

LSTM Derivations

Minh-Thang Luong
lmthang@stanford.edu

Abigail See
abisee@stanford.edu

1 LSTM Architecture

An LSTM block at layer $l \in \{1, \dots, L\}$ and time $t \in \{1, \dots, T\}$ consists of:

- The hidden state $\mathbf{h}_t^l \in \mathbb{R}^n$
- The memory cell $\mathbf{c}_t^l \in \mathbb{R}^n$
- The input gate $\mathbf{i}_t^l \in [0, 1]^n$
- The forget gate $\mathbf{f}_t^l \in [0, 1]^n$
- The output gate $\mathbf{o}_t^l \in [0, 1]^n$
- The input modulation gate $\hat{\mathbf{h}}_t^l \in [0, 1]^n$

We call n the LSTM block size.

2 Forward Propagation

We use the formulation of [Zaremba et al., 2014]. For a single LSTM block at layer l and time t , the new hidden state \mathbf{h}_t^l and memory cell \mathbf{c}_t^l are calculated from \mathbf{h}_t^{l-1} , \mathbf{h}_{t-1}^l and \mathbf{c}_{t-1}^l like so:

$$\begin{pmatrix} \mathbf{i}_t^l \\ \mathbf{f}_t^l \\ \mathbf{o}_t^l \\ \hat{\mathbf{h}}_t^l \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \text{tanh} \end{pmatrix} \mathbf{T}_{4n \times 2n} \begin{bmatrix} \mathbf{h}_t^{l-1} \\ \mathbf{h}_{t-1}^l \end{bmatrix} \quad (1)$$

$$\mathbf{c}_t^l = \mathbf{f}_t^l \circ \mathbf{c}_{t-1}^l + \mathbf{i}_t^l \circ \hat{\mathbf{h}}_t^l \quad (2)$$

$$\mathbf{h}_t^l = \mathbf{o}_t^l \circ \tanh(\mathbf{c}_t^l) \quad (3)$$

where sigm and tanh are applied element-wise, \circ denotes element-wise multiplication, and $\mathbf{T}_{4n \times 2n}$ is a $4n \times 2n$ matrix of weights that depends on l but not t .¹ If $l = 1$ then \mathbf{h}_t^{l-1} is the input vector x_t . If $t = 1$ then \mathbf{h}_{t-1}^l and \mathbf{c}_{t-1}^l are taken to be zero.

¹Note: Sometimes these equations are written omitting the superscript l and writing \mathbf{h}_t^{l-1} as x_t , but for the purposes of deriving the back-propagation equations, we need to refer to the layer l explicitly.

3 Backward Propagation

3.1 Definitions

In this section we define some additional notation that will help us to derive the necessary back-propagation equations.

Definition 1. Let U and V refer to the $n \times n$ weight matrices corresponding to the following portions of $\mathbf{T}_{4n \times 2n}$:

$$\mathbf{T}_{4n \times 2n} = \begin{bmatrix} U_i & V_i \\ U_f & V_f \\ U_o & V_o \\ U_{\hat{h}} & V_{\hat{h}} \end{bmatrix} \quad (4)$$

In particular, we will use the superscript l to denote these matrices used to calculate layer l in Equation (1).

Definition 2. For all $l \in \{1, \dots, L\}$ and $t \in \{1, \dots, T\}$, define the weighted inputs

$$\begin{aligned} z_i^l(t) &= U_i \mathbf{h}_t^{l-1} + V_i \mathbf{h}_{t-1}^l & z_f^l(t) &= U_f \mathbf{h}_t^{l-1} + V_f \mathbf{h}_{t-1}^l \\ z_o^l(t) &= U_o \mathbf{h}_t^{l-1} + V_o \mathbf{h}_{t-1}^l & z_{\hat{h}}^l(t) &= U_{\hat{h}} \mathbf{h}_t^{l-1} + V_{\hat{h}} \mathbf{h}_{t-1}^l \end{aligned}$$

so that

$$\begin{pmatrix} \mathbf{i}_t^l \\ \mathbf{f}_t^l \\ \mathbf{o}_t^l \\ \hat{\mathbf{h}}_t^l \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \text{tanh} \end{pmatrix} \begin{pmatrix} z_i^l(t) \\ z_f^l(t) \\ z_o^l(t) \\ z_{\hat{h}}^l(t) \end{pmatrix} \quad (5)$$

where sigm and tanh are applied element-wise. We call $z_i^l(t)$ the weighted input to the input gate \mathbf{i}_t^l .

Definition 3. Define the error of the input, forget, output and input modulation gates at layer l and time t to be

$$\begin{aligned} \delta_i^l(\mathbf{t}) &= \frac{\partial L}{\partial z_i^l(t)} & \delta_f^l(\mathbf{t}) &= \frac{\partial L}{\partial z_f^l(t)} \\ \delta_o^l(\mathbf{t}) &= \frac{\partial L}{\partial z_o^l(t)} & \delta_{\hat{h}}^l(\mathbf{t}) &= \frac{\partial L}{\partial z_{\hat{h}}^l(t)} \end{aligned}$$

where L is the loss function. Note: $\delta_i^l(\mathbf{t})$ is the partial derivative of L with respect to the weighted input $z_i^l(t)$, not \mathbf{i}_t^l .

Definition 4. Define the error of the hidden state and cell and layer l and time t to be

$$\delta_h^l(\mathbf{t}) = \frac{\partial L}{\partial \mathbf{h}_t^l} \quad \delta_c^l(\mathbf{t}) = \frac{\partial L}{\partial \mathbf{c}_t^l} \quad (6)$$

where L is the loss function.

3.2 Derivations

In this section we will derive expressions for $\delta_h^l(\mathbf{t})$, $\delta_c^l(\mathbf{t})$, $\delta_i^l(\mathbf{t})$, $\delta_f^l(\mathbf{t})$, $\delta_o^l(\mathbf{t})$, and $\delta_{\hat{h}}^l(\mathbf{t})$ in terms of the δ values for the $(l+1, t)$ and $(l, t+1)$ blocks. These expressions will enable us to do back-propagation through time and layers. If you are not interested in the derivations, skip ahead to Section 3.3 to see the final back-propagation equations.

Lemma 1. For $l \in \{1, \dots, L-1\}$ and $t \in \{1, \dots, T-1\}$,

$$\delta_h^l(\mathbf{t}) = \begin{bmatrix} U_i^\top & U_f^\top & U_o^\top & U_{\hat{h}}^\top & V_i^\top & V_f^\top & V_o^\top & V_{\hat{h}}^\top \end{bmatrix} \begin{bmatrix} \delta_i^l(\mathbf{t}+1) \\ \delta_f^l(\mathbf{t}+1) \\ \delta_o^l(\mathbf{t}+1) \\ \delta_{\hat{h}}^l(\mathbf{t}+1) \\ \delta_i^{l+1}(\mathbf{t}) \\ \delta_f^{l+1}(\mathbf{t}) \\ \delta_o^{l+1}(\mathbf{t}) \\ \delta_{\hat{h}}^{l+1}(\mathbf{t}) \end{bmatrix} \quad (7)$$

where each of the U and V matrices are with respect to layer l .

Note the left matrix in the multiplication has dimensions $n \times 8n$, the right matrix $8n \times n$, and $\delta_h^l(\mathbf{t})$ is $n \times 1$.

Proof. We prove this element-wise. For any $j = 1, \dots, n$:

$$\delta_h^l(\mathbf{t})_j = \frac{\partial L}{\partial (\mathbf{h}_t^l)_j} \quad (\text{definition of } \delta_h^l(\mathbf{t}))$$

Now, because \mathbf{h}_t^l affects $(z_i, z_f, z_o, z_{\hat{h}})$ for $(l, t+1)$ and $(l+1, t)$, we take the chain rule over these eight variables. Therefore the above equation can be written as

$$\begin{aligned} & \sum_{k=1}^n \left(\frac{\partial L}{\partial z_i^l(t+1)_k} \frac{\partial z_i^l(t+1)_k}{\partial (\mathbf{h}_t^l)_j} + \frac{\partial L}{\partial z_f^l(t+1)_k} \frac{\partial z_f^l(t+1)_k}{\partial (\mathbf{h}_t^l)_j} \right. \\ & + \frac{\partial L}{\partial z_o^l(t+1)_k} \frac{\partial z_o^l(t+1)_k}{\partial (\mathbf{h}_t^l)_j} + \frac{\partial L}{\partial z_{\hat{h}}^l(t+1)_k} \frac{\partial z_{\hat{h}}^l(t+1)_k}{\partial (\mathbf{h}_t^l)_j} \\ & + \frac{\partial L}{\partial z_i^{l+1}(t)_k} \frac{\partial z_i^{l+1}(t)_k}{\partial (\mathbf{h}_t^l)_j} + \frac{\partial L}{\partial z_f^{l+1}(t)_k} \frac{\partial z_f^{l+1}(t)_k}{\partial (\mathbf{h}_t^l)_j} \\ & \left. + \frac{\partial L}{\partial z_o^{l+1}(t)_k} \frac{\partial z_o^{l+1}(t)_k}{\partial (\mathbf{h}_t^l)_j} + \frac{\partial L}{\partial z_{\hat{h}}^{l+1}(t)_k} \frac{\partial z_{\hat{h}}^{l+1}(t)_k}{\partial (\mathbf{h}_t^l)_j} \right) \end{aligned}$$

First note that the first of each pair is some δ e.g.

$$\frac{\partial L}{\partial z_i^l(t+1)_k} = \delta_i^l(\mathbf{t}+1) \quad (\text{by definition})$$

The second of each pair can be evaluated like so:

$$\begin{aligned}
\frac{\partial z_i^l(t+1)_k}{\partial(\mathbf{h}_t^l)_j} &= \frac{\partial}{\partial(\mathbf{h}_t^l)_j} (U_i^l \mathbf{h}_{t+1}^l + V_i^l \mathbf{h}_t^l)_k && \text{(definition of } z_i^l(t+1)) \\
&= \frac{\partial}{\partial(\mathbf{h}_t^l)_j} \left(\sum_{m=1}^n (V_i^l)_{km} (\mathbf{h}_t^l)_m \right) && (U_i^l \mathbf{h}_{t+1}^l \text{ does not depend on } \mathbf{h}_t^l) \\
&= (V_i^l)_{kj} && \text{(expression equals 0 except when } m = j)
\end{aligned}$$

so the first of the eight sums can be written as

$$\sum_{k=1}^n \frac{\partial L}{\partial z_i^l(t+1)_k} \frac{\partial z_i^l(t+1)_k}{\partial(\mathbf{h}_t^l)_j} = \sum_{k=1}^n \delta_i^l(\mathbf{t}+1) (V_i^l)_{kj} = [(V_i^l)^\top \delta_i^l(\mathbf{t}+1)]_j \quad (8)$$

Finding similar expressions for the other seven sums, we obtain

$$\begin{aligned}
\delta_{\mathbf{h}}^l(\mathbf{t}) &= (U_i^l)^\top \delta_i^l(\mathbf{t}+1) + (U_f^l)^\top \delta_f^l(\mathbf{t}+1) + (U_o^l)^\top \delta_o^l(\mathbf{t}+1) + (U_h^l)^\top \delta_h^l(\mathbf{t}+1) \\
&+ (V_i^l)^\top \delta_i^{l+1}(\mathbf{t}) + (V_f^l)^\top \delta_f^{l+1}(\mathbf{t}) + (V_o^l)^\top \delta_o^{l+1}(\mathbf{t}) + (V_h^l)^\top \delta_h^{l+1}(\mathbf{t})
\end{aligned}$$

□

Lemma 2. For $l \in \{1, \dots, L\}$ and $t \in \{1, \dots, T-1\}$,

$$\delta_{\mathbf{c}}^l(\mathbf{t}) = \delta_{\mathbf{c}}^l(\mathbf{t}+1) \circ \mathbf{f}_{t+1}^l + \delta_{\mathbf{h}}^l(\mathbf{t}) \circ \boldsymbol{\sigma}_t^l \circ \tanh'(\mathbf{c}_t^l) \quad (9)$$

Proof. We prove this element-wise. For any $j = 1, \dots, n$:

$$\begin{aligned}
\delta_{\mathbf{c}}^l(\mathbf{t})_j &= \frac{\partial L}{\partial(\mathbf{c}_t^l)_j} && \text{(definition of } \delta_{\mathbf{c}}^l(\mathbf{t})) \\
&= \sum_{k=1}^n \frac{\partial L}{\partial(\mathbf{c}_{t+1}^l)_k} \frac{\partial(\mathbf{c}_{t+1}^l)_k}{\partial(\mathbf{c}_t^l)_j} + \sum_{k=1}^n \frac{\partial L}{\partial(\mathbf{h}_t^l)_k} \frac{\partial(\mathbf{h}_t^l)_k}{\partial(\mathbf{c}_t^l)_j} && \text{(chain rule)}
\end{aligned}$$

The second equality follows from the fact that \mathbf{c}_t^l affects \mathbf{h}_t^l and \mathbf{c}_{t+1}^l . For the first part of the expression, note that

$$\begin{aligned}
\frac{\partial(\mathbf{c}_{t+1}^l)_k}{\partial(\mathbf{c}_t^l)_j} &= \frac{\partial}{\partial(\mathbf{c}_t^l)_j} \left(\mathbf{f}_{t+1}^l \circ \mathbf{c}_t^l + \mathbf{i}_{t+1}^l \circ \hat{\mathbf{h}}_{t+1}^l \right)_k && \text{(by definition of } \mathbf{c}_{t+1}^l) \\
&= \begin{cases} \mathbf{f}_{t+1}^l & \text{if } k = j \\ 0 & \text{otherwise.} \end{cases}
\end{aligned}$$

For the second part, note that

$$\begin{aligned}
\frac{\partial(\mathbf{h}_t^l)_k}{\partial(\mathbf{c}_t^l)_j} &= \frac{\partial}{\partial(\mathbf{c}_t^l)_j} (\boldsymbol{\sigma}_t^l \circ \tanh(\mathbf{c}_t^l))_k && \text{(by definition of } \mathbf{h}_t^l) \\
&= \begin{cases} (\boldsymbol{\sigma}_t^l)_j \tanh'(\mathbf{c}_t^l)_j & \text{if } k = j \\ 0 & \text{otherwise.} \end{cases}
\end{aligned}$$

Combining the previous three equations and using the definitions of $\delta_c^l(\mathbf{t} + 1)$ and $\delta_h^l(\mathbf{t})$, we obtain

$$\delta_c^l(\mathbf{t})_j = \delta_c^l(\mathbf{t} + 1)_j (\mathbf{f}_{t+1}^l)_j + \delta_h^l(\mathbf{t})_j (\mathbf{o}_t^l)_j \tanh'(\mathbf{c}_t^l)_j \quad (10)$$

□

Lemma 3. For $l \in \{1, \dots, L\}$ and $t \in \{1, \dots, T\}$,

$$\delta_i^l(\mathbf{t}) = \delta_c^l(\mathbf{t}) \circ \text{sigm}'(z_i^l(t)) \circ \hat{\mathbf{h}}_t^l \quad (11)$$

Proof. We prove this element-wise. For any $j = 1, \dots, n$:

$$\begin{aligned} \delta_i^l(\mathbf{t})_j &= \frac{\partial L}{\partial z_i^l(t)_j} && \text{(definition of } \delta_i^l(\mathbf{t}) \text{)} \\ &= \sum_{k=1}^n \frac{\partial L}{\partial (\mathbf{c}_t^l)_k} \frac{\partial (\mathbf{c}_t^l)_k}{\partial z_i^l(t)_j} && \text{(chain rule)} \\ &= \sum_{k=1}^n \delta_c^l(\mathbf{t})_k \frac{\partial}{\partial z_i^l(t)_j} \left(\mathbf{f}_t^l \circ \mathbf{c}_{t-1}^l + \mathbf{i}_t^l \circ \hat{\mathbf{h}}_t^l \right)_k && \text{(definition of } \delta_c^l(\mathbf{t}) \text{ and } \mathbf{c}_t^l \text{)} \\ &= \delta_c^l(\mathbf{t})_j \frac{\partial}{\partial z_i^l(t)_j} \left(\mathbf{i}_t^l \circ \hat{\mathbf{h}}_t^l \right)_j && \text{(expression equals 0 except when } k = j \text{)} \\ &= \delta_c^l(\mathbf{t})_j \text{sigm}'(z_i^l(t))_j (\hat{\mathbf{h}}_t^l)_j && \text{(definition of } \mathbf{i}_t^l \text{ in terms of } z_i^l(t) \text{)} \end{aligned}$$

Note that for the second equality we took the chain rule with respect to the elements of \mathbf{c}_t^l , because \mathbf{i}_t^l affects \mathbf{c}_t^l . □

Lemma 4. For $l \in \{1, \dots, L\}$ and $t \in \{1, \dots, T\}$,

$$\delta_f^l(\mathbf{t}) = \delta_c^l(\mathbf{t}) \circ \text{sigm}'(z_f^l(t)) \circ \mathbf{c}_{t-1}^l \quad (12)$$

Proof. We prove this element-wise. For any $j = 1, \dots, n$:

$$\begin{aligned} \delta_f^l(\mathbf{t})_j &= \frac{\partial L}{\partial z_f^l(t)_j} && \text{(definition of } \delta_f^l(\mathbf{t}) \text{)} \\ &= \sum_{k=1}^n \frac{\partial L}{\partial (\mathbf{c}_t^l)_k} \frac{\partial (\mathbf{c}_t^l)_k}{\partial z_f^l(t)_j} && \text{(chain rule)} \\ &= \sum_{k=1}^n \delta_c^l(\mathbf{t})_k \frac{\partial}{\partial z_f^l(t)_j} \left(\mathbf{f}_t^l \circ \mathbf{c}_{t-1}^l + \mathbf{i}_t^l \circ \hat{\mathbf{h}}_t^l \right)_k && \text{(definition of } \delta_c^l(\mathbf{t}) \text{ and } \mathbf{c}_t^l \text{)} \\ &= \delta_c^l(\mathbf{t})_j \frac{\partial}{\partial z_f^l(t)_j} \left(\mathbf{f}_t^l \circ \mathbf{c}_{t-1}^l \right)_j && \text{(expression equals 0 except when } k = j \text{)} \\ &= \delta_c^l(\mathbf{t})_j \text{sigm}'(z_f^l(t))_j (\mathbf{c}_{t-1}^l)_j && \text{(definition of } \mathbf{f}_t^l \text{ in terms of } z_f^l(t) \text{)} \end{aligned}$$

Note that for the second equality we took the chain rule with respect to the elements of \mathbf{c}_t^l , because \mathbf{f}_t^l affects \mathbf{c}_t^l . □

Lemma 5. For $l \in \{1, \dots, L\}$ and $t \in \{1, \dots, T\}$,

$$\delta_o^l(\mathbf{t}) = \delta_h^l(\mathbf{t}) \circ \text{sigm}'(z_o^l(t)) \circ \tanh(\mathbf{c}_t^l) \quad (13)$$

Proof. We prove this element-wise. For any $j = 1, \dots, n$:

$$\begin{aligned} \delta_o^l(\mathbf{t})_j &= \frac{\partial L}{\partial z_o^l(t)_j} && \text{(definition of } \delta_o^l(\mathbf{t}) \text{)} \\ &= \sum_{k=1}^n \frac{\partial L}{\partial (\mathbf{h}_t^l)_k} \frac{\partial (\mathbf{h}_t^l)_k}{\partial z_o^l(t)_j} && \text{(chain rule)} \\ &= \sum_{k=1}^n \delta_h^l(\mathbf{t})_k \frac{\partial}{\partial z_o^l(t)_j} (\mathbf{o}_t^l \circ \tanh(\mathbf{c}_t^l))_k && \text{(definition of } \delta_h^l(\mathbf{t}) \text{ and } \mathbf{h}_t^l \text{)} \\ &= \delta_h^l(\mathbf{t})_j \frac{\partial}{\partial z_o^l(t)_j} (\mathbf{o}_t^l \circ \tanh(\mathbf{c}_t^l))_j && \text{(expression equals 0 except when } k = j \text{)} \\ &= \delta_h^l(\mathbf{t})_j \text{sigm}'(z_o^l(t))_j \tanh(\mathbf{c}_t^l)_j && \text{(definition of } \mathbf{o}_t^l \text{ in terms of } z_o^l(t) \text{)} \end{aligned}$$

Note that for the second equality we took the chain rule with respect to the elements of \mathbf{h}_t^l , because \mathbf{o}_t^l affects \mathbf{h}_t^l . \square

Lemma 6. For $l \in \{1, \dots, L\}$ and $t \in \{1, \dots, T\}$,

$$\delta_h^l(\mathbf{t}) = \delta_c^l(\mathbf{t}) \circ \mathbf{i}_t^l \circ \tanh'(z_h^l(t)) \quad (14)$$

Proof. We prove this element-wise. For any $j = 1, \dots, n$:

$$\begin{aligned} \delta_h^l(\mathbf{t})_j &= \frac{\partial L}{\partial z_h^l(t)_j} && \text{(definition of } \delta_h^l(\mathbf{t}) \text{)} \\ &= \sum_{k=1}^n \frac{\partial L}{\partial (\mathbf{c}_t^l)_k} \frac{\partial (\mathbf{c}_t^l)_k}{\partial z_h^l(t)_j} && \text{(chain rule)} \\ &= \sum_{k=1}^n \delta_c^l(\mathbf{t})_k \frac{\partial}{\partial z_h^l(t)_j} (\mathbf{f}_t^l \circ \mathbf{c}_{t-1}^l + \mathbf{i}_t^l \circ \hat{\mathbf{h}}_t^l)_k && \text{(definition of } \delta_c^l(\mathbf{t}) \text{ and } \mathbf{c}_t^l \text{)} \\ &= \delta_c^l(\mathbf{t})_j \frac{\partial}{\partial z_h^l(t)_j} (\mathbf{i}_t^l \circ \hat{\mathbf{h}}_t^l)_j && \text{(expression equals 0 except when } k = j \text{)} \\ &= \delta_c^l(\mathbf{t})_j (\mathbf{i}_t^l)_j \tanh'(z_h^l(t))_j && \text{(definition of } \hat{\mathbf{h}}_t^l \text{ in terms of } z_h^l(t) \text{)} \end{aligned}$$

Note that for the second equality we took the chain rule with respect to the elements of \mathbf{c}_t^l , because $\hat{\mathbf{h}}_t^l$ affects \mathbf{c}_t^l . \square

Lemma 7. For all $l \in \{1, \dots, L\}$,

$$\begin{aligned} \frac{\partial L}{\partial U_i^l} &= \sum_{t=1}^T (\mathbf{h}_t^{l-1}) (\boldsymbol{\delta}_i^l(\mathbf{t}))^\top & \frac{\partial L}{\partial V_i^l} &= \sum_{t=1}^T (\mathbf{h}_{t-1}^l) (\boldsymbol{\delta}_i^l(\mathbf{t}))^\top \\ \frac{\partial L}{\partial U_f^l} &= \sum_{t=1}^T (\mathbf{h}_t^{l-1}) (\boldsymbol{\delta}_f^l(\mathbf{t}))^\top & \frac{\partial L}{\partial V_f^l} &= \sum_{t=1}^T (\mathbf{h}_{t-1}^l) (\boldsymbol{\delta}_f^l(\mathbf{t}))^\top \\ \frac{\partial L}{\partial U_o^l} &= \sum_{t=1}^T (\mathbf{h}_t^{l-1}) (\boldsymbol{\delta}_o^l(\mathbf{t}))^\top & \frac{\partial L}{\partial V_o^l} &= \sum_{t=1}^T (\mathbf{h}_{t-1}^l) (\boldsymbol{\delta}_o^l(\mathbf{t}))^\top \\ \frac{\partial L}{\partial U_{\hat{h}}^l} &= \sum_{t=1}^T (\mathbf{h}_t^{l-1}) (\boldsymbol{\delta}_{\hat{h}}^l(\mathbf{t}))^\top & \frac{\partial L}{\partial V_{\hat{h}}^l} &= \sum_{t=1}^T (\mathbf{h}_{t-1}^l) (\boldsymbol{\delta}_{\hat{h}}^l(\mathbf{t}))^\top \end{aligned}$$

Proof. We will prove the identities for the input gate i only; the proofs for f , o and \hat{h} are identical. First recall Definition 2 for the weighted input:

$$z_i^l(t) = U_i \mathbf{h}_t^{l-1} + V_i \mathbf{h}_{t-1}^l$$

Now, for any $j, k \in \{1, \dots, n\}$, consider the effect of $(U_i^l)_{jk}$. It maps from the k th element of \mathbf{h}_t^{l-1} to the j th element of $z_i^l(t)$, for all t . Therefore applying the chain rule we obtain

$$\begin{aligned} \frac{\partial L}{\partial (U_i^l)_{jk}} &= \sum_{t=1}^T \frac{\partial L}{\partial z_i^l(t)_j} \frac{\partial z_i^l(t)_j}{\partial (U_i^l)_{jk}} && \text{(chain rule)} \\ &= \sum_{t=1}^T \boldsymbol{\delta}_i^l(\mathbf{t})_j (\mathbf{h}_t^{l-1})_k && \text{(definition of } \boldsymbol{\delta}_i^l(\mathbf{t}) \text{ and } z_i^l(t)) \end{aligned}$$

Therefore

$$\frac{\partial L}{\partial U_i^l} = \sum_{t=1}^T (\mathbf{h}_t^{l-1}) (\boldsymbol{\delta}_i^l(\mathbf{t}))^\top \quad (15)$$

The expression for $\partial L / \partial V_i^l$ is derived similarly, by noting that $(V_i^l)_{jk}$ maps from the k th element of \mathbf{h}_{t-1}^l to the j th element of $z_i^l(t)$. \square

Corollary 1. For all $l \in \{1, \dots, L\}$,

$$\begin{bmatrix} \partial L / \partial U_i^l & \partial L / \partial V_i^l \\ \partial L / \partial U_f^l & \partial L / \partial V_f^l \\ \partial L / \partial U_o^l & \partial L / \partial V_o^l \\ \partial L / \partial U_{\hat{h}}^l & \partial L / \partial V_{\hat{h}}^l \end{bmatrix} = \sum_{t=1}^T \begin{bmatrix} \boldsymbol{\delta}_i^l(\mathbf{t}) \\ \boldsymbol{\delta}_f^l(\mathbf{t}) \\ \boldsymbol{\delta}_o^l(\mathbf{t}) \\ \boldsymbol{\delta}_{\hat{h}}^l(\mathbf{t}) \end{bmatrix} \begin{bmatrix} \mathbf{h}_t^{l-1} & \mathbf{h}_{t-1}^l \end{bmatrix} \quad (16)$$

Proof. This is simply a rearrangement of Lemma 7. \square

3.3 Summary

Now we have derived all the necessary equations, we have an algorithm to calculate the necessary error values for each LSTM block, and thus calculate the derivative of the loss function with respect to our various weights.

For $l \in \{1, \dots, L-1\}$ and $t \in \{1, \dots, T\}$, we calculate $\delta_h^l(t)$ as follows:

$$\delta_h^l(t) = \begin{bmatrix} U_i^\top & U_f^\top & U_o^\top & U_{\hat{h}}^\top & V_i^\top & V_f^\top & V_o^\top & V_{\hat{h}}^\top \end{bmatrix} \begin{bmatrix} \delta_i^l(t+1) \\ \delta_f^l(t+1) \\ \delta_o^l(t+1) \\ \delta_{\hat{h}}^l(t+1) \\ \delta_i^{l+1}(t) \\ \delta_f^{l+1}(t) \\ \delta_o^{l+1}(t) \\ \delta_{\hat{h}}^{l+1}(t) \end{bmatrix}$$

$$\begin{aligned} \delta_c^l(t) &= \delta_c^l(t+1) \circ f_{t+1}^l + \delta_h^l(t) \circ \mathbf{o}_t^l \circ \tanh'(\mathbf{c}_t^l) \\ \delta_i^l(t) &= \delta_c^l(t) \circ \text{sigm}'(z_i^l(t)) \circ \hat{\mathbf{h}}_t^l \\ \delta_o^l(t) &= \delta_h^l(t) \circ \text{sigm}'(z_o^l(t)) \circ \tanh(\mathbf{c}_t^l) \\ \delta_f^l(t) &= \delta_c^l(t) \circ \text{sigm}'(z_f^l(t)) \circ \mathbf{c}_{t-1}^l \\ \delta_{\hat{h}}^l(t) &= \delta_c^l(t) \circ \mathbf{i}_t^l \circ \tanh'(z_{\hat{h}}^l(t)) \end{aligned}$$

Note: if $t = T$ then we take $\delta^l(t+1)$ to be zero for \mathbf{i} , \mathbf{f} , \mathbf{o} , $\hat{\mathbf{h}}$ and \mathbf{c} . **if $l = L$ how do we calculate $\delta_h^l(t)$?**

Once we have calculated the above error values for all l and t , we can calculate the derivative of the loss function with respect to our various weights. In particular, for $l \in \{1, \dots, L\}$:

$$\begin{bmatrix} \partial L / \partial U_i^l & \partial L / \partial V_i^l \\ \partial L / \partial U_f^l & \partial L / \partial V_f^l \\ \partial L / \partial U_o^l & \partial L / \partial V_o^l \\ \partial L / \partial U_{\hat{h}}^l & \partial L / \partial V_{\hat{h}}^l \end{bmatrix} = \sum_{t=1}^T \begin{bmatrix} \delta_i^l(t) \\ \delta_f^l(t) \\ \delta_o^l(t) \\ \delta_{\hat{h}}^l(t) \end{bmatrix} \begin{bmatrix} \mathbf{h}_t^{l-1} & \mathbf{h}_{t-1}^l \end{bmatrix}$$

We then use these derivatives to apply gradient descent to U^l and V^l .

4 Random

$$\begin{aligned} &\delta_{c^{(2)}} \\ &\delta_{h^{(2)}} \\ &\delta_{c^{(1)}} \\ &\delta_{h^{(1)}} \\ &\delta_c += \delta_h \mathbf{o}_t \tanh'(\mathbf{c}_t) \\ &\delta_c = \delta_c \circ \mathbf{f}_t \\ &\delta_h += \text{upper grad} \end{aligned}$$

5 Other Recurrent Units

Different recurrent units:

RNN

$$\mathbf{h}_t = \sigma \left(\mathbf{T}_{n \times 2n} \begin{bmatrix} \mathbf{x}_t \\ \mathbf{h}_{t-1} \end{bmatrix} \right) \quad (17)$$

$$\mathbf{T}_{n \times 2n} = [\mathbf{W}_{xh} \mathbf{W}_{hh}] \quad (18)$$

$$\mathbf{h}_t = \sigma (\mathbf{W}_{xh} \mathbf{x}_t + \mathbf{W}_{hh} \mathbf{h}_{t-1}) \quad (19)$$

$$\frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} = \text{diag} (\sigma'(\dots)) \mathbf{W}_{hh}^\top \quad (20)$$

$$\left\| \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \right\| \leq \|\text{diag} (\sigma'(\dots))\| \|\mathbf{W}_{hh}^\top\| \quad (21)$$

$$\leq \gamma \lambda_1 \quad (22)$$

$$\left\| \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-k}} \right\| \leq (\gamma \lambda_1)^k \rightarrow 0 \quad \text{if } \lambda_1 < \frac{1}{\gamma} \quad (23)$$

$$\frac{\partial \mathbf{c}_t}{\partial \mathbf{c}_{t-1}} = \mathbf{I} \quad (24)$$

GRU [Cho et al., 2014]

$$\begin{pmatrix} \mathbf{z}_t \\ \mathbf{r}_t \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \end{pmatrix} \mathbf{T}_{2n \times 2n} \begin{bmatrix} \mathbf{x}_t \\ \mathbf{h}_{t-1} \end{bmatrix} \quad (25)$$

$$\hat{\mathbf{h}}_t = \tanh(\mathbf{W} \mathbf{x}_t + \mathbf{r}_t \circ \mathbf{U} \mathbf{h}_{t-1}) \quad (26)$$

$$\mathbf{h}_t = \mathbf{z}_t \circ \mathbf{h}_{t-1} + (1 - \mathbf{z}_t) \circ \hat{\mathbf{h}}_t \quad (27)$$

My unit (maybe we should try to implement this!)

$$\begin{pmatrix} \mathbf{i}_t \\ \mathbf{f}_t \\ \hat{\mathbf{h}}_t \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{tanh} \end{pmatrix} \mathbf{T}_{3n \times 2n} \begin{bmatrix} \mathbf{x}_t \\ \mathbf{h}_{t-1} \end{bmatrix} \quad (28)$$

$$\mathbf{h}_t = \mathbf{f}_t \circ \mathbf{h}_{t-1} + \mathbf{i}_t \circ \hat{\mathbf{h}}_t \quad (29)$$

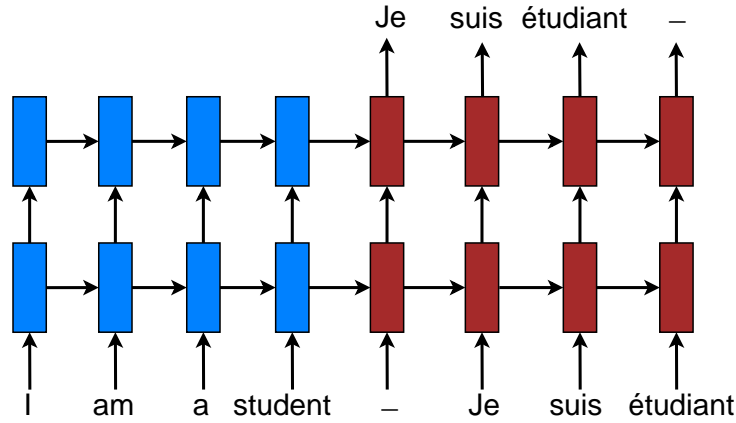


Figure 1: NMT

6 Neural Machine Translation

[Sutskever et al., 2014]

6.1 Attention

Content-based

$$\mathbf{a}_t = \text{Attend}(\mathbf{h}_{t-1}, \bar{\mathbf{h}}_{1..S}) \quad (30)$$

Location-based

$$\mathbf{a}_t = \text{Attend}(\mathbf{h}_{t-1}, \mathbf{a}_{t-1}) \quad (31)$$

Hybrid

$$\mathbf{a}_t = \text{Attend}(\mathbf{h}_{t-1}, \mathbf{a}_{t-1}, \bar{\mathbf{h}}_{1..S}) \quad (32)$$

7 Conclusion and Future Work

References

- [Cho et al., 2014] Cho, K., van Merriënboer, B., Gulcehre, C., Bougares, F., Schwenk, H., and Bengio, Y. (2014). Learning phrase representations using RNN encoder-decoder for statistical machine translation. In *EMNLP*.
- [Sutskever et al., 2014] Sutskever, I., Vinyals, O., and Le, Q. V. (2014). Sequence to sequence learning with neural networks. In *NIPS*.
- [Zaremba et al., 2014] Zaremba, W., Sutskever, I., and Vinyals, O. (2014). Recurrent neural network regularization. *CoRR*, abs/1409.2329.