

# 深度神经网络

深度神经网络(Deep Neural Networks, DNN)是深度学习的基础。假设 DNN 的总层数为  $L$ (包括输入层), 第  $l(l=1,2,\dots,L)$ 层的神经元个数为  $k_l$ , 第  $l$  层的输出向量为  $\mathbf{a}_{k_l}^l$ , 第  $l-1$  层到  $l$  层的权重矩阵为  $\mathbf{W}_{k_l \times k_{l-1}}^l$ , 第  $l$  层的偏倚向量为  $\mathbf{b}_{k_l}^l$ , 因此 DNN 的前向传播过程为:

(1).初始化  $\mathbf{a}_{k_1}^1 = \mathbf{x}$ ,  $\mathbf{x}$  表示输入层的样本向量;

(2).for  $l = 2$  to  $L$ , 计算:

$$\mathbf{a}_{k_l}^l = f(\mathbf{z}_{k_l}^l) = f(\mathbf{W}_{k_l \times k_{l-1}}^l \mathbf{a}_{k_{l-1}}^{l-1} + \mathbf{b}_{k_l}^l)$$

上式中, 最后输出的结果为  $\mathbf{a}_{k_L}^L$ ,  $f$  表示的是激活函数,  $\mathbf{z}_{k_l}^l$  表示第  $l$  层的输入向量。

有了 DNN 的前向传播过程, 可以通过通过反向传播算法来更新各层参数的值。假设损失函数为  $J$ , 根据输出层的激活函数和损失函数类型可以分为如下三种:

(1).输出层的激活函数为 *softmax*, 损失函数为交叉熵损失函数, 当  $l=L$  时, 则:

$$\delta_{k_L}^L = \frac{\partial J}{\partial \mathbf{z}_{k_L}^L} = \text{softmax}(\mathbf{z}_{k_L}^L) - \mathbf{y}_{k_L}$$

上式中,  $\mathbf{y}_{k_L}$  表示输出层的标签;

(2).输出层的损失函数为  $J = \frac{1}{2} \|\mathbf{a}_{k_L}^L - \mathbf{y}_{k_L}\|^2 = \frac{1}{2} \|f(\mathbf{z}_{k_L}^L) - \mathbf{y}_{k_L}\|^2$ , 当  $l=L$  时, 则:

$$\delta_{k_L}^L = \frac{\partial J}{\partial \mathbf{z}_{k_L}^L} = (f(\mathbf{z}_{k_L}^L) - \mathbf{y}_{k_L}) \odot f'(\mathbf{z}_{k_L}^L)$$

上式中,  $\mathbf{y}_{k_L}$  表示输出层的输出值;

(3).输出层只有一个神经元, 激活函数为 *sigmoid*, 损失函数为交叉熵损失函数, 当  $l=L$  时, 则:

$$\delta_{k_L}^L = \frac{\partial J}{\partial \mathbf{z}_{k_L}^L} = \text{sigmoid}(\mathbf{z}_{k_L}^L) - \mathbf{y}_{k_L}$$

上式中, 由于输出层只有一个神经元, 因此,  $\mathbf{y}_{k_L}$  表示 1 或者 0;

对于以上三种情况都有当  $l=2$  to  $L-1$  时, 则:

$$\delta_{k_l}^l = \frac{\partial J}{\partial \mathbf{z}_{k_l}^l} = \frac{\partial J}{\partial \mathbf{z}_{k_L}^L} \frac{\partial \mathbf{z}_{k_L}^L}{\partial \mathbf{z}_{k_{L-1}}^{L-1}} \dots \frac{\partial \mathbf{z}_{k_{l+2}}^{l+2}}{\partial \mathbf{z}_{k_{l+1}}^l} \frac{\partial \mathbf{z}_{k_{l+1}}^{l+1}}{\partial \mathbf{z}_{k_l}^l} = \delta_{k_{l+1}}^{l+1} \frac{\partial \mathbf{z}_{k_{l+1}}^{l+1}}{\partial \mathbf{z}_{k_l}^l}$$

又因为  $\mathbf{z}_{k_{l+1}}^{l+1} = \mathbf{W}_{k_{l+1} \times k_l}^{l+1} f(\mathbf{z}_{k_l}^l) + \mathbf{b}_{k_{l+1}}^{l+1}$ , 所以,

$$\delta_{k_l}^l = \delta_{k_{l+1}}^{l+1} \frac{\partial z_{k_{l+1}}^{l+1}}{f(z_{k_l}^l)} \odot \frac{f(z_{k_l}^l)}{\partial z_{k_l}^l} = (\mathbf{W}_{k_{l+1} \times k_l}^{l+1})^T \delta_{k_{l+1}}^{l+1} \odot \frac{f(z_{k_l}^l)}{\partial z_{k_l}^l} = ((\mathbf{W}_{k_{l+1} \times k_l}^{l+1})^T \delta_{k_{l+1}}^{l+1}) \odot f'(z_{k_l}^l)$$

因此，损失函数  $J$  对  $\mathbf{W}_{k_l \times k_{l-1}}^l$  和  $\mathbf{b}_{k_l}^l$  的偏导为

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{W}_{k_l \times k_{l-1}}^l} &= \frac{\partial J}{\partial \mathbf{z}_{k_l}^l} \frac{\partial \mathbf{z}_{k_l}^l}{\partial \mathbf{W}_{k_l \times k_{l-1}}^l} = \delta_{k_l}^l (\mathbf{a}_{k_{l-1}}^{l-1})^T \\ \frac{\partial J}{\partial \mathbf{b}_{k_l}^l} &= \frac{\partial J}{\partial \mathbf{z}_{k_l}^l} \frac{\partial \mathbf{z}_{k_l}^l}{\partial \mathbf{b}_{k_l}^l} = \delta_{k_l}^l \end{aligned}$$

最后，根据如下公式更新 RNN 模型中的参数：

$$\theta_{update} = \theta - \alpha \frac{\partial L}{\partial \theta}$$

上式中， $\theta$  代表的是模型参数  $\mathbf{W}_{k_l \times k_{l-1}}^l$  和  $\mathbf{b}_{k_l}^l$ ， $\alpha$  代表学习率， $\theta_{update}$  代表更新完以后的参数。

但是，在 DNN 中会出现梯度消失和梯度爆炸的现象。出现梯度消失，是因为更新参数  $\mathbf{W}_{k_l \times k_{l-1}}^l$  和  $\mathbf{b}_{k_l}^l$  时，都与前面的模型参数连乘，又因为激活函数的导数在 0 到 1 之间，经过多次迭代相乘，极容易出现梯度消失。出现梯度爆炸，是因为更新参数  $\mathbf{W}_{k_l \times k_{l-1}}^l$  和  $\mathbf{b}_{k_l}^l$  时，当激活函数的导数不是特别小同时模型参数特别大，就会出现梯度爆炸，但是梯度爆炸的概率比较小。

**注：**

1. 本文用  $\odot$  表示 Hadamard 积，对于两个维度相同的向量  $\mathbf{a} = [a_1, a_2, \dots, a_m]^T$  和  $\mathbf{b} = [b_1, b_2, \dots, b_m]^T$ ，则  $\mathbf{a} \odot \mathbf{b} = [a_1 b_1, a_2 b_2, \dots, a_m b_m]^T$ ；
2. 假设两个维度相同的向量  $\mathbf{y} = [y_1, y_2, \dots, y_m]^T$  和  $\mathbf{x} = [x_1, x_2, \dots, x_m]^T$  满足如下公式：

$$\mathbf{y} = f(\mathbf{x}) \Leftrightarrow \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} = \begin{bmatrix} f(x_1) \\ f(x_2) \\ \vdots \\ f(x_m) \end{bmatrix}$$

则：

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} f'(x_1) \\ f'(x_2) \\ \vdots \\ f'(x_m) \end{bmatrix}$$

3. 假设两个维度相同的向量  $\mathbf{y} = [y_1, y_2, \dots, y_m]^T$  和  $\mathbf{x} = [x_1, x_2, \dots, x_n]^T$  满足如下公式：

$$\mathbf{y} = \mathbf{A}_{m \times n} \mathbf{x}$$

则：

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \mathbf{A}_{m \times n}^T$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{A}_{m \times n}} = \mathbf{x}^T$$