

LSTM: From theory to practice

Adrian Mihai Iosif, Matei Macri, Tudor Berariu*

January 2016

1 Introduction

This document is both a resource for understanding the mathematics of LSTM (Section 2) and a tutorial for step by step implementation in Torch (Section 3).

The LSTM module was first introduced in [1] as a solution to the *vanishing gradient* problem that made training vanilla RNNs difficult.

2 Mathematical foundations

Notations and other conventions

In this document the following notations are adopted:

	D	memory cell dimension
	N	input vector size
g_ι	\mathbb{R}	input gate
g_ϕ	\mathbb{R}	forget gate
g_ω	\mathbb{R}	output gate
\mathbf{x}	\mathbb{R}^N	inputs
\mathbf{z}_c	\mathbb{R}^D	input vector (a better name, maybe?)
\mathbf{s}	\mathbb{R}^D	the actual memory of the cell
\mathbf{z}_h	\mathbb{R}^D	output vector
\mathbf{w}_ι	$\mathbb{R}^{(N+2D+1)}$	input gate parameters
\mathbf{w}_ϕ	$\mathbb{R}^{(N+2D+1)}$	forget gate parameters
\mathbf{w}_ω	$\mathbb{R}^{(N+2D+1)}$	output gate parameters
\mathbf{W}_c	$\mathbb{R}^{(N+D+1) \times D}$	input vector parameters
f_*		activation functions (e.g. logistic) - applied element-wise

Vertical concatenation. We use a simplified notation for the vertical concatenation of two vectors \mathbf{a} , \mathbf{b} : $[\mathbf{a}; \mathbf{b}]$ instead of $[\mathbf{a}^T; \mathbf{b}^T]^T$.

Element-wise multiplication. We use \odot to denote element-wise multiplication of equally sized tensors.

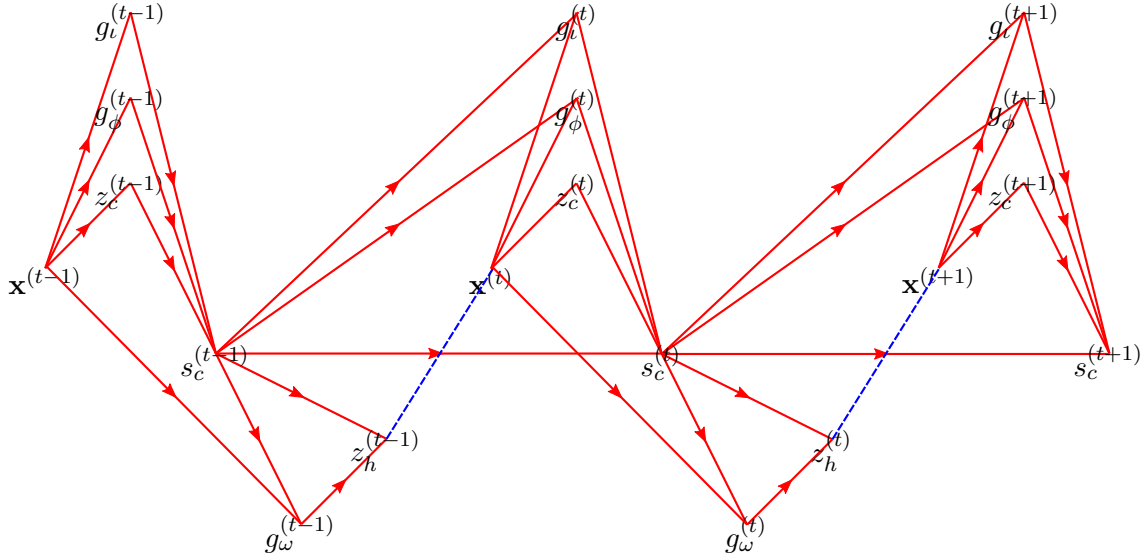
*This author wrote Section 2.

Unidimensional tensors. In this document we use the notation \mathbf{v} for a column vector and \mathbf{v}^T whenever a row vector is needed. However, for various Jacobians we use a notation such as $\mathbf{j} = \frac{\partial \alpha}{\partial \mathbf{v}}$ for a row vector (a $1 \times \text{size}(\mathbf{v})$ Jacobian). Vectors that are not defined as Jacobians are always columns.

2.1 The Forward Phase

In what follows we describe the general formulas used to compute the output vector of the LSTM cell.

The computational graph



Input gate and input value

$$a_l^{(t)} = \underbrace{\mathbf{w}_l^T}_{1 \times (N+2D+1)} \cdot \underbrace{\left[\mathbf{x}^{(t)}; \mathbf{z}_h^{(t-1)}; \mathbf{s}^{(t-1)}; 1 \right]}_{(N+2D+1)} \quad (1)$$

$$g_l^{(t)} = f_l \left(a_l^{(t)} \right) \quad (2)$$

$$\mathbf{a}_c^{(t)} = \underbrace{\mathbf{W}_c^T}_{D \times (N+D+1)} \cdot \underbrace{\left[\mathbf{x}^{(t)}; \mathbf{z}_h^{(t-1)}; 1 \right]}_{(N+D+1)} \quad (3)$$

$$\mathbf{z}_c^{(t)} = f_c \left(\mathbf{a}_c^{(t)} \right) \quad (4)$$

Forget gate

This gate actually acts as a *keep* gate.

$$a_\phi^{(t)} = \mathbf{w}_\phi^T \cdot [\mathbf{x}^{(t)}; \mathbf{z}_h^{(t-1)}; \mathbf{s}^{(t-1)}; 1] \quad (5)$$

$$g_\phi^{(t)} = f_\phi(a_\phi^{(t)}) \quad (6)$$

Cell value

$$\mathbf{s}^{(t)} = g_\phi^{(t)} \mathbf{s}^{(t-1)} + g_\iota^{(t)} \mathbf{z}_c^{(t)} \quad (7)$$

Output gate and output value

$$a_\omega^{(t)} = \mathbf{w}_\omega^T \cdot [\mathbf{x}^{(t)}; \mathbf{z}_h^{(t-1)}; \mathbf{s}^{(t)}; 1] \quad (8)$$

$$g_\omega^{(t)} = f_\omega(a_\omega^{(t)}) \quad (9)$$

$$\mathbf{z}_h^{(t)} = f_h(g_\omega^{(t)} f_s(\mathbf{s}^{(t)})) \quad (10)$$

2.2 The Backward Phase

In this subsection we present the exact form of the partial derivatives of some error function with respect to the parameters of the LSTM cell.

Notations for various partial derivatives

The partial derivatives of the error E with respect to the parameters:

$$\delta_\omega^{(t)} \stackrel{not.}{=} \frac{\partial E}{\partial \mathbf{w}_\omega} \quad \delta_\phi^{(t)} \stackrel{not.}{=} \frac{\partial E}{\partial \mathbf{w}_\phi} \quad \delta_\iota^{(t)} \stackrel{not.}{=} \frac{\partial E}{\partial \mathbf{w}_\iota} \quad (11)$$

$$\Delta_c^{(t)} \stackrel{not.}{=} \frac{\partial E}{\partial \mathbf{W}_c} \quad (12)$$

The partial derivatives of the error E with respect to the gates:

$$\delta_{g_\omega}^{(t)} \stackrel{not.}{=} \frac{\partial E}{\partial g_\omega^{(t)}} \quad \delta_{g_\phi}^{(t)} \stackrel{not.}{=} \frac{\partial E}{\partial g_\phi^{(t)}} \quad \delta_{g_\iota}^{(t)} \stackrel{not.}{=} \frac{\partial E}{\partial g_\iota^{(t)}} \quad (13)$$

The partial derivatives of the error E with respect to $\mathbf{x}^{(t)}$ and with respect to $\mathbf{z}_h^{(t-1)}$:

$$\underbrace{\delta_x^{(t)}}_{1 \times N} \stackrel{not.}{=} \underbrace{\frac{\partial E}{\partial \mathbf{z}_h^{(t)}}}_{1 \times D} \underbrace{\frac{\partial \mathbf{z}_h^{(t)}}{\partial \mathbf{x}^{(t)}}}_{D \times N} \quad (14)$$

$$\underbrace{\delta_{zz}^{(t) \rightarrow (t-1)}}_{1 \times D} \stackrel{not.}{=} \underbrace{\frac{\partial E}{\partial \mathbf{z}_h^{(t)}}}_{1 \times D} \underbrace{\frac{\partial \mathbf{z}_h^{(t)}}{\partial \mathbf{z}_h^{(t-1)}}}_{D \times D} \quad (15)$$

$$\delta_s^{(t)} \stackrel{not.}{=} \frac{\partial E}{\partial \mathbf{s}^{(t)}} \quad \delta_{\mathbf{z}_c}^{(t)} \stackrel{not.}{=} \frac{\partial E}{\partial \mathbf{z}_c^{(t)}} \quad \delta_{\mathbf{z}_h}^{(t)} \stackrel{not.}{=} \frac{\partial E}{\partial \mathbf{z}_h^{(t)}} \quad (16)$$

Inner loops

Before computing the needed gradients, let's take a closer look at $\frac{\partial E}{\partial \mathbf{s}^{(t)}}$. This gradient has two components. The first corresponds to the error flowing through $\mathbf{z}_h^{(t)}$ and the second corresponds to the inner loops of the LSTM (the connections to $g_l^{(t+1)}$, $g_\phi^{(t+1)}$, and $\mathbf{s}^{(t+1)}$).

$$\begin{aligned}
\delta_s^{(t)} \stackrel{not.}{=} \frac{\partial E}{\partial \mathbf{s}^{(t)}} &= \frac{\partial E}{\partial \mathbf{z}_h^{(t)}} \frac{\partial \mathbf{z}_h^{(t)}}{\partial \mathbf{s}^{(t)}} + \frac{\partial E}{\partial g_\phi^{(t+1)}} \frac{\partial g_\phi^{(t+1)}}{\partial \mathbf{s}^{(t)}} + \frac{\partial E}{\partial g_l^{(t+1)}} \frac{\partial g_l^{(t+1)}}{\partial \mathbf{s}^{(t)}} + g_\phi^{(t+1)} \delta_s^{(t+1)T} \\
&= \frac{\partial E}{\partial \mathbf{z}_h^{(t)}} \frac{\partial \mathbf{z}_h^{(t)}}{\partial \left(g_\omega^{(t)} f_s(\mathbf{s}^{(t)}) \right)} \frac{\partial \left(g_\omega^{(t)} f_s(\mathbf{s}^{(t)}) \right)}{\partial \mathbf{s}^{(t)}} \\
&\quad + \frac{\partial E}{\partial g_\phi^{(t+1)}} \frac{\partial g_\phi^{(t+1)}}{\partial \mathbf{s}^{(t)}} + \frac{\partial E}{\partial g_l^{(t+1)}} \frac{\partial g_l^{(t+1)}}{\partial \mathbf{s}^{(t)}} + g_\phi^{(t+1)} \delta_s^{(t+1)T} \\
&= \underbrace{\delta_{\mathbf{z}_h}^{(t)}}_{1 \times D} \underbrace{\text{diag} \left(f'_h \left(g_\omega^{(t)} f_s(\mathbf{s}^{(t)}) \right) \right)}_{D \times D} \underbrace{\left(g_\omega^{(t)} \text{diag} \left(f'_s(\mathbf{s}^{(t)}) \right) + f'_\omega(a_\omega) f_s(\mathbf{s}^{(t)}) \mathbf{w}_{\omega,s}^T \right)}_{D \times D} + \\
&\quad + \delta_{g_\phi}^{(t+1)} f'_\phi(a_\phi^{(t+1)}) \mathbf{w}_{\phi,s}^T + \delta_{g_l}^{(t+1)} f'_l(a_l^{(t+1)}) \mathbf{w}_{l,s}^T + g_\phi^{(t+1)} \delta_s^{(t+1)T} \\
&= \delta_s^{(t) \rightarrow (t)} + \delta_s^{(t+1) \rightarrow (t)}
\end{aligned} \tag{17}$$

It's easier to understand the above matriceal expression by observing an element of the $1 \times D$ row vector $\frac{\partial E}{\partial \mathbf{z}_h^{(t)}} \frac{\partial \mathbf{z}_h^{(t)}}{\partial \mathbf{s}^{(t)}}$:

$$\begin{aligned}
\delta_s^{(t) \rightarrow (t)}[i] &= \sum_{j=1} \frac{\partial E}{\partial \mathbf{z}_h^{(t)}[j]} \frac{\partial \mathbf{z}_h^{(t)}[j]}{\partial \mathbf{s}^{(t)}[i]} \\
&= \sum_{j=1} \left[\frac{\partial E}{\partial \mathbf{z}_h^{(t)}[j]} f'_h \left(g_\omega^{(t)} f_s(\mathbf{s}^{(t)}[j]) \right) \left(g_\omega^{(t)} \delta_{i,j} f'_s(\mathbf{s}^{(t)}[i]) + f_s(\mathbf{s}^{(t)}[j]) f'_\omega(a_\omega^{(t)}) \mathbf{w}_{\omega,s}[i] \right) \right]
\end{aligned} \tag{18}$$

In the above expression $\delta_{i,j}$ is the Kronecker-Delta function.

The objective of the backward phase

Given $\delta_{\mathbf{z}_h}^{(t)}$, and $\delta_s^{(t+1) \rightarrow (t)}$, the following derivatives need to be computed: $\delta_\omega^{(t)}$, $\delta_\phi^{(t)}$, $\delta_l^{(t)}$, $\Delta_c^{(t)}$ (for optimization purposes) $\delta_x^{(t)}$, $\delta_{zz}^{(t) \rightarrow (t-1)}$, and $\delta_s^{(t) \rightarrow (t-1)}$ (in order to compute other gradients back in the architecture).

Order of computation

Gradients need to be computed in the following order: $\delta_{g_\omega}^{(t)}$, $\delta_\omega^{(t)}$, $\delta_s^{(t) \rightarrow (t)}$, $\delta_s^{(t)}$, $\delta_{z_c}^{(t)}$, $\Delta_c^{(t)}$, $\delta_{g_l}^{(t)}$, $\delta_l^{(t)}$, $\delta_{g_\phi}^{(t)}$, $\delta_\phi^{(t)}$, $\delta_s^{(t) \rightarrow (t-1)}$, $\delta_x^{(t)}$, $\delta_{zz}^{(t) \rightarrow (t-1)}$.

The gradients with respect to the output gate and its parameters

$$\delta_{g_\omega}^{(t) \text{ not.}} \frac{\partial E}{\partial g_\omega^{(t)}} = \frac{\partial E}{\partial \mathbf{z}_h^{(t)}} \frac{\partial \mathbf{z}_h^{(t)}}{\partial g_\omega^{(t)}} = \underbrace{\delta_{\mathbf{z}_h}^{(t)}}_{1 \times D} \left(\underbrace{f'_h \left(g_\omega^{(t)} f_s \left(\mathbf{s}^{(t)} \right) \right)}_{D \times 1} \odot \underbrace{f_s \left(\mathbf{s}^{(t)} \right)}_{D \times 1} \right) \quad (19)$$

$$\delta_\omega^{(t) \text{ not.}} \frac{\partial E}{\partial \mathbf{w}_\omega^{(t)}} = \frac{\partial E}{\partial g_\omega^{(t)}} \frac{\partial g_\omega^{(t)}}{\partial \mathbf{w}_\omega^{(t)}} = \delta_{g_\omega}^{(t)} f'_\omega \left(a_\omega^{(t)} \right) \left[\mathbf{x}^{(t)}; \mathbf{z}_h^{(t-1)}; \mathbf{s}^{(t)}; 1 \right] \quad (20)$$

The gradients with respect to the memory cell

$$\delta_s^{(t) \rightarrow (t)} \stackrel{\text{not.}}{=} \frac{\partial E}{\partial \mathbf{z}_h^{(t)}} \frac{\partial \mathbf{z}_h^{(t)}}{\mathbf{s}^{(t)}} \quad (21)$$

$$= \delta_{\mathbf{z}_h}^{(t)} \text{diag} \left(f'_h \left(g_\omega^{(t)} f_s \left(\mathbf{s}^{(t)} \right) \right) \right) \left(g_\omega^{(t)} \text{diag} \left(f'_s \left(\mathbf{s}^{(t)} \right) \right) + f'_\omega \left(a_\omega \right) f_s \left(\mathbf{s}^{(t)} \right) \mathbf{w}_{\omega,s}^T \right) \quad (22)$$

$$= \left(\delta_{\mathbf{z}_h}^{(t)} \odot f'_h \left(g_\omega^{(t)} f_s \left(\mathbf{s}^{(t)} \right) \right)^T \right) \left(g_\omega^{(t)} \text{diag} \left(f'_s \left(\mathbf{s}^{(t)} \right) \right) + f'_\omega \left(a_\omega \right) f_s \left(\mathbf{s}^{(t)} \right) \mathbf{w}_{\omega,s}^T \right) \quad (23)$$

$$\delta_s^{(t)} \stackrel{\text{not.}}{=} \delta_s^{(t) \rightarrow (t)} + \delta_s^{(t+1) \rightarrow (t)} \quad (24)$$

The gradients with respect to the input value and its parameters

$$\delta_{\mathbf{z}_c}^{(t) \text{ not.}} \frac{\partial E}{\partial \mathbf{z}_c^{(t)}} = \frac{\partial E}{\partial \mathbf{s}^{(t)}} \frac{\partial \mathbf{s}^{(t)}}{\partial \mathbf{z}_c^{(t)}} = \delta_s^{(t)} \left(g_i^{(t)} \mathbf{I}_D \right) = g_i^{(t)} \delta_s^{(t)} \quad (25)$$

$$\Delta_c^{(t) \text{ not.}} \frac{\partial E}{\partial \mathbf{W}_c^{(t)}} = \frac{\partial E}{\partial \mathbf{z}_c^{(t)}} \frac{\partial \mathbf{z}_c^{(t)}}{\partial \mathbf{a}_c^{(t)}} \frac{\partial \mathbf{a}_c^{(t)}}{\partial \mathbf{W}_c^{(t)}} = \underbrace{\left(\delta_{\mathbf{z}_c}^{(t)T} \odot f'_c \left(\mathbf{a}_c \right) \right)}_{D \times 1} \underbrace{\left[\mathbf{x}^{(t)T}; \mathbf{s}^{(t-1)T}; 1 \right]}_{1 \times (N+D+1)} \quad (26)$$

The gradients with respect to the input gate and its parameters

$$\delta_{g_l}^{(t) \text{ not.}} \frac{\partial E}{\partial g_l^{(t)}} = \frac{\partial E}{\partial \mathbf{s}^{(t)}} \frac{\partial \mathbf{s}^{(t)}}{\partial g_l^{(t)}} = \delta_s^{(t)} \mathbf{z}_c^{(t)} \quad (27)$$

$$\delta_l^{(t) \text{ not.}} \frac{\partial E}{\partial \mathbf{w}_l^{(t)}} = \frac{\partial E}{\partial g_l^{(t)}} \frac{\partial g_l^{(t)}}{\partial \mathbf{w}_l^{(t)}} = \delta_{g_l}^{(t)} f'_l \left(a_l^{(t)} \right) \left[\mathbf{x}^{(t)T}; \mathbf{z}_h^{(t-1)T}; \mathbf{s}^{(t-1)T}; 1 \right] \quad (28)$$

The gradients with respect to the forget gate and its parameters

$$\delta_{g_\phi}^{(t) \text{ not.}} \frac{\partial E}{\partial g_\phi^{(t)}} = \frac{\partial E}{\partial \mathbf{s}^{(t)}} \frac{\partial \mathbf{s}^{(t)}}{\partial g_\phi^{(t)}} = \delta_s^{(t)} \mathbf{s}^{(t-1)} \quad (29)$$

$$\delta_\phi^{(t) \text{ not.}} \frac{\partial E}{\partial \mathbf{w}_\phi^{(t)}} = \frac{\partial E}{\partial g_\phi^{(t)}} \frac{\partial g_\phi^{(t)}}{\partial \mathbf{w}_\phi^{(t)}} = \delta_{g_\phi}^{(t)} f'_\phi \left(a_\phi^{(t)} \right) \left[\mathbf{x}^{(t)T}; \mathbf{z}_h^{(t-1)T}; \mathbf{s}^{(t-1)T}; 1 \right] \quad (30)$$

The gradients with respect to incoming values

$$\delta_s^{(t) \rightarrow (t-1)} = \delta_{g_\phi}^{(t)} f'_\phi \left(a_\phi^{(t)} \right) \mathbf{w}_{\phi,s}^T + \delta_{g_\iota}^{(t)} f'_\iota \left(a_\iota^{(t)} \right) \mathbf{w}_{\iota,s}^T + g_\phi^{(t)} \delta_s^{(t)T} \quad (31)$$

$$\begin{aligned} \delta_x^{(t)} &\stackrel{\text{not.}}{=} \frac{\partial E}{\partial \mathbf{z}_h^{(t)}} \frac{\partial \mathbf{z}_h^{(t)}}{\partial \mathbf{x}^{(t)}} \\ &= \frac{\partial E}{\partial \mathbf{z}_h^{(t)}} \left(\frac{\partial \mathbf{z}_h^{(t)}}{\partial g_\iota^{(t)}} \frac{\partial g_\iota^{(t)}}{\partial \mathbf{x}^{(t)}} + \frac{\partial \mathbf{z}_h^{(t)}}{\partial g_\phi^{(t)}} \frac{\partial g_\phi^{(t)}}{\partial \mathbf{x}^{(t)}} + \frac{\partial \mathbf{z}_h^{(t)}}{\partial g_\omega^{(t)}} \frac{\partial g_\omega^{(t)}}{\partial \mathbf{x}^{(t)}} + \frac{\partial \mathbf{z}_h^{(t)}}{\partial \mathbf{z}_c^{(t)}} \frac{\partial \mathbf{z}_c^{(t)}}{\partial \mathbf{x}^{(t)}} \right) \\ &= \delta_{g_\iota}^{(t)} f'_\iota \left(a_\iota^{(t)} \right) \mathbf{w}_{\iota,x}^T + \delta_{g_\phi}^{(t)} f'_\phi \left(a_\phi^{(t)} \right) \mathbf{w}_{\phi,x}^T + \delta_{g_\omega}^{(t)} f'_\omega \left(a_\omega^{(t)} \right) \mathbf{w}_{\omega,x}^T + \underbrace{\left(\delta_{\mathbf{z}_c}^{(t)} \odot f'_c \left(\mathbf{a}_c^{(t)} \right)^T \right)}_{1 \times D} \underbrace{\mathbf{W}_{c,x}^T}_{D \times N} \end{aligned} \quad (32)$$

$$\begin{aligned} \delta_{zz}^{(t) \rightarrow (t-1)} &\stackrel{\text{not.}}{=} \frac{\partial E}{\partial \mathbf{z}_h^{(t)}} \frac{\partial \mathbf{z}_h^{(t)}}{\partial \mathbf{z}_h^{(t-1)}} \\ &= \frac{\partial E}{\partial \mathbf{z}_h^{(t)}} \left(\frac{\partial \mathbf{z}_h^{(t)}}{\partial g_\iota^{(t)}} \frac{\partial g_\iota^{(t)}}{\partial \mathbf{z}_h^{(t-1)}} + \frac{\partial \mathbf{z}_h^{(t)}}{\partial g_\phi^{(t)}} \frac{\partial g_\phi^{(t)}}{\partial \mathbf{z}_h^{(t-1)}} + \frac{\partial \mathbf{z}_h^{(t)}}{\partial g_\omega^{(t)}} \frac{\partial g_\omega^{(t)}}{\partial \mathbf{z}_h^{(t-1)}} + \frac{\partial \mathbf{z}_h^{(t)}}{\partial \mathbf{z}_c^{(t)}} \frac{\partial \mathbf{z}_c^{(t)}}{\partial \mathbf{z}_h^{(t-1)}} \right) \\ &= \delta_{g_\iota}^{(t)} f'_\iota \left(a_\iota^{(t)} \right) \mathbf{w}_{\iota,z}^T + \delta_{g_\phi}^{(t)} f'_\phi \left(a_\phi^{(t)} \right) \mathbf{w}_{\phi,z}^T + \delta_{g_\omega}^{(t)} f'_\omega \left(a_\omega^{(t)} \right) \mathbf{w}_{\omega,z}^T + \\ &\quad + \underbrace{\left(\delta_{\mathbf{z}_c}^{(t)} \odot f'_c \left(\mathbf{a}_c^{(t)} \right)^T \right)}_{1 \times D} \underbrace{\mathbf{W}_{c,z}^T}_{D \times D} \end{aligned} \quad (33)$$

3 Torch Implementation

TODO

References

- [1] Sepp Hochreiter and Jürgen Schmidhuber. “Long short-term memory”. In: *Neural computation* 9.8 (1997), pp. 1735–1780.