

LSTM: From theory to practice

Adrian Mihai Iosif, Matei Macri, Tudor Berariu

January 2016

1 Introduction

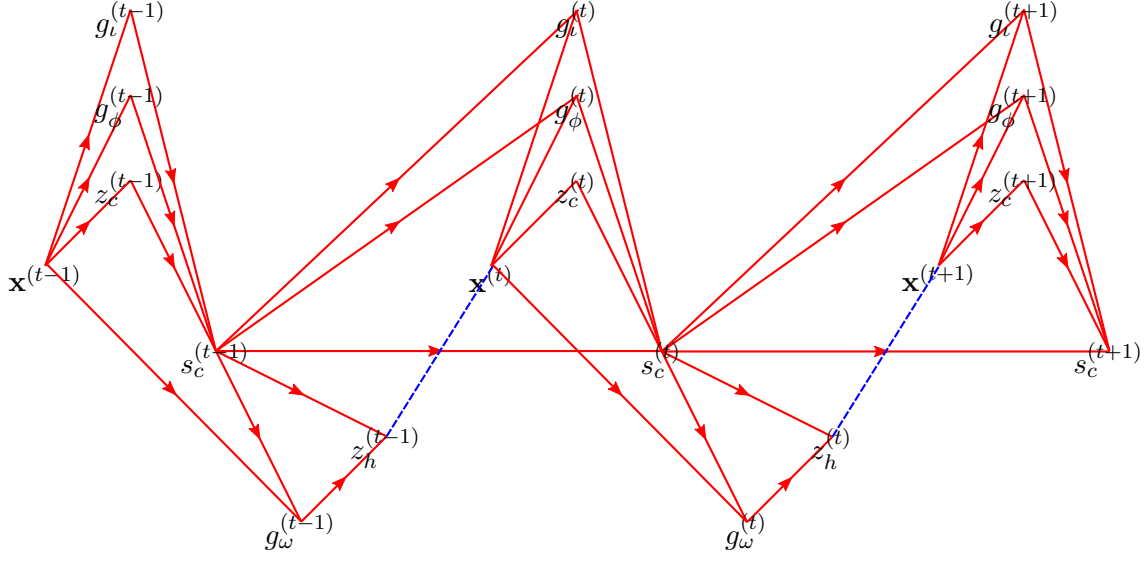
This document is a both a resource for understanding the mathematics of LSTM (Section 2) and a tutorial containing a step by step implementation in Torch (Section 3). The LSTM module was first introduced in [1]. It solves the *vanishing gradient* problem that makes training vanilla RNNs difficult.

2 Mathematical foundations

Conventions. We propose the following notation for the variables:

- g_* - gates
- z_c / z_h - input / output values
- \mathbf{w}_* - weights (parameters)
- f_* - activation functions (e.g. logistic)
- \mathbf{x} - inputs
- s - a scalar value (the actual memory of the LSTM cell)

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



2.1 The Forward Phase

Input gate and input value.

$$a_i^{(t)} = \mathbf{w}_i^T \cdot [\mathbf{x}^{(t)}; z_h^{(t-1)}; s^{(t-1)}; 1] \quad (1)$$

$$g_i^{(t)} = f_i(a_i^{(t)}) \quad (2)$$

$$a_c^{(t)} = \mathbf{w}_c^T \cdot [\mathbf{x}^{(t)}; z_h^{(t-1)}; 1] \quad (3)$$

$$z_c^{(t)} = f_c(a_c^{(t)}) \quad (4)$$

Forget gate. This is actually a *keep* gate:

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

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

Cell value.

$$s^{(t)} = g_\phi^{(t)} s^{(t-1)} + g_i^{(t)} z_c^{(t)} \quad (7)$$

Output gate and output value.

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

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

$$z_h^{(t)} = f_h(g_\omega^{(t)} s^{(t)}) \quad (10)$$

2.2 The Backward Phase

Notations. In what follows the following notations are used for various partial derivatives:

$$\begin{aligned}
\delta_x^{(t)} &\stackrel{\text{not.}}{=} \frac{\partial E}{\partial z_h^{(t)}} \frac{\partial z_h^{(t)}}{\partial \mathbf{x}^{(t)}}; & \delta_{zz}^{(t) \rightarrow (t-1)} &\stackrel{\text{not.}}{=} \frac{\partial E}{\partial z_h^{(t)}} \frac{\partial z_h^{(t)}}{\partial z_h^{(t-1)}}; \\
\delta_s^{(t)} &\stackrel{\text{not.}}{=} \frac{\partial E}{\partial s^{(t)}}; & \delta_c^{(t)} &\stackrel{\text{not.}}{=} \frac{\partial E}{\partial z_c^{(t)}}; & \delta_h^{(t)} &\stackrel{\text{not.}}{=} \frac{\partial E}{\partial z_h^{(t)}}; \\
\delta_{g_\omega}^{(t)} &\stackrel{\text{not.}}{=} \frac{\partial E}{\partial g_\omega^{(t)}}; & \delta_{g_\phi}^{(t)} &\stackrel{\text{not.}}{=} \frac{\partial E}{\partial g_\phi^{(t)}}; & \delta_{g_l}^{(t)} &\stackrel{\text{not.}}{=} \frac{\partial E}{\partial g_l^{(t)}} \\
\delta_\omega^{(t)} &\stackrel{\text{not.}}{=} \frac{\partial E}{\partial \mathbf{w}_\omega}; & \delta_\phi^{(t)} &\stackrel{\text{not.}}{=} \frac{\partial E}{\partial \mathbf{w}_\phi}; & \delta_l^{(t)} &\stackrel{\text{not.}}{=} \frac{\partial E}{\partial \mathbf{w}_l}
\end{aligned}$$

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

$$\begin{aligned}
\delta_s^{(t)} &\stackrel{\text{not.}}{=} \frac{\partial E}{\partial s^{(t)}} = \frac{\partial E}{\partial z_h^{(t)}} \frac{\partial z_h^{(t)}}{\partial s^{(t)}} + \frac{\partial E}{\partial g_\phi^{(t+1)}} \frac{\partial g_\phi^{(t+1)}}{\partial s^{(t)}} + \frac{\partial E}{\partial g_l^{(t+1)}} \frac{\partial g_l^{(t+1)}}{\partial s^{(t)}} + \frac{\partial E}{\partial s^{(t+1)}} \frac{\partial s^{(t+1)}}{\partial s^{(t)}} \\
&= \delta_h^{(t)} f'_\omega \left(g_\omega^{(t)} s^{(t)} \right) g_\omega^{(t)} + \delta_{g_\phi}^{(t+1)} f'_\phi \left(a_\phi^{(t+1)} \right) w_{\phi,s} + \delta_{g_l}^{(t+1)} f'_l \left(a_l^{(t+1)} \right) w_{l,s} + \delta_s^{(t+1)} g_\phi^{(t+1)} \\
&= \delta_s^{(t) \rightarrow (t)} + \delta_s^{(t+1) \rightarrow (t)}
\end{aligned} \tag{11}$$

Goal. Given $\delta_h^{(t)}$, and $\delta_s^{(t+1) \rightarrow (t)}$, the following derivatives need to be computed: $\frac{\partial E}{\partial W_*}$, $\delta_x^{(t)}$, $\delta_{zz}^{(t) \rightarrow (t-1)}$, and $\delta_s^{(t) \rightarrow (t-1)}$.

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

$$\delta_{g_\omega}^{(t)} \stackrel{\text{not.}}{=} \frac{\partial E}{\partial g_\omega^{(t)}} = \frac{\partial E}{\partial z_h^{(t)}} \frac{\partial z_h^{(t)}}{\partial g_\omega^{(t)}} = \delta_h^{(t)} f'_h \left(g_\omega^{(t)} s^{(t)} \right) s^{(t)} \quad (12)$$

$$\delta_\omega^{(t)} \stackrel{\text{not.}}{=} \frac{\partial E}{\partial \mathbf{w}_\omega^{(t)}} = \frac{\partial E}{\partial g_\omega^{(t)}} \frac{\partial g_\omega^{(t)}}{\partial \mathbf{w}_\omega^{(t)}} = \delta_{g_\omega}^{(t)} f'_\omega \left(a_\omega^{(t)} \right) \cdot \left[\mathbf{x}^{(t)}; z_h^{(t-1)}; s^{(t)}; 1 \right] \quad (13)$$

$$\delta_s^{(t) \rightarrow (t)} \stackrel{\text{not.}}{=} \frac{\partial E}{\partial z_h^{(t)}} \frac{\partial z_h^{(t)}}{s^{(t)}} = \delta_h^{(t)} f'_\omega \left(g_\omega^{(t)} s^{(t)} \right) g_\omega^{(t)} \quad (14)$$

$$\delta_s^{(t)} \stackrel{\text{not.}}{=} \delta_s^{(t) \rightarrow (t)} + \delta_s^{(t+1) \rightarrow (t)} \quad (15)$$

$$\delta_c^{(t)} \stackrel{\text{not.}}{=} \frac{\partial E}{\partial z_c^{(t)}} = \frac{\partial E}{\partial s^{(t)}} \frac{\partial s^{(t)}}{z_c^{(t)}} = \delta_s^{(t)} g_i^{(t)} \quad (16)$$

$$\delta_{g_i}^{(t)} \stackrel{\text{not.}}{=} \frac{\partial E}{\partial g_i^{(t)}} = \frac{\partial E}{\partial s^{(t)}} \frac{\partial s^{(t)}}{g_i^{(t)}} = \delta_s^{(t)} z_c^{(t)} \quad (17)$$

$$\delta_l^{(t)} \stackrel{\text{not.}}{=} \frac{\partial E}{\partial \mathbf{w}_l^{(t)}} = \frac{\partial E}{\partial g_l^{(t)}} \frac{\partial g_l^{(t)}}{\partial \mathbf{w}_l^{(t)}} = \delta_{g_\phi}^{(t)} f'_l \left(a_l^{(t)} \right) \cdot \left[\mathbf{x}^{(t)}; z_h^{(t-1)}; s^{(t-1)}; 1 \right] \quad (18)$$

$$\delta_{g_\phi}^{(t)} \stackrel{\text{not.}}{=} \frac{\partial E}{\partial g_\phi^{(t)}} = \frac{\partial E}{\partial s^{(t)}} \frac{\partial s^{(t)}}{g_\phi^{(t)}} = \delta_s^{(t)} s^{(t-1)} \quad (19)$$

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

$$\delta_s^{(t) \rightarrow (t-1)} = \delta_{g_\phi}^{(t)} f'_\phi \left(a_\phi^{(t)} \right) w_{\phi,s} + \delta_{g_l}^{(t)} f'_l \left(a_l^{(t)} \right) w_{l,s} + \delta_s^{(t)} g_\phi^{(t)} \quad (21)$$

$$\begin{aligned} \delta_x^{(t)} &\stackrel{\text{not.}}{=} \frac{\partial E}{\partial z_h^{(t)}} \frac{\partial z_h^{(t)}}{\partial \mathbf{x}^{(t)}} \\ &= \frac{\partial E}{\partial z_h^{(t)}} \left(\frac{\partial z_h^{(t)}}{\partial g_l^{(t)}} \frac{\partial g_l^{(t)}}{\partial \mathbf{x}^{(t)}} + \frac{\partial z_h^{(t)}}{\partial g_\phi^{(t)}} \frac{\partial g_\phi^{(t)}}{\partial \mathbf{x}^{(t)}} + \frac{\partial z_h^{(t)}}{\partial g_\omega^{(t)}} \frac{\partial g_\omega^{(t)}}{\partial \mathbf{x}^{(t)}} + \frac{\partial z_h^{(t)}}{\partial z_c^{(t)}} \frac{\partial z_c^{(t)}}{\partial \mathbf{x}^{(t)}} \right) \\ &= \delta_{g_l}^{(t)} f'_l \left(a_l^{(t)} \right) \mathbf{w}_{l,x} + \delta_{g_\phi}^{(t)} f'_\phi \left(a_\phi^{(t)} \right) \mathbf{w}_{\phi,x} + \delta_{g_\omega}^{(t)} f'_\omega \left(a_\omega^{(t)} \right) \mathbf{w}_{\omega,x} + \delta_c^{(t)} f'_c \left(a_c^{(t)} \right) \mathbf{w}_{c,x} \end{aligned} \quad (22)$$

$$\begin{aligned} \delta_x^{(t)} &\stackrel{\text{not.}}{=} \frac{\partial E}{\partial z_h^{(t)}} \frac{\partial z_h^{(t)}}{\partial z_h^{(t-1)}} \\ &= \frac{\partial E}{\partial z_h^{(t)}} \left(\frac{\partial z_h^{(t)}}{\partial g_l^{(t)}} \frac{\partial g_l^{(t)}}{\partial z_h^{(t-1)}} + \frac{\partial z_h^{(t)}}{\partial g_\phi^{(t)}} \frac{\partial g_\phi^{(t)}}{\partial z_h^{(t-1)}} + \frac{\partial z_h^{(t)}}{\partial g_\omega^{(t)}} \frac{\partial g_\omega^{(t)}}{\partial z_h^{(t-1)}} + \frac{\partial z_h^{(t)}}{\partial z_c^{(t)}} \frac{\partial z_c^{(t)}}{\partial z_h^{(t-1)}} \right) \\ &= \delta_{g_l}^{(t)} f'_l \left(a_l^{(t)} \right) \mathbf{w}_{l,z} + \delta_{g_\phi}^{(t)} f'_\phi \left(a_\phi^{(t)} \right) \mathbf{w}_{\phi,z} + \delta_{g_\omega}^{(t)} f'_\omega \left(a_\omega^{(t)} \right) \mathbf{w}_{\omega,z} + \delta_c^{(t)} f'_c \left(a_c^{(t)} \right) \mathbf{w}_{c,z} \end{aligned} \quad (23)$$

3 Torch Implementation

References

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