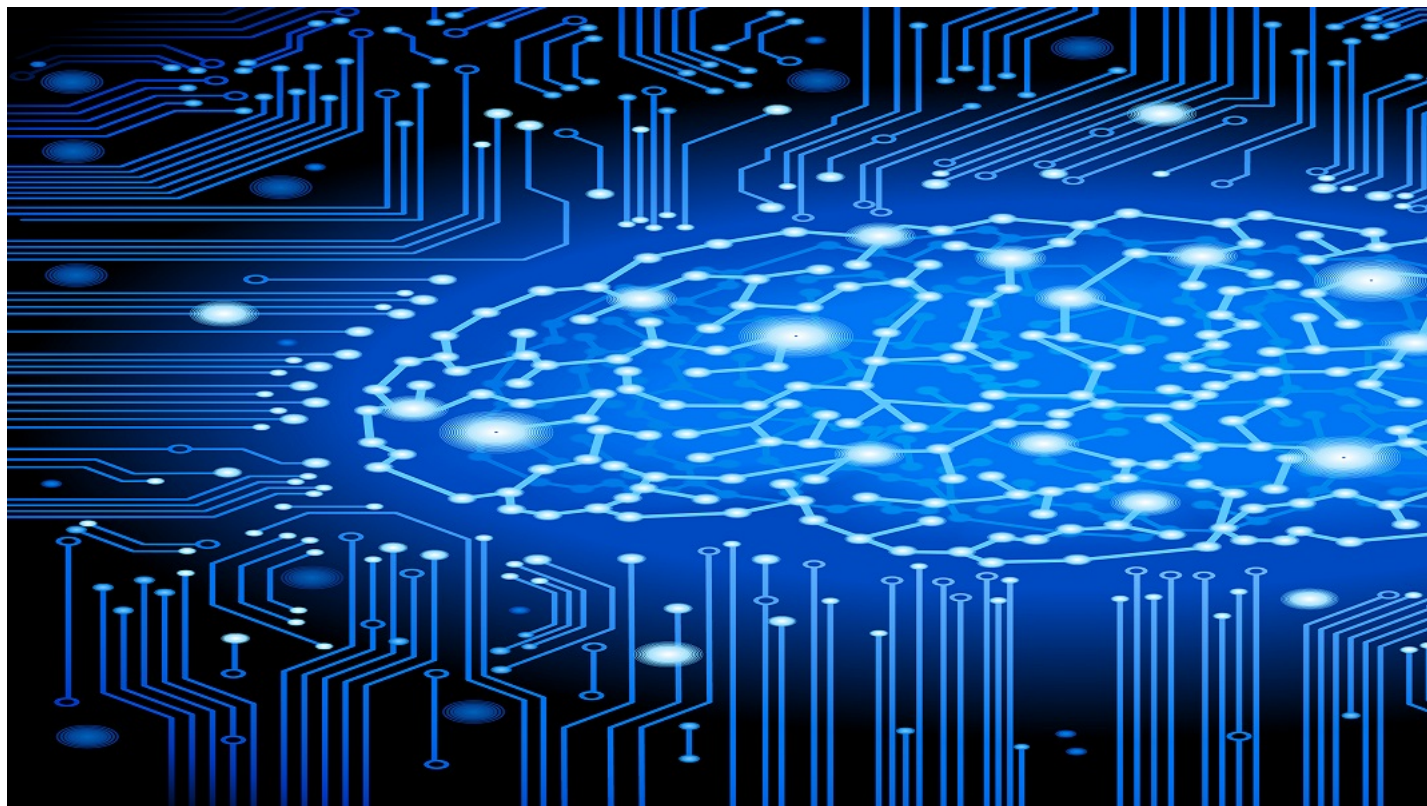


[\[关闭\]](#)

@hanbingtao 2017-03-01 01:08 字数 27318 阅读 20252

# 零基础入门深度学习(6) - 长短时记忆网络(LSTM)

机器学习 深度学习入门



无论即将到来的是大数据时代还是人工智能时代，亦或是传统行业使用人工智能在云上处理大数据的时代，作为一个有理想有追求的程序员，不懂深度学习（Deep Learning）这个超热的技术，会不会感觉马上就out了？现在救命稻草来了，《零基础入门深度学习》系列文章旨在讲帮助爱编程的你从零基础达到入门级水平。零基础意味着你不需要太多的数学知识，只要会写程序就行了，没错，这是专门为程序员写的文章。虽然文中会有很多公式你也许看不懂，但同时也会有更多的代码，程序员的你一定能看懂的（我周围是一群狂热的Clean Code程序员，所以我写的代码也不会很差）。

## 文章列表

[零基础入门深度学习\(1\) - 感知器](#)[零基础入门深度学习\(2\) - 线性单元和梯度下降](#)[零基础入门深度学习\(3\) - 神经网络和反向传播算法](#)[零基础入门深度学习\(4\) - 卷积神经网络](#)[零基础入门深度学习\(5\) - 循环神经网络](#)[零基础入门深度学习\(6\) - 长短时记忆网络\(LSTM\)](#)[零基础入门深度学习\(7\) - 递归神经网络](#)

## 往期回顾

在上一篇文章中，我们介绍了**循环神经网络**以及它的训练算法。我们也介绍了**循环神经网络**很难训练的原因，这导致了它在实际应用中，很难处理长距离的依赖。在本文中，我们将介绍一种改进之后的循环神经网络：**长短时记忆网络(Long Short Term Memory Network, LSTM)**，它成功的解决了原始循环神经网络的缺陷，成为当前最流行的RNN，在语音识别、图片描述、自然语言处理等许多领域中成功应用。但不幸的一面是，**LSTM**的结构很复杂，因此，我们需要花上一些力气，才能把LSTM以及它的训练算法弄明白。在搞清楚**LSTM**之后，我们再介绍一种**LSTM**的变体：**GRU (Gated Recurrent Unit)**。它的结构比**LSTM**简单，而效果却和**LSTM**一样好，因此，它正在逐渐流行起来。最后，我们仍然会动手实现一个**LSTM**。

## 长短时记忆网络是啥

我们首先了解一下长短时记忆网络产生的背景。回顾一下[零基础入门深度学习\(5\) - 循环神经网络](#)中推导的，误差项沿时间反向传播的公式：

$$\delta_k^T = \delta_i^T \prod_{i=k}^{t-1} \text{diag}[f'(\text{net}_i)] W \quad (1)$$

我们可以根据下面的不等式，来获取 $\delta_k^T$ 的模的上界（模可以看做对 $\delta_k^T$ 中每一项值的大小的度量）：

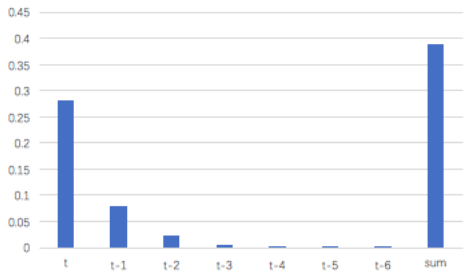
$$\begin{aligned}\|\delta_k^T\| &\leq \|\delta_t^T\| \prod_{i=k}^{t-1} \|diag[f'(\mathbf{net}_i)]\| \|W\| \\ &\leq \|\delta_t^T\| (\beta_f \beta_w)^{t-k}\end{aligned}\tag{2}$$
$$\tag{3}$$

我们可以看到，误差项 $\delta$ 从t时刻传递到k时刻，其值的上界是 $\beta_f \beta_w$ 的指数函数。 $\beta_f \beta_w$ 分别是对角矩阵 $diag[f'(\mathbf{net}_i)]$ 和矩阵W模的上界。显然，除非 $\beta_f \beta_w$ 乘积的值位于1附近，否则，当t-k很大时（也就是误差传递很多个时刻时），整个式子的值就会变得极小（当 $\beta_f \beta_w$ 乘积小于1）或者极大（当 $\beta_f \beta_w$ 乘积大于1），前者就是**梯度消失**，后者就是**梯度爆炸**。虽然科学家们搞出了很多技巧（比如怎样初始化权重），让 $\beta_f \beta_w$ 的值尽可能贴近于1，终究还是难以抵挡指数函数的威力。

**梯度消失**到底意味着什么？在[零基础入门深度学习\(5\) - 循环神经网络](#)中我们已证明，权重数组W最终的梯度是各个时刻的梯度之和，即：

$$\begin{aligned}\nabla_W E &= \sum_{k=1}^t \nabla_{W_k} E \\ &= \nabla_{W_t} E + \nabla_{W_{t-1}} E + \nabla_{W_{t-2}} E + \dots + \nabla_{W_1} E\end{aligned}\tag{4}$$
$$\tag{5}$$

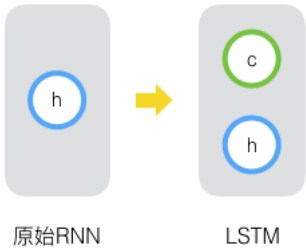
假设某轮训练中，各时刻的梯度以及最终的梯度之和如下图：



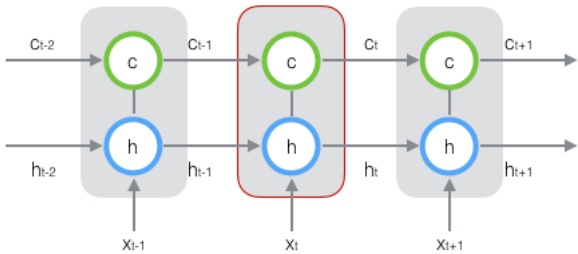
我们就可以看到，从上图的t-3时刻开始，梯度已经几乎减少到0了。那么，从这个时刻开始再往之前走，得到的梯度（几乎为零）就不会对最终的梯度值有任何贡献，这就相当于无论t-3时刻之前的网络状态h是什么，在训练中都不会对权重数组W的更新产生影响，也就是网络事实上已经忽略了t-3时刻之前的状态。这就是原始RNN无法处理长距离依赖的原因。

既然找到了问题的原因，那么我们就能解决它。从问题的定位到解决，科学家们大概花了7、8年时间。终于有一天，Hochreiter和Schmidhuber两位科学家发明出**长短时记忆网络**，一举解决这个问题。

其实，**长短时记忆网络**的思路比较简单。原始RNN的隐藏层只有一个状态，即h，它对于短期的输入非常敏感。那么，假如我们再增加一个状态，即c，让它来保存长期的状态，那么问题不就解决了吗？如下图所示：



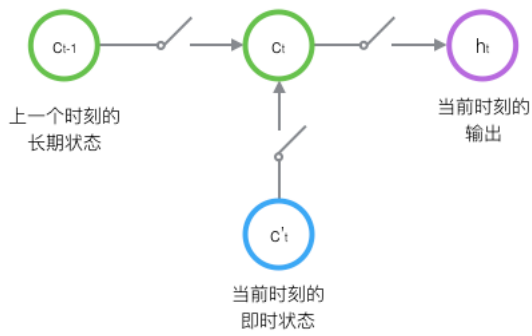
新增加的状态c，称为**单元状态(cell state)**。我们把上图按照时间维度展开：



上图仅仅是一个示意图，我们可以看出，在t时刻，LSTM的输入有三个：当前时刻网络的输入值 $\mathbf{x}_t$ 、上一时刻LSTM的输出值 $\mathbf{h}_{t-1}$ 、以及上一时刻的单元状态 $\mathbf{c}_{t-1}$ ；LSTM的输出有两个：当前时刻LSTM输出值 $\mathbf{h}_t$ 、和当前时刻的单元状态 $\mathbf{c}_t$ 。注意 $\mathbf{x}$ 、 $\mathbf{h}$ 、 $\mathbf{c}$ 都是**向量**。

LSTM的关键，就是怎样控制长期状态c。在这里，LSTM的思路是使用三个控制开关。第一个开关，负责控制继续保存长期状态c；第二个开关，负责控制把即时状态输入到长期状态c；第三个开关，负责控制是否把长期状态c作为当前的LSTM的输出。三个开关的作用如下图所示：

### 长期状态c的控制



接下来，我们要描述一下，输出h和单元状态c的具体计算方法。

## 长短时记忆网络的前向计算

前面描述的开关是怎样在算法中实现的呢？这就用到了门（gate）的概念。门实际上就是一层全连接层，它的输入是一个向量，输出是一个0到1之间的实数向量。假设W是门的权重向量，b是偏置项，那么门可以表示为：

$$g(\mathbf{x}) = \sigma(W\mathbf{x} + \mathbf{b})$$

门的使用，就是用门的输出向量按元素乘以我们需要控制的那个向量。因为门的输出是0到1之间的实数向量，那么，当门输出为0时，任何向量与之相乘都会得到0向量，这就相当于啥都不能通过；输出为1时，任何向量与之相乘都不会有任何改变，这就相当于啥都可以通过。因为 $\sigma$ （也就是sigmoid函数）的值域是(0,1)，所以门的状态都是半开半闭的。

LSTM用两个门来控制单元状态c的内容，一个是遗忘门（forget gate），它决定了上一时刻的单元状态 $\mathbf{c}_{t-1}$ 有多少保留到当前时刻 $\mathbf{c}_t$ ；另一个是输入门（input gate），它决定了当前时刻网络的输入 $\mathbf{x}_t$ 有多少保存到单元状态 $\mathbf{c}_t$ 。LSTM用输出门（output gate）来控制单元状态 $\mathbf{c}_t$ 有多少输出到LSTM的当前输出值 $\mathbf{h}_t$ 。

我们先来看一下遗忘门：

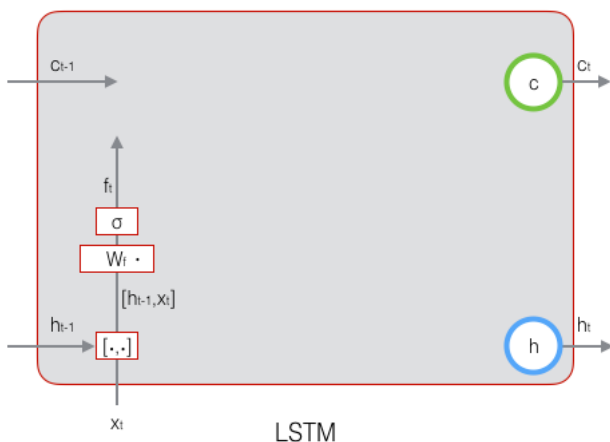
$$\mathbf{f}_t = \sigma(W_f \cdot [\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_f) \quad (\text{式1})$$

上式中， $W_f$ 是遗忘门的权重矩阵， $[\mathbf{h}_{t-1}, \mathbf{x}_t]$ 表示把两个向量连接成一个更长的向量， $\mathbf{b}_f$ 是遗忘门的偏置项， $\sigma$ 是sigmoid函数。如果输入的维度是 $d_x$ ，隐藏层的维度是 $d_h$ ，单元状态的维度是 $d_c$ （通常 $d_c = d_h$ ），则遗忘门的权重矩阵 $W_f$ 维度是 $d_c \times (d_h + d_x)$ 。事实上，权重矩阵 $W_f$ 都是两个矩阵拼接而成的：一个是 $W_{fh}$ ，它对应着输入项 $\mathbf{h}_{t-1}$ ，其维度为 $d_c \times d_h$ ；一个是 $W_{fx}$ ，它对应着输入项 $\mathbf{x}_t$ ，其维度为 $d_c \times d_x$ 。 $W_f$ 可以写为：

$$[W_f] \begin{bmatrix} \mathbf{h}_{t-1} \\ \mathbf{x}_t \end{bmatrix} = [W_{fh} \quad W_{fx}] \begin{bmatrix} \mathbf{h}_{t-1} \\ \mathbf{x}_t \end{bmatrix} \quad (6)$$

$$= W_{fh}\mathbf{h}_{t-1} + W_{fx}\mathbf{x}_t \quad (7)$$

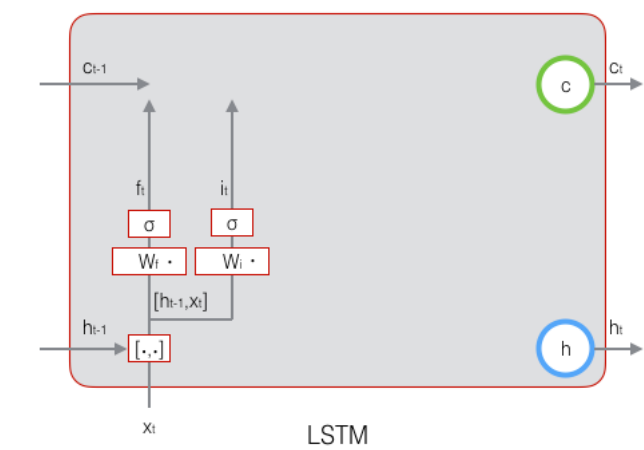
下图显示了遗忘门的计算：



接下来看看输入门：

$$\mathbf{i}_t = \sigma(W_i \cdot [\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_i) \quad (\text{式2})$$

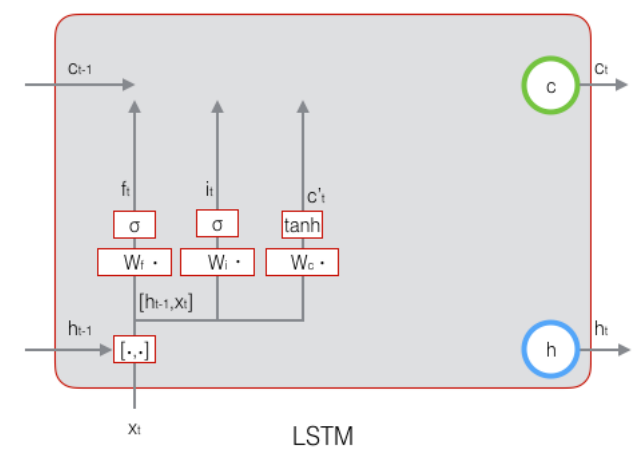
上式中， $W_i$ 是输入门的权重矩阵， $\mathbf{b}_i$ 是输入门的偏置项。下图表示了输入门的计算：



接下来，我们计算用于描述当前输入的单元状态 $\tilde{c}_t$ ，它是根据上一次的输出和本次输入来计算的：

$$\tilde{c}_t = \tanh(W_c \cdot [h_{t-1}, x_t] + b_c) \tag{式3}$$

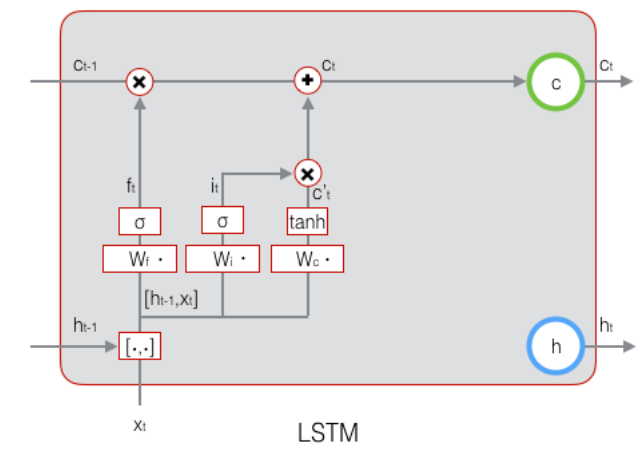
下图是 $\tilde{c}_t$ 的计算：



现在，我们计算当前时刻的单元状态 $c_t$ 。它是由上一次的单元状态 $c_{t-1}$ 按元素乘以遗忘门 $f_t$ ，再用当前输入的单元状态 $\tilde{c}_t$ 按元素乘以输入门 $i_t$ ，再将两个积加和产生的：

$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t \tag{式4}$$

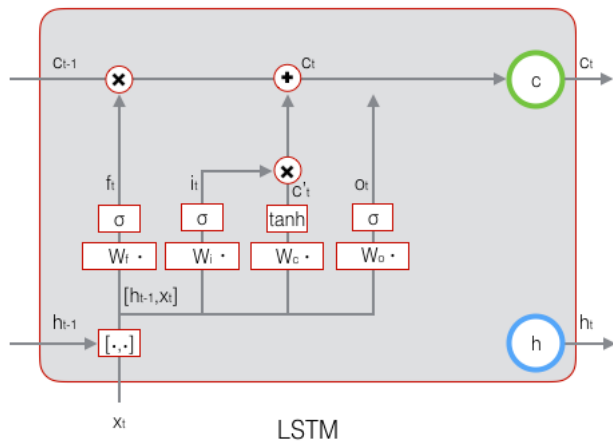
符号 $\circ$ 表示按元素乘。下图是 $c_t$ 的计算：



这样，我们就把LSTM关于当前的记忆 $\tilde{c}_t$ 和长期的记忆 $c_{t-1}$ 组合在一起，形成了新的单元状态 $c_t$ 。由于遗忘门的控制，它可以保存很久很久之前的信息，由于输入门的控制，它又可以避免当前无关紧要的内容进入记忆。下面，我们要看看输出门，它控制了长期记忆对当前输出的影响：

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o) \tag{式5}$$

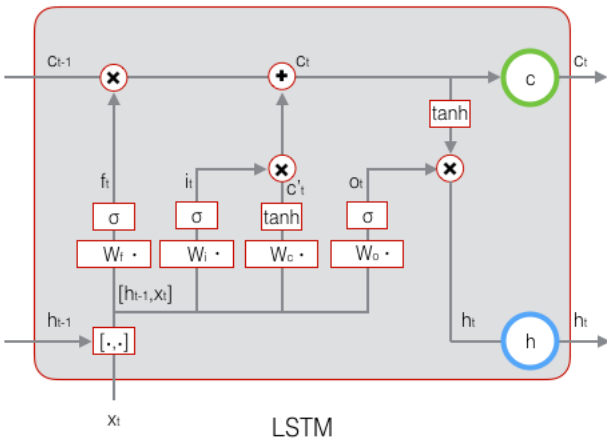
下图表示输出门的计算：



LSTM最终的输出，是由输出门和单元状态共同确定的：

$$\mathbf{h}_t = \mathbf{o}_t \circ \tanh(\mathbf{c}_t) \quad (\text{式6})$$

下图表示LSTM最终输出的计算：



式1到式6就是LSTM前向计算的全部公式。至此，我们就把LSTM前向计算讲完了。

## 长短时记忆网络的训练

熟悉我们这个系列文章的同学都清楚，训练部分往往比前向计算部分复杂多了。LSTM的前向计算都这么复杂，那么，可想而知，它的训练算法一定是非常非常复杂的。现在只有做几次深呼吸，再一头扎进公式海洋吧。

### LSTM训练算法框架

LSTM的训练算法仍然是反向传播算法，对于这个算法，我们已经非常熟悉了。主要有下面三个步骤：

1. 前向计算每个神经元的输出值，对于LSTM来说，即 $\mathbf{f}_t$ 、 $\mathbf{i}_t$ 、 $\mathbf{c}_t$ 、 $\mathbf{o}_t$ 、 $\mathbf{h}_t$ 五个向量的值。计算方法已经在上一节中描述过了。
2. 反向计算每个神经元的误差项 $\delta$ 值。与循环神经网络一样，LSTM误差项的反向传播也是包括两个方向：一个是沿时间的反向传播，即从当前t时刻开始，计算每个时刻的误差项；一个是将误差项向上一层传播。
3. 根据相应的误差项，计算每个权重的梯度。

### 关于公式和符号的说明

首先，我们对推导中用到的一些公式、符号做一下必要的说明。

接下来的推导中，我们设定gate的激活函数为sigmoid函数，输出的激活函数为tanh函数。他们的导数分别为：

$$\sigma(z) = y = \frac{1}{1 + e^{-z}} \quad (8)$$

$$\sigma'(z) = y(1 - y) \quad (9)$$

$$\tanh(z) = y = \frac{e^z - e^{-z}}{e^z + e^{-z}} \quad (10)$$

$$\tanh'(z) = 1 - y^2 \quad (11)$$

从上面可以看出，sigmoid和tanh函数的导数都是原函数的函数。这样，我们一旦计算原函数的值，就可以用它来计算出导数的值。

LSTM需要学习的参数共有8组，分别是：遗忘门的权重矩阵 $\mathbf{W}_f$ 和偏置项 $\mathbf{b}_f$ 、输入门的权重矩阵 $\mathbf{W}_i$ 和偏置项 $\mathbf{b}_i$ 、输出门的权重矩阵 $\mathbf{W}_o$ 和偏置项 $\mathbf{b}_o$ ，以及计算单元状态的权重矩阵 $\mathbf{W}_c$ 和偏置项 $\mathbf{b}_c$ 。因为权重矩阵的两部分在反向传播中使用不同的公式，因此在后续的推导中，权重矩阵 $\mathbf{W}_f$ 、 $\mathbf{W}_i$ 、 $\mathbf{W}_c$ 、 $\mathbf{W}_o$ 都将被写为分开的两个矩阵： $\mathbf{W}_{fh}$ 、 $\mathbf{W}_{fx}$ 、 $\mathbf{W}_{ih}$ 、 $\mathbf{W}_{ix}$ 、 $\mathbf{W}_{oh}$ 、 $\mathbf{W}_{ox}$ 、 $\mathbf{W}_{ch}$ 、 $\mathbf{W}_{cx}$ 。

我们解释一下按元素乘 $\circ$ 符号。当 $\circ$ 作用于两个**向量**时，运算如下：

$$\mathbf{a} \circ \mathbf{b} = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ \dots \\ a_n \end{bmatrix} \circ \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ \dots \\ b_n \end{bmatrix} = \begin{bmatrix} a_1 b_1 \\ a_2 b_2 \\ a_3 b_3 \\ \dots \\ a_n b_n \end{bmatrix}$$

当 $\circ$ 作用于一个**向量**和一个**矩阵**时，运算如下：

$$\mathbf{a} \circ X = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ \dots \\ a_n \end{bmatrix} \circ \begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots & x_{1n} \\ x_{21} & x_{22} & x_{23} & \dots & x_{2n} \\ x_{31} & x_{32} & x_{33} & \dots & x_{3n} \\ \dots & \dots & \dots & \dots & \dots \\ x_{n1} & x_{n2} & x_{n3} & \dots & x_{nn} \end{bmatrix} \quad (12)$$

$$= \begin{bmatrix} a_1 x_{11} & a_1 x_{12} & a_1 x_{13} & \dots & a_1 x_{1n} \\ a_2 x_{21} & a_2 x_{22} & a_2 x_{23} & \dots & a_2 x_{2n} \\ a_3 x_{31} & a_3 x_{32} & a_3 x_{33} & \dots & a_3 x_{3n} \\ \dots & \dots & \dots & \dots & \dots \\ a_n x_{n1} & a_n x_{n2} & a_n x_{n3} & \dots & a_n x_{nn} \end{bmatrix} \quad (13)$$

当 $\circ$ 作用于两个**矩阵**时，两个矩阵对应位置的元素相乘。按元素乘可以在某些情况下简化矩阵和向量运算。例如，当一个对角矩阵右乘一个矩阵时，相当于用对角矩阵的对角线组成的向量按元素乘那个矩阵：

$$\text{diag}[\mathbf{a}]X = \mathbf{a} \circ X$$

当一个行向量右乘一个对角矩阵时，相当于这个行向量按元素乘那个矩阵对角线组成的向量：

$$\mathbf{a}^T \text{diag}[\mathbf{b}] = \mathbf{a} \circ \mathbf{b}$$

上面这两点，在我们后续推导中会多次用到。

在t时刻，LSTM的输出值为 $\mathbf{h}_t$ 。我们定义t时刻的误差项 $\delta_t$ 为：

$$\delta_t \stackrel{\text{def}}{=} \frac{\partial E}{\partial \mathbf{h}_t}$$

注意，和前面几篇文章不同，我们这里假设误差项是损失函数对输出值的导数，而不是对加权输入 $\text{net}_t^j$ 的导数。因为LSTM有四个加权输入，分别对应 $\mathbf{f}_t$ 、 $\mathbf{i}_t$ 、 $\mathbf{c}_t$ 、 $\mathbf{o}_t$ ，我们希望往上一层传递一个误差项而不是四个。但我们仍然需要定义出这四个加权输入，以及他们对应的误差项。

$$\text{net}_{f,t} = W_f[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_f \quad (14)$$

$$= W_{fh} \mathbf{h}_{t-1} + W_{fx} \mathbf{x}_t + \mathbf{b}_f \quad (15)$$

$$\text{net}_{i,t} = W_i[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_i \quad (16)$$

$$= W_{ih} \mathbf{h}_{t-1} + W_{ix} \mathbf{x}_t + \mathbf{b}_i \quad (17)$$

$$\text{net}_{\tilde{c},t} = W_c[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_c \quad (18)$$

$$= W_{ch} \mathbf{h}_{t-1} + W_{cx} \mathbf{x}_t + \mathbf{b}_c \quad (19)$$

$$\text{net}_{o,t} = W_o[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_o \quad (20)$$

$$= W_{oh} \mathbf{h}_{t-1} + W_{ox} \mathbf{x}_t + \mathbf{b}_o \quad (21)$$

$$\delta_{f,t} \stackrel{\text{def}}{=} \frac{\partial E}{\partial \text{net}_{f,t}} \quad (22)$$

$$\delta_{i,t} \stackrel{\text{def}}{=} \frac{\partial E}{\partial \text{net}_{i,t}} \quad (23)$$

$$\delta_{\tilde{c},t} \stackrel{\text{def}}{=} \frac{\partial E}{\partial \text{net}_{\tilde{c},t}} \quad (24)$$

$$\delta_{o,t} \stackrel{\text{def}}{=} \frac{\partial E}{\partial \text{net}_{o,t}} \quad (25)$$

## 误差项沿时间的反向传递

沿时间反向传递误差项，就是要计算出t-1时刻的误差项 $\delta_{t-1}$ 。

$$\delta_{t-1}^T = \frac{\partial E}{\partial \mathbf{h}_{t-1}} \quad (26)$$

$$= \frac{\partial E}{\partial \mathbf{h}_t} \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \quad (27)$$

$$= \delta_t^T \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \quad (28)$$

我们知道， $\frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}}$ 是一个Jacobian矩阵。如果隐藏层h的维度是N的话，那么它就是一个 $N \times N$ 矩阵。为了求出它，我们列出 $\mathbf{h}_t$ 的计算公式，即前面的式6和式4：

$$\mathbf{h}_t = \mathbf{o}_t \circ \tanh(\mathbf{c}_t) \quad (29)$$

$$\mathbf{c}_t = \mathbf{f}_t \circ \mathbf{c}_{t-1} + \mathbf{i}_t \circ \tilde{\mathbf{c}}_t \quad (30)$$

显然， $\mathbf{o}_t$ 、 $\mathbf{f}_t$ 、 $\mathbf{i}_t$ 、 $\tilde{\mathbf{c}}_t$ 都是 $\mathbf{h}_{t-1}$ 的函数，那么，利用全导数公式可得：

$$\begin{aligned} \delta_t^T \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} &= \delta_t^T \frac{\partial \mathbf{h}_t}{\partial \mathbf{o}_t} \frac{\partial \mathbf{o}_t}{\partial \mathbf{net}_{o,t}} \frac{\partial \mathbf{net}_{o,t}}{\partial \mathbf{h}_{t-1}} + \delta_t^T \frac{\partial \mathbf{h}_t}{\partial \mathbf{c}_t} \frac{\partial \mathbf{c}_t}{\partial \mathbf{f}_t} \frac{\partial \mathbf{f}_t}{\partial \mathbf{net}_{f,t}} \frac{\partial \mathbf{net}_{f,t}}{\partial \mathbf{h}_{t-1}} + \delta_t^T \frac{\partial \mathbf{h}_t}{\partial \mathbf{c}_t} \frac{\partial \mathbf{c}_t}{\partial \mathbf{i}_t} \frac{\partial \mathbf{i}_t}{\partial \mathbf{net}_{i,t}} \frac{\partial \mathbf{net}_{i,t}}{\partial \mathbf{h}_{t-1}} + \delta_t^T \frac{\partial \mathbf{h}_t}{\partial \mathbf{c}_t} \frac{\partial \mathbf{c}_t}{\partial \tilde{\mathbf{c}}_t} \frac{\partial \tilde{\mathbf{c}}_t}{\partial \mathbf{net}_{\tilde{c},t}} \\ &= \delta_{o,t}^T \frac{\partial \mathbf{net}_{o,t}}{\partial \mathbf{h}_{t-1}} + \delta_{f,t}^T \frac{\partial \mathbf{net}_{f,t}}{\partial \mathbf{h}_{t-1}} + \delta_{i,t}^T \frac{\partial \mathbf{net}_{i,t}}{\partial \mathbf{h}_{t-1}} + \delta_{\tilde{c},t}^T \frac{\partial \mathbf{net}_{\tilde{c},t}}{\partial \mathbf{h}_{t-1}} \quad (\text{式7}) \end{aligned}$$

下面，我们要把式7中的每个偏导数都求出来。根据式6，我们可以求出：

$$\frac{\partial \mathbf{h}_t}{\partial \mathbf{o}_t} = \text{diag}[\tanh(\mathbf{c}_t)] \quad (33)$$

$$\frac{\partial \mathbf{h}_t}{\partial \mathbf{c}_t} = \text{diag}[\mathbf{o}_t \circ (1 - \tanh(\mathbf{c}_t)^2)] \quad (34)$$

根据式4，我们可以求出：

$$\frac{\partial \mathbf{c}_t}{\partial \mathbf{f}_t} = \text{diag}[\mathbf{c}_{t-1}] \quad (35)$$

$$\frac{\partial \mathbf{c}_t}{\partial \mathbf{i}_t} = \text{diag}[\tilde{\mathbf{c}}_t] \quad (36)$$

$$\frac{\partial \mathbf{c}_t}{\partial \tilde{\mathbf{c}}_t} = \text{diag}[\mathbf{i}_t] \quad (37)$$

因为：

$$\mathbf{o}_t = \sigma(\mathbf{net}_{o,t}) \quad (38)$$

$$\mathbf{net}_{o,t} = W_{oh} \mathbf{h}_{t-1} + W_{ox} \mathbf{x}_t + \mathbf{b}_o \quad (39)$$

$$\mathbf{f}_t = \sigma(\mathbf{net}_{f,t}) \quad (40)$$

$$\mathbf{net}_{f,t} = W_{fh} \mathbf{h}_{t-1} + W_{fx} \mathbf{x}_t + \mathbf{b}_f \quad (41)$$

$$\mathbf{i}_t = \sigma(\mathbf{net}_{i,t}) \quad (42)$$

$$\mathbf{net}_{i,t} = W_{ih} \mathbf{h}_{t-1} + W_{ix} \mathbf{x}_t + \mathbf{b}_i \quad (43)$$

$$\tilde{\mathbf{c}}_t = \tanh(\mathbf{net}_{\tilde{c},t}) \quad (44)$$

$$\mathbf{net}_{\tilde{c},t} = W_{ch} \mathbf{h}_{t-1} + W_{cx} \mathbf{x}_t + \mathbf{b}_c \quad (45)$$

$$\tilde{\mathbf{c}}_t = \tanh(\mathbf{net}_{\tilde{c},t}) \quad (46)$$

$$\mathbf{net}_{\tilde{c},t} = W_{ch} \mathbf{h}_{t-1} + W_{cx} \mathbf{x}_t + \mathbf{b}_c \quad (47)$$

$$\mathbf{net}_{\tilde{c},t} = W_{ch} \mathbf{h}_{t-1} + W_{cx} \mathbf{x}_t + \mathbf{b}_c \quad (48)$$

我们很容易得出：

$$\frac{\partial \mathbf{o}_t}{\partial \mathbf{net}_{o,t}} = \text{diag}[\mathbf{o}_t \circ (1 - \mathbf{o}_t)] \quad (49)$$

$$\frac{\partial \mathbf{net}_{o,t}}{\partial \mathbf{h}_{t-1}} = W_{oh} \quad (50)$$

$$\frac{\partial \mathbf{f}_t}{\partial \mathbf{net}_{f,t}} = \text{diag}[\mathbf{f}_t \circ (1 - \mathbf{f}_t)] \quad (51)$$

$$\frac{\partial \mathbf{net}_{f,t}}{\partial \mathbf{h}_{t-1}} = W_{fh} \quad (52)$$

$$\frac{\partial \mathbf{i}_t}{\partial \mathbf{net}_{i,t}} = \text{diag}[\mathbf{i}_t \circ (1 - \mathbf{i}_t)] \quad (53)$$

$$\frac{\partial \mathbf{net}_{i,t}}{\partial \mathbf{h}_{t-1}} = W_{ih} \quad (54)$$

$$\frac{\partial \tilde{\mathbf{c}}_t}{\partial \mathbf{net}_{\tilde{c},t}} = \text{diag}[1 - \tilde{\mathbf{c}}_t^2] \quad (55)$$

$$\frac{\partial \mathbf{net}_{\tilde{c},t}}{\partial \mathbf{h}_{t-1}} = W_{ch} \quad (56)$$

将上述偏导数带入到式7，我们得到：



$$\delta_{t-1} = \delta_{o,t}^T \frac{\partial \text{net}_{o,t}}{\partial \mathbf{h}_{t-1}} + \delta_{f,t}^T \frac{\partial \text{net}_{f,t}}{\partial \mathbf{h}_{t-1}} + \delta_{i,t}^T \frac{\partial \text{net}_{i,t}}{\partial \mathbf{h}_{t-1}} + \delta_{\tilde{c},t}^T \frac{\partial \text{net}_{\tilde{c},t}}{\partial \mathbf{h}_{t-1}} \quad (57)$$

$$= \delta_{o,t}^T W_{oh} + \delta_{f,t}^T W_{fh} + \delta_{i,t}^T W_{ih} + \delta_{\tilde{c},t}^T W_{ch} \quad (\text{式8}) \quad (58)$$

根据 $\delta_{o,t}$ 、 $\delta_{f,t}$ 、 $\delta_{i,t}$ 、 $\delta_{\tilde{c},t}$ 的定义，可知：

$$\delta_{o,t}^T = \delta_t^T \circ \tanh(\mathbf{c}_t) \circ \mathbf{o}_t \circ (1 - \mathbf{o}_t) \quad (\text{式9}) \quad (59)$$

$$\delta_{f,t}^T = \delta_t^T \circ \mathbf{o}_t \circ (1 - \tanh(\mathbf{c}_t)^2) \circ \mathbf{c}_{t-1} \circ \mathbf{f}_t \circ (1 - \mathbf{f}_t) \quad (\text{式10}) \quad (60)$$

$$\delta_{i,t}^T = \delta_t^T \circ \mathbf{o}_t \circ (1 - \tanh(\mathbf{c}_t)^2) \circ \tilde{\mathbf{c}}_t \circ \mathbf{i}_t \circ (1 - \mathbf{i}_t) \quad (\text{式11}) \quad (61)$$

$$\delta_{\tilde{c},t}^T = \delta_t^T \circ \mathbf{o}_t \circ (1 - \tanh(\mathbf{c}_t)^2) \circ \mathbf{i}_t \circ (1 - \tilde{\mathbf{c}}^2) \quad (\text{式12}) \quad (62)$$

式8到式12就是将误差沿时间反向传播一个时刻的公式。有了它，我们可以写出将误差项向前传递到任意k时刻的公式：

$$\delta_k^T = \prod_{j=k}^{t-1} \delta_{o,j}^T W_{oh} + \delta_{f,j}^T W_{fh} + \delta_{i,j}^T W_{ih} + \delta_{\tilde{c},j}^T W_{ch} \quad (\text{式13})$$

### 将误差项传递到上一层

我们假设当前为第l层，定义l-1层的误差项是误差函数对l-1层加权输入的导数，即：

$$\delta_t^{l-1} \stackrel{\text{def}}{=} \frac{\partial E}{\partial \text{net}_t^{l-1}}$$

本次LSTM的输入 $\mathbf{x}_t$ 由下面的公式计算：

$$\mathbf{x}_t^l = f^{l-1}(\text{net}_t^{l-1})$$

上式中， $f^{l-1}$ 表示第l-1层的激活函数。

因为 $\text{net}_{f,t}^l$ 、 $\text{net}_{i,t}^l$ 、 $\text{net}_{\tilde{c},t}^l$ 、 $\text{net}_{o,t}^l$ 都是 $\mathbf{x}_t$ 的函数， $\mathbf{x}_t$ 又是 $\text{net}_t^{l-1}$ 的函数，因此，要求出E对 $\text{net}_t^{l-1}$ 的导数，就需要使用全导数公式：

$$\begin{aligned} \frac{\partial E}{\partial \text{net}_t^{l-1}} &= \frac{\partial E}{\partial \text{net}_{f,t}^l} \frac{\partial \text{net}_{f,t}^l}{\partial \mathbf{x}_t^l} \frac{\partial \mathbf{x}_t^l}{\partial \text{net}_t^{l-1}} + \frac{\partial E}{\partial \text{net}_{i,t}^l} \frac{\partial \text{net}_{i,t}^l}{\partial \mathbf{x}_t^l} \frac{\partial \mathbf{x}_t^l}{\partial \text{net}_t^{l-1}} + \frac{\partial E}{\partial \text{net}_{\tilde{c},t}^l} \frac{\partial \text{net}_{\tilde{c},t}^l}{\partial \mathbf{x}_t^l} \frac{\partial \mathbf{x}_t^l}{\partial \text{net}_t^{l-1}} + \frac{\partial E}{\partial \text{net}_{o,t}^l} \frac{\partial \text{net}_{o,t}^l}{\partial \mathbf{x}_t^l} \frac{\partial \mathbf{x}_t^l}{\partial \text{net}_t^{l-1}} \\ &= \delta_{f,t}^T W_{fx} \circ f'(\text{net}_t^{l-1}) + \delta_{i,t}^T W_{ix} \circ f'(\text{net}_t^{l-1}) + \delta_{\tilde{c},t}^T W_{cx} \circ f'(\text{net}_t^{l-1}) + \delta_{o,t}^T W_{ox} \circ f'(\text{net}_t^{l-1}) \\ &= (\delta_{f,t}^T W_{fx} + \delta_{i,t}^T W_{ix} + \delta_{\tilde{c},t}^T W_{cx} + \delta_{o,t}^T W_{ox}) \circ f'(\text{net}_t^{l-1}) \end{aligned} \quad (64)$$

$$= (\delta_{f,t}^T W_{fx} + \delta_{i,t}^T W_{ix} + \delta_{\tilde{c},t}^T W_{cx} + \delta_{o,t}^T W_{ox}) \circ f'(\text{net}_t^{l-1}) \quad (\text{式14}) \quad (65)$$

式14就是将误差传递到上一层的公式。

### 权重梯度的计算

对于 $W_{fh}$ 、 $W_{ih}$ 、 $W_{ch}$ 、 $W_{oh}$ 的权重梯度，我们知道它的梯度是各个时刻梯度之和（证明过程请参考文章[零基础入门深度学习\(5\) - 循环神经网络](#)），我们首先求出它们在t时刻的梯度，然后再求出他们最终的梯度。

我们已经求得了误差项 $\delta_{o,t}$ 、 $\delta_{f,t}$ 、 $\delta_{i,t}$ 、 $\delta_{\tilde{c},t}$ ，很容易求出t时刻的 $W_{oh}$ 、的 $W_{ih}$ 、的 $W_{fh}$ 、的 $W_{ch}$ ：

$$\frac{\partial E}{\partial W_{oh,t}} = \frac{\partial E}{\partial \text{net}_{o,t}} \frac{\partial \text{net}_{o,t}}{\partial W_{oh,t}} \quad (66)$$

$$= \delta_{o,t} \mathbf{h}_{t-1}^T \quad (67)$$

$$= \delta_{o,t} \mathbf{h}_{t-1}^T \quad (68)$$

$$\frac{\partial E}{\partial W_{fh,t}} = \frac{\partial E}{\partial \text{net}_{f,t}} \frac{\partial \text{net}_{f,t}}{\partial W_{fh,t}} \quad (69)$$

$$= \delta_{f,t} \mathbf{h}_{t-1}^T \quad (70)$$

$$= \delta_{f,t} \mathbf{h}_{t-1}^T \quad (71)$$

$$\frac{\partial E}{\partial W_{ih,t}} = \frac{\partial E}{\partial \text{net}_{i,t}} \frac{\partial \text{net}_{i,t}}{\partial W_{ih,t}} \quad (72)$$

$$= \delta_{i,t} \mathbf{h}_{t-1}^T \quad (73)$$

$$= \delta_{i,t} \mathbf{h}_{t-1}^T \quad (74)$$

$$\frac{\partial E}{\partial W_{ch,t}} = \frac{\partial E}{\partial \text{net}_{\tilde{c},t}} \frac{\partial \text{net}_{\tilde{c},t}}{\partial W_{ch,t}} \quad (75)$$

$$= \delta_{\tilde{c},t} \mathbf{h}_{t-1}^T \quad (76)$$

将各个时刻的梯度加在一起，就能得到最终的梯度：



$$\frac{\partial E}{\partial W_{oh}} = \sum_{j=1}^t \delta_{o,j} \mathbf{h}_{j-1}^T \quad (77)$$

$$\frac{\partial E}{\partial W_{fh}} = \sum_{j=1}^t \delta_{f,j} \mathbf{h}_{j-1}^T \quad (78)$$

$$\frac{\partial E}{\partial W_{ih}} = \sum_{j=1}^t \delta_{i,j} \mathbf{h}_{j-1}^T \quad (79)$$

$$\frac{\partial E}{\partial W_{ch}} = \sum_{j=1}^t \delta_{\bar{c},j} \mathbf{h}_{j-1}^T \quad (80)$$

对于偏置项 $\mathbf{b}_f$ 、 $\mathbf{b}_i$ 、 $\mathbf{b}_c$ 、 $\mathbf{b}_o$ 的梯度，也是将各个时刻的梯度加在一起。下面是各个时刻的偏置项梯度：

$$\frac{\partial E}{\partial \mathbf{b}_{o,t}} = \frac{\partial E}{\partial \mathbf{net}_{o,t}} \frac{\partial \mathbf{net}_{o,t}}{\partial \mathbf{b}_{o,t}} \quad (81)$$

$$= \delta_{o,t} \quad (82)$$

$$\quad (83)$$

$$\frac{\partial E}{\partial \mathbf{b}_{f,t}} = \frac{\partial E}{\partial \mathbf{net}_{f,t}} \frac{\partial \mathbf{net}_{f,t}}{\partial \mathbf{b}_{f,t}} \quad (84)$$

$$= \delta_{f,t} \quad (85)$$

$$\quad (86)$$

$$\frac{\partial E}{\partial \mathbf{b}_{i,t}} = \frac{\partial E}{\partial \mathbf{net}_{i,t}} \frac{\partial \mathbf{net}_{i,t}}{\partial \mathbf{b}_{i,t}} \quad (87)$$

$$= \delta_{i,t} \quad (88)$$

$$\quad (89)$$

$$\frac{\partial E}{\partial \mathbf{b}_{c,t}} = \frac{\partial E}{\partial \mathbf{net}_{\bar{c},t}} \frac{\partial \mathbf{net}_{\bar{c},t}}{\partial \mathbf{b}_{c,t}} \quad (90)$$

$$= \delta_{\bar{c},t} \quad (91)$$

下面是最终的偏置项梯度，即将各个时刻的偏置项梯度加在一起：

$$\frac{\partial E}{\partial \mathbf{b}_o} = \sum_{j=1}^t \delta_{o,j} \quad (92)$$

$$\frac{\partial E}{\partial \mathbf{b}_i} = \sum_{j=1}^t \delta_{i,j} \quad (93)$$

$$\frac{\partial E}{\partial \mathbf{b}_f} = \sum_{j=1}^t \delta_{f,j} \quad (94)$$

$$\frac{\partial E}{\partial \mathbf{b}_c} = \sum_{j=1}^t \delta_{\bar{c},j} \quad (95)$$

对于 $W_{fx}$ 、 $W_{ix}$ 、 $W_{cx}$ 、 $W_{ox}$ 的权重梯度，只需要根据相应的误差项直接计算即可：

$$\frac{\partial E}{\partial W_{ox}} = \frac{\partial E}{\partial \mathbf{net}_{o,t}} \frac{\partial \mathbf{net}_{o,t}}{\partial W_{ox}} \quad (96)$$

$$= \delta_{o,t} \mathbf{x}_t^T \quad (97)$$

$$\quad (98)$$

$$\frac{\partial E}{\partial W_{fx}} = \frac{\partial E}{\partial \mathbf{net}_{f,t}} \frac{\partial \mathbf{net}_{f,t}}{\partial W_{fx}} \quad (99)$$

$$= \delta_{f,t} \mathbf{x}_t^T \quad (100)$$

$$\quad (101)$$

$$\frac{\partial E}{\partial W_{ix}} = \frac{\partial E}{\partial \mathbf{net}_{i,t}} \frac{\partial \mathbf{net}_{i,t}}{\partial W_{ix}} \quad (102)$$

$$= \delta_{i,t} \mathbf{x}_t^T \quad (103)$$

$$\quad (104)$$

$$\frac{\partial E}{\partial W_{cx}} = \frac{\partial E}{\partial \mathbf{net}_{\bar{c},t}} \frac{\partial \mathbf{net}_{\bar{c},t}}{\partial W_{cx}} \quad (105)$$

$$= \delta_{\bar{c},t} \mathbf{x}_t^T \quad (106)$$

以上就是LSTM的训练算法的全部公式。因为这里面存在很多重复的模式，仔细看看，会觉得并不是太复杂。

当然，LSTM存在着相当多的变体，读者可以在互联网上找到很多资料。因为大家已经熟悉了基本LSTM的算法，因此理解这些变体比较容易，因此本文就不再赘述了。

## 长短时记忆网络的实现

在下面的实现中，LstmLayer的参数包括输入维度、输出维度、隐藏层维度，单元状态维度等于隐藏层维度。gate的激活函数为sigmoid函数，输出的激活函数为tanh。

### 激活函数的实现

我们先实现两个激活函数：sigmoid和tanh。

```

1. class SigmoidActivator(object):
2.     def forward(self, weighted_input):
3.         return 1.0 / (1.0 + np.exp(-weighted_input))
4.
5.     def backward(self, output):
6.         return output * (1 - output)
7.
8.
9. class TanhActivator(object):
10.    def forward(self, weighted_input):
11.        return 2.0 / (1.0 + np.exp(-2 * weighted_input)) - 1.0
12.
13.    def backward(self, output):
14.        return 1 - output * output

```

### LSTM初始化

和前两篇文章代码架构一样，我们把LSTM的实现放在LstmLayer类中。

根据LSTM前向计算和方向传播算法，我们需要初始化一系列矩阵和向量。这些矩阵和向量有两类用途，一类是用于保存模型参数，例如 $W_f$ 、 $W_i$ 、 $W_o$ 、 $W_c$ 、 $b_f$ 、 $b_i$ 、 $b_o$ 、 $b_c$ ；另一类是保存各种中间计算结果，以便于反向传播算法使用，它们包括 $h_t$ 、 $f_t$ 、 $i_t$ 、 $o_t$ 、 $c_t$ 、 $\tilde{c}_t$ 、 $\delta_t$ 、 $\delta_{f,t}$ 、 $\delta_{i,t}$ 、 $\delta_{o,t}$ 、 $\delta_{\tilde{c},t}$ ，以及各个权重对应的梯度。

在构造函数的初始化中，只初始化了与forward计算相关的变量，与backward相关的变量没有初始化。这是因为构造LSTM对象的时候，我们还不知道它未来是用于训练（既有forward又有backward）还是推理（只有forward）。

```

1. class LstmLayer(object):
2.     def __init__(self, input_width, state_width,
3.                  learning_rate):
4.         self.input_width = input_width
5.         self.state_width = state_width
6.         self.learning_rate = learning_rate
7.         # 门的激活函数
8.         self.gate_activator = SigmoidActivator()
9.         # 输出的激活函数
10.        self.output_activator = TanhActivator()
11.        # 当前时刻初始化为t0
12.        self.times = 0
13.        # 各个时刻的单元状态向量c
14.        self.c_list = self.init_state_vec()
15.        # 各个时刻的输出向量h
16.        self.h_list = self.init_state_vec()
17.        # 各个时刻的遗忘门f
18.        self.f_list = self.init_state_vec()
19.        # 各个时刻的输入门i
20.        self.i_list = self.init_state_vec()
21.        # 各个时刻的输出门o
22.        self.o_list = self.init_state_vec()
23.        # 各个时刻的即时状态c~
24.        self.ct_list = self.init_state_vec()
25.        # 遗忘门权重矩阵Wfh, Wfx, 偏置项bf
26.        self.Wfh, self.Wfx, self.bf = (
27.            self.init_weight_mat())
28.        # 输入门权重矩阵Wih, Wix, 偏置项bi
29.        self.Wih, self.Wix, self.bi = (
30.            self.init_weight_mat())
31.        # 输出门权重矩阵Woh, Wox, 偏置项bo
32.        self.Woh, self.Wox, self.bo = (
33.            self.init_weight_mat())
34.        # 单元状态权重矩阵Wch, Wcx, 偏置项bc
35.        self.Wch, self.Wcx, self.bc = (
36.            self.init_weight_mat())
37.
38.    def init_state_vec(self):
39.        """
40.        初始化保存状态的向量
41.        """
42.        state_vec_list = []
43.        state_vec_list.append(np.zeros(
44.            (self.state_width, 1)))
45.        return state_vec_list
46.
47.    def init_weight_mat(self):
48.        """
49.        初始化权重矩阵
50.        """
51.        Wh = np.random.uniform(-1e-4, 1e-4,
52.                                (self.state_width, self.state_width))
53.        Wx = np.random.uniform(-1e-4, 1e-4,
54.                                (self.state_width, self.input_width))
55.        b = np.zeros((self.state_width, 1))

```

```
56.         return Wh, Wx, b
```

## 前向计算的实现

forward方法实现了LSTM的前向计算：

```
1.     def forward(self, x):
2.         '''
3.         根据式1-式6进行前向计算
4.         '''
5.         self.times += 1
6.         # 遗忘门
7.         fg = self.calc_gate(x, self.Wfx, self.Wfh,
8.                             self.bf, self.gate_activator)
9.         self.f_list.append(fg)
10.        # 输入门
11.        ig = self.calc_gate(x, self.Wix, self.Wih,
12.                            self.bi, self.gate_activator)
13.        self.i_list.append(ig)
14.        # 输出门
15.        og = self.calc_gate(x, self.Wox, self.Woh,
16.                            self.bo, self.gate_activator)
17.        self.o_list.append(og)
18.        # 即时状态
19.        ct = self.calc_gate(x, self.Wcx, self.Wch,
20.                            self.bc, self.output_activator)
21.        self.ct_list.append(ct)
22.        # 单元状态
23.        c = fg * self.c_list[self.times - 1] + ig * ct
24.        self.c_list.append(c)
25.        # 输出
26.        h = og * self.output_activator.forward(c)
27.        self.h_list.append(h)
28.
29.    def calc_gate(self, x, Wx, Wh, b, activator):
30.        '''
31.        计算门
32.        '''
33.        h = self.h_list[self.times - 1] # 上次的LSTM输出
34.        net = np.dot(Wh, h) + np.dot(Wx, x) + b
35.        gate = activator.forward(net)
36.        return gate
```

从上面的代码我们可以看到，门的计算都是相同的算法，而门和 $\tilde{c}_t$ 的计算仅仅是激活函数不同。因此我们提出了calc\_gate方法，这样减少了很多重复代码。

## 反向传播算法的实现

backward方法实现了LSTM的反向传播算法。需要注意的是，与backward相关的内部状态变量是在调用backward方法之后才初始化的。这种延迟初始化的一个好处是，如果LSTM只是用来推理，那么就不需要初始化这些变量，节省了很多内存。

```
1.     def backward(self, x, delta_h, activator):
2.         '''
3.         实现LSTM训练算法
4.         '''
5.         self.calc_delta(delta_h, activator)
6.         self.calc_gradient(x)
```

算法主要分成两个部分，一部分使计算误差项：

```
1.     def calc_delta(self, delta_h, activator):
2.         # 初始化各个时刻的误差项
3.         self.delta_h_list = self.init_delta() # 输出误差项
4.         self.delta_o_list = self.init_delta() # 输出门误差项
5.         self.delta_i_list = self.init_delta() # 输入门误差项
6.         self.delta_f_list = self.init_delta() # 遗忘门误差项
7.         self.delta_ct_list = self.init_delta() # 即时输出误差项
8.
9.         # 保存从上一层传递下来的当前时刻的误差项
10.        self.delta_h_list[-1] = delta_h
11.
12.        # 迭代计算每个时刻的误差项
13.        for k in range(self.times, 0, -1):
14.            self.calc_delta_k(k)
15.
16.    def init_delta(self):
17.        '''
18.        初始化误差项
19.        '''
20.        delta_list = []
21.        for i in range(self.times + 1):
22.            delta_list.append(np.zeros(
23.                (self.state_width, 1)))
24.        return delta_list
25.
26.    def calc_delta_k(self, k):
27.        '''
28.        根据k时刻的delta_h, 计算k时刻的delta_f、
29.        delta_i、delta_o、delta_ct, 以及k-1时刻的delta_h
30.        '''
31.        # 获得k时刻前向计算的值
32.        ig = self.i_list[k]
33.        og = self.o_list[k]
34.        fg = self.f_list[k]
35.        ct = self.ct_list[k]
```

```

36.         c = self.c_list[k]
37.         c_prev = self.c_list[k-1]
38.         tanh_c = self.output_activator.forward(c)
39.         delta_k = self.delta_h_list[k]
40.
41.         # 根据式9计算delta_o
42.         delta_o = (delta_k * tanh_c *
43.                    self.gate_activator.backward(og))
44.         delta_f = (delta_k * og *
45.                    (1 - tanh_c * tanh_c) * c_prev *
46.                    self.gate_activator.backward(fg))
47.         delta_i = (delta_k * og *
48.                    (1 - tanh_c * tanh_c) * ct *
49.                    self.gate_activator.backward(ig))
50.         delta_ct = (delta_k * og *
51.                     (1 - tanh_c * tanh_c) * ig *
52.                     self.output_activator.backward(ct))
53.         delta_h_prev = (
54.             np.dot(delta_o.transpose(), self.Woh) +
55.             np.dot(delta_i.transpose(), self.Wih) +
56.             np.dot(delta_f.transpose(), self.Wfh) +
57.             np.dot(delta_ct.transpose(), self.Wch)
58.         ).transpose()
59.
60.         # 保存全部delta值
61.         self.delta_h_list[k-1] = delta_h_prev
62.         self.delta_f_list[k] = delta_f
63.         self.delta_i_list[k] = delta_i
64.         self.delta_o_list[k] = delta_o
65.         self.delta_ct_list[k] = delta_ct

```

另一部分是计算梯度：

```

1.     def calc_gradient(self, x):
2.         # 初始化遗忘门权重梯度矩阵和偏置项
3.         self.Wfh_grad, self.Wfx_grad, self.bf_grad = (
4.             self.init_weight_gradient_mat())
5.         # 初始化输入门权重梯度矩阵和偏置项
6.         self.Wih_grad, self.Wix_grad, self.bi_grad = (
7.             self.init_weight_gradient_mat())
8.         # 初始化输出门权重梯度矩阵和偏置项
9.         self.Woh_grad, self.Wox_grad, self.bo_grad = (
10.            self.init_weight_gradient_mat())
11.        # 初始化单元状态权重梯度矩阵和偏置项
12.        self.Wch_grad, self.Wcx_grad, self.bc_grad = (
13.            self.init_weight_gradient_mat())
14.
15.        # 计算对上一次输出h的权重梯度
16.        for t in range(self.times, 0, -1):
17.            # 计算各个时刻的梯度
18.            (Wfh_grad, bf_grad,
19.             Wih_grad, bi_grad,
20.             Woh_grad, bo_grad,
21.             Wch_grad, bc_grad) = (
22.                self.calc_gradient_t(t))
23.            # 实际梯度是各时刻梯度之和
24.            self.Wfh_grad += Wfh_grad
25.            self.bf_grad += bf_grad
26.            self.Wih_grad += Wih_grad
27.            self.bi_grad += bi_grad
28.            self.Woh_grad += Woh_grad
29.            self.bo_grad += bo_grad
30.            self.Wch_grad += Wch_grad
31.            self.bc_grad += bc_grad
32.            print '-----%d-----' % t
33.            print Wfh_grad
34.            print self.Wfh_grad
35.
36.            # 计算对本次输入x的权重梯度
37.            xt = x.transpose()
38.            self.Wfx_grad = np.dot(self.delta_f_list[-1], xt)
39.            self.Wix_grad = np.dot(self.delta_i_list[-1], xt)
40.            self.Wox_grad = np.dot(self.delta_o_list[-1], xt)
41.            self.Wcx_grad = np.dot(self.delta_ct_list[-1], xt)
42.
43.        def init_weight_gradient_mat(self):
44.            """
45.            初始化权重矩阵
46.            """
47.            Wh_grad = np.zeros((self.state_width,
48.                                self.state_width))
49.            Wx_grad = np.zeros((self.state_width,
50.                                self.input_width))
51.            b_grad = np.zeros((self.state_width, 1))
52.            return Wh_grad, Wx_grad, b_grad
53.
54.        def calc_gradient_t(self, t):
55.            """
56.            计算每个时刻t权重的梯度
57.            """
58.            h_prev = self.h_list[t-1].transpose()
59.            Wfh_grad = np.dot(self.delta_f_list[t], h_prev)
60.            bf_grad = self.delta_f_list[t]
61.            Wih_grad = np.dot(self.delta_i_list[t], h_prev)
62.            bi_grad = self.delta_f_list[t]
63.            Woh_grad = np.dot(self.delta_o_list[t], h_prev)
64.            bo_grad = self.delta_f_list[t]
65.            Wch_grad = np.dot(self.delta_ct_list[t], h_prev)
66.            bc_grad = self.delta_ct_list[t]
67.            return Wfh_grad, bf_grad, Wih_grad, bi_grad, \
68.                Woh_grad, bo_grad, Wch_grad, bc_grad

```

## 梯度下降算法的实现

下面是用梯度下降算法来更新权重：

```

1. def update(self):
2.     """
3.     按照梯度下降，更新权重
4.     """
5.     self.Wfh -= self.learning_rate * self.Wfh_grad
6.     self.Wfx -= self.learning_rate * self.Whx_grad
7.     self.bf -= self.learning_rate * self.bf_grad
8.     self.Wih -= self.learning_rate * self.Whi_grad
9.     self.Wix -= self.learning_rate * self.Whi_grad
10.    self.bi -= self.learning_rate * self.bi_grad
11.    self.Woh -= self.learning_rate * self.Wof_grad
12.    self.Wox -= self.learning_rate * self.Wox_grad
13.    self.bo -= self.learning_rate * self.bo_grad
14.    self.Wch -= self.learning_rate * self.Wcf_grad
15.    self.Wcx -= self.learning_rate * self.Wcx_grad
16.    self.bc -= self.learning_rate * self.bc_grad

```

## 梯度检查的实现

和RecurrentLayer一样，为了支持梯度检查，我们需要支持重置内部状态：

```

1. def reset_state(self):
2.     # 当前时刻初始化为t0
3.     self.times = 0
4.     # 各个时刻的单元状态向量c
5.     self.c_list = self.init_state_vec()
6.     # 各个时刻的输出向量h
7.     self.h_list = self.init_state_vec()
8.     # 各个时刻的遗忘门f
9.     self.f_list = self.init_state_vec()
10.    # 各个时刻的输入门i
11.    self.i_list = self.init_state_vec()
12.    # 各个时刻的输出门o
13.    self.o_list = self.init_state_vec()
14.    # 各个时刻的即时状态c~
15.    self.ct_list = self.init_state_vec()

```

最后，是梯度检查的代码：

```

1. def data_set():
2.     x = [np.array([[1], [2], [3]]),
3.           np.array([[2], [3], [4]])]
4.     d = np.array([[1], [2]])
5.     return x, d
6.
7. def gradient_check():
8.     """
9.     梯度检查
10.    """
11.    # 设计一个误差函数，取所有节点输出项之和
12.    error_function = lambda o: o.sum()
13.
14.    lstm = LstmLayer(3, 2, 1e-3)
15.
16.    # 计算forward值
17.    x, d = data_set()
18.    lstm.forward(x[0])
19.    lstm.forward(x[1])
20.
21.    # 求取sensitivity map
22.    sensitivity_array = np.ones(lstm.h_list[-1].shape,
23.                                dtype=np.float64)
24.
25.    # 计算梯度
26.    lstm.backward(x[1], sensitivity_array, IdentityActivator())
27.
28.    # 检查梯度
29.    epsilon = 10e-4
30.    for i in range(lstm.Wfh.shape[0]):
31.        for j in range(lstm.Wfh.shape[1]):
32.            lstm.Wfh[i, j] += epsilon
33.            lstm.reset_state()
34.            lstm.forward(x[0])
35.            lstm.forward(x[1])
36.            err1 = error_function(lstm.h_list[-1])
37.            lstm.Wfh[i, j] -= 2*epsilon
38.            lstm.reset_state()
39.            lstm.forward(x[0])
40.            lstm.forward(x[1])
41.            err2 = error_function(lstm.h_list[-1])
42.            expect_grad = (err1 - err2) / (2 * epsilon)
43.            lstm.Wfh[i, j] += epsilon
44.            print 'weights(%d,%d): expected - actual %.4e - %.4e' % (
45.                i, j, expect_grad, lstm.Wfh_grad[i, j])
46.    return lstm

```

我们只对 $W_{fh}$ 做了检查，读者可以自行增加对其他梯度的检查。下面是某次梯度检查的结果：

```
weights(0,0): expected - actual 1.3908e-09 - 1.3909e-09
weights(0,1): expected - actual 1.4017e-09 - 1.4017e-09
weights(1,0): expected - actual 1.4017e-09 - 1.4017e-09
weights(1,1): expected - actual 1.4126e-09 - 1.4126e-09
```

## GRU

前面我们讲了一种普通的LSTM，事实上LSTM存在很多**变体**，许多论文中的LSTM都或多或少的不太一样。在众多的LSTM变体中，**GRU (Gated Recurrent Unit)**也许是最成功的一种。它对LSTM做了很多简化，同时却保持着和LSTM相同的效果。因此，GRU最近变得越来越流行。

GRU对LSTM做了两个大改动：

- 1. 将输入门、遗忘门、输出门变为两个门：更新门（Update Gate） $\mathbf{z}_t$ 和重置门（Reset Gate） $\mathbf{r}_t$ 。
- 2. 将单元状态与输出合并为一个状态： $\mathbf{h}$ 。

GRU的前向计算公式为：

$$\mathbf{z}_t = \sigma(W_z \cdot [\mathbf{h}_{t-1}, \mathbf{x}_t])$$

(107)

$$\mathbf{r}_t = \sigma(W_r \cdot [\mathbf{h}_{t-1}, \mathbf{x}_t])$$

(108)

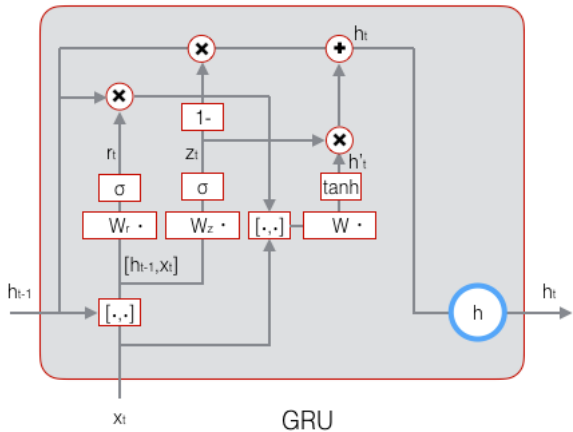
$$\tilde{\mathbf{h}}_t = \tanh(W \cdot [\mathbf{r}_t \circ \mathbf{h}_{t-1}, \mathbf{x}_t])$$

(109)

$$\mathbf{h} = (1 - \mathbf{z}_t) \circ \mathbf{h}_{t-1} + \mathbf{z}_t \circ \tilde{\mathbf{h}}_t$$

(110)

下图是GRU的示意图：



GRU的训练算法比LSTM简单一些，留给读者自行推导，本文就不再赘述了。

## 小结

至此，LSTM——也许是结构最复杂的一类神经网络——就讲完了，相信拿下几篇文章的读者们搞定这篇文章也不在话下吧！现在我们已经了解**循环神经网络**和它最流行的变体——**LSTM**，它们都可以用来处理序列。但是，有时候仅仅拥有处理序列的能力还不够，还需要处理比序列更为复杂的结构（比如树结构），这时候就需要用到另外一类网络：**递归神经网络(Recursive Neural Network)**，巧合的是，它的缩写也是**RNN**。在下一篇文章中，我们将介绍**递归神经网络**和它的训练算法。现在，漫长的烧脑暂告一段落，休息一下吧:)



## 参考资料

1. [CS224d: Deep Learning for Natural Language Processing](#)

2. [Understanding LSTM Networks](#)

3. [LSTM Forward and Backward Pass](#)

• 内容目录

- [零基础入门深度学习\(6\) - 长短时记忆网络\(LSTM\)](#)

■ [文章列表](#)

■ [往期回顾](#)

■ [长短时记忆网络是啥](#)

■ [长短时记忆网络的前向计算](#)

■ [长短时记忆网络的训练](#)

■ [LSTM训练算法框架](#)

■ [关于公式和符号的说明](#)

■ [误差项沿时间的反向传递](#)

■ [将误差项传递到上一层](#)

■ [权重梯度的计算](#)

■ [长短时记忆网络的实现](#)

■ [激活函数的实现](#)

■ [LSTM初始化](#)

■ [前向计算的实现](#)

■ [反向传播算法的实现](#)

■ [梯度下降算法的实现](#)

■ [梯度检查的实现](#)

■ [GRU](#)

■ [小结](#)

■ [参考资料](#)

•

•

○ 

■ [机器学习 7](#)

■ [零基础入门深度学习\(7\) - 递归神经网络](#)

■ [零基础入门深度学习\(6\) - 长短时记忆网络\(LSTM\)](#)

■ [零基础入门深度学习\(5\) - 循环神经网络](#)

■ [零基础入门深度学习\(4\) - 卷积神经网络](#)

■ [零基础入门深度学习\(3\) - 神经网络和反向传播算法](#)

■ [零基础入门深度学习\(2\) - 线性单元和梯度下降](#)

■ [零基础入门深度学习\(1\) - 感知器](#)

○ 

■ [深度学习入门 7](#)

■ [零基础入门深度学习\(7\) - 递归神经网络](#)

■ [零基础入门深度学习\(6\) - 长短时记忆网络\(LSTM\)](#)

■ [零基础入门深度学习\(5\) - 循环神经网络](#)

■ [零基础入门深度学习\(4\) - 卷积神经网络](#)

■ [零基础入门深度学习\(3\) - 神经网络和反向传播算法](#)

■ [零基础入门深度学习\(2\) - 线性单元和梯度下降](#)

■ [零基础入门深度学习\(1\) - 感知器](#)

○ 

搜索

○ 以下【标签】将用于标记这篇文稿：

•

•

•

•

○ [下载客户端](#)

○ [关注开发者](#)

○ [报告问题，建议](#)

○ [联系我们](#)

•

添加新批注

保存

取消

在作者公开此批注前，只有你和作者可见。

保存

取消

修改

保存

取消

删除

• 私有

• 公开

https://zybuluo.com/hanbingtao/note/581764

15/16



- 删除

查看更早的 5 条回复

回复批注



通知

取消 确认

- ☐
- ☐