

在学习神经网络之前，我们先来了解一下神经元是如何工作的：

生物上，神经元有着多样的形式，但是所有的神经元都是将电信号从一端传输到另一端，沿着轴突，将电信号从树突传到树突。继而，这些信号从一个神经元传到另外一个神经元。

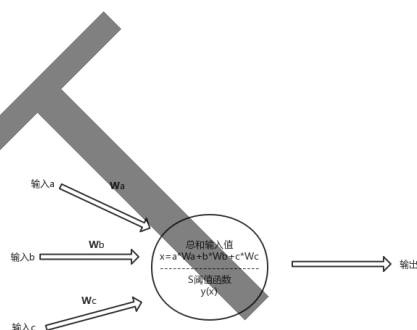
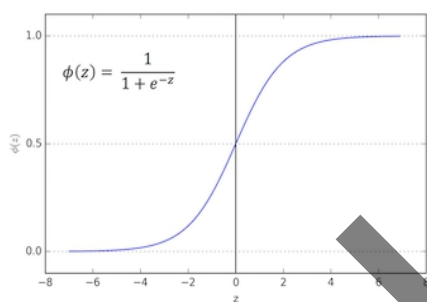
现在我们已经了解了，神经元是如何去工作的---->它接收了一个电输入，输出为另外一个电信号。映射到机器学习上就是，接收一个输入--->中间处理--->弹出一个输出。

同样的，在生物上，神经元不会立即反应，而是会抑制输出，直到输入增强，强大到可以触发输出。我们可以认为在产生输出之前，输入必须达到一个阈值。就比如往水杯里面倒水---水满了就会溢出。换言之，我们可以想象一下，用打火机去烧手指，如果你的火在手指下停留的时间短，你几乎不会有疼痛，但火稍微久一点，你就会感到被烧到了，

这就相当于是达到了阈值，你的神经元将输入的信号转换为疼痛感了。

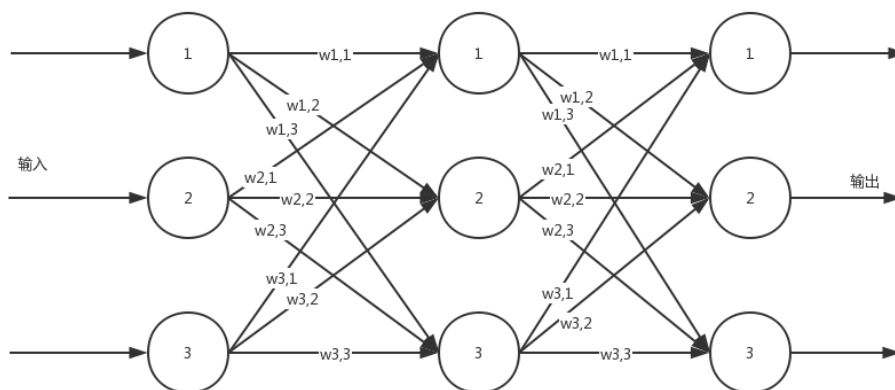
了解了阈值这个概念以后，我们就需要定义神经网络里面阈值的激活函数了，在数学上，有许多激活函数可以达到这样的效果，如简单的阶跃函数，以及 S 函数 (sigmoid function)---> $y=1/(1+e^{-x})$

如图我们可以这样理解，当输入较小的时候，输出为 0，一旦输入达到阈值，输出就一跃而起。就如神经元被激发了。



现在我们知道了一个神经元的输入输出了。

但是，实际中，每个神经元接受来自其之前多个神经元的输入，并当神经元被激发了，它也同时提供信号给更多的神经元。将这种理解映射到神经网络上，也就是将这种自然形态映射到我们人造的模型上：就是构建多层神经元，每一层中的神经元都在与其前后层的神经元互相连接。

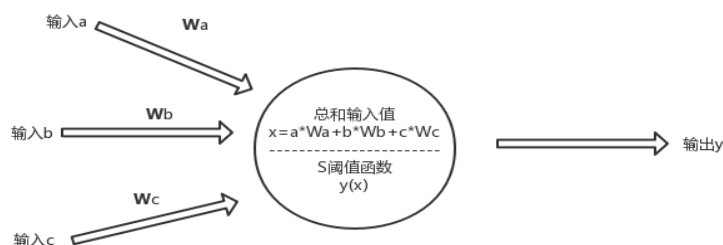


三层神经网络-自制

这里，我们需要解释  $w$  是什么？

$w$  就是我们常说的权重，较大的权重将放大信号，较小的权重将弱化信号。

加上权重的神经网络图：



这是我思考的一个问题是：对于这种简单的模型我们可以手写计算，但是对于更复杂的模型我们应该怎么办呢？

这就可以联系到我们学过的线代的知识，利用矩阵计算，而这也是更容易让计算机去表达的一种计算形式。

如：你有三个输入记为  $in1, in2, in3$ ，你的权重就会有  $3*3=9$  个。

要算出你的总输入值：

$$X = \begin{pmatrix} w_{1,1} & w_{2,1} & w_{3,1} \\ w_{1,2} & w_{2,2} & w_{3,2} \\ w_{1,3} & w_{2,3} & w_{3,3} \end{pmatrix} * \begin{pmatrix} in1 \\ in2 \\ in3 \end{pmatrix}$$

联系上面的三层神经网络图这是很清晰的，利用这种形式，我们将这模型推广可以写为  $X = W^T * I$ ,  $W$  为权重矩阵， $I$  为输入矩阵

最后输出则为： $O = \text{sigmoid}(X)$ ,  $O$  为输出矩阵。

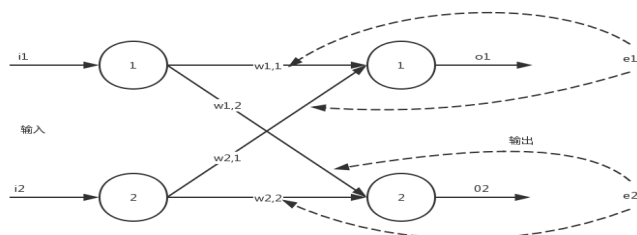
上面我们介绍了输入和输出两层。但隐藏在输入输出中还有若干的隐藏层我们也称其为中间层。

如定义： $I = \begin{pmatrix} 0.5 \\ 0.6 \\ 0.1 \end{pmatrix}$ ,  $W_{input\_hidden}^T = \begin{pmatrix} 0.1 & 0.7 & 0.8 \\ 0.3 & 0.6 & 0.9 \\ 0.6 & 0.5 & 0.1 \end{pmatrix}$ ,  $W_{in\_hidden}^T$  为输入层与

隐藏层之间的权重矩阵， $X_{hidden} = W_{in\_hidden}^T * I$ ,  $X_{hidden}$  就是我们所求的中间层。 $O_{hidden} = \text{sigmoid}(X_{hidden})$ ,  $O_{hidden}$  为中间层的输出。

假设只有一层中间层则， $X_{output} = W_{output\_hidden}^T * O_{hidden}$

看到这里，似乎这一切都完美的结束了。其实不然，我们还需要考虑学习来自多个节点的权重及更新权重等问题。



利用反向传播误差的方法，观察上图：

将第一个输出节点的误差记为  $e_1$ ，这个值由训练数据的所期望的输出值  $t_1$  与实际输出值  $o_1$  之间的差，即  $e_1 = t_1 - o_1$

从图中，按照所连链接的比例，也就是权重  $w_{1,1}$  与  $w_{2,1}$ ，对误差  $e_1$  进行了分割，即使用  $e_1$  的一部分更新  $w_{1,1}$  为  $\frac{w_{1,1}}{w_{1,1} + w_{2,1}}$ ，另外一部分更新  $w_{2,1}$  为

$$\frac{w_{2,1}}{w_{1,1} + w_{2,1}}$$

当有多个层的时候，我们需要考虑隐藏层。

$$error_{hidden} = W_{output\_hidden}^T * error_{output}$$

如何更新权重？这是我们需要思考的一个问题？

不想了，直接肝。

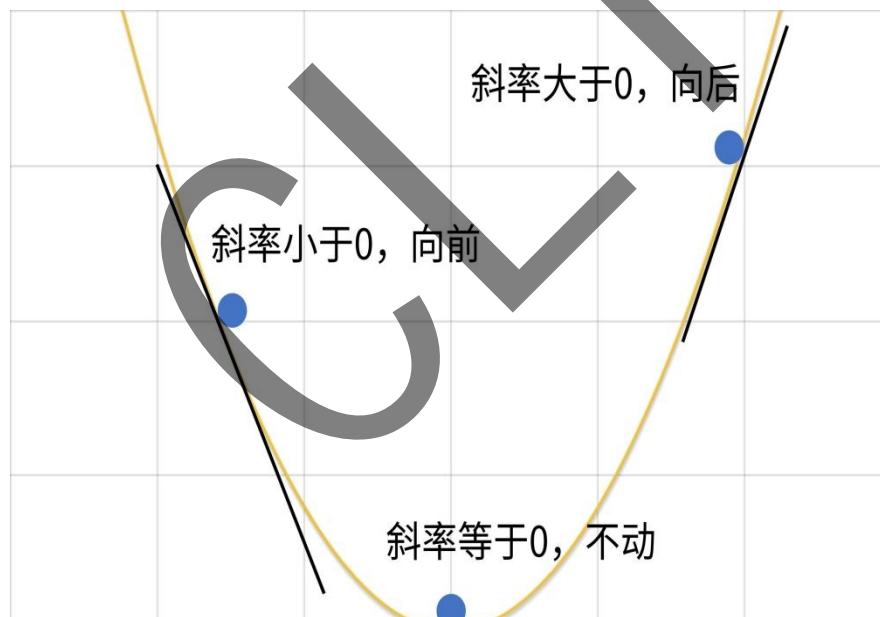
对于反向传播就是为了最小化 loss 求梯度的过程了。

梯度概念这里不做解释了。

假如有一个函数， $z = f(x, y)$ ，那么  $z$  的梯度为：

$$\frac{\partial z}{\partial x_i} + \frac{\partial z}{\partial y_j}$$

$z$  函数梯度的方向是  $z$  增加最快的方向。我们在深度学习里，需要降低 loss，因此我们是选择 loss 函数梯度的反方向



推荐一篇经典论文供大家研读：

[1] Hochreiter S, Schmidhuber J. Long Short-Term Memory[J]. Neural Computation, 1997, 9(8):1735-1780.

下面我们具体讲一些神经网络的模型：

### 一、感知器神经网络 (Perceptron neural network)

其算法框架为：

---

#### Perceptron Algorithm

1. 随机取  $w, b$
  2. 取一个训练样本  $(x, y)$ 
    - (1) 若  $w^T x + b > 0$  且  $y = -1$  则：  
 $w = w - x, b = b - 1$
    - (2) 若  $w^T x + b < 0$  且  $y = 1$  则：  
 $w = w + x, b = b + 1$
  3. 再取另一条件  $(x, y)$  回到 2
  4. 终止条件，直到所有输入输出对都不满足 2 中 (1)、(2) 之一，退出循环
- 

可以想一下：为什么  $w^T x + b > 0$ ，且  $y = -1$  则： $w = w - x, b = b - 1$

意思是当  $w^T x + b$  为一个正数的时候， $y$  是一个负数，这表明着  $w, b$  对  $x$  分类错误。

所以给出  $w = w - x, b = b - 1$ ，那为什么是  $w = w - x, b = b - 1$ ？

我给出一个简单的推导：

令  $w_{\text{新}} = w - x, b_{\text{新}} = b - 1$ ，带入  $w^T x + b$ ，有  $w_{\text{新}}^T x + b_{\text{新}} = (w - x)^T * x + (b - 1)$

化简得到  $w_{\text{新}}^T x + b_{\text{新}} = w^T x + b - (||x||^2 + 1)$ ，这样就浅显易懂了。

因为  $x$  模的平方是一个正数再加上 1 就为一个恒大于 1 的数，结合算法就是迭代一次就向负方向走之上 1 的数，不断迭代将分类错误的拉到正确的这边就结束迭代。

同理 (2)  $+b < 0$  且  $y = 1$  则： $w = w + x, b = b + 1$  的道理一样。

感知器是由两层神经网络构成，所以感知器模型公式可以写成：

$$y = f\left(\sum_{i=1}^n w_i x_i - \theta\right) = f(w^T x - \theta)$$

其中， $x$  为样本的特征向量，是感知器模型的输入； $w, \theta$  是感知器模型的参数， $w$  为权重， $\theta$  为阈值。假定  $f$  为阶跃函数，那么感知机模型的公式可进一步表示

为： $y = \text{sgn}(w^T x - \theta) = \begin{cases} 1, & w^T x - \theta \geq 0 \\ 0, & w^T x - \theta < 0 \end{cases}$ ， $\text{sgn}$  代表符号函数。

在  $n$  维超平面中，感知器模型将  $n$  维空间划分为  $w^T x - \theta \geq 0$  和  $w^T x - \theta < 0$  两个子空间，并输出相应的值 1、0。

## 二、误差逆传播 (BP)算法

适用：多层前馈神经网络、递归神经网络等

BP 算法是什么？（下面具体理论来源于周志华老师的西瓜书内容）

下面我们来看看 BP 算法究竟是什么样. 给定训练集  $D = \{(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_m, \mathbf{y}_m)\}$ ,  $\mathbf{x}_i \in \mathbb{R}^d$ ,  $\mathbf{y}_i \in \mathbb{R}^l$ , 即输入示例由  $d$  个属性描述, 输出  $l$  维实值向量. 为便于讨论, 图 5.7 给出了一个拥有  $d$  个输入神经元、 $l$  个输出神经元、 $q$  个隐层神经元的多层前馈网络结构, 其中输出层第  $j$  个神经元的阈值用  $\theta_j$  表示, 隐层第  $h$  个神经元的阈值用  $\gamma_h$  表示. 输入层第  $i$  个神经元与隐层第  $h$  个神经元之间的连接权为  $v_{ih}$ , 隐层第  $h$  个神经元与输出层第  $j$  个神经元之间的连接权为  $w_{hj}$ . 记隐层第  $h$  个神经元接收到的输入为  $\alpha_h = \sum_{i=1}^d v_{ih}x_i$ , 输出层第  $j$  个神经元接收到的输入为  $\beta_j = \sum_{h=1}^q w_{hj}b_h$ , 其中  $b_h$  为隐层第  $h$  个神经元的输出. 假设隐层和输出层神经元都使用图 5.2(b) 中的 Sigmoid 函数.

对训练例  $(\mathbf{x}_k, \mathbf{y}_k)$ , 假定神经网络的输出为  $\hat{\mathbf{y}}_k = (\hat{y}_1^k, \hat{y}_2^k, \dots, \hat{y}_l^k)$ , 即

$$\hat{y}_j^k = f(\beta_j - \theta_j), \quad (5.3)$$

则网络在  $(\mathbf{x}_k, \mathbf{y}_k)$  上的均方误差为

$$E_k = \frac{1}{2} \sum_{j=1}^l (\hat{y}_j^k - y_j^k)^2. \quad (5.4)$$

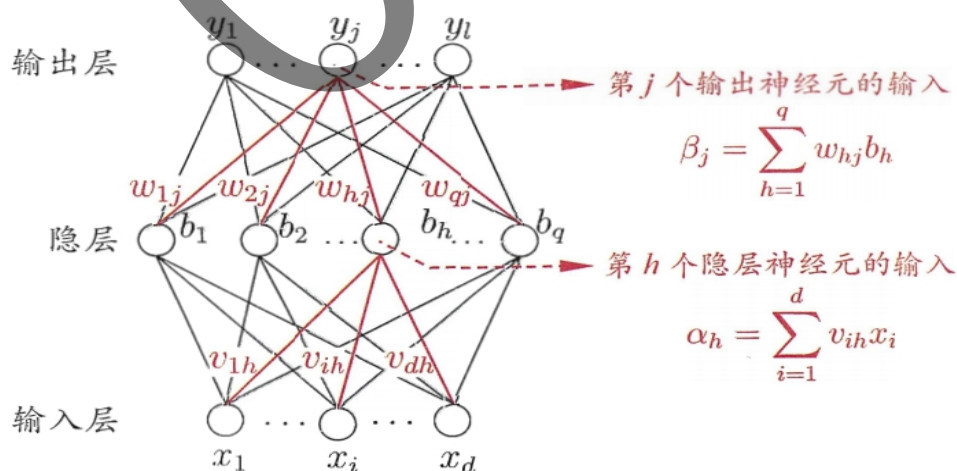


图 5.7 BP 网络及算法中的变量符号

下为对输出层权重、阈值变化量公式的推导：



均方误差:

$$E = \frac{1}{2} \sum_{j=1}^n (\hat{y}_j - y_j)^2 \quad \text{求导} \rightarrow \frac{\partial E}{\partial \hat{y}_j} = \hat{y}_j - y_j$$

第j个输出层神经元的输入:

$$B_j = \sum_{h=1}^n w_{hj} b_h \quad \frac{\partial B_j}{\partial w_{hj}} = b_h$$

又BP算法属于(梯度下降)策略,以目标的负梯度方向对参数进行调整. 对误差E, 给定学习率 $\eta$ , 有:

$$\Delta w_{hj} = -\eta \frac{\partial E}{\partial w_{hj}}$$

对w求导: 其中  $\frac{\partial E}{\partial w_{hj}} = \frac{\partial E}{\partial \hat{y}_j} \cdot \frac{\partial \hat{y}_j}{\partial B_j} \cdot \frac{\partial B_j}{\partial w_{hj}}$

由第j个输出层神经元的输出为:

$$\hat{y}_j = f(B_j - \theta_j) \quad \frac{\partial \hat{y}_j}{\partial B_j} = f'(B_j - \theta_j)$$

由此处f函数使用的sgn为激活函数: 有.

$$f'(x) = f(x)(1 - f(x))$$

$$\begin{aligned} i) \quad \frac{\partial \hat{y}_j}{\partial B_j} &= f'(B_j - \theta_j) = f(B_j - \theta_j)(1 - f(B_j - \theta_j)) \\ &= \hat{y}_j(1 - \hat{y}_j) \end{aligned}$$

$$\Delta w_{hj} = -\eta \cdot \hat{y}_j(1 - \hat{y}_j) \cdot b_h$$

↑  
输出层权重变化量.

输出层两值变化量:

$$\begin{aligned}
 \text{由 } \frac{\partial E}{\partial \theta_j} &= \frac{\partial E}{\partial \hat{y}_j} * \frac{\partial \hat{y}_j}{\partial \theta_j} = \frac{\partial E}{\partial \hat{y}_j} * \frac{\partial [f(B_j - \theta_j)]}{\partial \theta_j} \\
 &= \frac{\partial E}{\partial \hat{y}_j} * f'(B_j - \theta_j) * (-1) \\
 &= \frac{\partial E}{\partial \hat{y}_j} * f(B_j - \theta_j) * [1 - f(B_j - \theta_j)] * (-1) \\
 &= \frac{\partial E}{\partial \hat{y}_j} * \hat{y}_j (1 - \hat{y}_j) * (-1) \\
 &= \frac{\partial \left[ \frac{1}{2} \sum_{j=1}^L (\hat{y}_j - y_j)^2 \right]}{\partial \hat{y}_j} * \hat{y}_j (1 - \hat{y}_j) * (-1) \\
 &= \frac{1}{2} * 2 * (\hat{y}_j - y_j) * 1 * \hat{y}_j (1 - \hat{y}_j) * (-1) \\
 &= (y_j - \hat{y}_j) \hat{y}_j (1 - \hat{y}_j) \\
 &= g_j \\
 \text{所以 } \Delta \theta_j &= -\eta \frac{\partial E}{\partial \theta_j} = -\eta g_j \\
 &= -\eta \hat{y}_j (1 - \hat{y}_j) (y_j - \hat{y}_j)
 \end{aligned}$$

同理：通过对输出层公式的推导我们可以推出隐藏层的权重与阈值的变化量： $\Delta v_{ih} = \eta b_h (1 - b_h) \sum_{j=1}^l w_{hj} g_j x_i$

$$\Delta v_h = -\eta b_h (1 - b_h) \sum_{j=1}^l w_{hj} g_j$$

介绍到这里我们可以使用 python 对其进行实现