

# A Note on BPTT for LSTM LM

Tomonari MASADA @ Nagasaki University

January 15, 2015

## 1 Forward pass

$K$  is the vocabulary size.  $N$  is the number of hidden layers.  $D_n$  is the number of hidden units at the  $n$ th layer. The input sequence is  $\mathbf{X} \equiv (\mathbf{x}_1, \dots, \mathbf{x}_T)$ , where each  $\mathbf{x}_t \equiv (x_{t1}, \dots, x_{tK})^\top$  is a one-hot vector.

- $\iota$  : input gates
- $\phi$  : forget gates
- $\eta$  : cells
- $\omega$  : output gates
- $h$  : cell outputs

$f$  and  $g$  are sigmoid functions.

$$\iota_t^n = f(\mathbf{W}_{x\iota}^n \mathbf{x}_t + \mathbf{W}_{h-\iota}^n \mathbf{h}_t^{n-1} + \mathbf{W}_{hu}^n \mathbf{h}_{t-1}^n + \mathbf{w}_{\eta\iota}^n \odot \boldsymbol{\eta}_{t-1}^n + \mathbf{b}_\iota^n) \equiv f(\mathbf{a}_{\iota,t}^n) \quad (1)$$

$$\phi_t^n = f(\mathbf{W}_{x\phi}^n \mathbf{x}_t + \mathbf{W}_{h-\phi}^n \mathbf{h}_t^{n-1} + \mathbf{W}_{h\phi}^n \mathbf{h}_{t-1}^n + \mathbf{w}_{\eta\phi}^n \odot \boldsymbol{\eta}_{t-1}^n + \mathbf{b}_\phi^n) \equiv f(\mathbf{a}_{\phi,t}^n) \quad (2)$$

$$\eta_t^n = \phi_t^n \odot \boldsymbol{\eta}_{t-1}^n + \iota_t^n \odot g(\mathbf{W}_{x\eta}^n \mathbf{x}_t + \mathbf{W}_{h-\eta}^n \mathbf{h}_t^{n-1} + \mathbf{W}_{h\eta}^n \mathbf{h}_{t-1}^n + \mathbf{b}_\eta^n) \equiv \phi_t^n \odot \boldsymbol{\eta}_{t-1}^n + \iota_t^n \odot g(\mathbf{a}_{\eta,t}^n) \quad (3)$$

$$\omega_t^n = f(\mathbf{W}_{x\omega}^n \mathbf{x}_t + \mathbf{W}_{h-\omega}^n \mathbf{h}_t^{n-1} + \mathbf{W}_{h\omega}^n \mathbf{h}_{t-1}^n + \mathbf{w}_{\eta\omega}^n \odot \boldsymbol{\eta}_t^n + \mathbf{b}_\omega^n) \equiv f(\mathbf{a}_{\omega,t}^n) \quad (4)$$

$$\mathbf{h}_t^n = \omega_t^n \odot g(\boldsymbol{\eta}_t^n), \quad (5)$$

where  $\odot$  is the element-wise product. The superscript  $n$  means the  $n$ -th layer.

These can be rewritten for  $d = 1, \dots, D_n$  separately:

$$\iota_{t,d}^n = f(\mathbf{w}_{x\iota,d}^{n\top} \mathbf{x}_t + \mathbf{w}_{h-\iota,d}^{n\top} \mathbf{h}_t^{n-1} + \mathbf{w}_{hu,d}^{n\top} \mathbf{h}_{t-1}^n + \mathbf{w}_{\eta\iota,d}^n \eta_{t-1,d}^n + \mathbf{b}_{\iota,d}^n) \equiv f(a_{\iota,t,d}^n) \quad (6)$$

$$\phi_{t,d}^n = f(\mathbf{w}_{x\phi,d}^{n\top} \mathbf{x}_t + \mathbf{w}_{h-\phi,d}^{n\top} \mathbf{h}_t^{n-1} + \mathbf{w}_{h\phi,d}^{n\top} \mathbf{h}_{t-1}^n + \mathbf{w}_{\eta\phi,d}^n \eta_{t-1,d}^n + \mathbf{b}_{\phi,d}^n) \equiv f(a_{\phi,t,d}^n) \quad (7)$$

$$\eta_{t,d}^n = \phi_{t,d}^n \eta_{t-1,d}^n + \iota_{t,d}^n g(\mathbf{w}_{x\eta,d}^{n\top} \mathbf{x}_t + \mathbf{w}_{h-\eta,d}^{n\top} \mathbf{h}_t^{n-1} + \mathbf{w}_{h\eta,d}^{n\top} \mathbf{h}_{t-1}^n + \mathbf{b}_{\eta,d}^n) \equiv \phi_{t,d}^n \eta_{t-1,d}^n + \iota_{t,d}^n g(a_{\eta,t,d}^n) \quad (8)$$

$$\omega_{t,d}^n = f(\mathbf{w}_{x\omega,d}^{n\top} \mathbf{x}_t + \mathbf{w}_{h-\omega,d}^{n\top} \mathbf{h}_t^{n-1} + \mathbf{w}_{h\omega,d}^{n\top} \mathbf{h}_{t-1}^n + \mathbf{w}_{\eta\omega,d}^n \eta_{t,d}^n + \mathbf{b}_{\omega,d}^n) \equiv f(a_{\omega,t,d}^n) \quad (9)$$

$$h_{t,d}^n = \omega_{t,d}^n g(\eta_{t,d}^n). \quad (10)$$

Word probabilities for each  $t$  are obtained as follows:

$$\hat{\mathbf{y}}_t = \mathbf{b}_y + \sum_{n=1}^N \mathbf{W}_{hy}^n \mathbf{h}_t^n, \quad \mathbf{y}_t \equiv \frac{\exp(\hat{\mathbf{y}}_{t,k})}{\sum_{k'=1}^K \exp(\hat{\mathbf{y}}_{t,k'})}. \quad (11)$$

These can be rewritten for  $k = 1, \dots, K$  separately:

$$\hat{y}_{t,k} = b_{y,k} + \sum_{n=1}^N \sum_{d=1}^{D_n} w_{hy,kd}^n h_{t,d}^n, \quad y_{t,k} \equiv \frac{\exp(\hat{y}_{t,k})}{\sum_{k'=1}^K \exp(\hat{y}_{t,k'})}. \quad (12)$$

## 2 Negative log likelihood

$$\mathcal{L}(\mathbf{X}) = -\sum_{t=1}^T \log y_{t,x_{t+1}} = -\sum_{t=1}^T \log \frac{\exp(\hat{y}_{t,x_{t+1}})}{\sum_{k=1}^K \exp(\hat{y}_{t,k})} = -\sum_{t=1}^T \hat{y}_{t,x_{t+1}} + \sum_{t=1}^T \log \left\{ \sum_{k=1}^K \exp(\hat{y}_{t,k}) \right\} \quad (13)$$

## 3 Backward pass

When  $f(a) = \frac{1}{1+\exp(-a)}$ ,  $f'(a) = f(a)(1-f(a))$ . When  $f(a) = \tanh(a)$ ,  $f'(a) = 1 - f(a)^2$ .

### 3.1

$$\begin{aligned} \frac{\partial \mathcal{L}(\mathbf{X})}{\partial \hat{y}_{t,k}} &= -\frac{\partial \hat{y}_{t,x_{t+1}}}{\partial \hat{y}_{t,k}} + \frac{\partial}{\partial \hat{y}_{t,k}} \log \left\{ \sum_{\bar{k}=1}^K \exp(\hat{y}_{t,\bar{k}}) \right\} = -\frac{\partial \hat{y}_{t,x_{t+1}}}{\partial \hat{y}_{t,k}} + \frac{\sum_{\bar{k}=1}^K \exp(\hat{y}_{t,\bar{k}}) \frac{\partial \hat{y}_{t,\bar{k}}}{\partial \hat{y}_{t,k}}}{\sum_{k=1}^K \exp(\hat{y}_{t,k})} \\ &= -\left\{ \delta(x_{t+1} = k) - y_{t,k} \right\} \end{aligned} \quad (14)$$

For  $w_{hy,kd}^n$ ,

$$\frac{\partial \hat{y}_{t,k}}{\partial w_{hy,kd}^n} = \frac{\partial b_{y,k}}{\partial w_{hy,kd}^n} + \sum_{n=1}^N \sum_{\bar{d}=1}^{D_n} \left( h_{t,\bar{d}}^n \frac{\partial w_{hy,k\bar{d}}^n}{\partial w_{hy,kd}^n} + w_{hy,k\bar{d}}^n \frac{\partial h_{t,\bar{d}}^n}{\partial w_{hy,kd}^n} \right) = h_{t,d}^n \quad (15)$$

$$\therefore \frac{\partial \mathcal{L}(\mathbf{X})}{\partial w_{hy,kd}^n} = \sum_{t=1}^T \frac{\partial \mathcal{L}(\mathbf{X})}{\partial \hat{y}_{t,k}} \frac{\partial \hat{y}_{t,k}}{\partial w_{hy,kd}^n} = -\sum_{t=1}^T h_{t,d}^n \left\{ \delta(x_{t+1} = k) - y_{t,k} \right\} \quad (16)$$

### 3.2 Output errors

$$\begin{aligned} \frac{\partial \mathcal{L}(\mathbf{X})}{\partial h_{t,d}^n} &= \sum_{k=1}^K \frac{\partial \mathcal{L}(\mathbf{X})}{\partial \hat{y}_{t,k}} \frac{\partial \hat{y}_{t,k}}{\partial h_{t,d}^n} + \sum_{\bar{d}=1}^{D_n} \frac{\partial \mathcal{L}(\mathbf{X})}{\partial a_{\omega,t+1,\bar{d}}^n} \frac{\partial a_{\omega,t+1,\bar{d}}^n}{\partial h_{t,d}^n} + \sum_{\bar{d}=1}^{D_n} \frac{\partial \mathcal{L}(\mathbf{X})}{\partial a_{\eta,t+1,\bar{d}}^n} \frac{\partial a_{\eta,t+1,\bar{d}}^n}{\partial h_{t,d}^n} \\ &\quad + \sum_{\bar{d}=1}^{D_n} \frac{\partial \mathcal{L}(\mathbf{X})}{\partial a_{\phi,t+1,\bar{d}}^n} \frac{\partial a_{\phi,t+1,\bar{d}}^n}{\partial h_{t,d}^n} + \sum_{\bar{d}=1}^{D_n} \frac{\partial \mathcal{L}(\mathbf{X})}{\partial a_{\iota,t+1,\bar{d}}^n} \frac{\partial a_{\iota,t+1,\bar{d}}^n}{\partial h_{t,d}^n} \\ &\quad + \sum_{\bar{d}=1}^{D_{n+1}} \frac{\partial \mathcal{L}(\mathbf{X})}{\partial a_{\omega,t,\bar{d}}^n} \frac{\partial a_{\omega,t,\bar{d}}^n}{\partial h_{t,d}^n} + \sum_{\bar{d}=1}^{D_{n+1}} \frac{\partial \mathcal{L}(\mathbf{X})}{\partial a_{\eta,t,\bar{d}}^n} \frac{\partial a_{\eta,t,\bar{d}}^n}{\partial h_{t,d}^n} \\ &\quad + \sum_{\bar{d}=1}^{D_{n+1}} \frac{\partial \mathcal{L}(\mathbf{X})}{\partial a_{\phi,t,\bar{d}}^n} \frac{\partial a_{\phi,t,\bar{d}}^n}{\partial h_{t,d}^n} + \sum_{\bar{d}=1}^{D_{n+1}} \frac{\partial \mathcal{L}(\mathbf{X})}{\partial a_{\iota,t,\bar{d}}^n} \frac{\partial a_{\iota,t,\bar{d}}^n}{\partial h_{t,d}^n} \\ &= \sum_{k=1}^K w_{hy,kd}^n \frac{\partial \mathcal{L}(\mathbf{X})}{\partial \hat{y}_{t,k}} + \sum_{\bar{d}=1}^{D_n} w_{h\omega,\bar{d}d}^n \frac{\partial \mathcal{L}(\mathbf{X})}{\partial a_{\omega,t+1,\bar{d}}^n} + \sum_{\bar{d}=1}^{D_n} w_{h\eta,\bar{d}d}^n \frac{\partial \mathcal{L}(\mathbf{X})}{\partial a_{\eta,t+1,\bar{d}}^n} \\ &\quad + \sum_{\bar{d}=1}^{D_n} w_{h\phi,\bar{d}d}^n \frac{\partial \mathcal{L}(\mathbf{X})}{\partial a_{\phi,t+1,\bar{d}}^n} + \sum_{\bar{d}=1}^{D_n} w_{h\iota,\bar{d}d}^n \frac{\partial \mathcal{L}(\mathbf{X})}{\partial a_{\iota,t+1,\bar{d}}^n} \\ &\quad + \sum_{\bar{d}=1}^{D_{n+1}} w_{h-\omega,\bar{d}d}^{n+1} \frac{\partial \mathcal{L}(\mathbf{X})}{\partial a_{\omega,t,\bar{d}}^{n+1}} + \sum_{\bar{d}=1}^{D_{n+1}} w_{h-\eta,\bar{d}d}^{n+1} \frac{\partial \mathcal{L}(\mathbf{X})}{\partial a_{\eta,t,\bar{d}}^{n+1}} \\ &\quad + \sum_{\bar{d}=1}^{D_{n+1}} w_{h-\phi,\bar{d}d}^{n+1} \frac{\partial \mathcal{L}(\mathbf{X})}{\partial a_{\phi,t,\bar{d}}^{n+1}} + \sum_{\bar{d}=1}^{D_{n+1}} w_{h-\iota,\bar{d}d}^{n+1} \frac{\partial \mathcal{L}(\mathbf{X})}{\partial a_{\iota,t,\bar{d}}^{n+1}} \end{aligned} \quad (17)$$

### 3.3 Output gates

$$\frac{\partial \mathcal{L}(\mathbf{X})}{\partial a_{\omega,t,d}^n} = \frac{\partial \mathcal{L}(\mathbf{X})}{\partial h_{t,d}^n} \frac{\partial h_{t,d}^n}{\partial a_{\omega,t,d}^n} \quad (18)$$

$$\frac{\partial h_{t,d}^n}{\partial a_{\omega,t,d}^n} = g(\eta_{t,d}^n) f'(a_{\omega,t,d}^n) \quad (19)$$

### 3.4 States

$$\begin{aligned} \frac{\partial \mathcal{L}(\mathbf{X})}{\partial \eta_{t,d}^n} &= \frac{\partial \mathcal{L}(\mathbf{X})}{\partial h_{t,d}^n} \frac{\partial h_{t,d}^n}{\partial \eta_{t,d}^n} + \frac{\partial \mathcal{L}(\mathbf{X})}{\partial \eta_{t+1,d}^n} \frac{\partial \eta_{t+1,d}^n}{\partial \eta_{t,d}^n} + \frac{\partial \mathcal{L}(\mathbf{X})}{\partial a_{\phi,t+1,d}^n} \frac{\partial a_{\phi,t+1,d}^n}{\partial \eta_{t,d}^n} + \frac{\partial \mathcal{L}(\mathbf{X})}{\partial a_{\iota,t+1,d}^n} \frac{\partial a_{\iota,t+1,d}^n}{\partial \eta_{t,d}^n} \\ &= \frac{\partial \mathcal{L}(\mathbf{X})}{\partial h_{t,d}^n} \frac{\partial h_{t,d}^n}{\partial \eta_{t,d}^n} + \phi_{t+1,d}^n \frac{\partial \mathcal{L}(\mathbf{X})}{\partial \eta_{t+1,d}^n} + w_{\eta\phi,d} \frac{\partial \mathcal{L}(\mathbf{X})}{\partial a_{\phi,t+1,d}^n} + w_{\eta\iota,d} \frac{\partial \mathcal{L}(\mathbf{X})}{\partial a_{\iota,t+1,d}^n} \end{aligned} \quad (20)$$

$$\begin{aligned} \frac{\partial h_{t,d}^n}{\partial \eta_{t,d}^n} &= g(\eta_{t,d}^n) \frac{\partial \omega_{t,d}^n}{\partial \eta_{t,d}^n} + \omega_{t,d}^n g'(\eta_{t,d}^n) \\ &= w_{\eta\omega,d}^n f'(a_{\omega,t,d}^n) g(\eta_{t,d}^n) + \omega_{t,d}^n g'(\eta_{t,d}^n) \\ &= w_{\eta\omega,d}^n \frac{\partial h_{t,d}^n}{\partial a_{\omega,t,d}^n} + \omega_{t,d}^n g'(\eta_{t,d}^n) \end{aligned} \quad (21)$$

$$\begin{aligned} \therefore \frac{\partial \mathcal{L}(\mathbf{X})}{\partial \eta_{t,d}^n} &= \frac{\partial \mathcal{L}(\mathbf{X})}{\partial h_{t,d}^n} \left\{ w_{\eta\omega,d}^n \frac{\partial h_{t,d}^n}{\partial a_{\omega,t,d}^n} + \omega_{t,d}^n g'(\eta_{t,d}^n) \right\} + \phi_{t+1,d}^n \frac{\partial \mathcal{L}(\mathbf{X})}{\partial \eta_{t+1,d}^n} + w_{\eta\phi,d} \frac{\partial \mathcal{L}(\mathbf{X})}{\partial a_{\phi,t+1,d}^n} + w_{\eta\iota,d} \frac{\partial \mathcal{L}(\mathbf{X})}{\partial a_{\iota,t+1,d}^n} \\ &= w_{\eta\omega,d}^n \frac{\partial \mathcal{L}(\mathbf{X})}{\partial a_{\omega,t,d}^n} + \omega_{t,d}^n g'(\eta_{t,d}^n) \frac{\partial \mathcal{L}(\mathbf{X})}{\partial h_{t,d}^n} + \phi_{t+1,d}^n \frac{\partial \mathcal{L}(\mathbf{X})}{\partial \eta_{t+1,d}^n} + w_{\eta\phi,d} \frac{\partial \mathcal{L}(\mathbf{X})}{\partial a_{\phi,t+1,d}^n} + w_{\eta\iota,d} \frac{\partial \mathcal{L}(\mathbf{X})}{\partial a_{\iota,t+1,d}^n} \end{aligned} \quad (22)$$

### 3.5 Cells, forget gates, and input gates

$$\frac{\partial \mathcal{L}(\mathbf{X})}{\partial a_{\eta,t,d}^n} = \frac{\partial \eta_{t,d}^n}{\partial a_{\eta,t,d}^n} \frac{\partial \mathcal{L}(\mathbf{X})}{\partial \eta_{t,d}^n} = \iota_{t,d}^n g'(a_{\eta,t,d}^n) \frac{\partial \mathcal{L}(\mathbf{X})}{\partial \eta_{t,d}^n} \quad (23)$$

$$\frac{\partial \mathcal{L}(\mathbf{X})}{\partial a_{\phi,t,d}^n} = \frac{\partial \eta_{t,d}^n}{\partial a_{\phi,t,d}^n} \frac{\partial \mathcal{L}(\mathbf{X})}{\partial \eta_{t,d}^n} = f'(a_{\phi,t,d}^n) \eta_{t-1,d}^n \frac{\partial \mathcal{L}(\mathbf{X})}{\partial \eta_{t,d}^n} \quad (24)$$

$$\frac{\partial \mathcal{L}(\mathbf{X})}{\partial a_{\iota,t,d}^n} = \frac{\partial \eta_{t,d}^n}{\partial a_{\iota,t,d}^n} \frac{\partial \mathcal{L}(\mathbf{X})}{\partial \eta_{t,d}^n} = f'(a_{\iota,t,d}^n) g(\eta_{t,d}^n) \frac{\partial \mathcal{L}(\mathbf{X})}{\partial \eta_{t,d}^n} \quad (25)$$

## References

- [1] Thomas Breuel, Volkmar Frinken, and Marcus Liwicki. Long Short-Term Memory Managing Long-Term Dependencies within Sequences. *nntutorial-lstm.pdf* at <http://lstm.iupr.com/files>
- [2] Alex Graves and Jürgen Schmidhuber. Framewise phoneme classification with bidirectional LSTM and other neural network architectures. *Neural Networks* 18 (2005) 602–610.
- [3] Alex Graves. Generating Sequences With Recurrent Neural Networks. *arXiv preprint arXiv:1308.0850*, 2013.