# Gradient backpropagation through a long short-term memory (LSTM) cell

Mostafa Sadeghi

**Abstract**

In this brief document, I am going to derive the equations for gradient backpropagation through an LSTM cell. Although with the presence of many deep learning softwares there is no longer any need to compute gradients by hand, this would be a nice exercise on the backpropagation method. Note that this document is by no means a tutorial on LSTM; instead, it is only a guide on how to do backpropagation in LSTMs. So, it assumes the reader has a prior knowledge on RNN and LSTM.

**Index Terms**

LSTM, RNN, backpropagation, gradient computation, computational graphs

## I. INTRODUCTION

Long short-term memory (LSTM) networks have been proven to be much more effective than vanilla recurrent neural network (RNN) networks in training over long sequences. This is due to the fact that in contrast to RNNs, in LSTMs, there is a direct flow of the so-called state variables all across the time axis that prevents gradient vanishing or explosion. However, gradient computations for an LSTM cell are more involved than vanilla RNNs.

### A. LSTM: forward pass

The overall structure of an LSTM cell is shown in Fig. 1. "[1]Similar to the vanilla RNN, at each timestep we receive an input $\mathbf{x}_t \in \mathbb{R}^D$ and the previous hidden state $\mathbf{h}_{t-1} \in \mathbb{R}^H$; the LSTM also maintains an $H$-dimensional *cell state*, so we also receive the previous cell state $\mathbf{c}_{t-1} \in \mathbb{R}^H$. The learnable parameters

Electrical Engineering Department, Sharif University of Technology, Tehran, Iran (e-mail: m.saadeghii@gmail.com.)
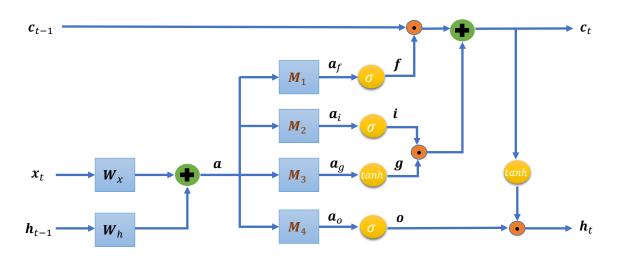
[1]http://cs231n.github.io/

Fig. 1: An LSTM cell. Square blocks denote matrix-vector multiplication, yellow and brown blocks denote elementwise operations, and the green blocks denote summation. Note that the matrices $\mathbf{M}_1, \ldots, \mathbf{M}_4$ are fixed and so, they are not trained. The bias vector has not been shown in this figure, because its gradient computation is very easy. For the definition of each variable, see the text.

of the LSTM are an *input-to-hidden* matrix $\mathbf{W}_x \in \mathbb{R}^{4H \times D}$, a *hidden-to-hidden* matrix $\mathbf{W}_h \in \mathbb{R}^{4H \times H}$ and a *bias vector* $\mathbf{b} \in \mathbb{R}^{4H}$.

At each timestep we first compute an *activation vector* $\mathbf{a} \in \mathbb{R}^{4H}$ as $\mathbf{a} = \mathbf{W}_x \mathbf{x}_t + \mathbf{W}_h \mathbf{h}_{t-1} + \mathbf{b}$. We then divide this into four vectors $\mathbf{a}_i, \mathbf{a}_f, \mathbf{a}_o, \mathbf{a}_g \in \mathbb{R}^H$ where $\mathbf{a}_i$ consists of the first $H$ elements of $\mathbf{a}$, $\mathbf{a}_f$ is the next $H$ elements of $\mathbf{a}$, etc. We then compute the *input gate* $\mathbf{i} \in \mathbb{R}^H$, *forget gate* $\mathbf{f} \in \mathbb{R}^H$, *output gate* $\mathbf{o} \in \mathbb{R}^H$ and *block input* $\mathbf{g} \in \mathbb{R}^H$ as

$$\mathbf{i} = \sigma(\mathbf{a}_i) \qquad \mathbf{f} = \sigma(\mathbf{a}_f) \qquad \mathbf{o} = \sigma(\mathbf{a}_o) \qquad \mathbf{g} = \tanh(\mathbf{a}_g)$$

where $\sigma$ is the sigmoid function and $\tanh$ is the hyperbolic tangent, both applied elementwise. Finally we compute the next cell state $\mathbf{c}_t$ and next hidden state $\mathbf{h}_t$ as

$$\mathbf{c}_t = f \odot \mathbf{c}_{t-1} + \mathbf{i} \odot \mathbf{g} \qquad \mathbf{h}_t = \mathbf{o} \odot \tanh(\mathbf{c}_t)$$

where $\odot$ is the elementwise product of vectors."

*B. LSTM: backward pass*

To make the computations easier, we define the following matrices: $\mathbf{M}_1, \ldots, \mathbf{M}_4$, where $\mathbf{M}_i$ is the $i$-th $H \times 4H$ submatrix of the $4H \times 4H$ identity matrix. For instance,

$$
\mathbf{M}_1 =
\begin{bmatrix}
1 & 0 & 0 & \ldots & 0 & 0 & 0 & \ldots & 0 \\
0 & 1 & 0 & \ldots & 0 & 0 & 0 & \ldots & 0 \\
\vdots & \vdots & \vdots & \ddots & \vdots & 0 & 0 & \ldots & 0 \\
0 & 0 & 0 & \ldots & 1 & 0 & 0 & \ldots & 0
\end{bmatrix}_{H \times 4H}
$$

It is then straightforward to verify the following

$$
\mathbf{a}_i = \mathbf{M}_1 \cdot \mathbf{a} \qquad \mathbf{a}_f = \mathbf{M}_2 \cdot \mathbf{a} \qquad \mathbf{a}_o = \mathbf{M}_3 \cdot \mathbf{a} \qquad \mathbf{a}_g = \mathbf{M}_4 \cdot \mathbf{a}
$$

Note that introducing these matrices is mainly for simplicity of the derivations, and it's not an efficient way for practical implementations. To derive the gradient expressions, we use graph computations by breaking the overall computations in an LSTM cell into simple operations. To this end, let's define the following additional intermediate variables:

$$
\mathbf{i}_g = \mathbf{i} \odot \mathbf{g} \qquad \mathbf{c}'_t = \tanh(\mathbf{c}_t) \qquad \mathbf{f}_c = \mathbf{f} \odot \mathbf{c}_{t-1}
$$

With these definitions, the cell outputs can be written as

$$
\mathbf{c}_t = \mathbf{f}_c + \mathbf{i}_g \qquad\qquad \mathbf{h}_t = \mathbf{o} \odot \mathbf{c}'_t
$$

Now, we are ready for gradient computations! Let's go ahead! Given the output gradients denoted $\mathbf{dh}_t$ and $\mathbf{dc}_t$ coming from the next LSTM cell, we want to compute the gradients of all the variables. In other words, we want to backpropagate the output gradients from the cell outputs all the way to the cell inputs. To do this, we can write (see also Fig. 2)

- $\mathbf{do} = \mathbf{dh}_t \odot \mathbf{c}'_t$                                       (Because $\mathbf{h}_t = \mathbf{o} \odot \mathbf{c}'_t$)
- $\mathbf{dc}'_t = \mathbf{dh}_t \odot \mathbf{o}$                                    (Because $\mathbf{h}_t = \mathbf{o} \odot \mathbf{c}'_t$)
- $\mathbf{di}_g = \mathbf{dc}_t + \dot{\tanh}(\mathbf{c}_t) \odot \mathbf{dc}'_t$     (Gradient comes from two sources to the top summation node; see Fig. 2)
- $\mathbf{df}_c = \mathbf{dc}_t + \dot{\tanh}(\mathbf{c}_t) \odot \mathbf{dc}'_t = \mathbf{dc}_t + (1 - \tanh^2(\mathbf{c}_t)) \odot \mathbf{dc}'_t = \mathbf{dc}_t + (1 - \mathbf{c}'^2_t) \odot \mathbf{dc}'_t$
- $\mathbf{di} = \mathbf{di}_g \odot \mathbf{g}$                                      (Because $\mathbf{i}_g = \mathbf{i} \odot \mathbf{g}$)
- $\mathbf{dg} = \mathbf{di}_g \odot \mathbf{i}$
- $\mathbf{df} = \mathbf{df}_c \odot \mathbf{c}_{t-1}$                               (Because $\mathbf{f}_c = \mathbf{f} \odot \mathbf{c}_{t-1}$)
- $\mathbf{da}_i = \dot{\sigma}(\mathbf{a}_i) \odot \mathbf{di} = \sigma(\mathbf{a}_i) \odot (1 - \sigma(\mathbf{a}_i)) \odot \mathbf{di} = \mathbf{i} \odot (1 - \mathbf{i}) \odot \mathbf{di}$     (Because $\mathbf{i} = \sigma(\mathbf{a}_i)$ and $\dot{\sigma} = \sigma(1 - \sigma)$)
- $\mathbf{da}_f = \dot{\sigma}(\mathbf{a}_f) \odot \mathbf{df} = \mathbf{f} \odot (1 - \mathbf{f}) \odot \mathbf{df}$
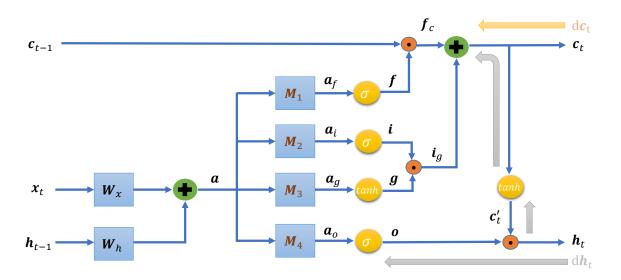
Fig. 2: The same as Fig. 1 but with illustration of gradient flows directions as well as indicating some intermediate variables.

- $d\mathbf{a}_o = \dot{\sigma}(\mathbf{a}_o) \odot d\mathbf{o} = \mathbf{o} \odot (1 - \mathbf{o}) \odot d\mathbf{o}$

- $d\mathbf{a}_g = \dot{\tanh}(\mathbf{a}_g) \odot d\mathbf{g} = (1 - \mathbf{g}^2) \odot d\mathbf{g}$

- $d\mathbf{a} = \mathbf{M}_1^T \cdot d\mathbf{a}_i + \mathbf{M}_2^T \cdot d\mathbf{a}_f + \mathbf{M}_3^T \cdot d\mathbf{a}_o + \mathbf{M}_4^T \cdot d\mathbf{a}_g$ (Because the gradient backpropagates through four paths to reach $\mathbf{a}$; see Fig. 2)

- $d\mathbf{x}_t = \mathbf{W}_x^T \cdot d\mathbf{a}$

- $d\mathbf{h}_{t-1} = \mathbf{W}_h^T \cdot d\mathbf{a}$

- $d\mathbf{W}_x = d\mathbf{a} \cdot \mathbf{x}_t^T$

- $d\mathbf{W}_h = d\mathbf{a} \cdot \mathbf{h}_{t-1}^T$

- $d\mathbf{b} = d\mathbf{a}$  (Bias gradient)

In the above, the power of 2 is applied elementwise. Note that the effect of $\mathbf{M}_1, \ldots, \mathbf{M}_4$ is to put $d\mathbf{a}_i, \ldots, d\mathbf{a}_g$ in their corresponding positions in a $4H$-long vector to create $d\mathbf{a}$. That is, they just concatenate $d\mathbf{a}_i, \ldots, d\mathbf{a}_g$ on top of each other.