

# DERIVING BACKPROPAGATION EQUATIONS FOR AN LSTM

JUNE 19, 2019

In this post I will derive the backpropagation equations for a LSTM cell in vectorised form. It assumes basic knowledge of LSTMs and backpropagation, which you can refresh at [Understanding LSTM Networks](#) and [A Quick Introduction to Backpropagation](#).

## DERIVATIONS

### FORWARD PROPAGATION

We will firstly remind ourselves of the forward propagation equations. The nomenclature followed is demonstrated in Figure 1. All equations correspond to one time step.

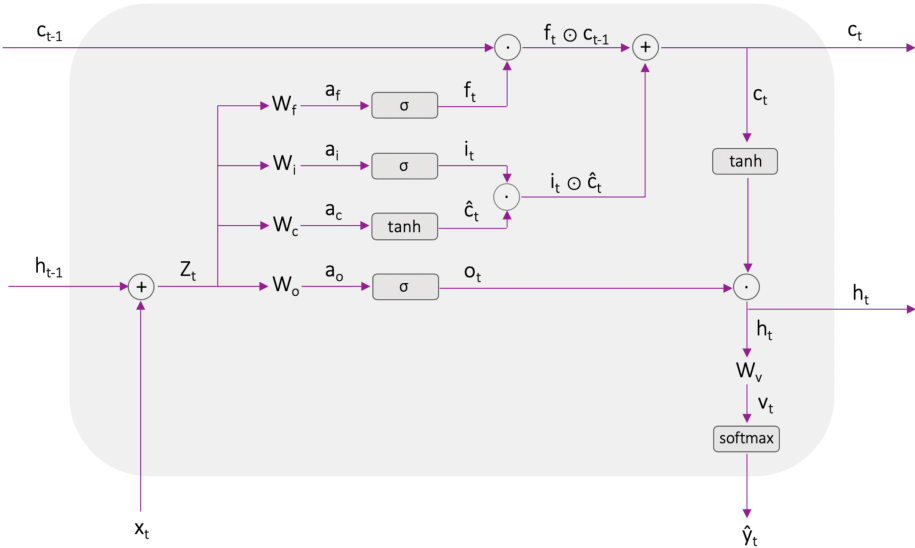


FIGURE 1: ARCHITECTURE OF A LSTM MEMORY CELL AT TIMESTEP T

$$h_{t-1} \in \mathbb{R}^{n_h}, \quad x_t \in \mathbb{R}^{n_x}$$
$$z_t = [h_{t-1}, x_t]$$

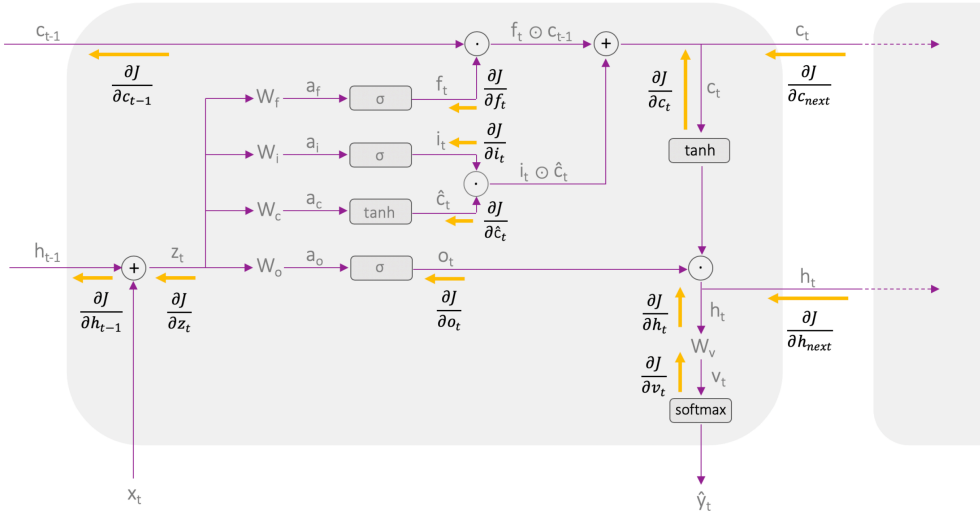
$$\begin{aligned} a_f &= W_f \cdot z_t + b_f, & f_t &= \sigma(a_f) \\ a_i &= W_i \cdot z_t + b_i, & i_t &= \sigma(a_i) \\ a_o &= W_o \cdot z_t + b_o, & o_t &= \sigma(a_o) \\ a_c &= W_c \cdot z_t + b_c, & \hat{c}_t &= \tanh(a_c) \end{aligned}$$

$$c\_t = i\_t \odot \hat{c\_t} + f\_t \odot c\_t - 1$$
$$h\_t = o\_t \odot \tanh(c\_t)$$

$$v\_t = W\_v \cdot h\_t + b\_v$$
$$\hat{y\_t} = softmax(v\_t)$$

## BACKWARD PROPAGATION

Backpropagation through a LSTM is not as straightforward as through other common Deep Learning architectures, due to the special way its underlying layers interact. Nonetheless, the approach is largely the same; identifying dependencies and recursively applying the chain rule.



**FIGURE 2: BACKPROPAGATION THROUGH A LSTM MEMORY CELL**

Cross-entropy loss with a softmax function are used at the output layer. The standard definition of the derivative of the cross-entropy loss ( $\frac{\partial J}{\partial v_t}$ ) is used directly; a detailed derivation can be found here.

## OUTPUT

$$\frac{\partial J}{\partial v\_t} = \hat{y\_t} - y\_t$$
$$\frac{\partial J}{\partial W\_v} = \frac{\partial J}{\partial v\_t} \cdot \frac{\partial v\_t}{\partial W\_v} \Rightarrow \frac{\partial J}{\partial W\_v} = \frac{\partial J}{\partial v\_t} \cdot h\_t^T$$
$$\frac{\partial J}{\partial b\_v} = \frac{\partial J}{\partial v\_t} \cdot \frac{\partial v\_t}{\partial b\_v} \Rightarrow \frac{\partial J}{\partial b\_v} = \frac{\partial J}{\partial v\_t}$$

## HIDDEN STATE

$$\frac{\partial J}{\partial h_t} = \frac{\partial J}{\partial v_t} \cdot \frac{\partial v_t}{\partial h_t} \Rightarrow \frac{\partial J}{\partial h_t} = W_{-v}^T \cdot \frac{\partial J}{\partial v_t}$$
$$\frac{\partial J}{\partial h_t} + = \frac{\partial J}{\partial h_{next}}$$

**OUTPUT GATE**

$$\frac{\partial J}{\partial o_t} = \frac{\partial J}{\partial h_t} \cdot \frac{\partial h_t}{\partial o_t} \Rightarrow \frac{\partial J}{\partial o_t} = \frac{\partial J}{\partial h_t} \odot \tanh(c_t)$$

$$\frac{\partial J}{\partial a_o} = \frac{\partial J}{\partial o_t} \cdot \frac{\partial o_t}{\partial a_o} \Rightarrow \frac{\partial J}{\partial a_o} = \frac{\partial J}{\partial h_t} \odot \tanh(c_t) \odot \frac{d(\sigma(a_o))}{da_o}$$
$$\Rightarrow \frac{\partial J}{\partial a_o} = \frac{\partial J}{\partial h_t} \odot \tanh(c_t) \odot \sigma(a_o)(1 - \sigma(a_o))$$
$$\Rightarrow \frac{\partial J}{\partial a_o} = \frac{\partial J}{\partial h_t} \odot \tanh(c_t) \odot o_t(1 - o_t)$$

$$\frac{\partial J}{\partial W_{-o}} = \frac{\partial J}{\partial a_o} \cdot \frac{\partial a_o}{\partial W_{-o}} \Rightarrow \frac{\partial J}{\partial W_{-o}} = \frac{\partial J}{\partial a_o} \cdot z_{-t}^T$$
$$\frac{\partial J}{\partial b_{-o}} = \frac{\partial J}{\partial a_o} \cdot \frac{\partial a_o}{\partial b_{-o}} \Rightarrow \frac{\partial J}{\partial b_{-o}} = \frac{\partial J}{\partial a_o}$$

**CELL STATE**

$$\frac{\partial J}{\partial c_t} = \frac{\partial J}{\partial h_t} \cdot \frac{\partial h_t}{\partial c_t} \Rightarrow \frac{\partial J}{\partial c_t} = \frac{\partial J}{\partial h_t} \odot o_t \odot (1 - \tanh(c_t)^2)$$

$$\frac{\partial J}{\partial c_t} + = \frac{\partial J}{\partial c_{next}}$$

$$\frac{\partial J}{\partial \hat{c}_t} = \frac{\partial J}{\partial c_t} \cdot \frac{\partial c_t}{\partial \hat{c}_t} \Rightarrow \frac{\partial J}{\partial \hat{c}_t} = \frac{\partial J}{\partial c_t} \odot i_t$$

$$\frac{\partial J}{\partial a_c} = \frac{\partial J}{\partial \hat{c}_t} \cdot \frac{\partial \hat{c}_t}{\partial a_c} \Rightarrow \frac{\partial J}{\partial a_c} = \frac{\partial J}{\partial c_t} \odot i_t \odot \frac{d(\tanh(a_c))}{da_c}$$
$$\Rightarrow \frac{\partial J}{\partial a_c} = \frac{\partial J}{\partial c_t} \odot i_t \odot (1 - \tanh(a_c)^2)$$
$$\Rightarrow \frac{\partial J}{\partial a_c} = \frac{\partial J}{\partial c_t} \odot i_t \odot (1 - \hat{c}_t^2)$$

$$\frac{\partial J}{\partial W_c} = \frac{\partial J}{\partial a_c} \cdot \frac{\partial a_c}{\partial W_c} \Rightarrow \frac{\partial J}{\partial W_c} = \frac{\partial J}{\partial a_c} \cdot z_t^T$$

$$\frac{\partial J}{\partial b_c} = \frac{\partial J}{\partial a_c} \cdot \frac{\partial a_c}{\partial b_c} \Rightarrow \frac{\partial J}{\partial b_c} = \frac{\partial J}{\partial a_c}$$

**INPUT GATE**

$$\frac{\partial J}{\partial i_t} = \frac{\partial J}{\partial c_t} \cdot \frac{\partial c_t}{\partial i_t} \Rightarrow \frac{\partial J}{\partial i_t} = \frac{\partial J}{\partial c_t} \odot \hat{c}_t$$

$$\begin{aligned} \frac{\partial J}{\partial a_i} &= \frac{\partial J}{\partial i_t} \cdot \frac{\partial i_t}{\partial a_i} \Rightarrow \frac{\partial J}{\partial a_i} = \frac{\partial J}{\partial c_t} \odot \hat{c}_t \odot \frac{d(\sigma(a_i))}{da_i} \\ \Rightarrow \frac{\partial J}{\partial a_i} &= \frac{\partial J}{\partial c_t} \odot \hat{c}_t \odot \sigma(a_i)(1 - \sigma(a_i)) \\ \Rightarrow \frac{\partial J}{\partial a_i} &= \frac{\partial J}{\partial c_t} \odot \hat{c}_t \odot i_t(1 - i_t) \end{aligned}$$

$$\frac{\partial J}{\partial W_i} = \frac{\partial J}{\partial a_i} \cdot \frac{\partial a_i}{\partial W_i} \Rightarrow \frac{\partial J}{\partial W_i} = \frac{\partial J}{\partial a_i} \cdot z_t^T$$

$$\frac{\partial J}{\partial b_i} = \frac{\partial J}{\partial a_i} \cdot \frac{\partial a_i}{\partial b_i} \Rightarrow \frac{\partial J}{\partial b_i} = \frac{\partial J}{\partial a_i}$$

**FORGET GATE**

$$\frac{\partial J}{\partial f_t} = \frac{\partial J}{\partial c_t} \cdot \frac{\partial c_t}{\partial f_t} \Rightarrow \frac{\partial J}{\partial f_t} = \frac{\partial J}{\partial c_t} \odot c_t - 1$$

$$\begin{aligned} \frac{\partial J}{\partial a_f} &= \frac{\partial J}{\partial f_t} \cdot \frac{\partial f_t}{\partial a_f} \Rightarrow \frac{\partial J}{\partial a_f} = \frac{\partial J}{\partial c_t} \odot c_t - 1 \odot \frac{d(\sigma(a_f))}{da_f} \\ \Rightarrow \frac{\partial J}{\partial a_f} &= \frac{\partial J}{\partial c_t} \odot c_t - 1 \odot \sigma(a_f)(1 - \sigma(a_f)) \\ \Rightarrow \frac{\partial J}{\partial a_f} &= \frac{\partial J}{\partial c_t} \odot c_t - 1 \odot f_t(1 - f_t) \end{aligned}$$

$$\frac{\partial J}{\partial W_{-f}} = \frac{\partial J}{\partial a_{-f}} \cdot \frac{\partial a_{-f}}{\partial W_{-f}} \Rightarrow \frac{\partial J}{\partial W_{-f}} = \frac{\partial J}{\partial a_{-f}} \cdot z_{-t}^T$$
$$\frac{\partial J}{\partial b_{-f}} = \frac{\partial J}{\partial a_{-f}} \cdot \frac{\partial a_{-f}}{\partial b_{-f}} \Rightarrow \frac{\partial J}{\partial b_{-f}} = \frac{\partial J}{\partial a_{-f}}$$

INPUT

$$\frac{\partial J}{\partial z_{-t}} = \frac{\partial J}{\partial a_{-f}} \cdot \frac{\partial a_{-f}}{\partial z_{-t}} + \frac{\partial J}{\partial a_{-i}} \cdot \frac{\partial a_{-i}}{\partial z_{-t}} + \frac{\partial J}{\partial a_{-o}} \cdot \frac{\partial a_{-o}}{\partial z_{-t}} + \frac{\partial J}{\partial a_{-c}} \cdot \frac{\partial a_{-c}}{\partial z_{-t}}$$
$$\Rightarrow \frac{\partial J}{\partial z_{-t}} = W_{-f}^T \cdot \frac{\partial J}{\partial a_{-f}} + W_{-i}^T \cdot \frac{\partial J}{\partial a_{-i}} + W_{-o}^T \cdot \frac{\partial J}{\partial a_{-o}} + W_{-c}^T \cdot \frac{\partial J}{\partial a_{-c}}$$

$$\frac{\partial J}{\partial h_{-t-1}} = \frac{\partial J}{\partial z_{-t}}[:, n_{-h}, :]$$
$$\frac{\partial J}{\partial c_{-t-1}} = \frac{\partial J}{\partial c_{-t}} \cdot \frac{\partial c_{-t}}{\partial c_{-t-1}} \Rightarrow \frac{\partial J}{\partial c_{-t-1}} = \frac{\partial J}{\partial c_{-t}} \odot f_{-t}$$

The above equations for forward propagation and back propagation will be calculated T times (number of time steps) in each training iteration. At the end of each training iteration, the weights will be updated using the accumulated cost gradient with respect to each weight for all time steps. Assuming Stochastic Gradient Descent, the update equations are the following:

$$\frac{\partial J}{\partial W_{-f}} = \sum_{-t}^T \frac{\partial J}{\partial W_{-f}^t}, \quad W_{-f+} = \alpha * \frac{\partial J}{\partial W_{-f}}$$
$$\frac{\partial J}{\partial W_{-i}} = \sum_{-t}^T \frac{\partial J}{\partial W_{-i}^t}, \quad W_{-i+} = \alpha * \frac{\partial J}{\partial W_{-i}}$$
$$\frac{\partial J}{\partial W_{-o}} = \sum_{-t}^T \frac{\partial J}{\partial W_{-o}^t}, \quad W_{-o+} = \alpha * \frac{\partial J}{\partial W_{-o}}$$
$$\frac{\partial J}{\partial W_{-c}} = \sum_{-t}^T \frac{\partial J}{\partial W_{-c}^t}, \quad W_{-c+} = \alpha * \frac{\partial J}{\partial W_{-c}}$$
$$\frac{\partial J}{\partial W_{-v}} = \sum_{-t}^T \frac{\partial J}{\partial W_{-v}^t}, \quad W_{-v+} = \alpha * \frac{\partial J}{\partial W_{-v}}$$

In the next post, we will implement the above equations using Numpy and train the resulting LSTM model on real data.