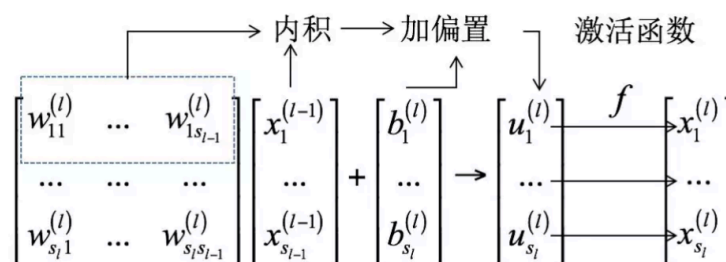


根据上面的结论可以方便的推导出神经网络的求导公式。假设神经网络有 n_l 层，第 l 层神经元个数为 s_l 。第 l 层从第 $l-1$ 层接收的输入向量为 $\mathbf{x}^{(l-1)}$ ，本层的权重矩阵为 $\mathbf{W}^{(l)}$ ，偏置向量为 $\mathbf{b}^{(l)}$ ，输出向量为 $\mathbf{x}^{(l)}$ 。该层的输出可以写成如下矩阵形式：

$$\begin{aligned}\mathbf{u}^{(l)} &= \mathbf{W}^{(l)} \mathbf{x}^{(l-1)} + \mathbf{b}^{(l)} \\ \mathbf{x}^{(l)} &= f(\mathbf{u}^{(l)})\end{aligned}$$

其中 $\mathbf{W}^{(l)}$ 是 $s_l \times s_{l-1}$ 的矩阵， $\mathbf{u}^{(l)}$ 和 $\mathbf{b}^{(l)}$ 是 s_l 维的向量。神经网络一个层实现的变换如下图所示：



如果将神经网络按照各个层展开，最后得到一个深层的复合函数，将其代入欧氏距离损失函数，依然是一个关于各个层的权重矩阵和偏置向量的复合函数：

$$L = \frac{1}{2} \left\| f\left(\mathbf{W}^{(L)} f\left(\mathbf{W}^{(L-1)} f\left(\mathbf{W}^{(L-2)} \left(\dots f\left(\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)}\right)\right) + \mathbf{b}^{(L-1)}\right) + \mathbf{b}^{(L-1)}\right) + \mathbf{b}^{(L)}\right) - \mathbf{y} \right\|^2$$

先从最外层算起

要计算某一层的权重矩阵和偏置向量的梯度，只能依赖于它紧贴着的外面那一层变量的梯度值，通过一次复合函数求导得到。

根据定义， $\mathbf{w}^{(l)}$ 和 $\mathbf{b}^{(l)}$ 是目标函数的自变量， $\mathbf{u}^{(l)}$ 和 $\mathbf{x}^{(l)}$ 可以看成是它们的函数。根据前面的结论，损失函数对权重矩阵的梯度为：

$$\nabla_{\mathbf{w}^{(l)}} L = \left(\nabla_{\mathbf{u}^{(l)}} L \right) \left(\mathbf{x}^{(l-1)} \right)^T$$

对偏置向量的梯度为：

$$\nabla_{\mathbf{b}^{(l)}} L = \nabla_{\mathbf{u}^{(l)}} L$$

现在的问题是，梯度 $\nabla_{\mathbf{u}^{(l)}} L$ 怎么计算？我们分两种情况讨论，如果第 l 层是输出层，在这里只考虑对单个样本的损失函数，根据上一节推导的结论，这个梯度为：

$$\nabla_{\mathbf{u}^{(l)}} L = \left(\nabla_{\mathbf{x}^{(l)}} L \right) \odot f'(\mathbf{u}^{(l)}) = (\mathbf{x}^{(l)} - \mathbf{y}) \odot f'(\mathbf{u}^{(l)})$$

这就是输出层的神经元输出值与期望值之间的误差。这样我们得到输出层权重的梯度为：

$$\nabla_{\mathbf{w}^{(l)}} L = (\mathbf{x}^{(l)} - \mathbf{y}) \odot f'(\mathbf{u}^{(l)}) \left(\mathbf{x}^{(l-1)} \right)^T$$

等号右边第一个乘法是向量对应元素乘；第二个乘法是矩阵乘，在这里是列向量与行向量的乘积，结果是一个矩阵，尺寸刚好和权重矩阵相同。损失函数对偏置项的梯度为：

$$\nabla_{\mathbf{b}^{(l)}} L = (\mathbf{x}^{(l)} - \mathbf{y}) \odot f'(\mathbf{u}^{(l)})$$

下面考虑第二种情况。如果第 l 层是隐含层，则有：

$$\mathbf{u}^{(l+1)} = \mathbf{W}^{(l+1)} \mathbf{x}^{(l)} + \mathbf{b}^{(l+1)} = \mathbf{W}^{(l+1)} f(\mathbf{u}^{(l)}) + \mathbf{b}^{(l+1)}$$

假设梯度 $\nabla_{\mathbf{u}^{(l+1)}} L$ 已经求出，根据前面的结论，有：

$$\nabla_{\mathbf{u}^{(l)}} L = \left(\nabla_{\mathbf{x}^{(l)}} L \right) \odot f'(\mathbf{u}^{(l)}) = \left(\left(\mathbf{W}^{(l+1)} \right)^T \nabla_{\mathbf{u}^{(l+1)}} L \right) \odot f'(\mathbf{u}^{(l)})$$

这是一个递推的关系，通过 $\nabla_{\mathbf{u}^{(l+1)}} L$ 可以计算出 $\nabla_{\mathbf{u}^{(l)}} L$ ，递推的终点是输出层，而输出层的梯度值我们之前已经算出。由于根据 $\nabla_{\mathbf{u}^{(l)}} L$ 可以计算出 $\nabla_{\mathbf{w}^{(l)}} L$ 和 $\nabla_{\mathbf{b}^{(l)}} L$ ，因此可以计算出任意层权重与偏置的梯度值。

为此我们定义误差项为损失函数对临时变量u的梯度：

$$\delta^{(l)} = \nabla_{w^{(l)}} L = \begin{cases} (x^{(l)} - y) \odot f'(u^{(l)}) & l = n_l \\ (W^{(l+1)})^T (\delta^{(l+1)}) \odot f'(u^{(l)}) & l \neq n_l \end{cases}$$

向量 $\delta^{(l)}$ 的尺寸和本层神经元的个数相同。这是一个递推的定义， $\delta^{(l)}$ 依赖于 $\delta^{(l+1)}$ ，递推的终点是输出层，它的误差项可以直接求出。

根据误差项可以方便的计算出对权重和偏置的偏导数。首先计算输出层的误差项，根据他得到权重和偏置项的梯度，这是起点；根据上面的递推公式，逐层向前，利用后一层的误差项计算出本层的误差项，从而得到本层权重和偏置项的梯度。

单个样本的反向传播算法在每次迭代时的流程为：

- 1.正向传播，利用当前权重和偏置值，计算每一层对输入样本的输出值
- 2.反向传播，对输出层的每一个节点计算其误差：

$$\delta^{(n_l)} = (x^{(n_l)} - y) \odot f'(u^{(n_l)})$$

- 3.对于 $l = n_l - 1, \dots, 2$ 的各层，计算第l层每个节点的误差：

$$\delta^{(l)} = (W^{(l+1)})^T \delta^{(l+1)} \odot f'(u^{(l)})$$

- 4.根据误差计算损失函数对权重的梯度值：

$$\nabla_{w^{(l)}} L = \delta^{(l)} (x^{(l-1)})^T$$

对偏置的梯度为：

$$\nabla_{b^{(l)}} L = \delta^{(l)}$$

- 5.用梯度下降法更新权重和偏置：

$$\begin{aligned} W^{(l)} &= W^{(l)} - \eta \nabla_{w^{(l)}} L \\ b^{(l)} &= b^{(l)} - \eta \nabla_{b^{(l)}} L \end{aligned}$$

实现时需要在正向传播时记住每一层的输入向量 $x^{(l-1)}$ ，本层的激活函数导数值 $f'(u^{(l)})$ 。

神经网络的训练算法可以总结为：

复合函数求导+梯度下降法

训练算法有两个版本：批量模式和单样本模式。批量模式每次梯度下降法迭代时对所有样本计算损失函数值，计算出对这些样本的总误差，然后用梯度下降法更新参数；单样本模式是每次对一个样本进行前向传播，计算对该样本的误差，然后更新参数，它可以天然的支持增量学习，即动态的加入新的训练样本进行训练。

在数学中，向量一般是列向量，但在编程语言中，向量一般按行存储，即是行向量，因此实现时计算公式略有不同，需要进行转置。正向传播时的计算公式为：

$$\mathbf{u}^{(l)} = \mathbf{x}^{(l-1)} \mathbf{W}^{(l)} + \mathbf{b}^{(l)}$$

感兴趣的读者可以阅读各个开源库，对比它们的计算公式。反向传播时的计算公式为：

$$\delta^{(l)} = \delta^{(l+1)} \left(\mathbf{W}^{(l+1)} \right)^T \odot f' \left(\mathbf{u}^{(l)} \right)$$

对权重矩阵的计算公式为：

$$\nabla_{\mathbf{W}} L = \left(\mathbf{x}^{(l-1)} \right)^T \delta^{(l)}$$

这些向量都是行向量。



