

AI从入门到放弃：BP神经网络算法推导及代码实现笔记

Author: 拉轰的二哈 Date: 2018.06.08

一. 前言:

- 作为AI入门小白，参考了一些文章，想记点笔记加深印象，发出来是给有需求的童鞋学习共勉，大神轻拍！
- 【毒鸡汤】：算法这东西，读完之后的状态多半是-->“我是谁，我在哪？”没事的，吭哧吭哧学总能学会，毕竟还有千千万万个算法等着你。

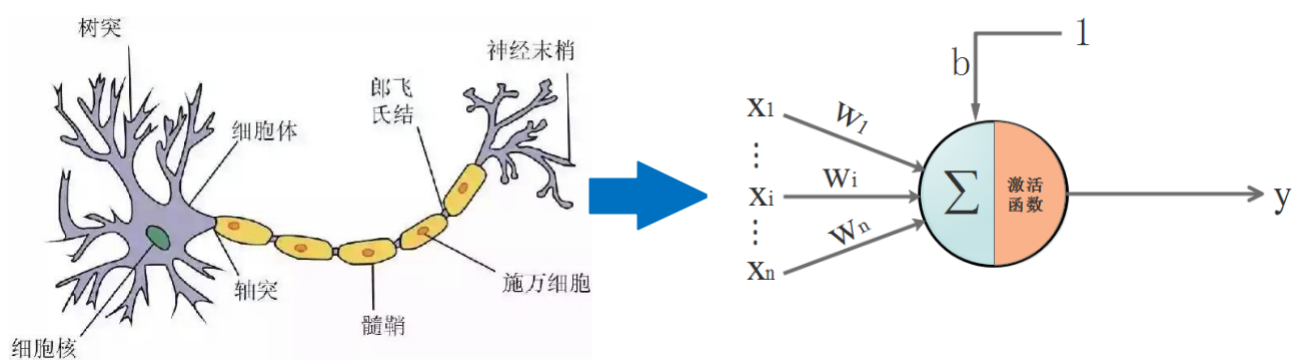
搞事！搞事！搞事



- 本文货很干，堪比沙哈拉大沙漠，自己挑的文章，含着泪也要读完！

二. 科普:

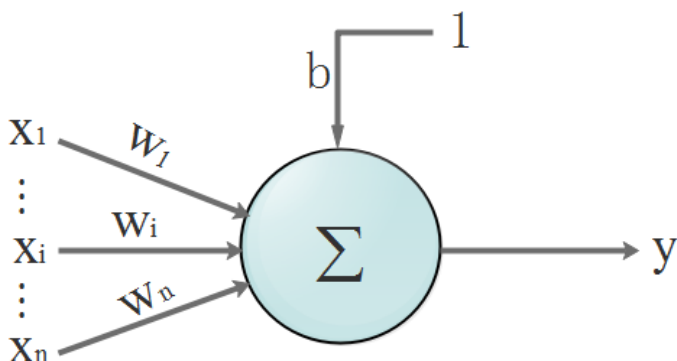
- 生物上的神经元就是接收四面八方的刺激（输入），然后做出反应（输出），给它一点☀就灿烂。
- 仿生嘛，于是喜欢放飞自我的 [某些人](#) 就提出了人工神经网络。一切的基础-->人工神经单元，看图：



三. 通往沙漠的入口: 神经元是什么, 有什么用:

开始前, 需要搞清楚一个很重要的问题: 人工神经网络里的神经元是什么, 有什么用。只有弄清楚这个问题, 你才知道你在哪里, 在做什么, 要往哪里去。

首先, 回顾一下神经元的结构, 看下图, 我们先忽略激活函数不管:



- 输入:

$$x_1, x_2, \dots, x_n \quad (3.1)$$

- 输出:

$$y \quad (3.2)$$

- 输入和输出的关系 (函数):

$$\begin{aligned} y &= (x_1 * w_1 + x_2 * w_2 + \dots + x_n * w_n) + b \\ &= \sum_{i=1}^n x_i * w_i + b \end{aligned} \quad (3.3)$$

其中, $w_i, i = 1, n$ 为权重 (weight), 待会就知道这货是什么了

没错, 开始晒公式了! 我们的数据都是离散的, 为了看得更清楚点, 所以换个表达方式, 把离散的数据写成向量。该不会忘了向量是啥吧? 回头致电问候一下当年的体育老师!

- 改写输入:

$$x = [x_1, x_2, \dots, x_n]^T \quad (3.4)$$

加 T 转置后, x 相当于一个 n 行 1 列的矩阵

- 改写权重:

$$w = [w_1, w_2, \dots, w_n] \quad (3.5)$$

- 那么输出就写成了:

$$\begin{aligned} y &= [w_1, w_2, \dots, w_n] \cdot [x_1, x_2, \dots, x_n]^T + b \\ &= wx + b \end{aligned}$$

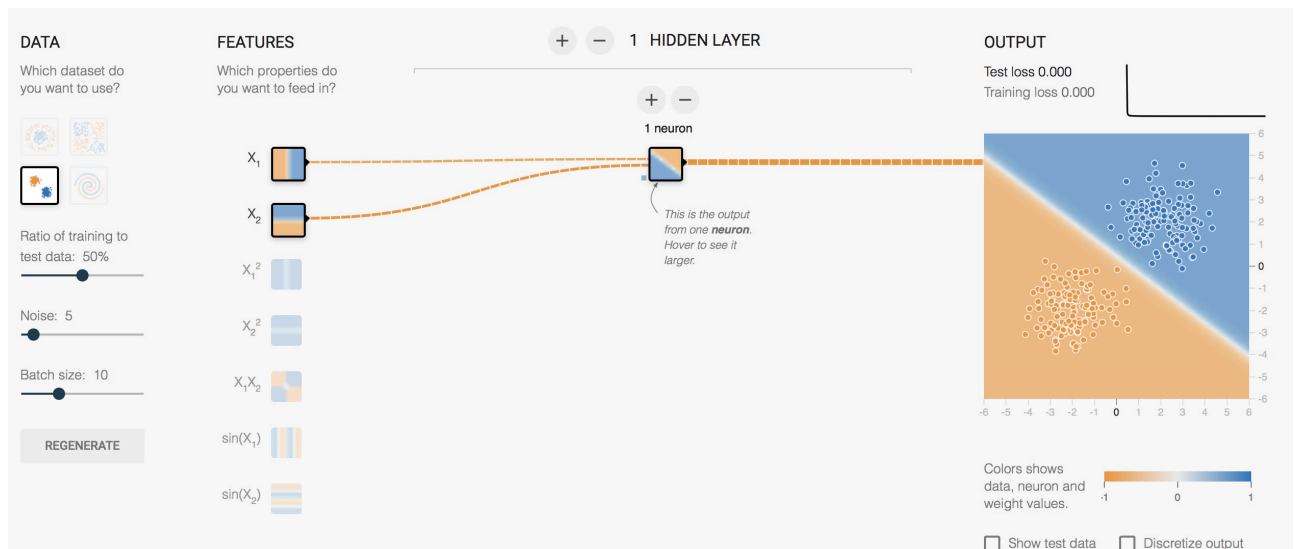
这是什么, 我们换个字母, 把 w 换成 k , 可以看到 w 就是直线的斜率啊啊啊:

$$= kx + b \quad (3.6)$$

现在回答问题刚才的问题:

- 一个神经元是什么：参照式（1.6），从函数图像角度看，这就是一根直线。
- 一个神经元有什么用：要说明用途就要给出一个应用场景：分类。一个神经元就是一条直线，相当于楚河汉界，可以把红棋绿棋分隔开，此时它就是个分类器。所以，在线性场景下，单个神经元能达到分类的作用，它总能学习到一条合适的直线，将两类元素区分出来。

先睹为快，看效果图，自己可以去玩：[传送门](#)



对上面的图简单说明一下：

- (x_1, x_2) 对于神经元的输入都是 x , 而对我们而言，这数据就是意义上的点的坐标，我们习惯写成 (x, y) 。

又要划重点了：

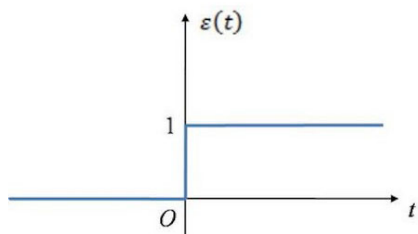


我们需要对神经元的输出做判定，那么就需要有判定规则，通过判定规则后我们才能拿到我们想要的结果，这个规则就是：

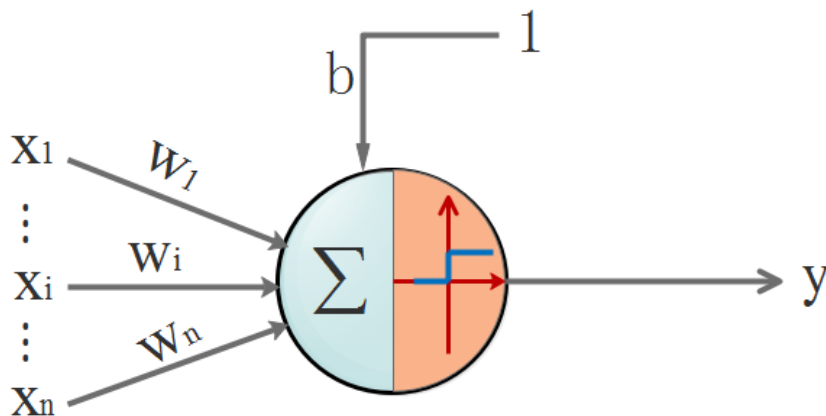
1. 假设，0代表红点，1代表蓝点（这些数据都是事先标定好的，在监督学习下，神经元会知道点是什么颜色并以这个已知结果作为标杆进行学习）
2. 当神经元输出小于等于 0 时，最终结果输出为 0，这是个红点
3. 当神经元输出大于 1 时，最终结果输出为 1，这是个蓝点

上面提到的规则让我闻到了激活函数的味道！（这里只是线性场景，虽然不合适，但是简单起见，使用了单位阶跃函数来描述激活函数的功能）当 $x \leq 0$ 时， $y = 0$ ；当 $x > 0$ 时， $y = 1$

这是阶跃函数的长相：



此时神经元的长相：



四. 茫茫大漠第一步：激活函数是什么，有什么用

从上面的例子，其实已经说明了激活函数的作用；但是，我们通常面临的问题，不是简单的线性问题，不能用单位阶跃函数作为激活函数，原因是：

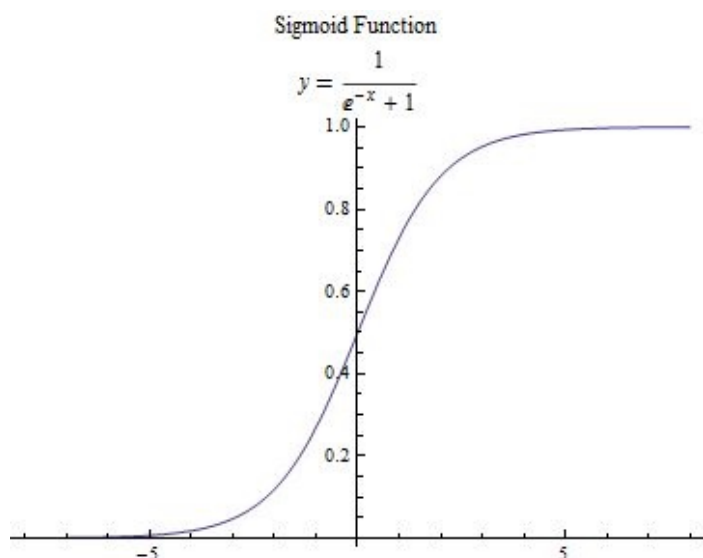
阶跃函数在 $x=0$ 时不连续，即不可导，在非 0 处导数为 0 。用人话说就是它具备输出限定在 $[0,1]$ ，但是它不具备丝滑的特性，这个特性很重要。并且在非 0 处导数为 0 ，也就是硬饱和，压根儿就没梯度可言，梯度也很重要，梯度意味着在神经元传播间是有反应的，而不是“死”了的。

接下来说明下，激活函数所具备的特性有什么，只挑重要的几点特性讲：

- ☑ 非线性：即导数不是常数，不然就退化成直线。对于一些画一条直线仍然无法分开的问题，非线性可以把直线掰弯，自从变弯以后，就能包罗万象了。
- ☑ 几乎处处可导：也就是具备“丝滑的特性”，不要应激过度，要做正常人。数学上，处处可导为后面降到的后向传播算法（BP算法）提供了核心条件
- ☑ 输出范围有限：一般是限定在 $[0,1]$ ，有限的输出范围使得神经元对于一些比较大的输入也会比较稳定。
- ☑ 非饱和性：饱和就是指，当输入比较大的时候，输出几乎没变化了，那么会导致梯度消失！什么是梯度消失：就是你天天给女生送花，一开始妹纸还惊喜，到后来直接麻木没反应了。梯度消失带来的负面影响就是会限制了神经网络表达能力，词穷的感觉你有过么。sigmoid, tanh函数都是软饱和的，阶跃函数是硬饱和。软是指输入趋于无穷大的时候输出无限接近上线，硬是指像阶跃函数那样，输入非 0 输出就已经始终都是上限值。数学表示我就懒得写了，[传送门在此](#)，里面有写到。如果激活函数是饱和的，带来的缺陷就是系统迭代更新变慢，系统收敛就慢，当然这是可以有办法弥补的，一种方法是使用交叉熵函数作为损失函数，这里不多说。ReLU是非饱和的，亲测效果挺不错，所以这货最近挺火的。
- ☑ 单调性：即导数符号不变。导出要么一直大于 0 ，要么一直小于 0 ，不要上蹿下跳。导数符号不变，让神经网络训练容易收敛。

这里只说我们用到的激活函数：

$$\text{Sigmoid函数: } y = \frac{1}{e^{-x} + 1} \quad (4.1)$$



求一下它的导数把，因为后面讲bp算法会直接套用它：

先祭出大杀器，高中数学之复合函数求导法则：

$$\text{法则1: } [u(x) \pm v(x)]' = u'(x) \pm v'(x) \quad (4.2)$$

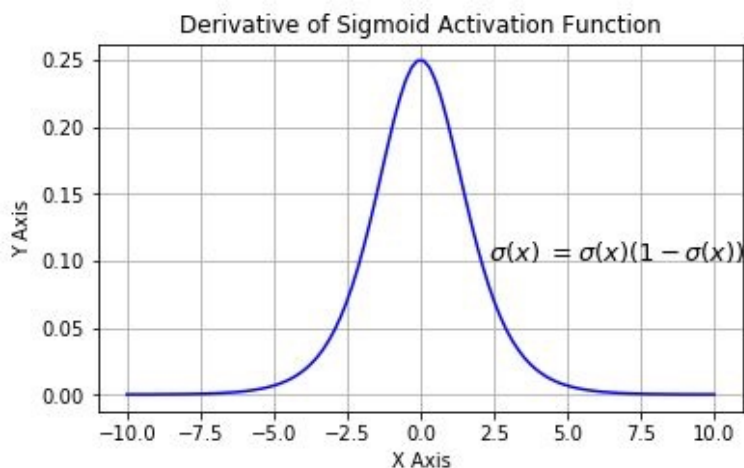
$$\text{法则2: } [u(x) * v(x)]' = u'(x) * v(x) + u(x) * v'(x) \quad (4.3)$$

$$\text{法则3: } \left(\frac{u}{v}\right)' = \frac{u'v - uv'}{v^2} \quad (4.4)$$

开始算算：

$$\begin{aligned} y' &= \left(\frac{1}{e^{-x} + 1}\right)' \\ &= \left(\frac{u}{v}\right)' \quad (\text{这里定: } u = 1, v = e^{-x} + 1) \\ &= \frac{u'v - uv'}{v^2} \\ &= \frac{1' * (e^{-x} + 1) - 1 * (e^{-x} + 1)'}{(e^{-x} + 1)^2} \\ &= \frac{e^{-x}}{(e^{-x} + 1)^2} \\ &= \frac{1}{1 + e^{-x}} * \frac{e^{-x}}{1 + e^{-x}} \\ &= \frac{1}{1 + e^{-x}} * \frac{1 + e^{-x} - 1}{1 + e^{-x}} \\ &= \frac{1}{1 + e^{-x}} * \left(1 - \frac{1}{1 + e^{-x}}\right), \text{ 因为: } y = \frac{1}{e^{-x} + 1}, \text{ 所以有:} \\ &= y * (1 - y) \quad (4.5) \end{aligned}$$

它的导数图像：



五. 沙漠中心的风暴：BP(Back Propagation)算法

1. 神经网络的结构

经过上面的介绍，单个神经元不足以让人心动，唯有组成网络。神经网络是一种分层结构，一般由输入层，隐藏层，输出层组成。所以神经网络至少有3层，隐藏层多于1，总层数大于3的就是我们所说的深度学习了。

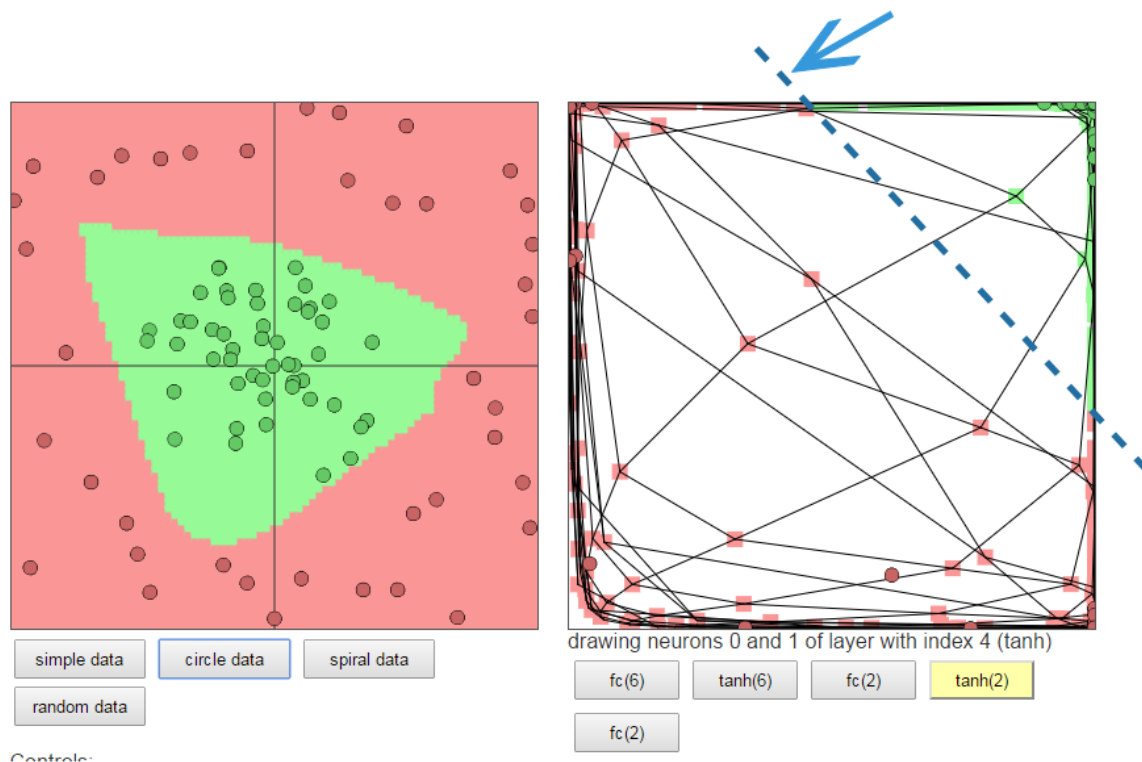
- 输入层：就是接收原始数据，然后往隐层送
- 输出层：神经网络的决策输出
- 隐藏层：该层可以说是神经网络的关键，相当于对数据做一次特征提取。隐藏层的意义，是把前一层的向量变成新的向量。就是坐标变换，说人话就是把数据做平移，旋转，伸缩，扭曲，让数据变得线性可分。可能这个不那么好理解，举个栗子：

下面的图左侧是原始数据，中间很多绿点，外围是很多红点，如果你是神经网络，你会怎么做呢？

一种做法：把左图的平面看成一块布，把它缝合成一个闭合的包包（相当于数据变换到了一个3维坐标空间），然后把有绿色点的部分撸到顶部（伸缩和扭曲），然后外围的红色点自然在另一端了，要是姿势还不够帅，就挪挪位置（平移）。这时候干脆利落的砍一刀，绿点红点就彻底区分开了。

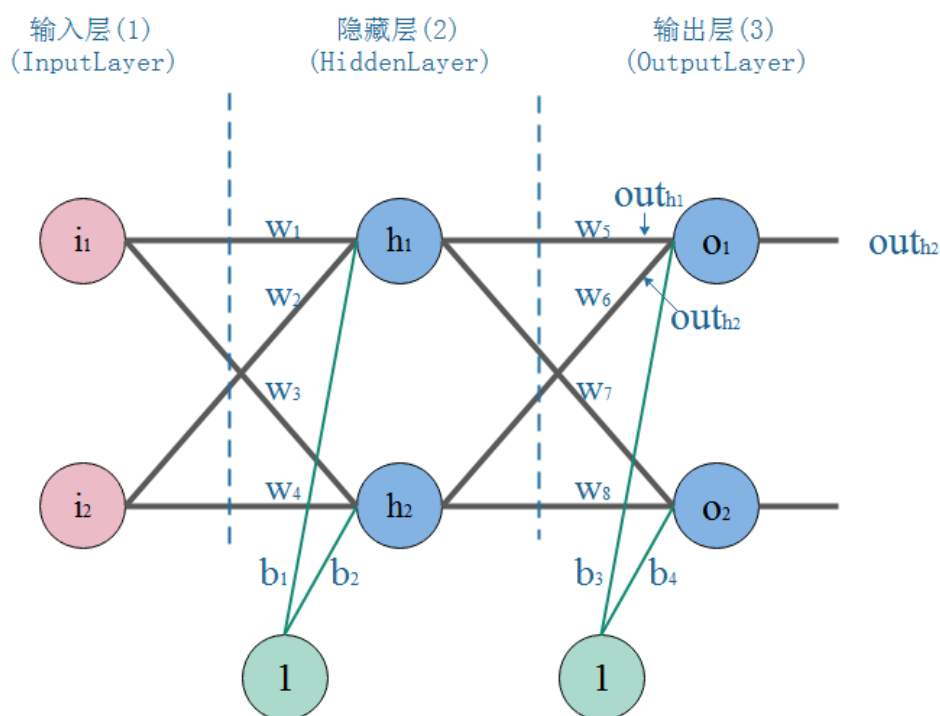
重要的东西再说一遍：神经网络换着坐标空间玩数据，根据需要，可降维，可升维，可大，可小，可圆可扁，就是这么“无敌”

这个也可以自己去玩玩，直观的感受一下：[传送门](#)



2.正反向传播过程

看图，这是一个典型的三层神经网络结构，第一层是输入层，第二层是隐藏层，第三层是输出层。PS:不同的应用场景，神经网络的结构要有针对性的设计，这里仅仅是为了推导算法和计算方便才采用这个简单的结构



我们以战士打靶，目标是训练战士能命中靶心成为神枪手作为场景：

那么我们手里有这样一些数据：一堆枪摆放的位置(x,y)，以及射击结果，命中靶心和不命中靶心。

我们的目标是：训练出一个神经网络模型，输入一个点的坐标（射击姿势），它就告诉你这个点是什么结果（是否命中）。

我们的方法是：训练一个能根据误差不断自我调整的模型，训练模型的步骤是：

- 正向传播：把点的坐标数据输入神经网络，然后开始一层一层的传播下去，直到输出层输出结果。
- 反向传播(BP)：就好比战士去靶场打靶，枪的摆放位置（输入），和靶心🎯（期望的输出）是已知。战士（神经网络）一开始的时候是这样做的，随便开一枪（w, b参数初始化称随机值），观察结果（这时候相当于进行了一次正向传播）。然后发现，偏离靶心左边，应该往右点儿打。所以战士开始根据偏离靶心的距离（误差，也称损失）调整了射击方向往右一点（这时，完成了一次反向传播）
- 当完成了一次正反向传播，也就完成了一次神经网络的训练迭代，反复调整射击角度（反复迭代），误差越来越小，战士打得越来越准，神枪手模型也就诞生了。

3.BP算法推导和计算

- 参数初始化：

输入： $i_1 = 0.1, i_2 = 0.2$

输出： $O_1 = 0.01, O_2 = 0.99$ ，相当于标定了

权重： $w_1 = 0.1, w_2 = 0.2, w_3 = 0.3, w_4 = 0.4$
 $w_5 = 0.5, w_6 = 0.6, w_7 = 0.7, w_8 = 0.8$

偏置： $b_1 = 0.55, b_2 = 0.56, b_3 = 0.66, b_4 = 0.67$

- 正向传播：

1.输入层-->隐层：

- 计算隐层神经元 h1 的输入加权和：

$$\begin{aligned} in_{h_1} &= w_1 * i_1 + w_2 * i_2 + 1 * b_1 \quad (5.1) \\ &= 0.1 * 0.1 + 0.2 * 0.2 + 1 * 0.55 \\ &= 0.6 \end{aligned}$$

- 计算隐层神经元 h1 的输出,需要通过激活函数Sigmoid，记得吧？：

$$\begin{aligned} out_{h_1} &= \frac{1}{e^{-in_{h_1}} + 1} \quad (5.2) \\ &= \frac{1}{e^{-0.6} + 1} \\ &= 0.6456563062 \end{aligned}$$

- 同理，可以算出隐层神经元 h2 的输出：

$$out_{h_2} = 0.6592603884$$

2.隐层-->输出层：

- 计算输出层神经元O1的输入加权和：

$$\begin{aligned} in_{O_1} &= w_5 * out_{h_1} + w_6 * out_{h_2} + 1 * b_3 \quad (5.3) \\ &= 0.5 * 0.6456563062 + 0.6 * 0.6592603884 + 1 * 0.66 \\ &= 1.3783843861 \end{aligned}$$

- 计算隐层神经元O1的输出：

$$\begin{aligned}
 out_{O_1} &= \frac{1}{e^{-in_{O_1}} + 1} \quad (5.4) \\
 &= \frac{1}{e^{-1.3783843861} + 1} \\
 &= 0.7987314002
 \end{aligned}$$

- 同理，可以算出隐层神经元O2的输出：

$$out_{O_2} = 0.8374488853$$

正向传播结束，我们看看输出层的输出结果：[0.7987314002, 0.8374488853]，但是我们希望它能输出 [0.01, 0.99]，所以明显的差太远了，这个时候我们就需要利用反向传播，更新权值w，然后重新计算输出

• 反向传播：

1. 计算输出误差：

$$\begin{aligned}
 E_{total} &= \sum_{i=1}^2 E_{out_{O_i}} \quad (5.5) \\
 &= E_{out_{O_1}} + E_{out_{O_2}} \\
 &= \frac{1}{2} (expected_{out_1} - out_{O_1})^2 + \frac{1}{2} (expected_{out_2} - out_{O_2})^2 \\
 &= \frac{1}{2} (O_1 - out_{O_1})^2 + \frac{1}{2} (O_2 - out_{O_2})^2 \\
 &= \frac{1}{2} * (0.01 - 0.7987314002)^2 + \frac{1}{2} * (0.99 - 0.8374488853)^2 \\
 &= 0.0116359213 + 0.3110486109 \\
 &= 0.3226845322 \\
 &\text{其中：} E_{out_{O_1}} = 0.0116359213, E_{out_{O_2}} = 0.3110486109
 \end{aligned}$$

PS: 这里我要说的是，用这个作为误差的计算，因为它简单，实际上用的时候效果不咋滴。【原因上文我提过了】：我在第四章说激活函数作用时，提到激活函数应具备“非饱和性”。如果激活函数是饱和的，带来的缺陷就是系统迭代更新变慢，系统收敛就慢，当然这是可以有办法弥补的，一种方法是使用交叉熵函数作为损失函数。

交叉熵做为代价函数能达到上面说的优化系统收敛下欧工，是因为它在计算误差对输入的梯度时，抵消了激活函数的导数项，从而避免了因为激活函数的“饱和性”给系统带来的负面影响。如果项了解更详细的证明可以点 --> [传送门](#)

$$E_{total} = \frac{1}{m} \sum_{i=1}^m (O \cdot \log out_O + (1 - O) \cdot \log(1 - out_O))$$

对输出的偏导数：

$$\frac{\partial E_{total}}{\partial out_O} = \frac{1}{m} \sum_{i=1}^m \left(\frac{O}{out_O} - \frac{1 - O}{1 - out_O} \right)$$

2. 隐层-->输出层的权值及偏置b的更新：

- 先放出链式求导法则：

假设 y 是 u 的函数，而 u 是 x 的函数： $y = f(u)$, $u = g(x)$

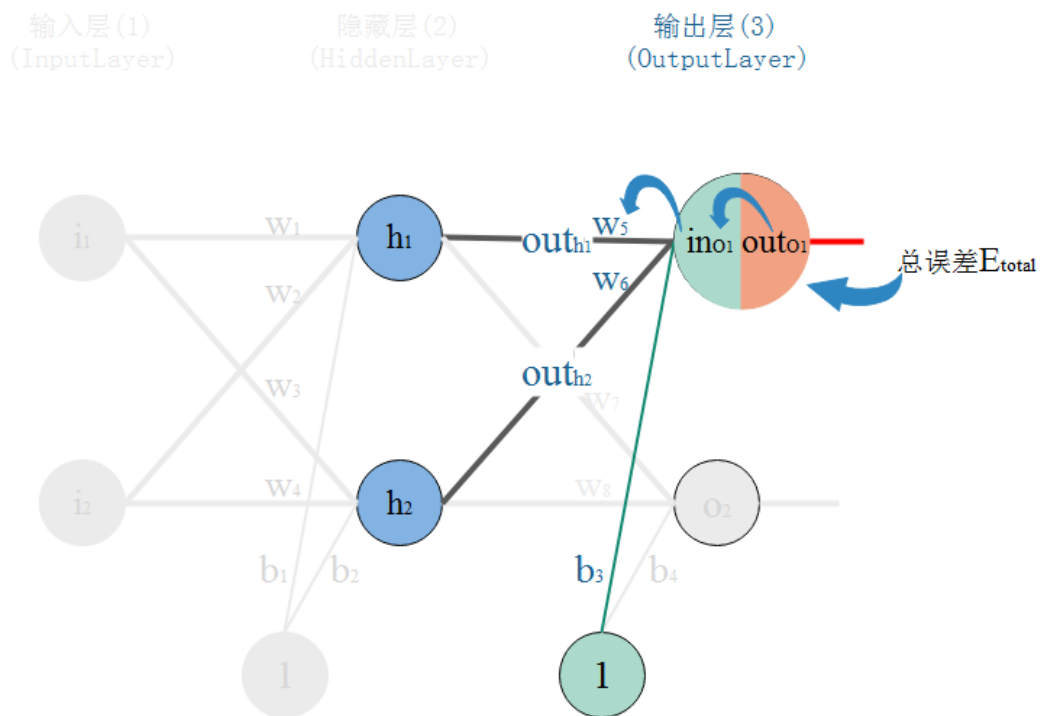
那么对应的复合函数就是： $y = f(g(x))$

那么 y 对 x 的导数则有： $\frac{dy}{dx} = \frac{dy}{du} \cdot \frac{du}{dx}$

◦ 以更新 w_5 举例：

我们知道，权重 w 的大小能直接影响输出， w 不合适那么会使得输出误差。要想直到某一个 w 值对误差影响的程度，可以用误差对该 w 的变化率来表达。如果 w 的一点点变动，就会导致误差增大很多，说明这个 w 对误差影响的程度就更大，也就是说，误差对该 w 的变化率越高。而误差对 w 的变化率就是误差对 w 的偏导。

所以，看下图，总误差的大小首先受输出层神经元 O_1 的输出影响，继续反推， O_1 的输出受它自己的输入的影响，而它自己的输入会受到 w_5 的影响。这就是连锁反应，从结果找根因。



那么，根据链式法则则有：

$$\frac{\partial E_{total}}{\partial w_5} = \frac{\partial E_{total}}{\partial out_{O_1}} \cdot \frac{\partial out_{O_1}}{\partial in_{O_1}} \cdot \frac{\partial in_{O_1}}{\partial w_5} \quad (5.6)$$

现在挨个计算：

$$\begin{aligned} \text{因为：} E_{total} &= \frac{1}{2} (O_1 - out_{O_1})^2 + \frac{1}{2} (O_2 - out_{O_2})^2 \\ \text{则有：} \frac{\partial E_{total}}{\partial out_{O_1}} &= \frac{\partial (\frac{1}{2} (O_1 - out_{O_1})^2 + \frac{1}{2} (O_2 - out_{O_2})^2)}{\partial out_{O_1}} \end{aligned} \quad (5.7)$$

$$\begin{aligned} &= 2 * \frac{1}{2} (O_1 - out_{O_1})^{2-1} * (0 - 1) + 0 \\ &= -(O_1 - out_{O_1}) \\ &= -(0.01 - 0.7987314002) \\ &= 0.7887314002 \end{aligned} \quad (5.8)$$

激活函数的导数看公式（4.5）：

$$out_{O_1} = \frac{1}{e^{-in_{O_1}} + 1} \quad (5.8)$$

$$\begin{aligned} \frac{\partial out_{O_1}}{\partial in_{O_1}} &= \frac{\partial \left(\frac{1}{e^{-in_{O_1}} + 1} \right)}{\partial in_{O_1}} \\ &= out_{O_1} (1 - out_{O_1}) \quad (5.9) \\ &= 0.7987314002 * (1 - 0.7987314002) \\ &= 0.1607595505 \end{aligned}$$

$$in_{O_1} = w_5 * out_{h_1} + w_6 * out_{h_2} + 1 * b_3 \quad (5.10)$$

$$\begin{aligned} \frac{\partial in_{O_1}}{\partial w_5} &= \frac{\partial (w_5 * out_{h_1} + w_6 * out_{h_2} + 1 * b_3)}{\partial w_5} \quad (5.11) \\ &= 1 * w_5^{(1-1)} * out_{h_1} + 0 + 0 \\ &= out_{h_1} \\ &= 0.6456563062 \end{aligned}$$

所以：

$$\begin{aligned} \frac{\partial E_{total}}{\partial w_5} &= \frac{\partial E_{total}}{\partial out_{O_1}} \cdot \frac{\partial out_{O_1}}{\partial in_{O_1}} \cdot \frac{\partial in_{O_1}}{\partial w_5} \quad (5.12) \\ &= 0.7887314002 * 0.1607595505 * 0.6456563062 \\ &= 0.0818667051 \end{aligned}$$

我们归纳一下式子：

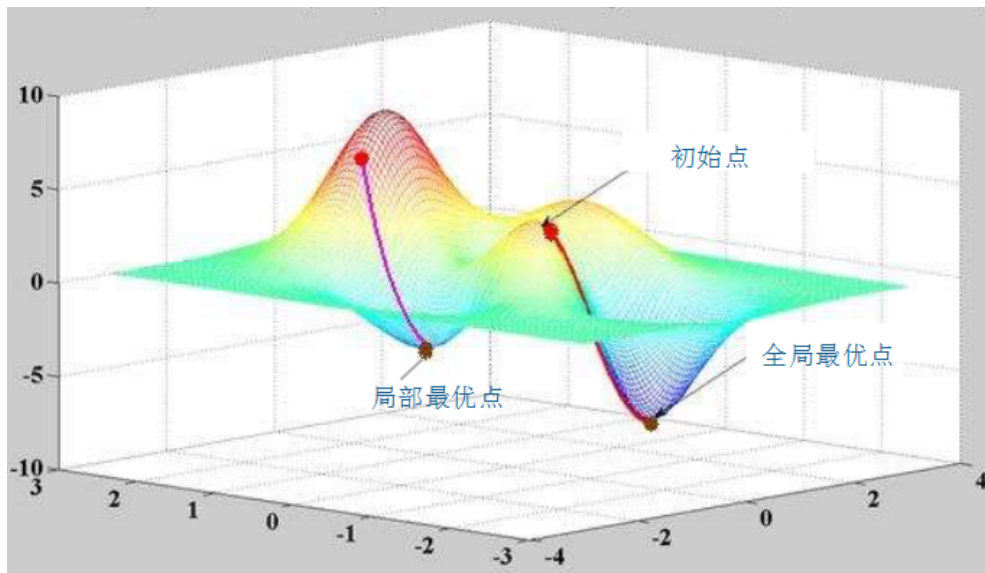
$$\begin{aligned} \frac{\partial E_{total}}{\partial w_5} &= \frac{\partial E_{total}}{\partial out_{O_1}} \cdot \frac{\partial out_{O_1}}{\partial in_{O_1}} \cdot \frac{\partial in_{O_1}}{\partial w_5} \\ &= -(O_1 - out_{O_1}) \cdot out_{O_1} \cdot (1 - out_{O_1}) \cdot out_{h_1} \quad (5.13) \\ &= \sigma_{O_1} \cdot out_{h_1} \\ \text{其中, } \sigma_{O_1} &= -(O_1 - out_{O_1}) \cdot out_{O_1} \cdot (1 - out_{O_1}) \end{aligned}$$

同理，更新输出层偏置b如下：

$$\begin{aligned} in_{O_1} &= w_5 * out_{h_1} + w_6 * out_{h_2} + 1 * b_3 \\ \frac{\partial in_{O_1}}{\partial b_3} &= \frac{\partial (w_5 * out_{h_1} + w_6 * out_{h_2} + 1 * b_3)}{\partial b_3} \quad (5.14) \\ &= 0 + 0 + b_3^{(1-1)} \\ &= 1 \\ \frac{\partial E_{total}}{\partial b_3} &= \frac{\partial E_{total}}{\partial out_{O_1}} \cdot \frac{\partial out_{O_1}}{\partial in_{O_1}} \cdot \frac{\partial in_{O_1}}{\partial b_3} \\ &= -(O_1 - out_{O_1}) \cdot out_{O_1} \cdot (1 - out_{O_1}) \cdot 1 \quad (5.15) \\ &= \sigma_{O_1} \\ \text{其中, } \sigma_{O_1} &= -(O_1 - out_{O_1}) \cdot out_{O_1} \cdot (1 - out_{O_1}) \end{aligned}$$

有个学习率的东西，学习率取个0.5。关于学习率，不能过高也不能过低。因为训练神经网络系统的过程，就是通过不断的迭代，找到让系统输出误差最小的参数的过程。每一次迭代都经过反向传播进行梯度下降，然而误差空间不是一个滑梯，一降到底，常规情况下就像坑洼的山地。学习率太小，那就很容易陷入局部最优，就是你认为的最低点并不是整个空间的最低点。如果学习率太高，那系统可能难以收敛，会在一个地方上串下跳，无法对准目标（目标是指误差空间的最低点），可以看图：

xy轴是权值w平面，z轴是输出总误差。整个误差曲面可以看到两个明显的低点，显然右边最低，属于全局最优。而左边的是次低，从局部范围看，属于局部最优。而图中，在给定初始点的情况下，标出的两条抵达低点的路线，已经是很理想情况的梯度下降路径。



现在可以更新w5的值了,就设定学习率为0.5吧:

$$\begin{aligned} w_5^+ &= w_5 - \alpha \cdot \frac{\partial E_{total}}{\partial w_5} \quad (5.16) \\ &= 0.5 - 0.5 * 0.0818667051 \\ &= 0.45906664745 \end{aligned}$$

归纳一下输出层w更新的公式:

$$\begin{aligned} w_O^+ &= w_O - \alpha \cdot (-out_O \cdot (1 - out_O) \cdot (O - out_O) \cdot out_h) \\ &= w_O + \alpha \cdot (O - out_O) \cdot out_O \cdot (1 - out_O) \cdot out_h \quad (5.17) \end{aligned}$$

同理可以计算出w6, w7, w8的更新值, 懒癌晚期, 懒得算了...

$$\begin{aligned} w_6^+ &= \dots \\ w_7^+ &= \dots \\ w_8^+ &= \dots \end{aligned}$$

同理更新偏置b:

$$b^+ = b_O - \alpha \cdot \frac{\partial E_{total}}{\partial b_O} \quad (5.18)$$

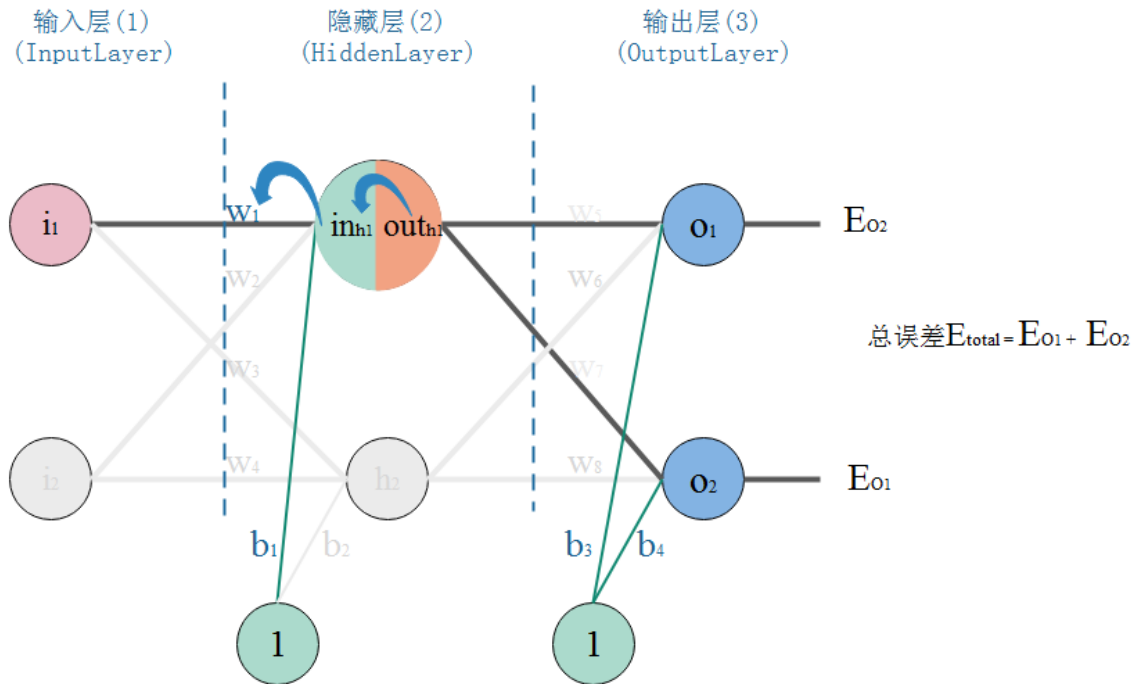
归纳一下输出层w更新的公式:

$$\begin{aligned}
 b_O^+ &= b_O - \alpha \cdot (-out_O \cdot (1 - out_O) \cdot (O - out_O)) \\
 &= b_O + \alpha \cdot (O - out_O) \cdot out_O \cdot (1 - out_O) \quad (5.19)
 \end{aligned}$$

3. 输入层-->隐层的权值及偏置**b**更新:

◦ 以更新**w1**为例:

仔细观察，我们在求w5的更新，误差反向传递路径输出层-->隐层，即out(O1)-->in(O1)-->w5，总误差只有一根线能传回来。但是求w1时，误差反向传递路径是隐藏层-->输入层，但是隐藏层的神经元是有2根线的，所以总误差沿着2个路径回来，也就是说，计算偏导时，要分开来算。看图：



那么，现在开始算总误差对**w1**的偏导：

$$\frac{\partial E_{total}}{\partial w_1} = \frac{\partial E_{total}}{\partial out_{h_1}} \cdot \frac{\partial out_{h_1}}{\partial in_{h_1}} \cdot \frac{\partial in_{h_1}}{\partial w_1} \quad (5.20)$$

$$= \left(\frac{\partial E_{O_1}}{\partial out_{h_1}} + \frac{\partial E_{O_2}}{\partial out_{h_1}} \right) \cdot \frac{\partial out_{h_1}}{\partial in_{h_1}} \cdot \frac{\partial in_{h_1}}{\partial w_1} \quad (5.21)$$

3.1 现在先算：

$$\frac{\partial E_{total}}{\partial out_{h_1}}$$

也就是：

$$\frac{\partial E_{total}}{\partial out_{h_1}} = \frac{\partial E_{O_1}}{\partial out_{h_1}} + \frac{\partial E_{O_2}}{\partial out_{h_1}} \quad (5.22)$$

挨个算：

$$\frac{\partial E_{O_1}}{\partial out_{h_1}} = \frac{\partial E_{O_1}}{\partial in_{O_1}} \cdot \frac{\partial in_{O_1}}{\partial out_{h_1}} \quad (5.23)$$

- 计算左边部分，参考式子(5.7), (5.8), (5.9):

$$\begin{aligned}
 \frac{\partial E_{O_1}}{\partial in_{O_1}} &= \frac{\partial E_{O_1}}{\partial out_{O_1}} \cdot \frac{\partial out_{O_1}}{\partial in_{O_1}} & (5.24) \\
 &= \frac{\partial \left(\frac{1}{2} (O_1 - out_{O_1})^2 \right)}{\partial out_{O_1}} \cdot \frac{\partial out_{O_1}}{\partial in_{O_1}} \\
 &= - (O_1 - out_{O_1}) \cdot \frac{\partial out_{O_1}}{\partial in_{O_1}} \\
 &= 0.7987314002 * 0.1607595505 \\
 &= 0.1284037009
 \end{aligned}$$

- 计算右边部分:

$$in_{O_1} = w_5 * out_{h_1} + w_6 * out_{h_2} + 1 * b_3$$

$$\begin{aligned}
 \frac{\partial in_{O_1}}{\partial out_{h_1}} &= \frac{\partial (w_5 * out_{h_1} + w_6 * out_{h_2} + 1 * b_3)}{\partial out_{h_1}} \\
 &= w_5 * out_{h_1}^{(1-1)} + 0 + 0 \\
 &= w_5 & (5.25) \\
 &= 0.5
 \end{aligned}$$

所以:

$$\begin{aligned}
 \frac{\partial E_{O_1}}{\partial out_{h_1}} &= \frac{\partial E_{O_1}}{\partial in_{O_1}} \cdot \frac{\partial in_{O_1}}{\partial out_{h_1}} & (5.26) \\
 &= 0.1284037009 * 0.5 \\
 &= 0.06420185045
 \end{aligned}$$

同理 ~ 同理 ~ :

$$\begin{aligned}
 \frac{\partial E_{O_2}}{\partial out_{h_1}} &= \frac{\partial E_{O_2}}{\partial in_{O_2}} \cdot \frac{\partial in_{O_2}}{\partial out_{h_1}} \\
 &= - (O_2 - out_{O_2}) \cdot \frac{\partial out_{O_2}}{\partial in_{O_2}} \cdot \frac{\partial in_{O_2}}{\partial out_{h_1}} \\
 &= - (O_2 - out_{O_2}) \cdot out_{O_2} (1 - out_{O_2}) \cdot w_7 & (5.27) \\
 &= -(0.99 - 0.8374488853) * 0.8374488853 * (1 - 0.8374488853) * 0.7 \\
 &= -0.0145365614
 \end{aligned}$$

所以(3.19)的值为:

$$\begin{aligned}
 \frac{\partial E_{total}}{\partial out_{h_1}} &= \frac{\partial E_{O_1}}{\partial out_{h_1}} + \frac{\partial E_{O_2}}{\partial out_{h_1}} \\
 &= 0.06420185045 + (-0.0145365614) \\
 &= 0.04966528905
 \end{aligned}$$

3.2 然后算:

$$\frac{\partial out_{h_1}}{\partial in_{h_1}}$$

$$\begin{aligned}\therefore out_{h_1} &= \frac{1}{e^{-in_{h_1}} + 1} \\ \therefore \frac{\partial out_{h_1}}{\partial in_{h_1}} &= \frac{\partial \left(\frac{1}{e^{-in_{h_1}} + 1} \right)}{\partial in_{h_1}} \\ &= out_{h_1} (1 - out_{h_1}) \\ &= 0.6456563062 * (1 - 0.6456563062) \\ &= 0.2287842405\end{aligned}$$

3.3 最后算：

$$\begin{aligned}\frac{\partial in_{h_1}}{\partial w_1} &= \frac{\partial (w_1 * i_1 + w_2 * i_2 + 1 * b)}{\partial w_1} \\ &= w_1^{(1-1)} * i_1 + 0 + 0 = i_1 = 0.1\end{aligned}$$

最后，将**3**者相乘，就算出：

$$\begin{aligned}\frac{\partial E_{total}}{\partial w_1} &= \frac{\partial E_{total}}{\partial out_{h_1}} \cdot \frac{\partial out_{h_1}}{\partial in_{h_1}} \cdot \frac{\partial in_{h_1}}{\partial w_1} \\ &= 0.04966528905 * 0.2287842405 * 0.1 \\ &= 0.0011362635\end{aligned}$$

我们归纳一下式子：

$$\begin{aligned}\frac{\partial E_{total}}{\partial w_1} &= \frac{\partial E_{total}}{\partial out_{h_1}} \cdot \frac{\partial out_{h_1}}{\partial in_{h_1}} \cdot \frac{\partial in_{h_1}}{\partial w_1} \\ &= \left(\frac{\partial E_{O_1}}{\partial out_{h_1}} + \frac{\partial E_{O_2}}{\partial out_{h_1}} \right) \cdot \frac{\partial out_{h_1}}{\partial in_{h_1}} \cdot \frac{\partial in_{h_1}}{\partial w_1} \\ &= \left(\sum_O \frac{\partial E_O}{\partial out_O} \cdot \frac{\partial out_O}{\partial in_O} \cdot \frac{\partial in_O}{\partial out_h} \right) \cdot \frac{\partial out_{h_1}}{\partial in_{h_1}} \cdot \frac{\partial in_{h_1}}{\partial w_1} \\ &= \left(\sum_O \sigma_O w_O \right) \cdot out_{h_1} (1 - out_{h_1}) \cdot i_1 \quad (5.28) \\ &= \sigma_{h_1} \cdot i_1\end{aligned}$$

$$\text{其中, } \sigma_{h_1} = \left(\sum_O \sigma_O w_O \right) \cdot out_{h_1} (1 - out_{h_1})$$

σ_O 看作输出层的误差量，然后该误差量和 w 相乘，相当于通过 w 以传播了过来；
如果是深层网络，隐藏层数量 > 1 ，那么公式中的 σ_O 写成 σ_h ， w_O 写成 w_h

现在，可以更新 **w1** 的值了：

$$\begin{aligned}w_1^+ &= w_1 - \alpha \cdot \frac{\partial E_{total}}{\partial w_1} \quad (5.29) \\ &= 0.1 - 0.1 * 0.0011362635 \\ &= 0.0998863737\end{aligned}$$

归纳一下，隐藏层 w 更新的公式：

$$\begin{aligned}
 w_h^+ &= w_h - \alpha \cdot \frac{\partial E_{total}}{\partial w} \\
 &= w_h + \alpha \cdot \left(- \sum_O \sigma_O w_O \right) \cdot out_h (1 - out_h) \cdot i \quad (5.30)
 \end{aligned}$$

如果隐藏层数量 > 1:

$$\begin{aligned}
 w_h^+ &= w_h - \alpha \cdot \frac{\partial E_{total}}{\partial w_h} \\
 &= w_h + \alpha \cdot \left(- \sum_{hh} \sigma_{hh} w_{hh} \right) \cdot out_h (1 - out_h) \cdot in_h \quad (5.31)
 \end{aligned}$$

hh 代表当前隐藏层的下一个隐藏层，有点拗口；深层网络，计算的式子就是递归计算的了。

同理可以计算出 w_2 ， w_3 ， w_4 的更新值，懒...

$$w_2^+ = \dots$$

$$w_3^+ = \dots$$

$$w_4^+ = \dots$$

同理，隐藏层偏置 b 的更新：

$$\frac{\partial E_{total}}{\partial b_h} = \left(\sum_h \sigma_h w_h \right) \cdot out_h (1 - out_h) \quad (5.32)$$

$$\begin{aligned}
 b_h^+ &= b_h - \alpha \cdot \frac{\partial E_{total}}{\partial b_h} \\
 &= b_h + \alpha \cdot \left(- \sum_h \sigma_h w_h \right) \cdot out_h (1 - out_h) \quad (5.33)
 \end{aligned}$$

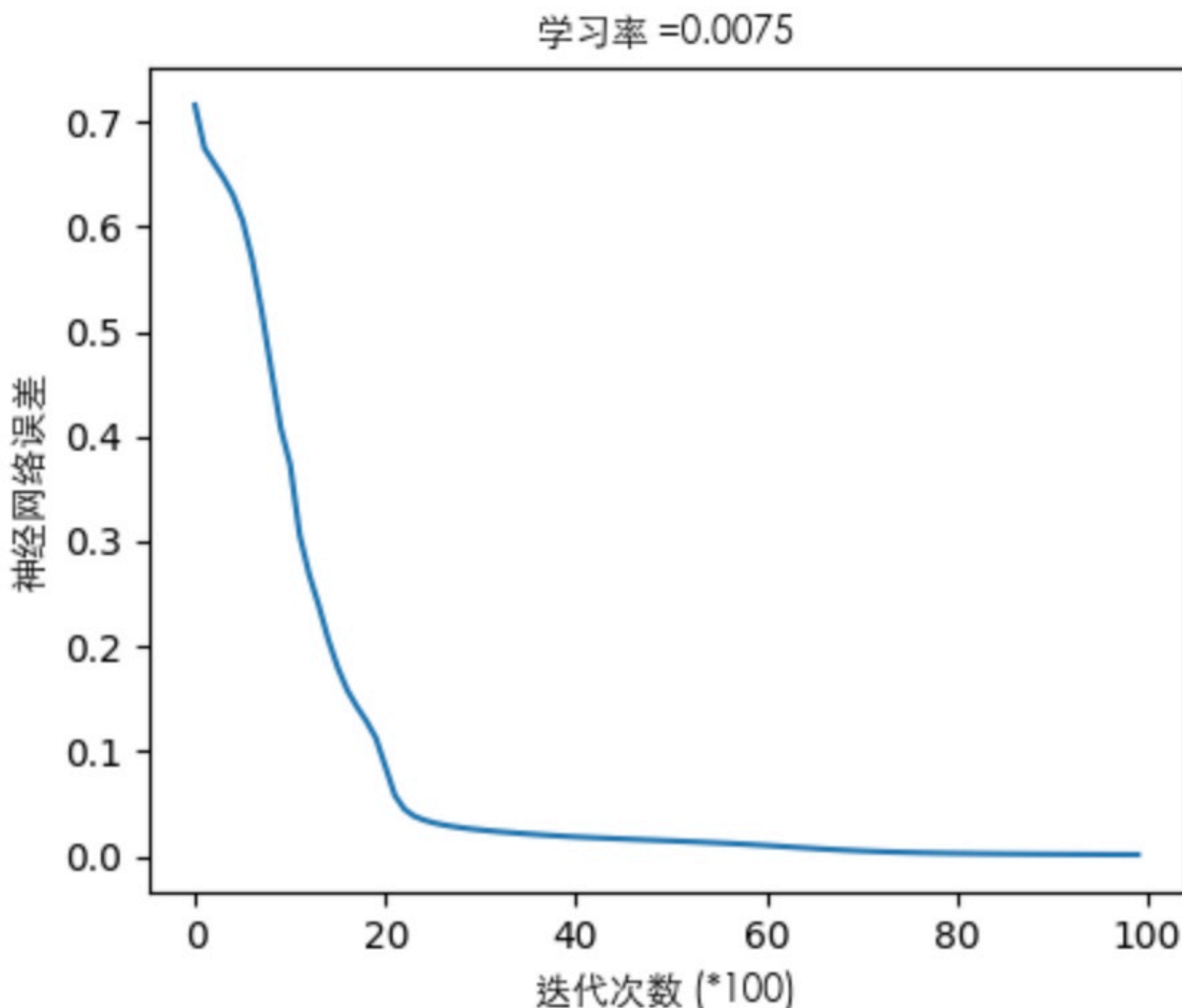
如果隐藏层数量 > 1:

$$\begin{aligned}
 b_h^+ &= b_h - \alpha \cdot \frac{\partial E_{total}}{\partial b_h} \\
 &= b_h + \alpha \cdot \left(- \sum_{hh} \sigma_{hh} w_{hh} \right) \cdot out_h (1 - out_h) \quad (5.34)
 \end{aligned}$$

同上 hh 代表当前隐藏层的下一个隐藏层。

4.结论：

我们通过亲力亲为的计算，走过了正向传播，也体会了反向传播，完成了一次训练(迭代)。随着迭代加深，输出层的误差会越来越小，专业点说就是系统趋于收敛。来一张系统误差随迭代次数变化的图来表明我刚才说描述：



六. 沙漠的绿洲：代码实现

1. 代码代码！

其实已经有很多机器学习的框架可以很简单的实现神经网络。但是我们的目标是：在看懂算法之后，我们是否能照着算法的整个过程，去实现一遍，可以加深对算法原理的理解，以及对算法实现思路的理解。顺便说打个call，numpy这个库，你值得拥有！

- 代码实现如下。代码里已经做了尽量啰嗦的注释，关键实现的地方对标了公式的编号，如果看不明白的地方多回来啃一下算法推导。对应代码也传到了github上。
- 代码能自己定义神经网络的结构，支持深度网络。代码实现了对红蓝颜色的点做分类的模型训练，通过3层网络结构，改变隐藏层的神经元个数，通过图形显示隐藏层神经元数量对问题的解释能力。
- 代码中还实现了不同激活函数。隐藏层可以根据需要换着激活函数玩，输出层一般就用sigmoid，当然想换也随你喜欢～

```
1 #coding:utf-8
2 import h5py
3 import sklearn.datasets
4 import sklearn.linear_model
5 import matplotlib
6 import matplotlib.font_manager as fm
```

```

7 import matplotlib.pyplot as plt
8 import numpy as np
9
10 np.random.seed(1)
11
12 font = fm.FontProperties(fname='/System/Library/Fonts/STHeiti Light.ttc')
13 matplotlib.rcParams['figure.figsize'] = (10.0, 8.0)
14
15 def sigmoid(input_sum):
16     """
17     函数:
18         激活函数Sigmoid
19     输入:
20         input_sum: 输入, 即神经元的加权和
21     返回:
22         output: 激活后的输出
23         input_sum: 把输入缓存起来返回
24     """
25     output = 1.0/(1+np.exp(-input_sum))
26     return output, input_sum
27
28
29 def sigmoid_back_propagation(derror_wrt_output, input_sum):
30     """
31     函数:
32         误差关于神经元输入的偏导:  $dE/dIn = dE/dOut * dOut/dIn$  参照式 (5.6)
33         其中:  $dOut/dIn$  就是激活函数的导数  $dy=y(1-y)$ , 见式 (5.9)
34          $dE/dOut$  误差对神经元输出的偏导, 见式 (5.8)
35     输入:
36         derror_wrt_output: 误差关于神经元输出的偏导:  $dE/dy_j = 1/2(d(expect\_to\_output - output)**2/doutput) = -(expect\_to\_output - output)$ 
37         input_sum: 输入加权和
38     返回:
39         derror_wrt_dinputs: 误差关于输入的偏导, 见式 (5.13)
40     """
41     output = 1.0/(1 + np.exp(- input_sum))
42     doutput_wrt_dinput = output * (1 - output)
43     derror_wrt_dinput = derror_wrt_output * doutput_wrt_dinput
44
45     return derror_wrt_dinput
46
47
48 def relu(input_sum):
49     """
50     函数:
51         激活函数ReLU
52     输入:
53         input_sum: 输入, 即神经元的加权和
54     返回:
55         outputs: 激活后的输出
56         input_sum: 把输入缓存起来返回
57     """
58     output = np.maximum(0, input_sum)
59     return output, input_sum
60
61

```

```

62 def relu_back_propagation(derror_wrt_output, input_sum):
63     """
64     函数:
65         误差关于神经元输入的偏导:  $dE / dIn = dE / dOut * dOut / dIn$ 
66         其中:  $dOut / dIn$  就是激活函数的导数
67          $dE / dOut$  误差对神经元输出的偏导
68     输入:
69         derror_wrt_output: 误差关于神经元输出的偏导
70         input_sum: 输入加权和
71     返回:
72         derror_wrt_dinputs: 误差关于输入的偏导
73     """
74     derror_wrt_dinputs = np.array(derror_wrt_output, copy=True)
75     derror_wrt_dinputs[input_sum <= 0] = 0
76
77     return derror_wrt_dinputs
78
79
80 def tanh(input_sum):
81     """
82     函数:
83         激活函数 tanh
84     输入:
85         input_sum: 输入, 即神经元的加权和
86     返回:
87         output: 激活后的输出
88         input_sum: 把输入缓存起来返回
89     """
90     output = np.tanh(input_sum)
91     return output, input_sum
92
93
94 def tanh_back_propagation(derror_wrt_output, input_sum):
95     """
96     函数:
97         误差关于神经元输入的偏导:  $dE / dIn = dE / dOut * dOut / dIn$ 
98         其中:  $dOut / dIn$  就是激活函数的导数  $\tanh'(x) = 1 - x^2$ 
99          $dE / dOut$  误差对神经元输出的偏导
100     输入:
101         derror_wrt_output: 误差关于神经元输出的偏导:  $dE / dy_j = 1/2(d(\text{expect\_to\_output} - \text{output})^2 / d\text{output}) = -(\text{expect\_to\_output} - \text{output})$ 
102         input_sum: 输入加权和
103     返回:
104         derror_wrt_dinputs: 误差关于输入的偏导
105     """
106     output = np.tanh(input_sum)
107     doutput_wrt_dinput = 1 - np.power(output, 2)
108     derror_wrt_dinput = derror_wrt_output * doutput_wrt_dinput
109
110     return derror_wrt_dinput
111
112
113 def activated(activation_choose, input):
114     """把正向激活包装一下"""
115     if activation_choose == "sigmoid":
116         return sigmoid(input)

```

```

117     elif activation_choose == "relu":
118         return relu(input)
119     elif activation_choose == "tanh":
120         return tanh(input)
121
122     return sigmoid(input)
123
124 def activated_back_propagation(activation_choose, derror_wrt_output, output):
125     """包装反向激活传播"""
126     if activation_choose == "sigmoid":
127         return sigmoid_back_propagation(derror_wrt_output, output)
128     elif activation_choose == "relu":
129         return relu_back_propagation(derror_wrt_output, output)
130     elif activation_choose == "tanh":
131         return tanh_back_propagation(derror_wrt_output, output)
132
133     return sigmoid_back_propagation(derror_wrt_output, output)
134
135 class NeuralNetwork:
136     def __init__(self, layers_strcuture, print_cost = False):
137         self.layers_strcuture = layers_strcuture
138         self.layers_num = len(layers_strcuture)
139
140         # 除掉输入层的网络层数，因为其他层才是真正的神经元层
141         self.param_layers_num = self.layers_num - 1
142
143         self.learning_rate = 0.0618
144         self.num_iterations = 2000
145         self.x = None
146         self.y = None
147         self.w = dict()
148         self.b = dict()
149         self.costs = []
150         self.print_cost = print_cost
151
152         self.init_w_and_b()
153
154     def set_learning_rate(self, learning_rate):
155         """设置学习率"""
156         self.learning_rate = learning_rate
157
158     def set_num_iterations(self, num_iterations):
159         """设置迭代次数"""
160         self.num_iterations = num_iterations
161
162     def set_xy(self, input, expected_output):
163         """设置神经网络的输入和期望的输出"""
164         self.x = input
165         self.y = expected_output
166
167     def init_w_and_b(self):
168         """
169         函数：
170             初始化神经网络所有参数
171         输入：
172             layers_strcuture：神经网络的结构，例如[2,4,3,1]，4层结构：

```

```

173         第0层输入层接收2个数据，第1层隐藏层4个神经元，第2层隐藏层3个神经元，第3层输出层1个神经元
174 返回：神经网络各层参数的索引表，用来定位权值  $w_i$  和偏置  $b_i$ ， $i$ 为网络层编号
175     """
176     np.random.seed(3)
177
178     # 当前神经元层的权值为  $n_i \times n_{(i-1)}$ 的矩阵， $i$ 为网络层编号， $n$ 为下标 $i$ 代表的网络层的节点个数
179     # 例如[2,4,3,1]，4层结构：第0层输入层为2，那么第1层隐藏层神经元个数为4
180     # 那么第1层的权值 $w$ 是一个  $4 \times 2$  的矩阵，如：
181     #     w1 = array([ [-0.96927756, -0.59273074],
182     #                  [ 0.58227367,  0.45993021],
183     #                  [-0.02270222,  0.13577601],
184     #                  [-0.07912066, -1.49802751] ])
185     # 当前层的偏置一般给0就行，偏置是个 $1 \times n_i$ 的矩阵， $n_i$ 为第 $i$ 层的节点个数，例如第1层为4个节点，那么：
186     #     b1 = array([ 0.,  0.,  0.,  0.])
187
188     for l in range(1, self.layers_num):
189         self.w["w" + str(l)] = np.random.randn(self.layers_structure[l],
self.layers_structure[l-1])/np.sqrt(self.layers_structure[l-1])
190         self.b["b" + str(l)] = np.zeros((self.layers_structure[l], 1))
191
192     return self.w, self.b
193
194     def layer_activation_forward(self, x, w, b, activation_choose):
195         """
196         函数：
197             网络层的正向传播
198         输入：
199             x：当前网络层输入（即上一层的输出），一般是所有训练数据，即输入矩阵
200             w：当前网络层的权值矩阵
201             b：当前网络层的偏置矩阵
202             activation_choose：选择激活函数 "sigmoid", "relu", "tanh"
203         返回：
204             output：网络层的激活输出
205             cache：缓存该网络层的信息，供后续使用：(x, w, b, input_sum) -> cache
206         """
207
208         # 对输入求加权和，见式 (5.1)
209         input_sum = np.dot(w, x) + b
210
211         # 对输入加权和进行激活输出
212         output, _ = activated(activation_choose, input_sum)
213
214         return output, (x, w, b, input_sum)
215
216     def forward_propagation(self, x):
217         """
218         函数：
219             神经网络的正向传播
220         输入：
221
222         返回：
223             output：正向传播完成后的输出层的输出
224             caches：正向传播过程中缓存每一个网络层的信息：(x, w, b, input_sum),... -> caches
225         """
226         caches = []
227

```

```

228         #作为输入层, 输出 = 输入
229         output_prev = x
230
231         #第0层为输入层, 只负责观察到输入的数据, 并不需要处理, 正向传播从第1层开始, 一直到输出层输出为止
232         # range(1, n) => [1, 2, ..., n-1]
233         L = self.param_layers_num
234         for l in range(1, L):
235             # 当前网络层的输入来自前一层的输出
236             input_cur = output_prev
237             output_prev, cache = self.layer_activation_forward(input_cur, self.w["w"+
str(l)], self.b["b" + str(l)], "tanh")
238             caches.append(cache)
239
240             output, cache = self.layer_activation_forward(output_prev, self.w["w" + str(L)],
self.b["b" + str(L)], "sigmoid")
241             caches.append(cache)
242
243         return output, caches
244
245     def show_caches(self, caches):
246         """显示网络层的缓存参数信息"""
247         i = 1
248         for cache in caches:
249             print("%dth Layer" % i)
250             print(" input: %s" % cache[0])
251             print(" w: %s" % cache[1])
252             print(" b: %s" % cache[2])
253             print(" input_sum: %s" % cache[3])
254             print("-----")
255             i += 1
256
257     def compute_error(self, output):
258         """
259         函数:
260             计算档次迭代的输出总误差
261         输入:
262
263         返回:
264
265         """
266
267         m = self.y.shape[1]
268
269         # 计算误差, 见式(5.5):  $E = \sum 1/2 (\text{期望输出} - \text{实际输出})^2$ 
270         # error = np.sum(0.5 * (self.y - output) ** 2) / m
271
272         # 交叉熵作为误差函数
273         error = -np.sum(np.multiply(np.log(output), self.y) + np.multiply(np.log(1 -
output), 1 - self.y)) / m
274         error = np.squeeze(error)
275
276         return error
277
278     def layer_activation_backward(self, derror_wrt_output, cache, activation_choose):
279         """
280         函数:

```

```

281         网络层的反向传播
282         输入:
283             derror_wrt_output: 误差关于输出的偏导
284             cache: 网络层的缓存信息 (x, w, b, input_sum)
285             activation_choose: 选择激活函数 "sigmoid", "relu", "tanh"
286         返回: 梯度信息, 即
287             derror_wrt_output_prev: 反向传播到上一层的误差关于输出的梯度
288             derror_wrt_dw: 误差关于权值的梯度
289             derror_wrt_db: 误差关于偏置的梯度
290         """
291         input, w, b, input_sum = cache
292         output_prev = input      # 上一层的输出 = 当前层的输入; 注意是'输入'不是输入的加权和 (
input_sum)
293         m = output_prev.shape[1]      # m是输入的样本数量, 我们要取均值, 所以下面的求值要除以m
294
295         # 实现式 (5.13) -> 误差关于权值w的偏导数
296         derror_wrt_dinput = activated_back_propagation(activation_choose, derror_wrt_output,
input_sum)
297         derror_wrt_dw = np.dot(derror_wrt_dinput, output_prev.T) / m
298
299         # 实现式 (5.32) -> 误差关于偏置b的偏导数
300         derror_wrt_db = np.sum(derror_wrt_dinput, axis=1, keepdims=True)/m
301
302         # 为反向传播到上一层提供误差传递, 见式 (5.28) 的  $(\sum \delta \cdot w)$  部分
303         derror_wrt_output_prev = np.dot(w.T, derror_wrt_dinput)
304
305         return derror_wrt_output_prev, derror_wrt_dw, derror_wrt_db
306
307     def back_propagation(self, output, caches):
308         """
309         函数:
310             神经网络的反向传播
311         输入:
312             output: 神经网络输
313             caches: 所有网络层 (输入层不算) 的缓存参数信息 [(x, w, b, input_sum), ...]
314         返回:
315             grads: 返回当前迭代的梯度信息
316         """
317         grads = {}
318         L = self.param_layers_num #
319         output = output.reshape(output.shape) # 把输出层输出输出重构成和期望输出一样的结构
320
321         expected_output = self.y
322
323         # 见式(5.8)
324         #derror_wrt_output = -(expected_output - output)
325
326         # 交叉熵作为误差函数
327         derror_wrt_output = - (np.divide(expected_output, output) - np.divide(1 -
expected_output, 1 - output))
328
329         # 反向传播: 输出层 -> 隐藏层, 得到梯度: 见式(5.8), (5.13), (5.15)
330         current_cache = caches[L - 1] # 取最后一层, 即输出层的参数信息
331         grads["derror_wrt_output" + str(L)], grads["derror_wrt_dw" + str(L)],
grads["derror_wrt_db" + str(L)] = \
332             self.layer_activation_backward(derror_wrt_output, current_cache, "sigmoid")

```

```

333
334     # 反向传播：隐藏层 -> 隐藏层，得到梯度：见式 (5.28)的 $(\Sigma \delta \cdot w)$ ，(5.28)，(5.32)
335     for l in reversed(range(L - 1)):
336         current_cache = caches[l]
337         derror_wrt_output_prev_temp, derror_wrt_dw_temp, derror_wrt_db_temp = \
338             self.layer_activation_backward(grads["derror_wrt_output" + str(l + 2)],
current_cache, "tanh")
339
340         grads["derror_wrt_output" + str(l + 1)] = derror_wrt_output_prev_temp
341         grads["derror_wrt_dw" + str(l + 1)] = derror_wrt_dw_temp
342         grads["derror_wrt_db" + str(l + 1)] = derror_wrt_db_temp
343
344     return grads
345
346     def update_w_and_b(self, grads):
347         """
348         函数：
349             根据梯度信息更新w, b
350         输入：
351             grads：当前迭代的梯度信息
352         返回：
353
354         """
355         # 权值w和偏置b的更新，见式：(5.16),(5.18)
356         for l in range(self.param_layers_num):
357             self.w["w" + str(l + 1)] = self.w["w" + str(l + 1)] - self.learning_rate *
grads["derror_wrt_dw" + str(l + 1)]
358             self.b["b" + str(l + 1)] = self.b["b" + str(l + 1)] - self.learning_rate *
grads["derror_wrt_db" + str(l + 1)]
359
360     def tarin_modle(self):
361         """训练神经网络模型"""
362
363         np.random.seed(5)
364         for i in range(0, self.num_iterations):
365             # 正向传播，得到网络输出，以及每一层的参数信息
366             output, caches = self.forward_propagation(self.x)
367
368             # 计算网络输出误差
369             cost = self.compute_error(output)
370
371             # 反向传播，得到梯度信息
372             grads = self.back_propagation(output, caches)
373
374             # 根据梯度信息，更新权值w和偏置b
375             self.update_w_and_b(grads)
376
377             # 当次迭代结束，打印误差信息
378             if self.print_cost and i % 1000 == 0:
379                 print ("Cost after iteration %i: %f" % (i, cost))
380             if self.print_cost and i % 1000 == 0:
381                 self.costs.append(cost)
382
383         # 模型训练完后显示误差曲线
384         if False:
385             plt.plot(np.squeeze(self.costs))

```



```

386         plt.ylabel(u'神经网络误差', fontproperties = font)
387         plt.xlabel(u'迭代次数 (*100)', fontproperties = font)
388         plt.title(u"学习率 =" + str(self.learning_rate), fontproperties = font)
389         plt.show()
390
391         return self.w, self.b
392
393     def predict_by_modle(self, x):
394         """使用训练好的模型（即最后求得w, b参数）来决策输入的样本的结果"""
395         output, _ = self.forward_propagation(x.T)
396         output = output.T
397         result = output / np.sum(output, axis=1, keepdims=True)
398         return np.argmax(result, axis=1)
399
400 def plot_decision_boundary(xy, colors, pred_func):
401     # xy是坐标点的集合，把集合的范围算出来
402     # 加减0.5相当于扩大画布的范围，不然画出来的图坐标点会落在图的边缘，逼死强迫症患者
403     x_min, x_max = xy[:, 0].min() - 0.5, xy[:, 0].max() + 0.5
404     y_min, y_max = xy[:, 1].min() - 0.5, xy[:, 1].max() + 0.5
405
406     # 以h为分辨率，生成采样点的网格，就像一张网覆盖所有颜色点
407     h = .01
408     xx, yy = np.meshgrid(np.arange(x_min, x_max, h), np.arange(y_min, y_max, h))
409
410     # 把网格点集合作为输入到模型，也就是预测这个采样点是什么颜色的点，从而得到一个决策面
411     Z = pred_func(np.c_[xx.ravel(), yy.ravel()])
412     Z = Z.reshape(xx.shape)
413
414     # 利用等高线，把预测的结果画出来，效果上就是画出红蓝点的分界线
415     plt.contourf(xx, yy, Z, cmap=plt.cm.Spectral)
416
417     # 训练用的红蓝点点也画出来
418     plt.scatter(xy[:, 0], xy[:, 1], c=colors, marker='o', cmap=plt.cm.Spectral,
419               edgecolors='black')
420
421 if __name__ == "__main__":
422     plt.figure(figsize=(16, 32))
423
424     # 用sklearn的数据样本集，产生2种颜色的坐标点，noise是噪声系数，噪声越大，2种颜色的点分布越凌乱
425     xy, colors = sklearn.datasets.make_moons(60, noise=1.0)
426
427     # 因为点的颜色是1bit，我们设计一个神经网络，输出层有2个神经元。
428     # 标定输出[1,0]为红色点，输出[0,1]为蓝色点
429     expect_output = []
430     for c in colors:
431         if c == 1:
432             expect_output.append([0,1])
433         else:
434             expect_output.append([1,0])
435
436     expect_output = np.array(expect_output).T
437
438     # 设计3层网络，改变隐藏层神经元的个数，观察神经网络分类红蓝点的效果
439     hidden_layer_neuron_num_list = [1,2,4,10,20,50]
440     for i, hidden_layer_neuron_num in enumerate(hidden_layer_neuron_num_list):

```

```

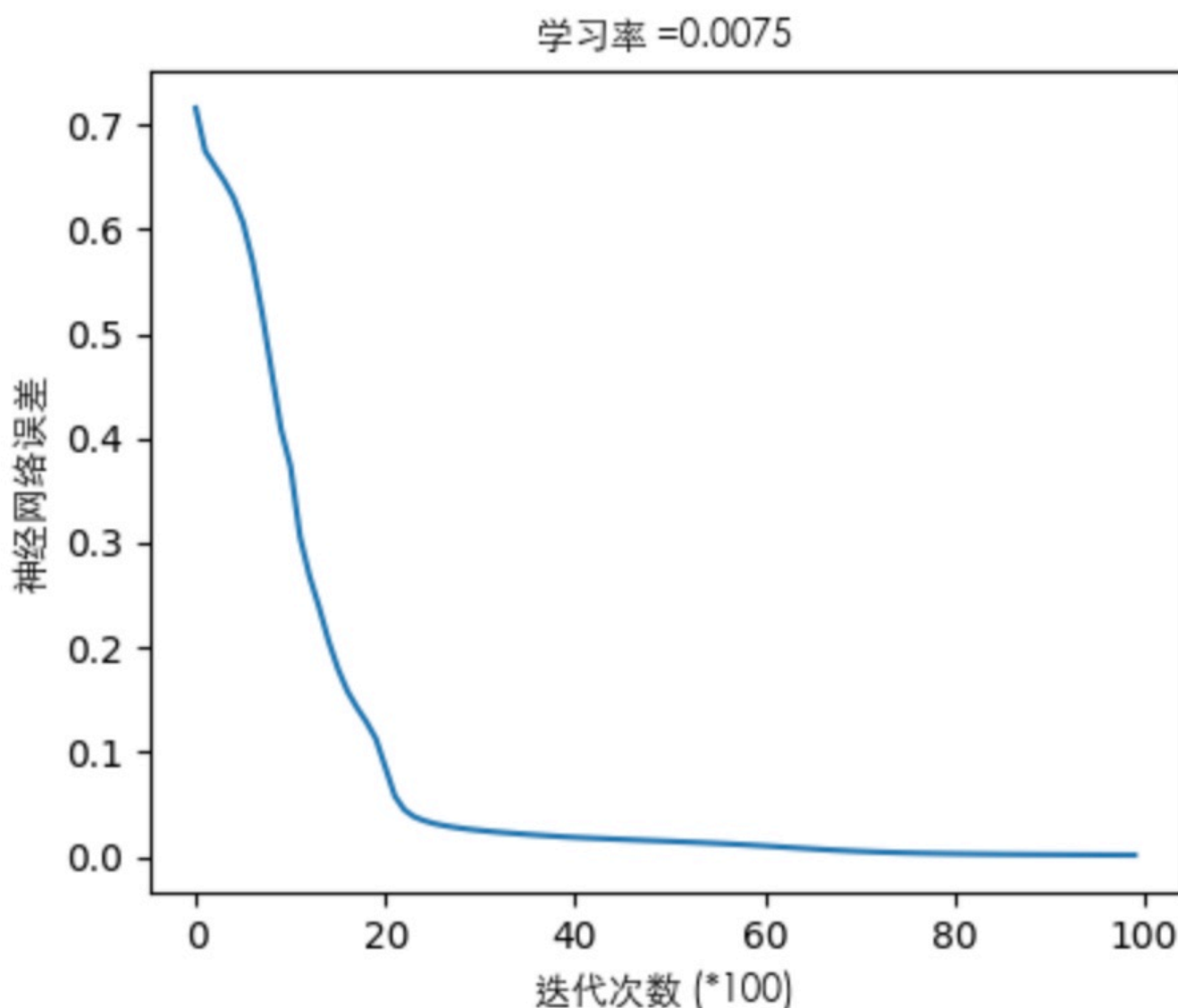
441 plt.subplot(5, 2, i + 1)
442 plt.title(u'隐藏层神经元数量: %d' % hidden_layer_neuron_num, fontproperties = font)
443
444 nn = NeuralNetwork([2, hidden_layer_neuron_num, 2], True)
445
446 # 输出和输入层都是2个节点，所以输入和输出的数据集都要是 nx2的矩阵
447 nn.set_xy(xy.T, expect_output)
448 nn.set_num_iterations(30000)
449 nn.set_learning_rate(0.1)
450 w, b = nn.tarin_modle()
451 plot_decision_boundary(xy, colors, nn.predict_by_modle)
452
453 plt.show()
454

```

2. 晒图晒图！

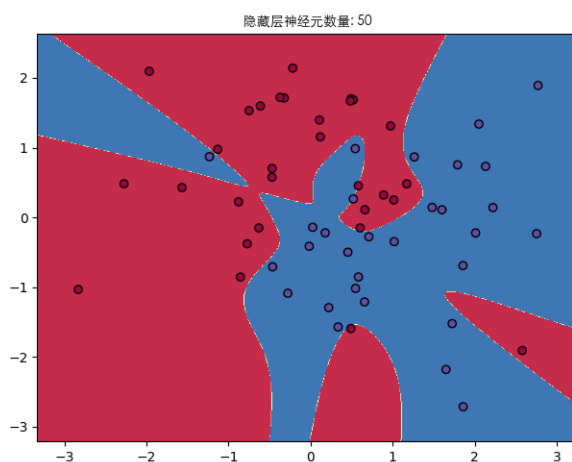
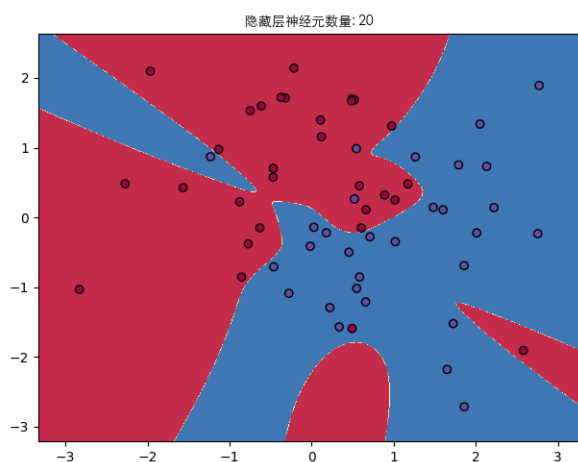
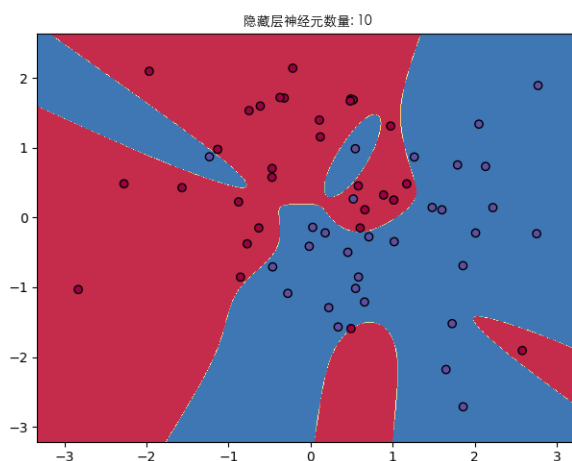
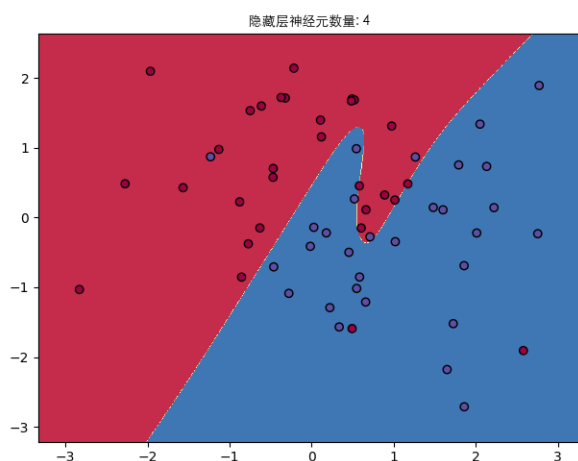
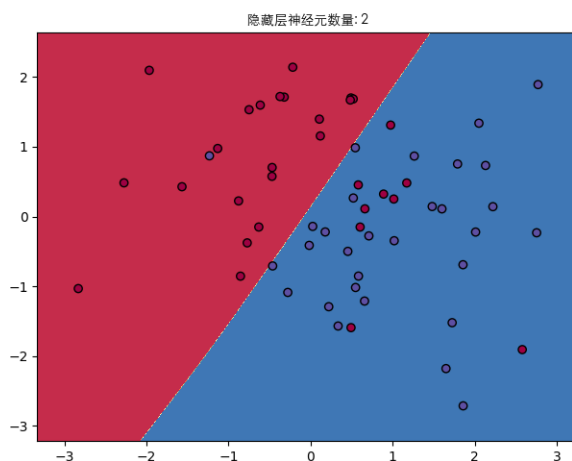
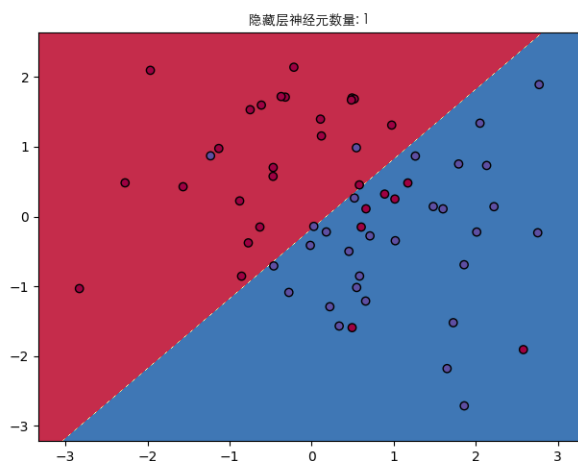
关于误差曲线(这里只举其中一个栗子):

- 通过看误差曲线，可以从一定程度上判定网络的效果，模型训练是否能收敛，收敛程度如何，都可以从误差曲线对梯度下降的过程能见一二。



3层网络的结构下，隐藏层只有一层，看图说明一下隐藏层神经元个数变化对神经网络表达能力的影响：

- 当隐藏层只有1个神经元：就像文章刚开始说的，一个神经元，就是个线性分类器，表达能力就一条直线而已，见式 (3.6)
- 2个神经元：线开始有点弯曲了，但是这次结果一点都不明显，尴尬。但从原理上神经网络开始具备了非线性表达能力
- 随着隐藏层神经元个数不断增加，神经网络表达能力越来越强，分类的效果越来越好。当然也不是神经元越多越好，可以开始考虑深度网络是不是效果更好一些。



7. 没有结局

记住一点，bp神经网络是其他各种神经网络中最简单的一种。只有学会了它，才能以此为基础展开对其他更复杂

的神经网络的学习。

虽然推导了并实现了算法，但是仍然是有很多疑问，这里就作为抛砖引玉吧：

- 神经网络的结构，即几层网络，输入输出怎么设计才最有效？
- 数学理论证明，三层的神经网络就能够以任意精度逼近任何非线性连续函数。那么为什么还需要有深度网络？
- 在不同应用场合下，激活函数怎么选择？
- 学习率怎么怎么选择？
- 训练次数设定多少训练出的模型效果更好？

AI，从入门到放弃，首篇结束。