:

$$\delta_l' = (z' - z) f'(\text{net}_l) = (1 - 3) \times f'(3) = -2 \times 1 = -2$$

$$\delta_i' = f'(\text{net}_i) \sum_l \delta_l' T_{li}$$

$$\delta_1' = f'(\text{net}_1) \delta_l' T_{l1} = f'(0) \times (-2) \times 1 = 0 \times (-2) \times 1 = 0$$

$$\delta_2' = f'(\text{net}_2) \delta_l' T_{l2} = f'(2) \times (-2) \times 1 = 1 \times (-2) \times 1 = -2$$

$$T_1(k+1) = T_1(k) + \Delta T_1 = T_1(k) + \eta \delta_l' y_1 = 1 + 1 \times (-2) \times 1 = -1$$

$$T_2(k+1) = T_2(k) + \Delta T_2 = T_2(k) + \eta \delta_l' y_2 = 1 + 1 \times (-2) \times 2 = -3$$

$$\begin{aligned} W_{11}(k+1) &= W_{11}(k) + \Delta W_{11} \\ &= W_{11}(k) + \eta \delta_1' x_1 = 0 + 1 \times 0 \times 1 = 0 \end{aligned}$$

$$\begin{aligned} W_{12}(k+1) &= W_{12}(k) + \Delta W_{12} \\ &= W_{12}(k) + \eta \delta_1' x_2 = 2 + 1 \times 0 \times 0 = 2 \end{aligned}$$

$$\begin{aligned} W_{21}(k+1) &= W_{21}(k) + \Delta W_{21} \\ &= W_{21}(k) + \eta \delta_2' x_1 = 2 + 1 \times (-2) \times 1 = 0 \end{aligned}$$

$$\begin{aligned} W_{22}(k+1) &= W_{22}(k) + \Delta W_{22} \\ &= W_{22}(k) + \eta \delta_2' x_2 = 1 + 1 \times (-2) \times 0 = 1 \end{aligned}$$

3. 第  $k+1$  次学习的正向过程如下:

$$y_i = f(\sum_j w_{ij} x_j - \theta_i) = f(net_i)$$

$$y_1 = f(\sum_{j=1}^2 w_{1j} x_j) = f(0 \times 1 + 2 \times 0) = f(0) = 1$$

$$y_2 = f(\sum_{j=1}^2 w_{2j} x_j) = f(0 \times 1 + 1 \times 0) = f(0) = 1$$

$$O_i = f(\sum_i T_i y_i - \theta_i) = f(net_i)$$

$$z = f(\sum_{i=1}^2 T_i y_i) = f(1 \times (-1) + 1 \times (-3)) = f(-4) = 1$$

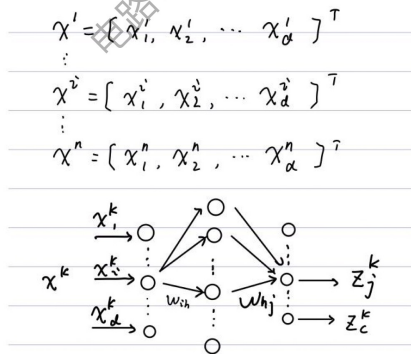
$$E = \frac{1}{2}(1-1)^2 = 0$$

5. 给定  $d$  维空间中的  $n$  个样本  $x_i, i = 1, \dots, n$ , 已知它们分别属于  $c$  个不同的类别。现在拟利用这些样本来训练一个三层前向神经网络（即包含一个输入层，一个隐含层和一个输出层）。假定采用如下平方损失函数作为该网络的目标函数：

$$E(w) = \sum_{k=1}^n \sum_{j=1}^c (t_j^k - z_j^k)^2$$

这里， $t_j^k$  表示样本  $x_k$  在输出层第  $j$  个结点的期望输出值（即该值已知，由样本  $x_k$  的类别标签来决定）， $z_j^k$  表示样本  $x_k$  在输出层第  $j$  个结点的实际输出值（即通过网络计算所得的输出值）， $w$  记录所有待学习的网络参数，包含输入层至隐含层的各个权重  $w_{ih}$  以及隐层至输出层的各个权重  $w_{hj}$ 。请结合上述三层前向神经网络，分别写出  $w_{ih}$  和  $w_{hj}$  的更新公式。（学习率为 1，此网络不包含激活函数）（15'）

5. 思路：题目太抽象，可以画个图辅助理解。另外要注意辨别  $n$  和  $d$  分别都是什么。



设隐层共  $H$  个神经元

$$\frac{\partial E}{\partial w_{hj}} = \frac{\partial \sum_{k=1}^n \sum_{j=1}^c (t_j^k - z_j^k)^2}{\partial w_{hj}} = \sum_{k=1}^n \sum_{j=1}^c [-2(t_j^k - z_j^k)] \cdot \frac{\partial z_j^k}{\partial w_{hj}}$$

$$= \sum_{k=1}^n \sum_{j=1}^c [-2(t_j^k - z_j^k)] \cdot \frac{\partial \sum_{h=1}^H w_{hj} (\sum_{i=1}^d x_i^k w_{ih})}{\partial w_{hj}}$$

$$= \sum_{k=1}^n \sum_{j=1}^c [-2(t_j^k - z_j^k)] \cdot \frac{\partial (w_{hj} \sum_{i=1}^d x_i^k w_{ih} + w_{j2} \sum_{i=1}^d x_i^k w_{i2} + \dots + w_{jH} \sum_{i=1}^d x_i^k w_{iH} + \dots)}{\partial w_{hj}}$$

$$= -2 \sum_{k=1}^n \sum_{j=1}^c (t_j^k - z_j^k) \cdot \sum_{i=1}^d x_i^k w_{ih}$$

$$\frac{\partial E}{\partial w_{ih}} = \sum_{k=1}^n \sum_{j=1}^c [-2(t_j^k - z_j^k)] \cdot \frac{\partial \sum_{h=1}^H w_{hj} (\sum_{i=1}^d x_i^k w_{ih})}{\partial w_{ih}}$$

$$= \sum_{k=1}^n \sum_{j=1}^c [-2(t_j^k - z_j^k)] \cdot \frac{\partial (w_{hj} \sum_{i=1}^d x_i^k w_{ih} + w_{j2} \sum_{i=1}^d x_i^k w_{i2} + \dots + w_{jH} \sum_{i=1}^d x_i^k w_{iH} + \dots)}{\partial w_{ih}}$$

$$= \sum_{k=1}^n \sum_{j=1}^c [-2(t_j^k - z_j^k)] \cdot \frac{\partial w_{hj} \sum_{i=1}^d x_i^k w_{ih}}{\partial w_{ih}}$$

$$= \sum_{k=1}^n \sum_{j=1}^c [-2(t_j^k - z_j^k)] \cdot \frac{\partial w_{hj} (x_1^k w_{1h} + x_2^k w_{2h} + \dots + x_i^k w_{ih} + \dots)}{\partial w_{ih}}$$

$$= \sum_{k=1}^n \sum_{j=1}^c [-2(t_j^k - z_j^k)] \cdot w_{hj} \cdot x_i^k$$

$$\begin{cases} w_{hj} = w_{hj} - \frac{\partial E}{\partial w_{hj}} \\ w_{ih} = w_{ih} - \frac{\partial E}{\partial w_{ih}} \end{cases}$$

# 深度前馈网络

## 推导

假设神经网络(NN)总共有  $L$  层

当第  $L-1$  层时, 权重求导

$$\frac{\partial J}{\partial W_{ij}^{L-1}} = \frac{\partial J}{\partial z_i^L} \frac{\partial z_i^L}{\partial W_{ij}^{L-1}} = \delta_i^L a_j^{L-1}$$

$$\delta_i^L = \frac{\partial J}{\partial z_i^L} = \frac{\partial}{\partial z_i^L} \sum_{i=1}^{S_L} \frac{1}{2} \|y_i - f(z_i^L)\|^2 = -(y_i - f(z_i^L)) f'(z_i^L)$$

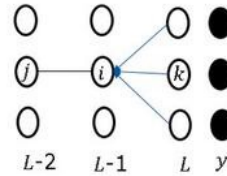
当第  $L-2$  层时, 权重求导

$$\frac{\partial J}{\partial W_{ij}^{L-2}} = \frac{\partial J}{\partial z_i^{L-1}} \frac{\partial z_i^{L-1}}{\partial W_{ij}^{L-2}} = \delta_i^{L-1} a_j^{L-2}$$

$$\begin{aligned} \delta_i^{L-1} &= \frac{\partial J}{\partial z_i^{L-1}} = \frac{\partial}{\partial z_i^{L-1}} \sum_{k=1}^{S_L} \frac{1}{2} \|y_k - f(z_k^L)\|^2 = \sum_{k=1}^{S_L} -(y_k - f(z_k^L)) f'(z_k^L) \frac{\partial z_k^L}{\partial z_i^{L-1}} \\ &= \sum_{k=1}^{S_L} \delta_k^L \cdot w_{ki}^{L-1} f'(z_i^{L-1}) \\ &= (\sum_{k=1}^{S_L} \delta_k^L w_{ki}^{L-1}) f'(z_i^{L-1}) \end{aligned}$$

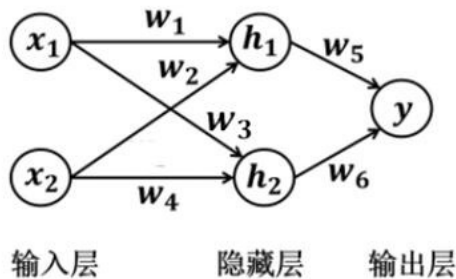
$$z_k^L = W_{ki}^{L-1} a_i^{L-1} + b_k^{L-1}$$

$$a_i^{L-1} = f(z_i^{L-1})$$



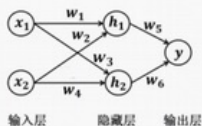
19

4 下图神经网络的隐藏层神经元采用 ReLU 激励函数, 输出层神经元无激励函数, 假设网络损失函数为  $\frac{1}{2}(y_{\text{预测}} - y_{\text{真实}})^2$ , 参数  $w_1, w_2, w_3, w_4, w_5, w_6$  的初始化 1, -2, -1, 2,  $\frac{1}{2}$ , -1, 学习率为 0.1, 只有一个样本 (1, 1), 其标签值是 1, 请问网络经过前馈运算、反向传播、再前馈运算, 损失值是多少? (此题 10 分)



## 手动计算BP神经网络

4 下图神经网络的隐藏层神经元采用 ReLU 激励函数, 输出层神经元无激励函数, 假设网络损失函数为  $\frac{1}{2}(y_{\text{预测}} - y_{\text{真实}})^2$ , 参数  $w_1, w_2, w_3, w_4, w_5, w_6$  的初始化 1, -2, -1, 2,  $\frac{1}{2}$ , -1, 学习率为 0.1, 只有一个样本 (1, 1), 其标签值是 1, 请问网络经过前馈运算、反向传播、再前馈运算, 损失值是多少? (此题 10 分)

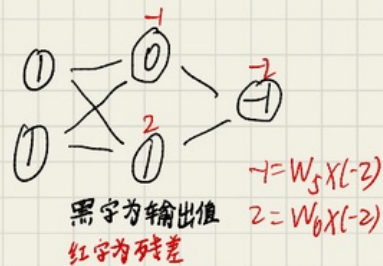


首先进行前馈传播:

$$h_1 = \text{relu}(x_1 w_1 + x_2 w_2) = \text{relu}(1 - 2) = 0$$

$$h_2 = \text{relu}(x_1 w_3 + x_2 w_4) = \text{relu}(-1 + 2) = 1$$

$$J_1 = \frac{1}{2} (y - \hat{y})^2 = \frac{1}{2} (1 - 1)^2 = 0$$



接下来进行反向传播, 此处的  $\Delta W_i$  为后层误差  $\times$  前层输出值  $\times f'(z_i)$

更新  $W_i$ :

$$W_i' = W_i - \text{learning\_rate} \times \Delta W_i$$

$$\Delta W_1 = (-1) \times 1 \times 0 = 0$$

$$W_1' = 1 \quad W_2' = -2 \quad W_3' = -1.2$$

$$\Delta W_2 = (-1) \times 1 \times 0 = 0$$

$$W_4' = 1.8 \quad W_5' = 0.5 \quad W_6' = -0.8$$

$$\Delta W_3 = 2 \times 1 \times 1 = 2$$

$$\Delta W_4 = 2 - 0.2 = 1.8$$

$$\Delta W_5 = (-1) \times 0 \times 1 = 0 \quad (\text{output 无激励函数})$$

$$W_6 = (-2) \times 1 \times 1 = -2$$

再次前向传播

$$h_1' = \text{relu}(x_1 w_1' + x_2 w_2') = \text{relu}(-1) = 0$$

$$h_2' = \text{relu}(x_1 w_3' + x_2 w_4') = \text{relu}(-1.2 + 1.8) = \text{relu}(0.6) = 0.6$$

$$\hat{y}_2 = h_1' w_5' + h_2' w_6' = 0 + 0.6 \times (-0.8) = -0.48$$

$$\Rightarrow J_2 = \frac{1}{2} (\hat{y}_2 - y)^2 = \frac{1369}{1250} = 1.0952$$



4.2 如采用  $\tanh()$  作为多层感知器中隐节点的激活函数，试推导 BP 算法，并讨论为什么多层感知器一般不常用  $\tanh()$  作为激活函数。

参考答案：

(1) 证明：（延续课件表述）

$$\text{for node } j, \quad \text{net}_j = \sum_i w_{ij} O_i, \quad O_j = f(\text{net}_j)$$

这里目标函数设为平方误差

$$E = \frac{1}{2} \sum_j (y_j - \hat{y}_j)^2$$

误差梯度

$$\frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial \text{net}_j} \frac{\partial \text{net}_j}{\partial w_{ij}} = \delta_j O_i, \quad \delta_j = \frac{\partial E}{\partial \text{net}_j}$$

对于输出节点

$$O_j = \hat{y}_j$$
$$\delta_j = \frac{\partial E}{\partial \hat{y}_j} \frac{\partial \hat{y}_j}{\partial \text{net}_j} = -(y_j - \hat{y}_j) f'(\text{net}_j)$$

又因为这里

$$f(x) = \tanh x$$
$$f'(x) = \frac{4e^{-2x}}{(1+e^{-2x})^2} = \frac{1}{(\frac{e^x + e^{-x}}{2})^2} = \frac{1}{\cosh^2 x}$$

则

$$\delta_j = \frac{\partial E}{\partial \hat{y}_j} \frac{\partial \hat{y}_j}{\partial \text{net}_j} = -(y_j - \hat{y}_j) f'(\text{net}_j) = -(y_j - \hat{y}_j) \frac{1}{\cosh^2(\text{net}_j)}$$

对于隐层节点

$$\delta_j = \frac{\partial E}{\partial \text{net}_j} = \sum_k \frac{\partial E}{\partial \text{net}_k} \frac{\partial \text{net}_k}{\partial O_j} \frac{\partial O_j}{\partial \text{net}_j} = \sum_k \delta_k w_{jk} f'(\text{net}_j) = \sum_k \delta_k w_{jk} \frac{1}{\cosh^2(\text{net}_j)}$$

权值学习

$$w_{ij}(t+1) = w_{ij}(t) + \Delta w_{ij}(t)$$
$$\Delta w_{ij}(t) = -\eta \delta_j(t) O_i(t)$$

其中， $\eta$  为学习步长

(2)  $\tanh()$  作为激活函数的缺点

多层感知器一般不常用  $\tanh()$  作为激活函数主要有两个原因，第一可能出现梯度消失的问题：由于激活函数在饱和区导数接近 0，导致在后向传递求导的链式法则下，小数相乘结果接近于 0，故不适合深层网络；第二个原因是  $\tanh()$  涉及到幂运算，速度较慢。