

(5条消息)深度学习（一）：DNN前向传播算法和反向传播算法_anshuai_aw1的博客-CSDN博客_前向传播和反向传播

 blog.csdn.net/anshuai_aw1/article/details/84615935

一、深度神经网络（DNN）模型

深度神经网络（Deep Neural Networks，以下简称DNN）是深度学习的基础，而理解DNN，首先我们要理解DNN模型，下面我们就对DNN的模型与前向传播算法做一个总结。

1.1 从感知机到神经网络

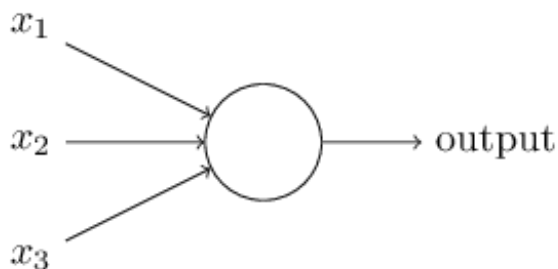
感知机的模型大家都比较熟悉，它是一个有若干输入和一个输出的模型，如下图：

输出和输入之间学习到一个线性关系，得到中间输出结果：

$$z = \sum_{i=1}^m w_i x_i + b$$

接着是一个神经元激活函数：

$$\text{sign}(z) = \begin{cases} -1 & z < 0 \\ 0 & z = 0 \\ 1 & z > 0 \end{cases}$$

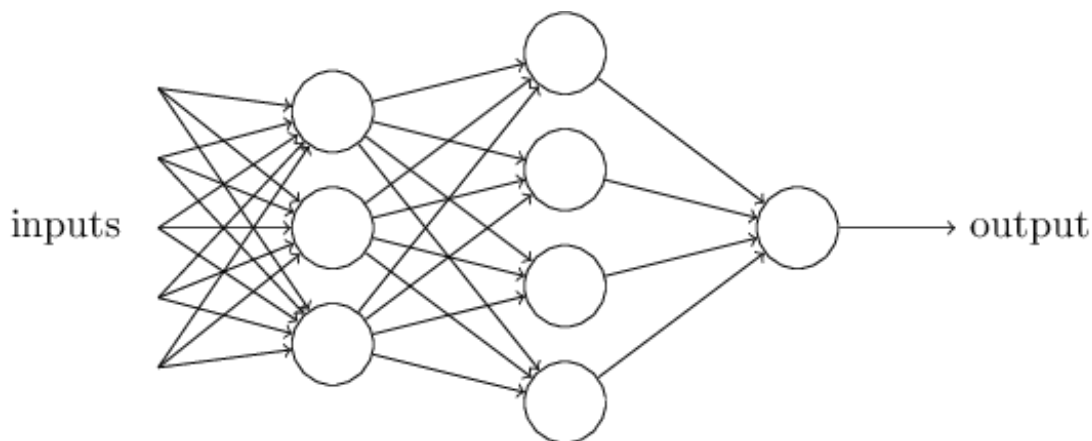


从而得到我们想要的输出结果1或者-1。

这个模型只能用于二元分类，且无法学习比较复杂的非线性模型，因此在工业界无法使用。

而神经网络则在感知机的模型上做了扩展，总结下主要有三点：

1) 加入了隐藏层，隐藏层可以有多层，增强模型的表达能力，如下图实例，当然增加了这么多隐藏层模型的复杂度也增加了好多。

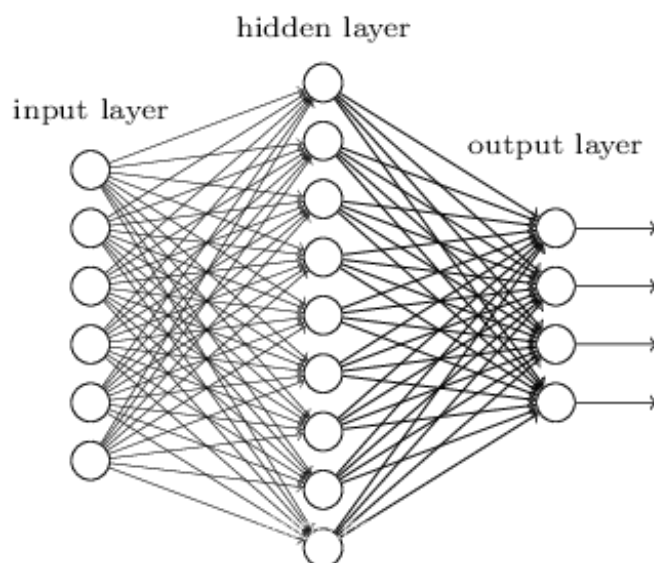


2) 输出层的神经元也可以不止一个输出，可以有多个输出，这样模型可以灵活的应用于分类回归，以及其他的机器学习领域比如降维和聚类等。多个神经元输出的输出层对应的一个实例如下图，输出层现在有4个神经元了。

3) 对激活函数做扩展，感知机的激活函数是 $\text{sign}(z)$ ，虽然简单但是处理能力有限，因此神经网络中一般使用的其他的激活函数，比如我们在逻辑回归里面使用过的Sigmoid函数，即：

$$f(z) = \frac{1}{1 + e^{-z}}$$

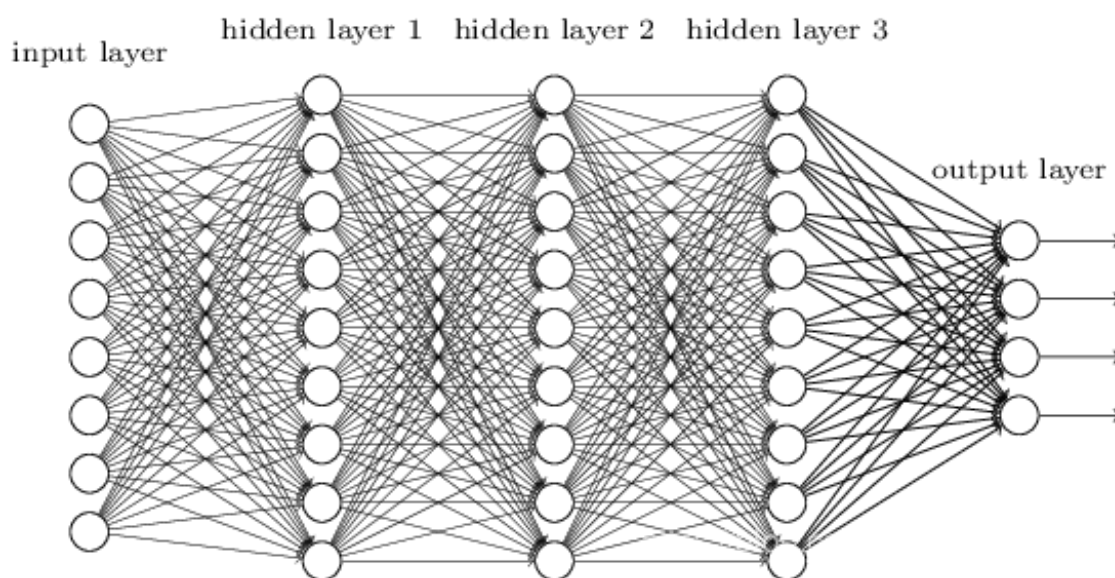
还有后来出现的tanx, softmax,和ReLU等。通过使用不同的激活函数，神经网络的表达能力进一步增强。对于各种常用的激活函数，我随后会整理一份资料。



1.2 DNN的基本结构

上一节我们了解了神经网络基于感知机的扩展，而DNN可以理解为有很多隐藏层的神经网络。这个很多其实也没有什么度量标准，多层神经网络和深度神经网络DNN其实也是指的一个东西，当然，DNN有时也叫做多层感知机（Multi-Layer perceptron,MLP），名字实在是多。后面我们讲到的神经网络都默认为DNN。

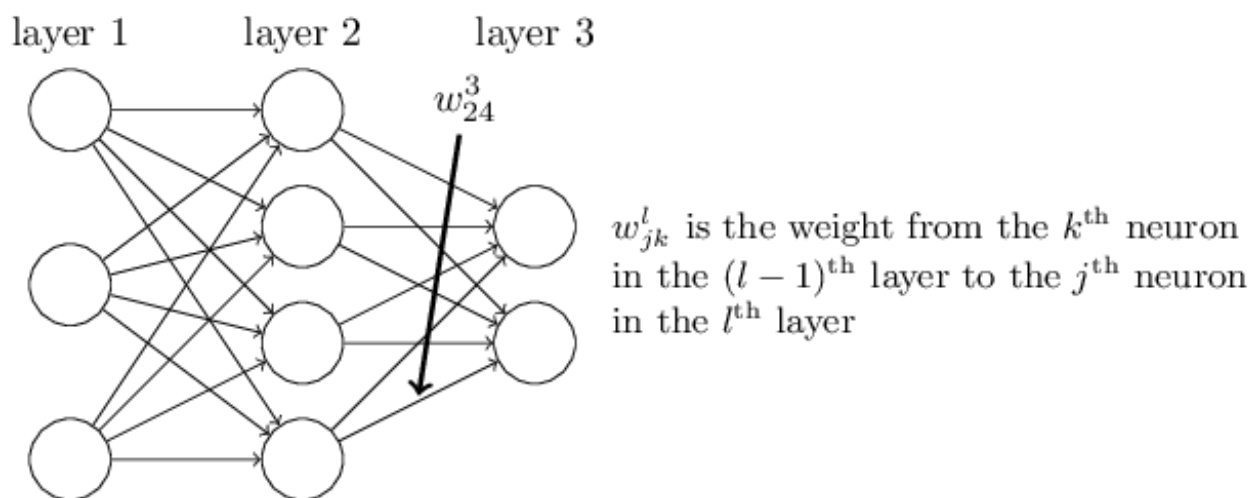
从DNN按不同层的位置划分，DNN内部的神经网络层可以分为三类，输入层，隐藏层和输出层，如下图示例，一般来说第一层是输入层，最后一层是输出层，而中间的层数都是隐藏层。



层与层之间是全连接的，也就是说，第 i 层的任意一个神经元一定与第 $i+1$ 层的任意一个神经元相连。虽然DNN看起来很复杂，但是从小的局部模型来说，还是和感知机一样，即一个线性关系 $z = \sum w_i x_i + b$ 加上一个激活函数 $\sigma(z)$ 。

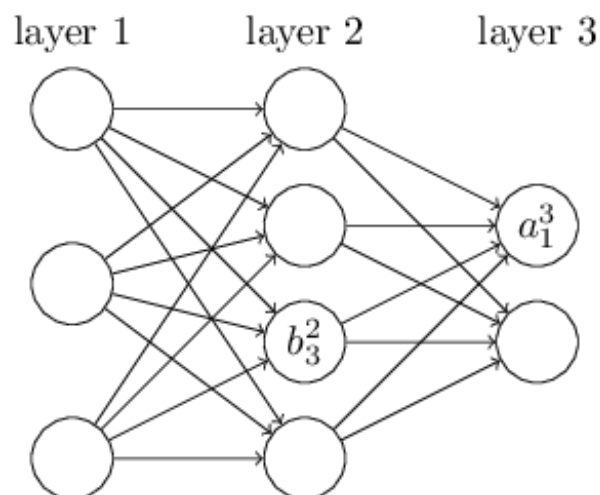
由于DNN层数多，则我们的线性关系系数 w 和偏倚 b 的数量也就是很多了。具体的参数在DNN是如何定义的呢？

首先我们来看看线性关系系数 w 的定义。以下图一个三层的DNN为例，第二层的第4个神经元到第三层的第2个神经元的线性系数定义为 w_{24}^3 。上标3代表线性系数 w 所在的层数，而下标对应的是输出的第三层索引2和输入的第二层索引4。你也许会问，为什么不是 w_{423} ，而是 w_{24}^3 呢？这主要是为了便于模型用于矩阵表示运算，如果是 w_{423} 每次进行矩阵运算是 $wTx+b$ ，需要进行转置。将输出的索引放在前面的话，则线性运算不用转置，即直接为 $wx+b$ 。总结下，第 $l-1$ 层的第 k 个神经元到第 l 层的第 j 个神经元的线性系数定义为 w_{jk}^l 。注意，输入层是没有 w 参数的。



再来看看偏倚 b 的定义。还是以这个三层的DNN为例，第二层的第3个神经元对应的偏倚定义为 b_{32} 。其中，上标2代表所在的层数，下标3代表偏倚所在的神经元的索引。输入层是没有偏倚参数 b 的。

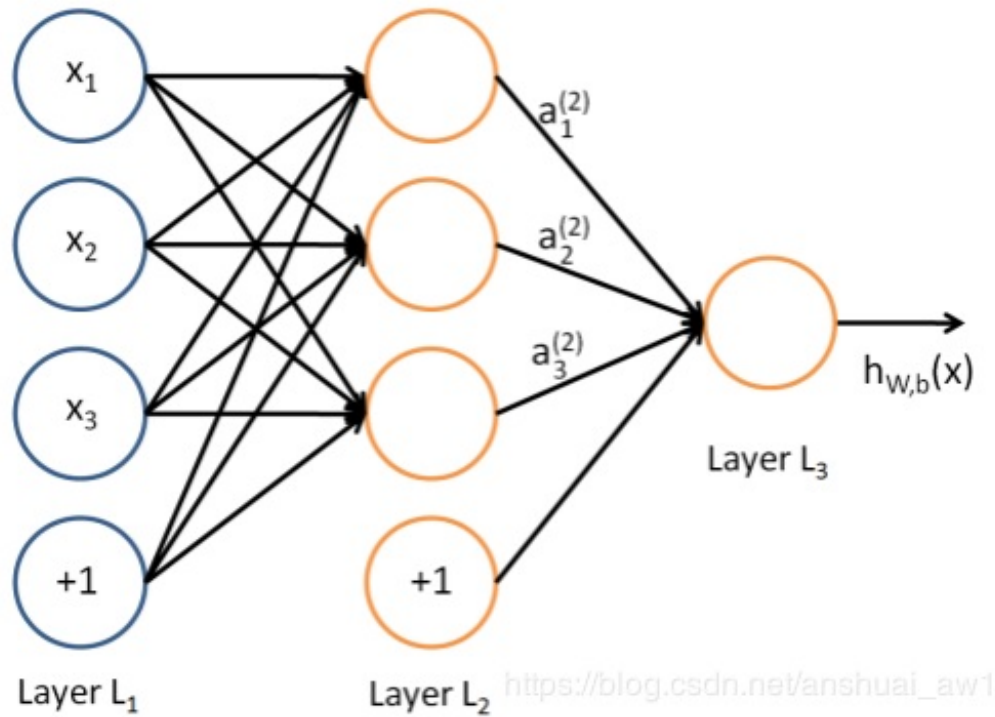
同样的道理，对于神经元的激活值而言，第3层的第1个神经元的激活值应该表示为 a_1^3 。



二、DNN前向传播算法

2.1 DNN前向传播算法数学原理

在上一节，我们已经介绍了DNN各层线性关系系数 w 和偏倚 b 的定义。假设我们选择的激活函数是 $\sigma(z)$ ，隐藏层和输出层的输出值为 a ，则对于下图的三层DNN，利用和感知机一样的思路，我们可以利用上一层的输出计算下一层的输出，也就是所谓的DNN前向传播算法。



对于第二层的输出 a_{12} , a_{22} , a_{32} ，我们有：

$$a_{12} = \sigma(z_{12}) = \sigma(w_{112} x_1 + w_{122} x_2 + w_{132} x_3 + b_{12})$$

$$a_{22} = \sigma(z_{22}) = \sigma(w_{212} x_1 + w_{222} x_2 + w_{232} x_3 + b_{22})$$

$$a_{32} = \sigma(z_{32}) = \sigma(w_{312} x_1 + w_{322} x_2 + w_{332} x_3 + b_{32})$$

对于第三层的输出 a_{13} ，我们有：

$$a_{13} = \sigma(z_{13}) = \sigma(w_{113} a_{12} + w_{123} a_{22} + w_{133} a_{32} + b_{13})$$

将上面的例子一般化，假设第 $l-1$ 层共有 m 个神经元，则对于第 l 层的第 j 个神经元的输出 a_{jl} ，我们有：

$$a_{jl} = \sigma(z_{jl}) = \sigma\left(\sum_{k=1}^m w_{jkl} a_{kl-1} + b_{jl}\right)$$

其中，如果 $l=2$ ，则对于的 a_{k1} 即为输入层的 x_k 。

从上面可以看出，使用代数法一个个的表示输出比较复杂，而如果使用矩阵法则比较的简洁。假设第 $l-1$ 层共有 m 个神经元，而第 l 层共有 n 个神经元，则第 l 层的线性系数 w 组成了一个 $n \times m$ 的矩阵 W_l ，第 l 层的偏倚 b 组成了一个 $n \times 1$ 的向量 b_l ，第 $l-1$ 层的输出 a 组成了一个 $m \times 1$ 的

向量 a_{l-1} ，第 l 层的未激活前线性输出 z 组成了一个 $n \times 1$ 的向量 z_l ，第 l 层的输出 a 组成了一个 $n \times 1$ 的向量 a_l 。则用矩阵法表示，第 l 层的输出为：

$$a_l = \sigma(z_l) = \sigma(W_l a_{l-1} + b_l)$$

这个表示方法简洁漂亮，后面我们的讨论都会基于上面的这个矩阵法表示来。所以，应该时刻记住我们符号的含义，否则在后面推导反向传播公式时会比较懵。

2.2 DNN前向传播算法

有了上一节的数学推导，DNN的前向传播算法也就不难了。所谓的DNN的前向传播算法也就是利用我们的若干个权重系数矩阵 W 和偏倚向量 b 来和输入值向量 x 进行一系列线性运算和激活运算，从输入层开始，一层层的向后计算，一直到运算到输出层，得到输出结果为止。

输入：总层数 L ，所有隐藏层和输出层对应的矩阵 W ，偏倚向量 b ，输入值向量 x

输出：输出层的输出 a_L

- 1) 初始化 $a_1 = x$
- 2) for $l=2$ to L , 计算：

$$a_l = \sigma(z_l) = \sigma(W_l a_{l-1} + b_l)$$

最后的结果即为输出 a_L 。

2.3 DNN前向传播算法小结

单独看DNN前向传播算法，似乎没有什么大用处，而且这一大堆的矩阵 W ，偏倚向量 b 对应的参数怎么获得呢？怎么得到最优的矩阵 W ，偏倚向量 b 呢？这个我们在下一章讲DNN的反向传播算法时再讲。而理解反向传播算法的前提就是理解DNN的模型与前向传播算法。这也是我们先讲前向传播算法的原因。

三、DNN反向传播算法

3.1 DNN反向传播算法要解决的问题

在了解DNN的反向传播算法（Back Propagation，BP）前，我们先要知道DNN反向传播算法要解决的问题，也就是说，什么时候我们需要这个反向传播算法？

回到我们监督学习的一般问题，假设我们有 m 个训练样

本： $\{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$ ，其中 x 为输入向量，特征维度为 n_{in} ，而 y 为输出向量，特征维度为 n_{out} 。我们需要利用这 m 个样本训练出一个模型，当有一个新的测试样本 $(x_{test}, ?)$ 来到时，我们可以预测 y_{test} 向量的输出。

如果我们采用DNN的模型，即我们使输入层有 n_{in} 个神经元，而输出层有 n_{out} 个神经元。再加上一些含有若干神经元的隐藏层。此时我们需要找到合适的所有隐藏层和输出层对应的线性系数矩阵 W 和偏倚向量 b ，让所有的训练样本输入计算出的输出尽可能的等于或很接近样本输出。怎么找到合适的参数呢？

如果大家对传统的机器学习的算法优化过程熟悉的话，这里就很容易联想到我们可以用一个合适的损失函数来度量训练样本的输出损失，接着对这个损失函数进行优化求最小化的极值，对应的一系列线性系数矩阵 W 和偏倚向量 b 即为我们的最终结果。在DNN中，损失函数优化极值求解的过程最常见的一般是通过梯度下降法来一步步迭代完成的，当然也可以是其他的迭代方法比如牛顿法与拟牛顿法。

对DNN的损失函数用梯度下降法进行迭代优化求极小值的过程即为我们的反向传播算法。

3.2 DNN反向传播算法的基本思路

在进行DNN反向传播算法前，我们需要选择一个损失函数，来度量训练样本计算出的输出和真实的训练样本输出之间的损失。你也许会问：训练样本计算出的输出是怎么得来的？这个输出是随机选择一系列 W, b ，用前向传播算法计算出来的。即通过一系列的计算： $a_l = \sigma(z_l) = \sigma(Wl_{l-1} + b_l)$ 。计算到输出层第 L 层对应的 a_L 即为前向传播算法计算出来的输出。

回到损失函数，DNN可选择的损失函数有不少，为了专注算法，这里我们使用最常见的均方差来度量损失。当然，针对不同的任务，可以选择不同的损失函数。即对于每个样本，我们期望最小化下式：

$$J(W, b, x, y) = \frac{1}{2} \|a_L - y\|_2^2$$

其中， a_L 和 y 为特征维度为 n_{out} 的向量，而 $\|S\|_2$ 为 S 的 L_2 范数。

损失函数有了，现在我们开始用梯度下降法迭代求解每一层的 W, b 。

注：以下是BP算法推导的过程，是本文最核心，也是神经网络最基本的公式推导。

思路第一步：

首先是输出层第 L 层。注意到输出层的 W, b 满足下式：

$$a_L = \sigma(z_L) = \sigma(WL_{L-1} + b_L) \quad (1)$$

这样对于输出层的参数，我们的损失函数变为：

$$J(W, b, x, y) = \frac{1}{2} \|a_L - y\|_2^2 = \frac{1}{2} \|\sigma(WL_{L-1} + b_L) - y\|_2^2 \quad (2)$$

这样求解 W, b 的梯度就简单了：

$$\frac{\partial W_L}{\partial J(W, b, x, y)} = \frac{\partial z_L}{\partial J(W, b, x, y)} \frac{\partial W_L}{\partial z_L} = \frac{\partial a_L}{\partial J(W, b, x, y)} \frac{\partial z_L}{\partial a_L} \frac{\partial W_L}{\partial z_L} = (a_L - y) \odot \sigma'(z_L) (a_L - 1) T(3)$$

$$\partial b^L \partial J(W, b, x, y) = \partial z^L \partial J(W, b, x, y) \quad \partial b^L \partial z^L = \partial a^L \partial J(W, b, x, y) \quad \partial z^L \partial a^L \quad \partial b^L \partial z^L = (a^{L-y}) \odot \sigma'(z^L) \quad (4)$$

注意上式中有一个符号 \odot ,它代表Hadamard积,对于两个维度相同的向量A $(a_1, a_2, \dots, a_n)^T$ 和B $(b_1, b_2, \dots, b_n)^T$,则 $A \odot B = (a_1 b_1, a_2 b_2, \dots, a_n b_n)^T$ 。

对于公式 (3) 和 (4), 我在这里多解释一下为什么是这样:

对于公式 (3): 前两项之所以是Hadamard积的形式, 是因为 $\partial a^L \partial J(W, b, x, y) \quad \partial z^L \partial a^L$ 都是针对同一层的神经元。如果我们考虑对于L层的第j个神经元, 即 $\partial a_j^L \partial J(W, b, x, y) \quad \sigma'(z_j^L)$, 那么整合这一层的神经元, 自然是 $(a^{L-y}) \odot \sigma'(z^L)$ 这样Hadamard积的形式。那么 $(a^{L-1})^T$ 为什么在 (3) 中的最后呢? 这涉及到矩阵求导的知识, 不是我们本文的重点。在这里, 用到的知识是: 如果

$$Y = WX + B$$

那么

$$\partial W \partial C = \partial Y \partial C^T X^T$$

这样, 即可推出公式 (3)。公式 (4) 与公式 (3) 类似。

思路第二步:

我们注意到在求解输出层的W,b的时候, 有公共的部分 $\partial z^L \partial J(W, b, x, y)$, 因此我们可以把公共的部分即对 z^L 先算出来, 记为:

$$\delta^L = \partial z^L \partial J(W, b, x, y) = (a^{L-y}) \odot \sigma'(z^L) \quad (5)$$

根据公式 (3) (4), 我们可以把输出层的梯度算出来, 那么如何计算上一层 $L-1$ 层的梯度, 上上层 $L-2$ 层的梯度呢? 这里我们需要一步步的递推, 注意到对于第1层的未激活输出 z^1 , 它的梯度可以表示为:

$$\delta^1 = \partial z^1 \partial J(W, b, x, y) = \partial z^L \partial J(W, b, x, y) \quad \partial z^{L-1} \partial z^L \quad \partial z^{L-2} \partial z^{L-1} \dots \partial z^1 \partial z^{l+1} \quad (6)$$

如果我们可以依次计算出第1层的 δ^1 , 则该层的 W^1, b^1 很容易计算? 为什么呢? 注意到根据前向传播算法, 我们有:

$$z^l = W^l a^{l-1} + b^l \quad (7)$$

所以根据上式我们可以很方便的计算出第1层的 W^1, b^1 的梯度如下:

$$\partial W^1 \partial J(W, b, x, y) = \partial z^1 \partial J(W, b, x, y) \quad \partial W^1 \partial z^1 = \delta^1 (a^{l-1})^T \quad (8)$$

$$\partial b^1 \partial J(W, b, x, y) = \partial z^1 \partial J(W, b, x, y) \quad \partial b^1 \partial z^1 = \delta^1 \quad (9)$$

思路第三步:

那么现在问题的关键就是要求出 δ^1 了。这里我们用数学归纳法, 第L层的 δ^L 上面我们已经求出, 假设第 $l+1$ 层的 δ^{l+1} 已经求出来了, 那么我们如何求出第1层的 δ^1 呢? 我们注意到:

$$\delta^1 = \partial z^1 \partial J(W, b, x, y) = \partial z^{l+1} \partial J(W, b, x, y) \quad \partial z^1 \partial z^{l+1} = \delta^{l+1} \partial z^1 \partial z^{l+1} \quad (10)$$

可见，用归纳法递推 δ_{l+1} 和 δ_l 的关键在于求解 $\partial z_l \partial z_{l+1}$ 。

而 z_{l+1} 和 z_l 的关系其实很容易找出：

$$z_{l+1} = W_{l+1}a_l + b_{l+1} = W_{l+1}\sigma(z_l) + b_{l+1} \quad (11)$$

这样很容易求出：

$$\partial z_l \partial z_{l+1} = (W_{l+1})^T \odot n_{l+1} \underbrace{\quad \quad \quad}_{(\sigma'(z_l), \dots, \sigma'(z_l))} \quad (12)$$

公式（12）的意思就是 $(W_{l+1})^T$ 的每一列都Hadamard积 $\sigma'(z_l)$ 。

将公式（12）代入公式（10）我们得到：

$$\delta_l = \delta_{l+1} \partial z_l \partial z_{l+1} = ((W_{l+1})^T \delta_{l+1}) \odot \sigma'(z_l) \quad (13)$$

公式（13）的推导过程为：

$$\delta_{l+1} \partial z_l \partial z_{l+1} = \partial z_l \partial (\delta_{l+1})^T z_{l+1} = \partial z_l \partial (\delta_{l+1})^T (W_{l+1} \sigma(z_l) + b_{l+1}) = \partial z_l \partial (\delta_{l+1})^T W_{l+1} \sigma(z_l) = ((\delta_{l+1})^T W_{l+1})^T \odot \sigma'(z_l) = ((W_{l+1})^T \delta_{l+1}) \odot \sigma'(z_l)$$

现在我们得到了 δ_l 的递推关系式，只要求出了某一层的 δ_l ，求解 W_l, b_l 的对应梯度就很简单。

总结：

其实，对于更新每一层的 W_l, b_l 的对应梯度，我们仔细观察整个过程，发现只需要四个公式就可以完整地进行更新。这就是著名的反向传播的四个公式，即公式(5)、(13)、(8)、(9)。我们稍加改动，使其可以应用到多种损失函数，即：

$$\delta_L = \partial a_L \partial J \odot \sigma'(z_L) \quad (BP1) \quad \delta_l = (W_{l+1})^T \delta_{l+1} \odot \sigma'(z_l) \quad (BP2) \quad \partial W_l \partial J = \delta_l (a_{l-1})^T \quad (BP3) \quad \partial b_l \partial J = \delta_l \quad (BP4)$$

3.3 DNN反向传播算法过程

现在我们总结下DNN反向传播算法的过程。由于梯度下降法有批量（Batch），小批量（mini-Batch），随机三个变种，为了简化描述，这里我们以最基本的批量梯度下降法为例来描述反向传播算法。实际上在业界使用最多的是mini-Batch的梯度下降法。不过区别仅仅在于迭代时训练样本的选择而已。

输入：总层数 L ，以及各隐藏层与输出层的神经元个数，激活函数，损失函数，迭代步长 α ，最大迭代次数 MAX 与停止迭代阈值 ϵ ，输入的 m 个训练样本 $\{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$

输出：各隐藏层与输出层的线性关系系数矩阵 W 和偏倚向量 b

1) 初始化各隐藏层与输出层的线性关系系数矩阵 W 和偏倚向量 b 的值为一个随机值。

2) for iter to 1 to MAX ：

2-1) for $i=1$ to m ：

a) 将DNN输入 a_1 设置为 x_i

b) for $l=2$ to L , 进行前向传播算法计算 $a_{i,l}=\sigma(z_{i,l})=\sigma(W_{l,i}a_{i,l-1}+b_l)$

c) 通过损失函数计算输出层的 $\delta_{i,L}$ (**BP1**)

d) for $l=L-1$ to 2 , 进行反向传播算法计算 $\delta_{i,l}=(W_{l+1})^T\delta_{i,l+1}\odot\sigma'(z_{i,l})$ (**BP2**)

2-2) for $l=2$ to L , 更新第 l 层的 W_l, b_l :

$W_l = W_l - \alpha i = 1 \sum_m \delta_{i,l}(a_{i,l-1})^T$ (BP3)

$b_l = b_l - \alpha i = 1 \sum_m \delta_{i,l}$ (BP4)

2-3) 如果所有 W, b 的变化值都小于停止迭代阈值 ϵ , 则跳出迭代循环到步骤3。

3) 输出各隐藏层与输出层的线性关系系数矩阵 W 和偏倚向量 b 。

注：上述的算法其实看起来有些复杂了，而实际上根据机器学习里普通的批梯度下降算法，我们在求出 m 个样本的 $\delta_{i,L}$ 后，求平均得到 δ_L ，然后反向继续更新 δ_l 。

3.4 DNN反向传播算法小结

有了DNN反向传播算法，我们就可以很方便的用DNN的模型去解决各种监督学习的分类回归问题。当然DNN的参数众多，矩阵运算量也很大，直接使用会有各种各样的问题。有哪些问题以及如何尝试解决这些问题并优化DNN模型与算法，我们会在其它博客里讲。