

# 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 ??) and a tutorial for step by step implementation in Torch (Section ??).

The LSTM module was first introduced in [hochreiter1997long] 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
$g_\iota$	$\mathbb{R}$	input gate
$g_\phi$	$\mathbb{R}$	forget gate
$g_\omega$	$\mathbb{R}$	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
$\mathbf{w}_\iota$	$\mathbb{R}^{(N+2D+1)}$	input gate parameters
$\mathbf{w}_\phi$	$\mathbb{R}^{(N+2D+1)}$	forget gate parameters
$\mathbf{w}_\omega$	$\mathbb{R}^{(N+2D+1)}$	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.

---

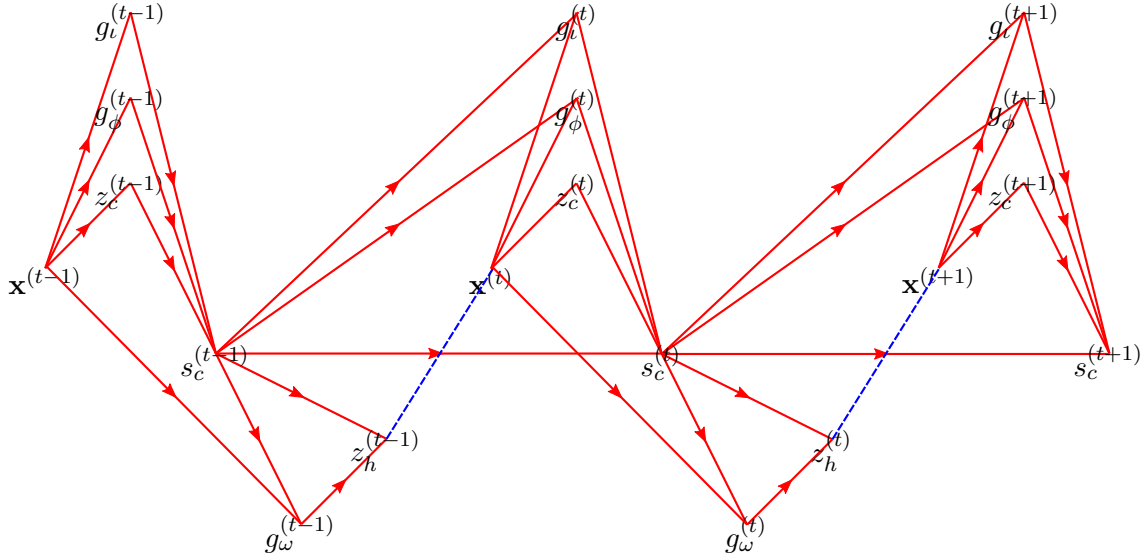
\*This author wrote Section ??.

**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

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

$$g_t^{(t)} = f_t \left( a_t^{(t)} \right) \quad (2)$$

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

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

### 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**

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

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

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

## 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  $\mathbf{z}_h^{(t-1)}$ :

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

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

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

## 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  $g_\iota^{(t+1)}$ ,  $g_\phi^{(t+1)}$ , and  $\mathbf{s}^{(t+1)}$ ).

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

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

## 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:  $\delta_\omega^{(t)}$ ,  $\delta_\phi^{(t)}$ ,  $\delta_\iota^{(t)}$ ,  $\Delta_c^{(t)}$  (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:  $\delta_{g_\omega}^{(t)}$ ,  $\delta_\omega^{(t)}$ ,  $\delta_s^{(t) \rightarrow (t)}$ ,  $\delta_s^{(t)}$ ,  $\delta_{z_c}^{(t)}$ ,  $\Delta_c^{(t)}$ ,  $\delta_{g_\iota}^{(t)}$ ,  $\delta_\iota^{(t)}$ ,  $\delta_{g_\phi}^{(t)}$ ,  $\delta_\phi^{(t)}$ ,  $\delta_s^{(t) \rightarrow (t-1)}$ ,  $\delta_x^{(t)}$ ,  $\delta_{zz}^{(t) \rightarrow (t-1)}$ .

## The gradients with respect to the output gate and its parameters

$$\delta_{g_\omega}^{(t)} \stackrel{\text{not.}}{=} \frac{\partial E}{\partial g_\omega^{(t)}} = \frac{\partial E}{\partial \mathbf{z}_h^{(t)}} \frac{\partial \mathbf{z}_h^{(t)}}{\partial g_\omega^{(t)}} = \delta_h^{(t)} \left( f'_h \left( g_\omega^{(t)} \mathbf{s}^{(t)} \right) \odot \mathbf{s}^{(t)} \right) \tag{19}$$

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

The gradients with respect to the memory cell

$$\delta_s^{(t) \rightarrow (t)} \stackrel{\text{not.}}{=} \frac{\partial E}{\partial \mathbf{z}_h^{(t)}} \frac{\partial \mathbf{z}_h^{(t)}}{\partial \mathbf{s}^{(t)}} = \delta_h^{(t)} \text{diag} \left( f'_\omega \left( g_\omega^{(t)} \mathbf{s}^{(t)} \right) \right) g_\omega^{(t)} \quad (21)$$

$$= g_\omega^{(t)} \delta_h^{(t)} \odot f'_\omega \left( g_\omega^{(t)} \mathbf{s}^{(t)} \right)^\text{T} \quad (22)$$

$$\delta_s^{(t)} \stackrel{\text{not.}}{=} \delta_s^{(t) \rightarrow (t)} + \delta_s^{(t+1) \rightarrow (t)} \quad (23)$$

The gradients with respect to the input value and its parameters

$$\delta_{\mathbf{z}_c}^{(t)} \stackrel{\text{not.}}{=} \frac{\partial E}{\partial \mathbf{z}_c^{(t)}} = \frac{\partial E}{\partial \mathbf{s}^{(t)}} \frac{\partial \mathbf{s}^{(t)}}{\partial \mathbf{z}_c^{(t)}} = \delta_s^{(t)} \left( g_i^{(t)} \mathbf{I}_D \right) = g_i^{(t)} \delta_s^{(t)} \quad (24)$$

$$\Delta_c^{(t)} \stackrel{\text{not.}}{=} \frac{\partial E}{\partial \mathbf{W}_c^{(t)}} = \frac{\partial E}{\partial \mathbf{z}_c^{(t)}} \frac{\partial \mathbf{z}_c^{(t)}}{\partial \mathbf{a}_c^{(t)}} \frac{\partial \mathbf{a}_c^{(t)}}{\partial \mathbf{W}_c^{(t)}} = \underbrace{\left( \delta_{\mathbf{z}_c}^{(t)\text{T}} \odot f'_c(\mathbf{a}_c) \right)}_{D \times 1} \underbrace{\left[ \mathbf{x}^{(t)\text{T}}; \mathbf{s}^{(t-1)\text{T}}; 1 \right]}_{1 \times (N+D+1)} \quad (25)$$

The gradients with respect to the input gate and its parameters

$$\delta_{g_i}^{(t)} \stackrel{\text{not.}}{=} \frac{\partial E}{\partial g_i^{(t)}} = \frac{\partial E}{\partial \mathbf{s}^{(t)}} \frac{\partial \mathbf{s}^{(t)}}{\partial g_i^{(t)}} = \delta_s^{(t)} \mathbf{z}_c^{(t)} \quad (26)$$

$$\delta_{\mathbf{z}_c}^{(t)} \stackrel{\text{not.}}{=} \frac{\partial E}{\partial \mathbf{w}_i^{(t)}} = \frac{\partial E}{\partial g_i^{(t)}} \frac{\partial g_i^{(t)}}{\partial \mathbf{w}_i^{(t)}} = \delta_{g_i}^{(t)} f'_i \left( a_i^{(t)} \right) \left[ \mathbf{x}^{(t)\text{T}}; \mathbf{z}_h^{(t-1)\text{T}}; \mathbf{s}^{(t-1)\text{T}}; 1 \right] \quad (27)$$

The gradients with respect to the forget gate and its parameters

$$\delta_{g_\phi}^{(t)} \stackrel{\text{not.}}{=} \frac{\partial E}{\partial g_\phi^{(t)}} = \frac{\partial E}{\partial \mathbf{s}^{(t)}} \frac{\partial \mathbf{s}^{(t)}}{\partial g_\phi^{(t)}} = \delta_s^{(t)} \mathbf{s}^{(t-1)} \quad (28)$$

$$\delta_\phi^{(t)} \stackrel{\text{not.}}{=} \frac{\partial E}{\partial \mathbf{w}_\phi^{(t)}} = \frac{\partial E}{\partial g_\phi^{(t)}} \frac{\partial g_\phi^{(t)}}{\partial \mathbf{w}_\phi^{(t)}} = \delta_{g_\phi}^{(t)} f'_\phi \left( a_\phi^{(t)} \right) \left[ \mathbf{x}^{(t)\text{T}}; \mathbf{z}_h^{(t-1)\text{T}}; \mathbf{s}^{(t-1)\text{T}}; 1 \right] \quad (29)$$

The gradients with respect to incoming values

$$\delta_s^{(t) \rightarrow (t-1)} = \delta_{g_\phi}^{(t)} f'_\phi \left( a_\phi^{(t)} \right) \mathbf{w}_{\phi,s}^\text{T} + \delta_{g_i}^{(t)} f'_i \left( a_i^{(t)} \right) \mathbf{w}_{i,s}^\text{T} + g_\phi^{(t)} \delta_s^{(t)\text{T}} \quad (30)$$

$$\begin{aligned} \delta_x^{(t)} &\stackrel{\text{not.}}{=} \frac{\partial E}{\partial \mathbf{z}_h^{(t)}} \frac{\partial \mathbf{z}_h^{(t)}}{\partial \mathbf{x}^{(t)}} \\ &= \frac{\partial E}{\partial \mathbf{z}_h^{(t)}} \left( \frac{\partial \mathbf{z}_h^{(t)}}{\partial g_i^{(t)}} \frac{\partial g_i^{(t)}}{\partial \mathbf{x}^{(t)}} + \frac{\partial \mathbf{z}_h^{(t)}}{\partial g_\phi^{(t)}} \frac{\partial g_\phi^{(t)}}{\partial \mathbf{x}^{(t)}} + \frac{\partial \mathbf{z}_h^{(t)}}{\partial g_\omega^{(t)}} \frac{\partial g_\omega^{(t)}}{\partial \mathbf{x}^{(t)}} + \frac{\partial \mathbf{z}_h^{(t)}}{\partial \mathbf{z}_c^{(t)}} \frac{\partial \mathbf{z}_c^{(t)}}{\partial \mathbf{x}^{(t)}} \right) \\ &= \delta_{g_i}^{(t)} f'_i \left( a_i^{(t)} \right) \mathbf{w}_{i,x}^\text{T} + \delta_{g_\phi}^{(t)} f'_\phi \left( a_\phi^{(t)} \right) \mathbf{w}_{\phi,x}^\text{T} + \delta_{g_\omega}^{(t)} f'_\omega \left( a_\omega^{(t)} \right) \mathbf{w}_{\omega,x}^\text{T} + \underbrace{\left( \delta_{\mathbf{z}_c}^{(t)} \odot f'_c \left( \mathbf{a}_c^{(t)} \right)^\text{T} \right)}_{1 \times D} \underbrace{\mathbf{W}_{c,x}^\text{T}}_{D \times N} \end{aligned} \quad (31)$$

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

### 3 Torch Implementation

TODO