

Deep Recurrent Neural Network

Alex Adamson
aadamson@stanford.edu

May 13, 2015

1 Network specification

We borrow our specification from Irsoy and Cardie, 2014.

For $i > 1$, we have

$$\vec{h}_t^{(i)} = f(\vec{W}_t^{(i)} \vec{h}_t^{(i-1)} + \vec{W}_t^{(i)} \overleftarrow{h}_t^{(i-1)} + \vec{V}^{(i)} \vec{h}_{t-1}^{(i)} + \vec{b}^{(i)}) \quad (1)$$

$$\overleftarrow{h}_t^{(i)} = f(\overleftarrow{W}_t^{(i)} \vec{h}_t^{(i-1)} + \overleftarrow{W}_t^{(i)} \overleftarrow{h}_t^{(i-1)} + \overleftarrow{V}^{(i)} \overleftarrow{h}_{t+1}^{(i)} + \overleftarrow{b}^{(i)}) \quad (2)$$

and for $i = 1$ we have

$$\vec{h}_t^{(1)} = f(\vec{W}^{(1)} x_t + \vec{V}^{(1)} \vec{h}_{t-1}^{(1)} + \vec{b}^{(1)}) \quad (3)$$

$$\overleftarrow{h}_t^{(1)} = f(\overleftarrow{W}^{(1)} x_t + \overleftarrow{V}^{(1)} \overleftarrow{h}_{t+1}^{(1)} + \overleftarrow{b}^{(1)}) \quad (4)$$

We only connect the last layer (which we denote layer L) to the output layer:

$$y_t = g(\underline{U} \vec{h}_t^{(L)} + \overleftarrow{U} \overleftarrow{h}_t^{(L)} + c) \quad (5)$$

2 Backpropagation derivation

For convenience, we use error vector notation when deriving the backpropagation updates.

For our loss function, we choose categorical cross-entropy.

Let δ_t^y be the error vector propagated by the softmax unit at timestep t .

Let $\vec{\delta}_t^{(i)}$ be the error vector propagated by the forward hidden unit in layer i at timestep t .

Let $\overleftarrow{\delta}_t^{(i)}$ be the error vector propagated by the backward hidden unit in layer i at timestep t .

Let f^* be the function such that $f^*(f(x)) = f'(x)$ where f is as above.

Then:

$$\delta_t^y = \hat{y}_t - y_t \quad (6)$$

$$\vec{\delta}_t^{(i)} = f^*(\vec{h}_t^{(i)}) \circ ((\vec{W}_t^{(i+1)})^T \vec{\delta}_t^{(i+1)} + (\overleftarrow{W}_t^{(i+1)})^T \overleftarrow{\delta}_t^{(i+1)} + (\vec{V}^{(i)})^T \vec{\delta}_{t+1}^{(i)}) \quad (7)$$

$$\overleftarrow{\delta}_t^{(i)} = f^*(\overleftarrow{h}_t^{(i)}) \circ ((\overleftarrow{W}_t^{(i+1)})^T \vec{\delta}_t^{(i+1)} + (\overleftarrow{W}_t^{(i+1)})^T \overleftarrow{\delta}_t^{(i+1)} + (\overleftarrow{V}^{(i)})^T \overleftarrow{\delta}_{t-1}^{(i)}) \quad (8)$$

With these in hand, we can find the actual updates:

$$\frac{\partial J}{\partial \vec{U}} = \delta^y (\vec{h}^{(L)})^T \quad (9)$$

$$\frac{\partial J}{\partial c} = \delta^y \cdot \mathbf{1} \quad (10)$$

$$\frac{\partial J}{\partial \vec{W}^{(i)}} = \vec{\delta}^{(i)} (\vec{h}^{(i-1)})^T \quad (11)$$

$$\frac{\partial J}{\partial \overleftarrow{W}^{(i)}} = \vec{\delta}^{(i)} (\overleftarrow{h}^{(i-1)})^T \quad (12)$$

$$\frac{\partial J}{\partial \vec{b}^{(i)}} = \vec{\delta}^{(i)} \cdot \mathbf{1} \quad (13)$$

$$\frac{\partial J}{\partial \overleftarrow{W}^{(i)}} = \overleftarrow{\delta}^{(i)} (\vec{h}^{(i-1)})^T \quad (14)$$

$$\frac{\partial J}{\partial \overleftarrow{W}^{(i)}} = \overleftarrow{\delta}^{(i)} (\overleftarrow{h}^{(i-1)})^T \quad (15)$$

$$\frac{\partial J}{\partial \overleftarrow{b}^{(i)}} = \overleftarrow{\delta}^{(i)} \cdot \mathbf{1} \quad (16)$$

$$\frac{\partial J}{\partial \vec{V}^{(i)}} = \sum_{t=1}^T \vec{\delta}_t^{(i)} (\vec{h}_{t-1}^{(i)})^T \quad (17)$$

$$\frac{\partial J}{\partial \overleftarrow{V}^{(i)}} = \sum_{t=1}^T \overleftarrow{\delta}_t^{(i)} (\overleftarrow{h}_{t+1}^{(i)})^T \quad (18)$$

$$\frac{\partial J}{\partial \vec{W}^{(1)}} = \vec{\delta}^{(1)} x_t^T \quad (19)$$

$$\frac{\partial J}{\partial \overleftarrow{W}^{(1)}} = \overleftarrow{\delta}^{(1)} x_t^T \quad (20)$$