

# 反向传播（Backpropagation）算法的数学原理

## 简单介绍

### 第一章 引言

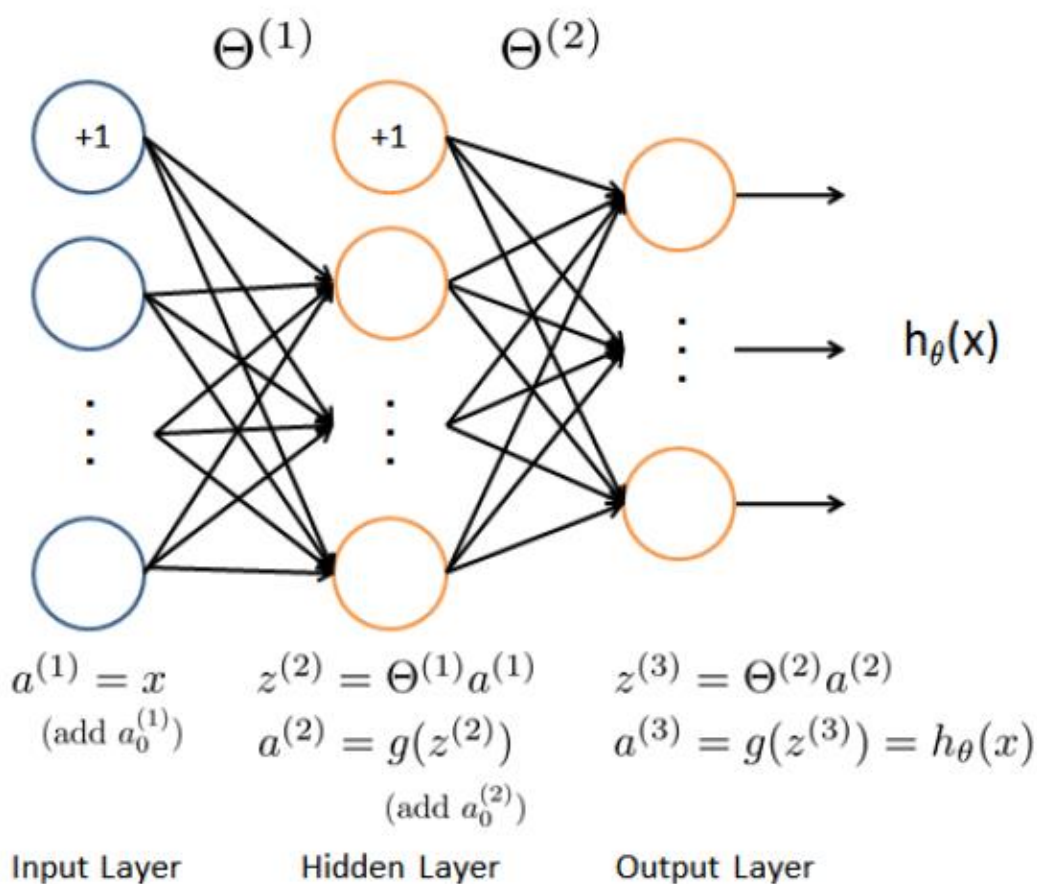


图 1：一个简单的人工神经网络

如图 1 所示，这是一个简单的人工神经网络，分为输入层，隐藏层和输出层。其中  $\theta(1), \theta(2)$  分别是第 1 层，第 2 层向前传递数据时的线性组合参数，将  $\theta(1), \theta(2)$  统称为  $\theta$ 。 $z^{(2)}, z^{(3)}$  是线性组合的输出，作为函数  $g(z)$  的输入， $g(z)$  是一个数学变换（激活函数），最后一层的输出，即  $a(3)$  或  $h_{\theta}(x)$ ，即人工神经网络的运算结果。

在训练时，将已经标记过的数据（标签可以表示为  $y$ ）输入网络，最开始的时候网络输出结果与标签结果  $y$  往往相距甚远，然后用优化算法（例如梯度下降法）不断修正  $\theta$ ，使网络的计算结果越来越接近标准标签结果  $y$ ，直至达到要求的精度。

我们首先定义一个 代价方程 (Cost Function)：

$$J(\theta) = [\text{计算结果与标准答案之间的差距}]$$

我们暂时不关心 Cost Function 的具体形式，只需要理解它是对计算误差的一种衡量，并且它是  $\theta$  的函数。我们的目标就是使计算误差尽量小。以梯度下降法为代表的很多优化算法可以用来解决上述的最小化问题。在优化过程中，这些算法通常会用到  $J(\theta)$  对  $\theta$  中各元素的偏导数（通常  $\theta$  为一个矩阵），反向传播算法就包含计算偏导数的来更新  $\theta$ 。

需要注意到的一点是，在这个简单的人工神经网络中，除了最后一层以外，每一层都要在自己的输出向量前附加一个值为 +1 的项（称为 Bias Term）。可以借用二维坐标系中的一条直线： $ax + b = 0$ ，Bias term 与对应参数相乘后的值就是直线的截距  $b$ 。

## 第二章 数学原理

### 2.1 定义

$L$  —— 神经网络的层数（从 1 开始编号，例如 layer1, layer2...）

$S_l$  —— 第  $l$  层所包含的单元数（不计入 Bias term）

$(x, y)$  —— 一条训练数据

$z_i^{(l)}$  —— 第  $l$  层中第  $i$  个单元的输入（ $i = 1, 2, \dots, S_l$ ）

$a_i^{(l)}$  —— 第  $l$  层中第  $i$  个单元的输出（ $i = 0, 1, 2, \dots, S_l$ ）， $a_0^{(l)} = +1$  即 bias term;

$$a_i^{(l)} = g(z_i^{(l)}) = \frac{1}{1 + e^{-z_i^{(l)}}}$$

$\Theta_{ij}^{(l)}$  —— 对应  $a_i^{(l)}$  的参数，用于计算  $z_i^{(l+1)}$ ；其中：

$$j = 0, 1, 2, \dots, S_l; i = 1, 2, \dots, S_{l+1}$$

$h_{\Theta}(x)$  —— 基于  $\Theta$ ，给定  $x$  后，人工神经网络的输出

$J(\Theta)$  —— 代价方程（Cost Function）：

$$J(\Theta) = -y \log(h_{\Theta}(x)) - (1 - y) \log(1 - h_{\Theta}(x))$$

$\delta_i^{(l)}$  —— 第  $l$  层第  $i$  个单元的计算误差：

$$\delta_i^{(l)} = \frac{\partial}{\partial z_i^{(l)}} J(\Theta) \quad (i = 1, 2, \dots, S_l)$$

### 2.2 数学推导

目标：对于所有的  $l, i, j$ ，计算  $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$

过程：利用链式求导法则

$$\frac{\partial J}{\partial \Theta_{ij}^{(l)}} = \underbrace{\frac{\partial J}{\partial a_i^{(l+1)}}}_{\textcircled{1}} \cdot \underbrace{\frac{\partial a_i^{(l+1)}}{\partial z_i^{(l+1)}}}_{\textcircled{2}} \cdot \underbrace{\frac{\partial z_i^{(l+1)}}{\partial \Theta_{ij}^{(l)}}}_{\textcircled{3}}$$

$$\begin{aligned} \textcircled{1} &= \frac{\partial J}{\partial a_i^{(l+1)}} = \frac{\partial J}{\partial z_1^{(l+2)}} \cdot \frac{\partial z_1^{(l+2)}}{\partial a_i^{(l+1)}} + \frac{\partial J}{\partial z_2^{(l+2)}} \cdot \frac{\partial z_2^{(l+2)}}{\partial a_i^{(l+1)}} + \dots + \frac{\partial J}{\partial z_{s_{l+2}}^{(l+2)}} \cdot \frac{\partial z_{s_{l+2}}^{(l+2)}}{\partial a_i^{(l+1)}} \\ &= \sum_{j=1}^{s_{l+2}} \frac{\partial J}{\partial z_j^{(l+2)}} \cdot \frac{\partial z_j^{(l+2)}}{\partial a_i^{(l+1)}} \\ &= \sum_{j=1}^{s_{l+2}} \delta_j^{(l+2)} \cdot \frac{\partial z_j^{(l+2)}}{\partial a_i^{(l+1)}} \\ &= \sum_{j=1}^{s_{l+2}} \delta_j^{(l+2)} \cdot \frac{\partial}{\partial a_i^{(l+1)}} \left( \Theta_{j0}^{(l+1)} a_0^{(l+1)} + \Theta_{j1}^{(l+1)} a_1^{(l+1)} + \dots + \Theta_{js_{l+1}}^{(l+1)} a_{s_{l+1}}^{(l+1)} \right) \\ &= \sum_{j=1}^{s_{l+2}} \delta_j^{(l+2)} \cdot \Theta_{ji}^{(l+1)} \end{aligned}$$

( $i = 1, 2, \dots, S_l$ ; 不包括 bias term)

$$\begin{aligned} \textcircled{2} &= \frac{\partial a_i^{(l+1)}}{\partial z_i^{(l+1)}} \\ &= \frac{\partial}{\partial z_i^{(l+1)}} g(z_i^{(l+1)}) \\ &= g(z_i^{(l+1)}) \cdot (1 - g(z_i^{(l+1)})) \\ &= a_i^{(l+1)} (1 - a_i^{(l+1)}) \end{aligned}$$

( $i = 1, 2, \dots, S_l$ ; 不包括 bias term)

$$\begin{aligned} \textcircled{3} &= \frac{\partial z_i^{(l+1)}}{\partial \Theta_{ij}^{(l)}} \\ &= \frac{\partial}{\partial \Theta_{ij}^{(l)}} \left( \Theta_{i0}^{(l)} a_0^{(l)} + \Theta_{i1}^{(l)} a_1^{(l)} + \dots + \Theta_{is_l}^{(l)} a_{s_l}^{(l)} \right) \\ &= a_j^{(l)} \end{aligned}$$

( $j = 0, 1, 2, \dots, S_l$ ; 包括 bias term)

$$\begin{aligned}
\therefore \frac{\partial J}{\partial \Theta_{ij}^{(l)}} &= \textcircled{1} \times \textcircled{2} \times \textcircled{3} \\
&= \sum_{k=1}^{S_{l+2}} \delta_k^{(l+2)} \Theta_{ki}^{(l+1)} \cdot a_i^{(l+1)} (1 - a_i^{(l+1)}) \cdot a_j^{(l)}
\end{aligned}$$

根据  $\delta$  的定义：

$$\begin{aligned}
\delta_i^{(l+1)} &= \textcircled{1} \times \textcircled{2} = \sum_{k=1}^{S_{l+2}} \delta_k^{(l+2)} \Theta_{ki}^{(l+1)} \cdot a_i^{(l+1)} (1 - a_i^{(l+1)}) \\
&(i = 1, 2, \dots, S_{l+1})
\end{aligned}$$

将第  $(l+1)$  层的所有单元误差表示为  $\delta^{(l+1)}$ ，则  $\delta_i^{(l+1)}$  是其第  $i$  个元素：

$$\therefore \delta^{(l+1)} = (\Theta^{(l+1)})^T \delta^{(l+2)} .* a^{(l+1)} .* (1 - a^{(l+1)}) \quad (4)$$

( $.*$ 表示矩阵对应元素之间两两相乘)

$$\therefore \frac{\partial J}{\partial \Theta^{(l)}} = \delta^{(l+1)} (a^{(l)})^T \quad (5)$$

式 (4) (5) 即 Prof. Ng 所给出的反向传播算法的公式

## 2.3 补充说明：

1. 公式 (4) 不能用来计算  $\delta^{(L)}$ ，因为  $\delta^{(L+1)}$  未定义。 $\delta^{(L)}$  的计算如下：

$$\begin{aligned}
J(\Theta) &= -y \log(h_{\Theta}(x)) - (1 - y) \log(1 - h_{\Theta}(x)) \\
&= -y \log(a^{(L)}) - (1 - y) \log(1 - a^{(L)}) \\
\delta_i^{(L)} &= \frac{\partial J}{\partial z_i^{(L)}} = \frac{\partial J}{\partial a_i^{(L)}} \cdot \frac{\partial a_i^{(L)}}{\partial z_i^{(L)}} \\
&= \frac{\partial}{\partial a_i^{(L)}} [-y_i \log(a_i^{(L)}) - (1 - y_i) \log(1 - a_i^{(L)})] \cdot \frac{\partial}{\partial z_i^{(L)}} g(\partial z_i^{(L)}) \\
&= a_i^{(L)} - y_i
\end{aligned}$$

$$\therefore S^{(L)} = a^{(L)} - y$$

2. 除了计算  $\delta^{(L)}$ ，其它运算都未用到  $J(\Theta)$  的具体表达式；
3. 由于  $z_0^{(l)}$  未定义，所以  $\delta_0^{(l)} = \frac{\partial J}{\partial z_0^{(l)}}$  未定义，所以不要试图计算 bias term 项的误差。

## 参考文献

[1]word:

[https://github.com/BasicCoder/NNBlog1/tree/master/%E5%8F%8D%E5%90%91%E4%BC%A0%E6%92%AD%E7%AE%97%E6%B3%95\(BP\)%E7%9A%84%E6%95%B0%E5%AD%A6%E5%8E%9F%E7%90%86%E7%AE%80%E5%8D%95%E4%BB%8B%E7%BB%8D](https://github.com/BasicCoder/NNBlog1/tree/master/%E5%8F%8D%E5%90%91%E4%BC%A0%E6%92%AD%E7%AE%97%E6%B3%95(BP)%E7%9A%84%E6%95%B0%E5%AD%A6%E5%8E%9F%E7%90%86%E7%AE%80%E5%8D%95%E4%BB%8B%E7%BB%8D)

[2]PDF: