

神经网络反向传播原理

1. 背景

遇到难题时，逃避往往很容易。神经网络中的反向传播就属于其中之一。Tensorflow给逃避提供了很好的条件。日子越久，觉得自己好像什么也不会。现在我要正视这个问题，来和我一起彻底弄清楚反向传播是怎么回事吧。公式推导的过程中最让人疑惑的可能不是内容本身，而是看起来吓人的公式。这篇文章我想尽一起办法对初学者友好。如果你觉得有些太过简单，可自行跳过，专注于你关注的部分。

2. 网络学习的目标

今天研究的神经网络是最基本的全连接网络，除了输入层神经元，其他层所有神经元都使用的是Sigmoid激活函数。在神经网络中，我们用损失函数(cost function/loss function)来表示神经网络在给定数据集上面的表现。通常来说，我们希望损失函数越小越好，而这个损失函数越变越小的过程就是学习(learning)的过程。在手写字符识别的例子中，在单个训练样本(图像-标签对)的损失函数，即预测输出和标签之间的差距，我们可以这样定义：

$$C(\mathbf{w}, \mathbf{b}) = \frac{1}{2} \|\mathbf{y} - \mathbf{o}\|^2 = \frac{1}{2} \sum_{i=0}^9 (\mathbf{y}_i - \mathbf{o}_i)^2 \quad (41)$$

一张图片大小为28x28，我们将其展开为一个一维的向量 \mathbf{x} ，那么这个向量有784个分量。这张图片对应的标签同样用一个一维的列向量 \mathbf{y} 来表示，由于我们最终分了10类，因此标签向量有10个分量。同理神经网络的输出也可以表示成一个有10个分量的向量 \mathbf{o} 。神经网络当中所有的权重可以构成一个向量 \mathbf{w} ，同样地所有偏置也可以构成一个向量 \mathbf{b} 。输入一张图片 \mathbf{x} ，经过由 \mathbf{w} 和 \mathbf{b} 所确定的神经网络的前向传播之后，我们可以得到对应输出 $\mathbf{o}_{\mathbf{w}, \mathbf{b}}(\mathbf{x})$ (简称为 \mathbf{o})。多个训练样本的损失函数就是单个损失函数(记为 C^i)的叠加：

$$L(\mathbf{w}, \mathbf{b}) = \sum_{i=1}^N C^i(\mathbf{w}, \mathbf{b}) \quad (42)$$

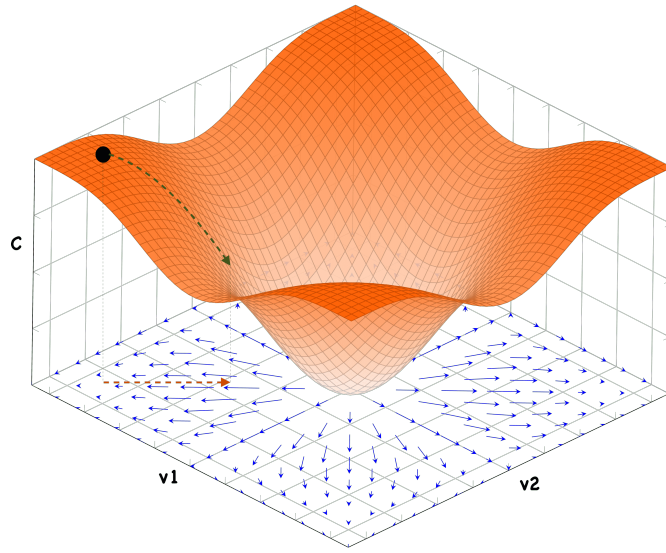
神经网络学习的过程就是不断调整 \mathbf{w}, \mathbf{b} 让 $L(\mathbf{w}, \mathbf{b})$ 逐渐变小的过程。

3. 参数更新的策略

我们现在来考虑，如何才能使得单个训练样本的损失函数 C 越来越小。我们将当前时刻的损失函数记为 C ，下个时刻的损失函数记为 C' ，那么损失的变化量可以表示成：

$$\Delta C = C' - C \quad (43)$$

要让损失函数越来越小，我们只要想办法让 ΔC 每次都为负就可以了。接下来的问题就变成了：如何才能使 ΔC 每次都为负？为了说清楚，我们现在考虑一个简单的情况，假设损失函数 $C(\mathbf{v})$ 只包含两个参数 v_1, v_2 ，即 $C(\mathbf{v}) = C(v_1, v_2)$ 。



根据微积分的知识，损失函数的变化量可以表示成：

$$\Delta C \approx \frac{\partial C}{\partial v_1} \Delta v_1 + \frac{\partial C}{\partial v_2} \Delta v_2 = \left[\frac{\partial C}{\partial v_1}, \frac{\partial C}{\partial v_2} \right] \cdot [\Delta v_1, \Delta v_2]^T \quad (44)$$

其中符号 \cdot 表示的是向量的内积。为了方便记 $\nabla C = \left[\frac{\partial C}{\partial v_1}, \frac{\partial C}{\partial v_2} \right]$, $\Delta v = [\Delta v_1, \Delta v_2]^T$, 于是上式可以写成：

$$\Delta C \approx \nabla C \cdot \Delta v \quad (45)$$

回到我们的问题：如何才能使 ΔC 每次都为负？于是有聪明人想出了这样一个办法。让参数 v 按照下面的方法更新，新的参数

$$v' = v - \eta \nabla C \quad (46)$$

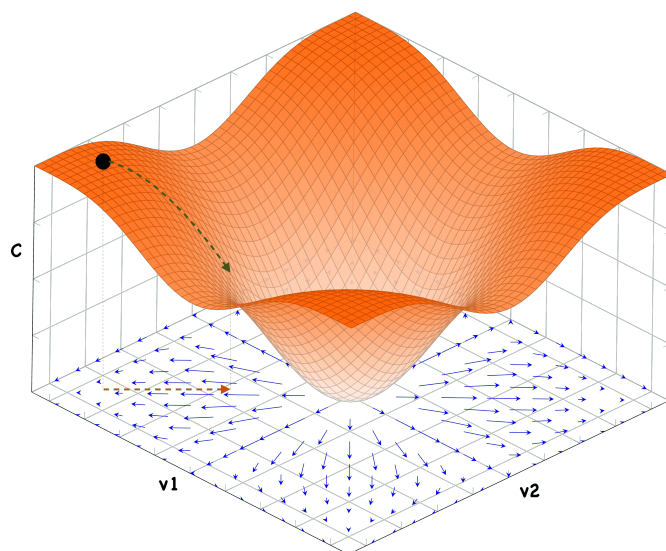
其中 η 是个正数。你应该会问为什么要这样更新呢？把等式右边的 v 移到左边就有

$$v' - v = \Delta v = -\eta \nabla C \quad (47)$$

将这个式子带入到损失函数的改变量中就得到

$$\Delta C \approx \nabla C \cdot (-\eta \nabla C) = -\eta (\nabla C)^2 \quad (48)$$

由于我们规定 $\eta > 0$, 且 $(\nabla C)^2 \geq 0$, 因此就有 $\Delta C < 0$ 。于是我们就找到一种参数更新策略使得每次更新都会使得损失函数减小。



我们把 $\nabla C = \left[\frac{\partial C}{\partial v_1}, \frac{\partial C}{\partial v_2} \right]$ 叫做损失函数 C 在 (v_1, v_2) 处的**梯度**。上图中用 v_1 - v_2 平面上的蓝色箭头表示梯度的大小和方向。由于在上面的策略当中，参数 v 的变化方向(红色箭头)和梯度的方向正好相反，即使得梯度减小。因此上面的参数更新策略被叫做**梯度下降方法**。

具体来说，对于上述例子中的每个参数的更新策略，可以由 $v' = v - \eta \nabla C$ 得到:

$$\begin{aligned} v_1 &\rightarrow v'_1 = v_1 - \eta \frac{\partial C}{\partial v_1} \\ v_2 &\rightarrow v'_2 = v_2 - \eta \frac{\partial C}{\partial v_2} \end{aligned} \quad (49)$$

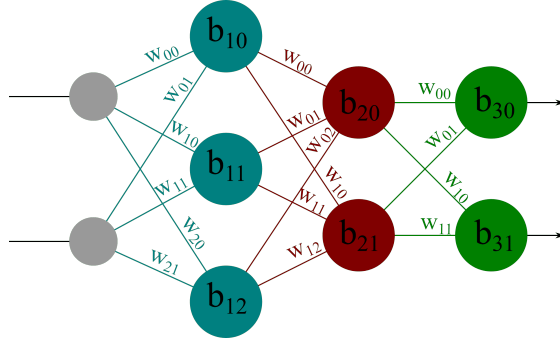
于是从二元函数推广，可以得到神经网络中任意一个权重和偏置的更新策略，

$$\begin{aligned} w_k &\rightarrow w'_k = w_k - \eta \frac{\partial C}{\partial w_k} \\ b_l &\rightarrow b'_l = b_l - \eta \frac{\partial C}{\partial b_l} \end{aligned} \quad (50)$$

到这里我们就找到了一种神经网络参数更新的策略，每一次网络参数的更新都可以让损失函数的值减小，从而使得网络的预测的值越来越接近训练数据的标签。要实现这个策略，我们只需要计算损失函数对所有的权重和偏置的偏微分, 即梯度的分量就可以了。

4. 符号标记

为了说清楚，这里我们使用下面这个简单的神经网络来做演示。这个神经网络有4层，每层的神经元个数分别是[2, 3, 2, 2]。由于输入层的神经元没有 **加权求和** 的操作，我这里把输入层叫做第0层。随后的1, 2, 3层分别用不同颜色加以区分。



为了方便后续说明，我们用

- 矩阵 \mathbf{W}^i 记录所有连接到第 i 层神经元对应的权重。例如这里的连接到第1层神经元的权重可以表示为：

$$\mathbf{W}^1 = \begin{bmatrix} w_{00} & w_{01} \\ w_{10} & w_{11} \\ w_{20} & w_{21} \end{bmatrix} \quad (51)$$

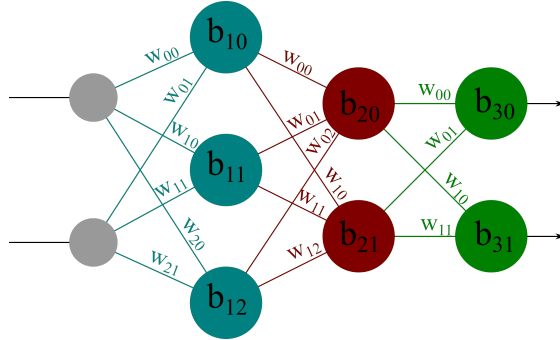
- 列向量 \mathbf{b}_i 记录第 i 层神经元所有的偏置值。例如这里的第1层神经元的所有偏置可以表示成：

$$\mathbf{b}_1 = \begin{bmatrix} b_{10} \\ b_{11} \\ b_{12} \end{bmatrix} \quad (52)$$

- 列向量 \mathbf{z}_i 记录第 i 层所有神经元的加权求和值；
- 列向量 \mathbf{a}_i 记录第 i 层所有神经元的输出值。

5. 计算梯度分量的方法

接下来要做的就是找到一个方法来计算损失函数 C 对任意一个权重和偏置的偏微分。



对于非输入层上的任意一个神经元的输出值 a 都可以表示成：

$$a = \sigma(z) \quad (53)$$

其中 z 表示这个神经元加权求和的结果， σ 表示Sigmoid函数。具体来说，对于第1层第0个神经元的输出值为

$$a_{10} = \sigma(z_{10}) \quad (54)$$

其中 $z_{10} = w_{00}a_{00} + w_{01}a_{01} + b_{10}$ 。可见任意一个神经元的输出 a 是中间值 z 的函数，而 z 又是 w 和 b 的函数。

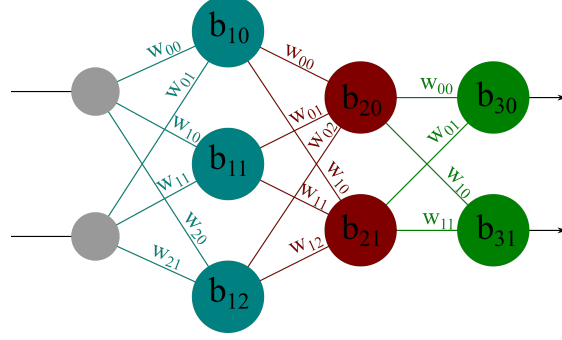
损失函数 C 是是权重 w 和偏置 b 的函数，我们同样也可以认为 C 是 z 的函数。

损失函数对第1层的其中一个权重 w_{00} 的偏导数，据链式求导法则可以得到：

$$\frac{\partial C}{\partial w_{00}} = \frac{\partial C}{\partial z_{10}} \frac{\partial z_{10}}{\partial w_{00}} \quad (55)$$

$$\frac{\partial C}{\partial \mathbf{W}^i} = \frac{\partial C}{\partial \mathbf{z}_i} \frac{\partial \mathbf{z}_i}{\partial \mathbf{W}^i} \quad (56)$$

5.1 计算 $\frac{\partial z}{\partial w}$



由于激活 $a_{10} = \sigma(z_{10})$ 的中间变量 $z_{10} = w_{00}a_{00} + w_{01}a_{01} + b_{10}$ ，因此有

$$\frac{\partial z_{10}}{\partial w_{00}} = a_{00} \quad (57)$$

即为与该神经元相连的上一层神经元的输出值。

推广1

推广得到 第一层所有的权重值 \mathbf{W}^1 对应的 $\frac{\partial z}{\partial w}$ 为:

$$\frac{\partial \mathbf{z}_1}{\partial \mathbf{W}^1} = \begin{bmatrix} \frac{\partial z_{10}}{\partial w_{00}} & \frac{\partial z_{10}}{\partial w_{01}} \\ \frac{\partial z_{11}}{\partial w_{10}} & \frac{\partial z_{11}}{\partial w_{11}} \\ \frac{\partial z_{12}}{\partial w_{20}} & \frac{\partial z_{12}}{\partial w_{21}} \end{bmatrix} = \begin{bmatrix} a_{00} & a_{01} \\ a_{00} & a_{01} \\ a_{00} & a_{01} \end{bmatrix} \equiv \mathbf{A}^0 \quad (58)$$

注意到，这里我们定义了一个矩阵 \mathbf{A}^0 表示第0层对应的这样一个矩阵。矩阵的特点是每一行都是该层的输出值 \mathbf{a}_0^\top 。

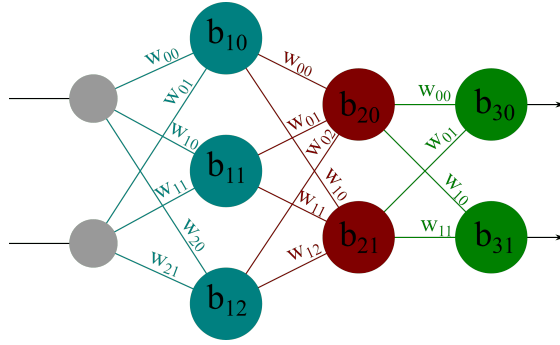
继续推广 得到第 i 层的所有权重 \mathbf{W}^i 对应的 $\frac{\partial z}{\partial w}$ 值为上一层第 $h = i - 1$ 层对应的 \mathbf{A}^h ，即

$$\frac{\partial \mathbf{z}_i}{\partial \mathbf{W}^i} = \mathbf{A}^h \quad (59)$$

而 \mathbf{A}^h 利用前向传播可以容易地计算出来。于是 $\frac{\partial z}{\partial w}$ 就计算完成了。

$$\frac{\partial C}{\partial \mathbf{W}^i} = \frac{\partial C}{\partial \mathbf{z}_i} \mathbf{A}^h \quad (60)$$

5.2 计算 $\frac{\partial C}{\partial z}$



$$\frac{\partial C}{\partial z_{10}} = \frac{\partial C}{\partial a_{10}} \frac{\partial a_{10}}{\partial z_{10}} \quad (61)$$

5.2.1 算 $\frac{\partial a}{\partial z}$

因为 $a_{10} = \sigma(z_{10})$, 因此对上式中第二个乘积项可以写成

$$\frac{\partial a_{10}}{\partial z_{10}} = \sigma'(z_{10}) \quad (62)$$

推广2

推广得到 这一层上所有神经元对应的 $\frac{\partial a}{\partial z}$ 为:

$$\frac{\partial \mathbf{a}_1}{\partial \mathbf{z}_1} = \begin{bmatrix} \frac{\partial a_{10}}{\partial z_{10}} \\ \frac{\partial a_{11}}{\partial z_{11}} \\ \frac{\partial a_{12}}{\partial z_{12}} \end{bmatrix} = \begin{bmatrix} \sigma'(z_{10}) \\ \sigma'(z_{11}) \\ \sigma'(z_{12}) \end{bmatrix} \equiv \sigma' \left(\begin{bmatrix} z_{10} \\ z_{11} \\ z_{12} \end{bmatrix} \right) = \sigma'(\mathbf{z}_1) \quad (63)$$

注意到这里定义了一个针对向量的函数操作, 即对某一个向量的操作就是对向量中每一个元素做这样的操作。

继续推广 第 i 层所有 $\frac{\partial a}{\partial z}$ 值为 $\sigma'(\mathbf{z}_i)$ 。即

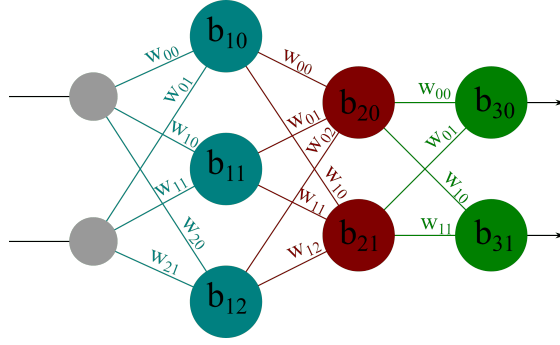
$$\frac{\partial \mathbf{a}_i}{\partial \mathbf{z}_i} = \sigma'(\mathbf{z}_i) \quad (64)$$

\mathbf{z}_i 在前向传播的时候可以容易地计算出来, 而 $\sigma'(z)$ 也是知道的。于是 $\frac{\partial a}{\partial z}$ 就计算出来了。

5.2.2 算 $\frac{\partial C}{\partial a}$

还剩下一项 $\frac{\partial C}{\partial a_{10}}$, (Todo, 添加一段解释)

$$\frac{\partial C}{\partial a_{10}} = \frac{\partial C}{\partial z_{20}} \frac{\partial z_{20}}{\partial a_{10}} + \frac{\partial C}{\partial z_{21}} \frac{\partial z_{21}}{\partial a_{10}} \quad (65)$$



由于 $z_{20} = w_{00}a_{10} + w_{01}a_{11} + w_{02}a_{12}$, 因此 $\frac{\partial z_{20}}{\partial a_{10}} = w_{00}$, 同理 $\frac{\partial z_{21}}{\partial a_{10}} = w_{10}$, 于是上式变成:

$$\frac{\partial C}{\partial a_{10}} = w_{00} \frac{\partial C}{\partial z_{20}} + w_{10} \frac{\partial C}{\partial z_{21}} = [w_{00}, w_{10}] \begin{bmatrix} \frac{\partial C}{\partial z_{20}} \\ \frac{\partial C}{\partial z_{21}} \end{bmatrix} \quad (66)$$

推广3

推广得到第1层上所有神经元的输出 \mathbf{a}_1 对应的 $\frac{\partial C}{\partial \mathbf{a}_1}$ 为:

$$\frac{\partial C}{\partial \mathbf{a}_1} = \begin{bmatrix} \frac{\partial C}{\partial a_{10}} \\ \frac{\partial C}{\partial a_{11}} \\ \frac{\partial C}{\partial a_{12}} \end{bmatrix} = \begin{bmatrix} w_{00} & w_{10} \\ w_{01} & w_{11} \\ w_{02} & w_{12} \end{bmatrix} \begin{bmatrix} \frac{\partial C}{\partial z_{20}} \\ \frac{\partial C}{\partial z_{21}} \end{bmatrix} = \begin{bmatrix} w_{00} & w_{01} & w_{02} \\ w_{10} & w_{11} & w_{12} \end{bmatrix}^T \begin{bmatrix} \frac{\partial C}{\partial z_{20}} \\ \frac{\partial C}{\partial z_{21}} \end{bmatrix} = \mathbf{W}^{2T} \frac{\partial C}{\partial \mathbf{z}_2} \quad (67)$$

继续推广 第 i 层神经元的输出 \mathbf{a}_i 对应的 $\frac{\partial C}{\partial \mathbf{a}_i}$ 为

$$\frac{\partial C}{\partial \mathbf{a}_i} = \mathbf{W}^{jT} \frac{\partial C}{\partial \mathbf{z}_j} \quad (68)$$

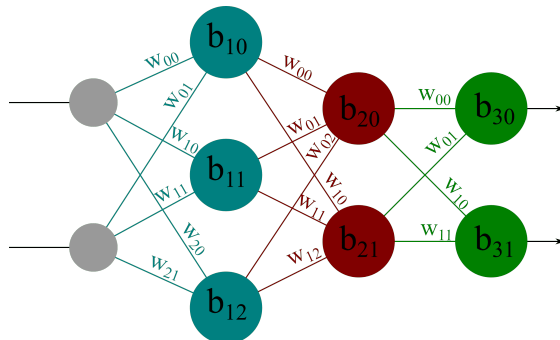
5.3 整理

$$\frac{\partial C}{\partial \mathbf{W}^i} = \frac{\partial C}{\partial \mathbf{z}_i} \mathbf{A}^h \quad (69)$$

$$\frac{\partial C}{\partial \mathbf{z}_i} = \left(\mathbf{W}^{jT} \frac{\partial C}{\partial \mathbf{z}_j} \right) \odot \sigma'(\mathbf{z}_i) \quad (70)$$

只要算出每个神经元的 $\frac{\partial C}{\partial \mathbf{z}}$ 即可算出梯度的改变量, 先把倒数第一层的 $\frac{\partial C}{\partial \mathbf{z}_{-1}}$ 算出, 一层层就可往前算出了。

5.4 计算 $\frac{\partial C}{\partial \mathbf{z}_{-1}}$



输出层的神经元的输出 \mathbf{a}_3 这里用 \mathbf{o} 来表示，单个训练数据的损失函数是：

$$C(\mathbf{w}, \mathbf{b}) = \frac{1}{2} \|\mathbf{y} - \mathbf{o}\|^2 = \frac{1}{2} \sum_{i=0}^1 (y_i - o_i)^2 = \frac{1}{2} (y_0^2 - 2y_0o_0 + o_0^2 + y_1^2 - 2y_1o_1 + o_1^2) \quad (71)$$

$$o_0 = \sigma(z_{30}) \quad (72)$$

$$\frac{\partial C}{\partial z_{30}} = \frac{\partial C}{\partial o_0} \frac{\partial o_0}{\partial z_{30}} = (o_0 - y_0) \sigma'(z_{30}) \quad (73)$$

同理

$$\frac{\partial C}{\partial z_{31}} = \frac{\partial C}{\partial o_1} \frac{\partial o_1}{\partial z_{31}} = (o_1 - y_1) \sigma'(z_{31}) \quad (74)$$

$$\frac{\partial C}{\partial \mathbf{z}_{-1}} = (\mathbf{o} - \mathbf{y}) \sigma'(\mathbf{z}_{-1}) \quad (75)$$

5.5 总结

对于权重

$$\frac{\partial C}{\partial \mathbf{W}^i} = \frac{\partial C}{\partial \mathbf{z}_i} \mathbf{A}^h \quad (76)$$

$$\frac{\partial C}{\partial \mathbf{z}_i} = \left(\mathbf{W}^{j\top} \frac{\partial C}{\partial \mathbf{z}_j} \right) \odot \sigma'(\mathbf{z}_i) \quad (77)$$

$$\frac{\partial C}{\partial \mathbf{z}_{-1}} = (\mathbf{o} - \mathbf{y}) \sigma'(\mathbf{z}_{-1}) \quad (78)$$

类似的, 损失函数对偏置的偏导可以容易地得到

$$\frac{\partial C}{\partial \mathbf{b}_i} = \frac{\partial C}{\partial \mathbf{z}_i} \quad (79)$$

6. 计算多个训练样本的梯度分量

对每一个训练样本得到的梯度分量求和，即可得到总的梯度分量：

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{W}^i} &= \sum_{n=1}^N \frac{\partial C^n}{\partial \mathbf{W}^i} \\ \frac{\partial L}{\partial \mathbf{b}_i} &= \sum_{n=1}^N \frac{\partial C^n}{\partial \mathbf{b}_i} \end{aligned} \quad (80)$$

7. 结语

现在来看，理解反向传播算法要从不同的维度的考虑，当考虑了对单个参数的导数以后，还得同时要考虑对这一层中所有参数的导数，还得通过迭代的方式得到对所有层的参数的倒数。另外在推导的过程中引入的数学记号也好像让理解变得困难了。事实却并不是这样。